

# Comparative Analysis of K-Means Method and Naïve Bayes Method for Brute Force Attack Visualization

Deris Stiawan

Computer Engineering Dept., Faculty of Comp. Science  
Sriwijaya University, Palembang, Indonesia  
deris@unsri.ac.id

Sari Sandra

Computer Engineering Dept., Faculty of Comp. Science,  
Sriwijaya University, Palembang, Indonesia  
sarischandrasalsabilla@gmail.com

Esam Alzahrani

College of Computer Science & IT  
Albaha University, Albaha, Saudi Arabia  
easz000@hotmail.com

Rahmat Budiarto

College of Computer Science & IT  
Albaha University, Albaha, Saudi Arabia  
rahmat@bu.edu.sa

**Abstract**— This paper presents 2-Dimensional visualization to categorize packets of network traffic into normal data pattern and attack data pattern based on the patterns resulted by a brute force attack. Two clustering methods: K-Means and Naïve Bayes methods are used to produce the data to be visualized. Experiments using ISCX and DARPA dataset were conducted. Brute force assaults on some service protocols. This paper focuses on SSH service for ISCX dataset and TELNET service for DARPA dataset. Visual analysis of the experimental results show a better results in term of accuracy by reducing false alarms.

**Keywords**— visualization, ISCX dataset, DARPA dataset, brute force, K-Means method and Naïve Bayes method

## I. INTRODUCTION

One of techniques, that most common used by the attackers is brute force attack with attack's percentage reached 25% of the total attacks[1], [2], [3]. In brute force attack, attacker attempts to login using SSH and TELNET protocol to reveal login password [4, 5]. This protocol enables data exchange between two network devices, which are widely used on Linux and other Unix-based systems.

In general brute force attack produces a data traffic that contains attack data pattern and normal data pattern. K-Means and Naïve Bayes methods are tools for data mining for intrusion detection system [6]. Experiments to categorize such type of data can be done off-line using data from existing datasets with the final purpose as a basis of providing visual results against attacks exist in the dataset.

This paper attempts to identify the patterns of attack data patterns as well as normal data patterns of brute force attack by utilizing recorded brute force attack data from two datasets: ISCX and DARPA. K-Means and Naïve Bayes methods are chosen with the rational of simplicity of the methods.

## II. RELATED WORKS

Authors in [4, 5], discuss about intrusion detection dataset using K-Means algorithm. This research tries to clustering datasets into normal category and attack category i.e. DOS,

Probe, R2L and U2R. Nevertheless, the research only uses NSL-KDD dataset for clustering as such no other benchmarking being done.

Besides K-Means algorithm, there is Naïve Bayes algorithm that can be used in identifying intrusion detection, such as in research work by [7]. The research discusses about the use of K-Means and Naïve Bayes algorithm to overcome false alarm using ISCX dataset.

Research in [8], discusses about automatic detection of attack using Parallel Coordinates Attack Visualization (PCAV). The proposed work detects Internet attack on large scale, such as Internet worms, DDOS and network scanning.

Furthermore, research by [9], discusses detection of attacks without monitoring the anomalies. The research uses KDD Cup 1999 dataset with K-Means clustering algorithm. The work utilizes cluster 3.0. tool and TreeView visualization tool.

## III. PREPARE YOUR PAPER BEFORE STYLING

The use of K-Means and Naïve Bayes is expected to provide better accuracy in attack detection and further gives a visual picture in categorizing a brute force attack. Moreover, this research uses software as additional tools to support the experiments. The following specifications of software requirements described in Table 1.

As shown in Table 1 that the software used are Snort and Visual Studio. Snort is used as Network Intrusion Detection System (NIDS) in detecting brute force attack in SSH service and TELNET service, while Visual Studio is used as a system of Attack Pattern and Normal Pattern and also used for visualization of data patterns. An overview of visualization system of brute force attack is shown in Fig. 1. The system consists of 4 components.

TABLE I. SOFTWARE REQUIREMENTS SPECIFICATION

Systems	Tools	Description
NIDS SSH and TELNET Attack Pattern and Normal Pattern Visualization	Snort Visual Studio Visual Studio	Version 2.9.8.0 2012 2012

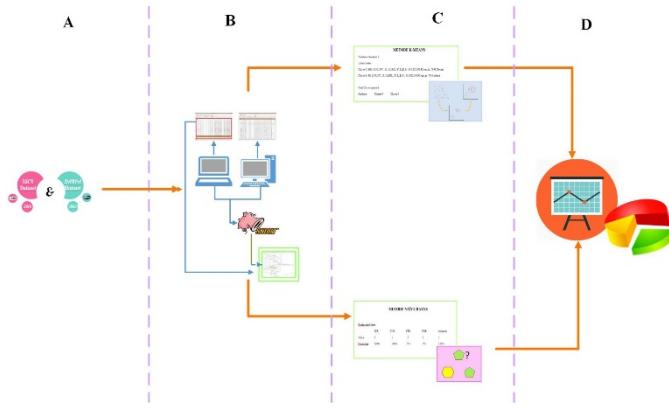


Fig. 1. Overview of brute force attack visualization: (A) dataset, (B) attack pattern and normal pattern, (C) K-Means and Naïve Bayes method and (D) visualization

The following sections explain each component.

#### A. Dataset

The dataset used in this research are ISCX dataset and DARPA dataset in CSV (Comma Separated Values) format. ISCX dataset is a dataset that used to capture network traffic, developed by Faculty of computer science, University of New Brunswick [10]. ISCX dataset simulates infiltrating the network from the inside, HTTP denial of service, distributed denial of service an IRC Botnet and brute force attack scenario SSH taken during 11-June 17, 2010.

The experiments in this paper focus only on one scenario of ISCX datasets i.e. Brute force attack on June 17, 2010 on the secure shell (SSH) service which consists of 20 features with 5540 packets. SSH provides fairly safe service remote login and also authentication and authorization systems, so to access the service, a system needs to login can cause attacks with brute force techniques occur. SSH service is part of TCP protocol.

DARPA dataset is collected in 1998 and 1999 by Information Systems Technology group of MIT Lincoln Laboratory, under Defense Advanced Research Projects Agency (DARPA) and Air Force Research Laboratory (AFRL/SNHS) sponsorship [11]. DARPA is created to simulate the traffic at U.S. air force bases to evaluate intrusion detection system. In DARPA dataset, brute force attack occurs at telecommunication network (TELNET) service which consists of 42 features with 1234 packets. Telnet does not use the security mechanisms such authentication system and encryption techniques and also transfer the data in plain-text, so that information becomes a major threat in network.

#### B. Attack Pattern and Normal Pattern

Attack pattern and normal pattern on this research use a matching program to define the brute force pattern. Results of the matching program are attack and normal dominant packets in dataset which will become pre-process on further research. Fig. 2 shows the flowchart of attack pattern and normal

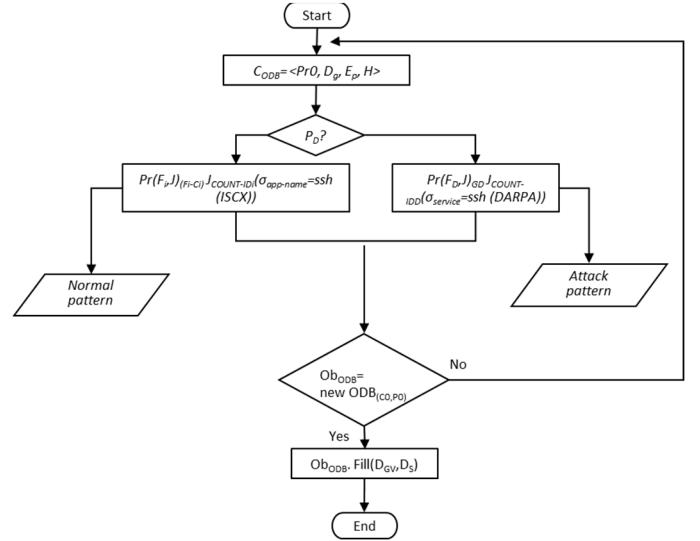


Fig. 2. Flowchart of Attack Pattern and Normal Pattern

pattern process. It shows some steps that occur in process of attack pattern and normal pattern identification. The first step is connecting to the database, where each dataset included in the program will connect to the database. Furthermore, the dataset will be through the second step (process dataset step) to get the desired output based on features in the dataset. When processing data is completed, then the result of the second step will be inserted into the datagridview in the program.

The Snort tool is also used in this stage to prove that there is brute force attack on ISCX and DARPA dataset. Snort can work in four modes that are sniffer, packet logger, Network Intrusion Detection System (NIDS) and Intrusion Prevention System (IPS). Snort NIDS mode is used in this experiment, with the setup of various rules based on rule options to detect the attack so that can distinguish normal packets from attack packets.

Besides matching programs and Snort tools, in this stage, system traceroutes to know the route through by attacker to reach servers in doing the attack. Traceroute is also useful in proving that there are certain client and server host on ISCX and DARPA dataset.

#### C. K-Means Method and Naïve Bayes Method

The K-Means method is one of non-hierarchy data clustering method that grouping data in one or more clusters. On K-Means method, distance of each data to each cluster is calculated using Euclidean Distance as in (1).

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} ; i = 1, 2, 3, \dots n \quad (1)$$

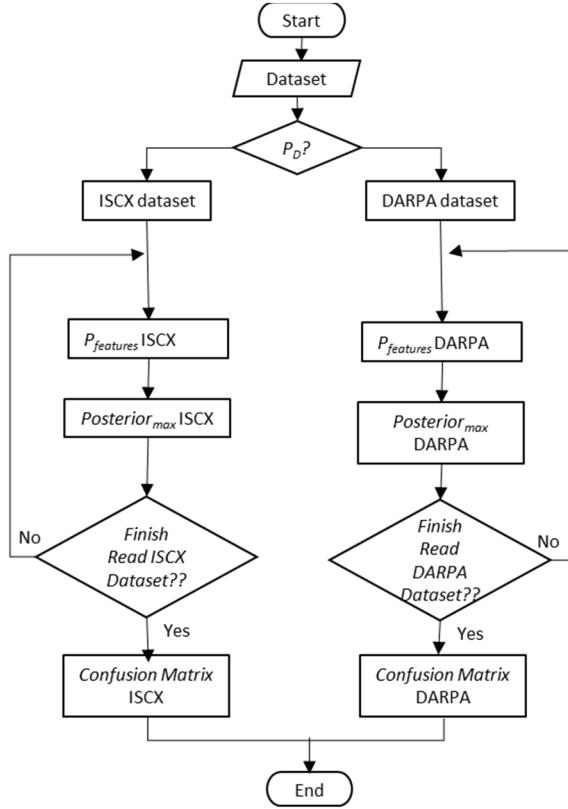


Fig. 3. Flowchart of K-Means method

In this experiment, the number of cluster ( $k$ ) is two, to distinguish attack packet and normal packets categories. Fig. 3 shows the flowchart of K-Means method consists of five steps.

Step 1 in K-Means method process is to determine the number of cluster ( $k$ ), the next step is initialization of centroid ( $Ce_o$ ). Initialization of centroid is determined based on attack pattern and normal pattern of pre-processing results in each dataset. Step 3 is a process of Euclidean distance ( $d$ ) to determine the nearest distance from each data packet to the centroid. Furthermore, the dataset will through Step 4, iteration cluster (CI) step. Cluster iteration step is performed if the result of the clusters have not been consistent, so the system do recalculation using new centroid value ( $Ce_n$ ) from membership of cluster. Last step of K-Means method is processing the confusion matrix to provide accuracy results.

Naïve Bayes method is also used in this experiment, to classify the data packet in pattern recognition. Naïve Bayes provides speed and high accuracy when applied into a large data. Naïve Bayes based on Bayes theorem as shown in (2).

$$P(H | X_i) = \frac{P(H) P(X_i | H)}{P(X_i)} \quad (2)$$

Naïve Bayes method categorizes data into specific categories based on the highest posterior probability  $P(H | X_i)$ . Classification of data will occur if and only if the posterior probability data  $H$  based on condition  $X_i$   $\{P(H | X_i)\}$  is smaller than posterior probability data  $H$  based on  $\{P(H | X_j)\}$  as shown in (3). Fig. 4 shows the flowchart of Naïve Bayes method.

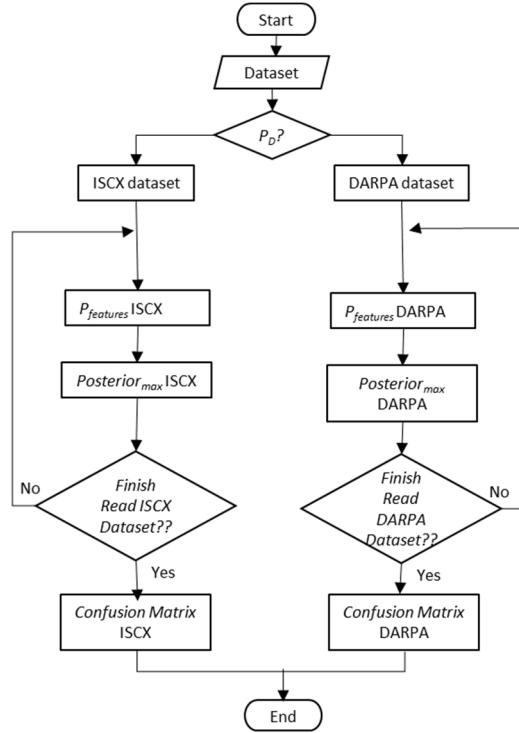


Fig. 4. Flowchart of Naïve Bayes method

$$P(H | X_i) > P(H | X_j) \quad (3)$$

where,  $j \geq 1$  dan  $j \neq i$

In Fig. 4 it is explained that there are three main steps in the process of Naïve Bayes method. Step 1 in Naïve Bayes method process is probability of features ( $P_{features}$ ), the next step is posterior process ( $P_{max}$ ) and the last step is process of the confusion matrix.

#### D. Visualization

The visualization component illustrates a pattern of attack and normal package in ISCX and DARPA dataset with parallel coordinate-visualization design. Visualization of parallel coordinate-describes information in two dimensional (2D).

## IV. RESULTS AND DISCUSSION

In the beginning of the experiment, validation process is performed to each dataset (ISCX and DARPA dataset) for detecting attack data pattern and normal data pattern which form a data correlation to data packet with IDS engine and traceroute. The results of the validation process of each dataset are shown in Fig. 5 and Fig. 6.

Fig. 7 shows dominant attack pattern and normal pattern on ISCX dataset with 1498 dominant value and 17 data packets. Dominant attack pattern has the following feature values.

"TotalSourceBytes" features value is 1274,  
 "TotalDestinationBytes" features value is 2343,  
 "TotalDestinationPackets" features value is 11,  
 "TotalSourcePackets" features value is 10 with "Direction"

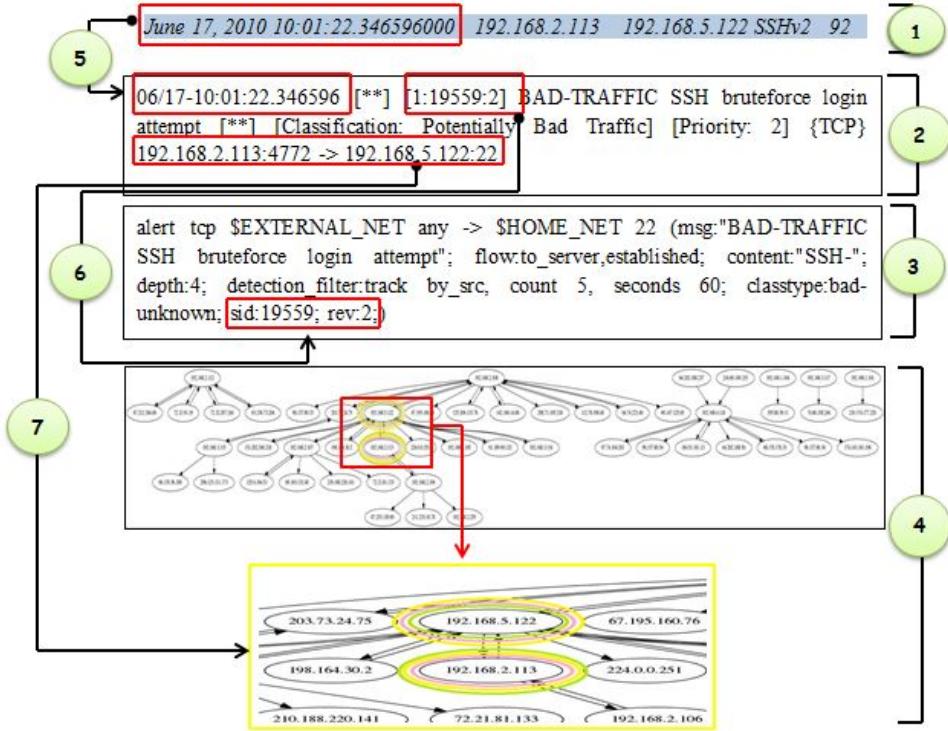


Fig. 5. Flowchart of ISCX Dataset

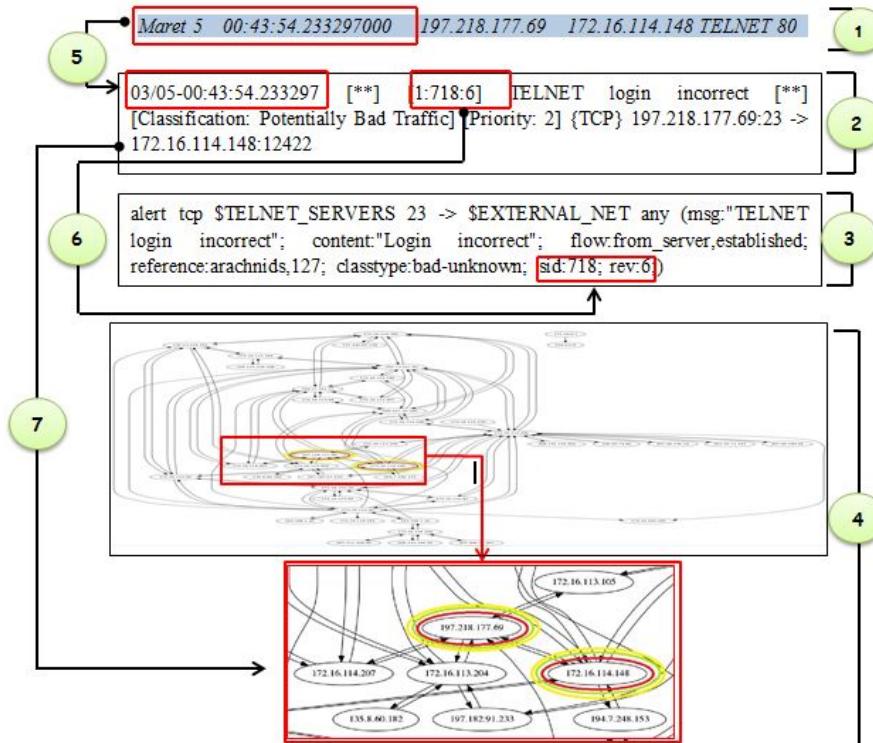


Fig. 6. Correlation of DARPA Dataset

features is "R2L", "SourceTCPFlagsDescription" and "DestinationTCPFlagsDescription" features are F, S, P,

Additionally, on "tags" attack known that dominant source is 131.202.243.90 with destination server is 192.168.5.122, where IP server is main server of ISCX IP, responsible for providing e-mail services so allowed brute force attack occurred in that IP server. Meanwhile, normal pattern in Fig. 7 has dominant data with the following features value.

"TotalSourceBytes" features value is 1724, "TotalDestinationBytes" features value is 6414, "TotalDestinationPackets" features value is 42, "TotalSourcePackets" features value is 15, "Direction" features is L2L, "SourceTCPFlagsDescription" features is S, R, P, A and "Destination-TCPFlagsDescription" features is S, P, A with IP source 192.168.4.120 and 192.168.5.122 for destination IP.

Fig. 8 shows attack pattern and normal pattern on DARPA dataset. Dominant attack pattern with the "Label" guess password is the data with "Protocol" features: TCP, "Service" features: TELNET, "Flags" features: RSTO, "DurationBytes" features value: 179 and "Num\_Failed\_logins" features minimal value is one. This dominant data, will be the pattern of a TELNET brute force attack with dominant number of data reaches 45 rows. Whereas, dominant normal pattern has dominant value amounted to 13. The data has "Protocol" features: TCP, "Service" features: TELNET, "Flags" features: S1, "DurationBytes" features value: 2832 and "Num\_Failed\_logins" features minimal value is zero. The next experiment step is the step of performing K-Means method and Naïve Bayes method in the dataset. The results of experiment of using K-Means and Naïve Bayes method of confusion matrix, are shown in Table 2.

From Table 2, it can be seen that the results of packet categorization using Naïve Bayes method get greater accuracy than K-Means method. Naïve Bayes method accuracy on ISCX and DARPA dataset are 99.68% and 98.779% respectively, whereas, the K-Means method on ISCX and DARPA dataset obtained 95.46% and 73.60%, respectively. Naïve Bayes method get false alarms lower than K-Means method.

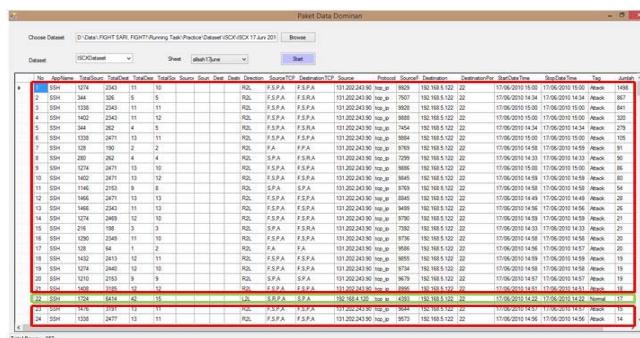


Fig. 7. Attack Pattern and Normal Pattern in ISCX Dataset

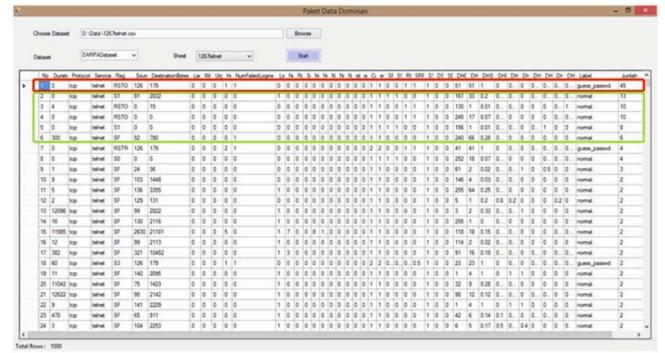


Fig. 8. Attack Pattern and Normal Pattern DARPA Dataset

TABLE II. SOFTWARE REQUIREMENTS SPECIFICATION

Methods	Dataset	Results			
		TP	TN	FP	Accuracy (%)
K-Means	ISCX	5197	92	245	6 99.68
	DARPA	52	857	325	1 73.60
Naïve Bayes	ISCX	5185	337	18	0 99.68
	DARPA	51	1169	13	2 98.779

The establishment of parallel coordinate visualization is implemented into normal pattern and attack pattern programs on ISCX and DARPA dataset. Screenshots of parallel coordinate visualization application are shown in Figure 9.

Figure 9 shows parallel coordinate visualization. (A) is parallel coordinate visualization of ISCX dataset by using K-Means method. The red line is attack by category in cluster 0, while the green line is cluster 1 that is normal. (B) is a DARPA dataset visualization using Naïve Bayes method with red line is category 0 (attack) and the green line is normal contains into category of cluster 1. (C) gives the results of ISCX datasets parallel coordinate visualization by using Naïve Bayes methods. The green line is data packet in normal category, while red line is attack data packet. (D) shows result of DARPA dataset visualization by using Naïve Bayes method.

## V. CONCLUSION AND FUTURE WORK

Based on the experiments that have been conducted, it is concluded that a brute force attack on ISCX dataset in SSH services form an attack pattern where one IP source focuses on attack one server, with the port of destination to be exploited is port 22. While, brute force attack on the TELNET service obtained a form of attack pattern where attackers IP experience minimal one failed login, with destination port to be exploited is port 23 and the destination bytes is 179 bytes.

K-Means method and Naïve Bayes method are implemented on dataset for categorizing several attack data packets or normal data packets. The results of both method implementations provide two dimensional (2D) visualization

that are scatter plots or parallel coordinate visualization, with good accuracy of categorization.

K-Means method and Naïve Bayes method are implemented in the ISCX dataset reached accuracy results up to 95.46% and 99.68%, respectively, while at DARPA dataset reached accuracy result up to 73.60% and 98.79%, respectively. Future research, include real time visualization with addition of attacks types such as SQL injection, probe, Internet worms and network scanning.

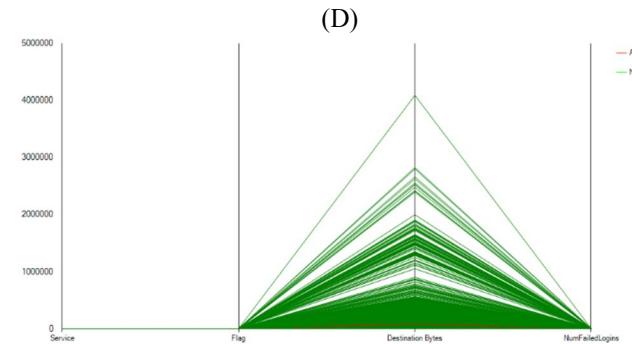
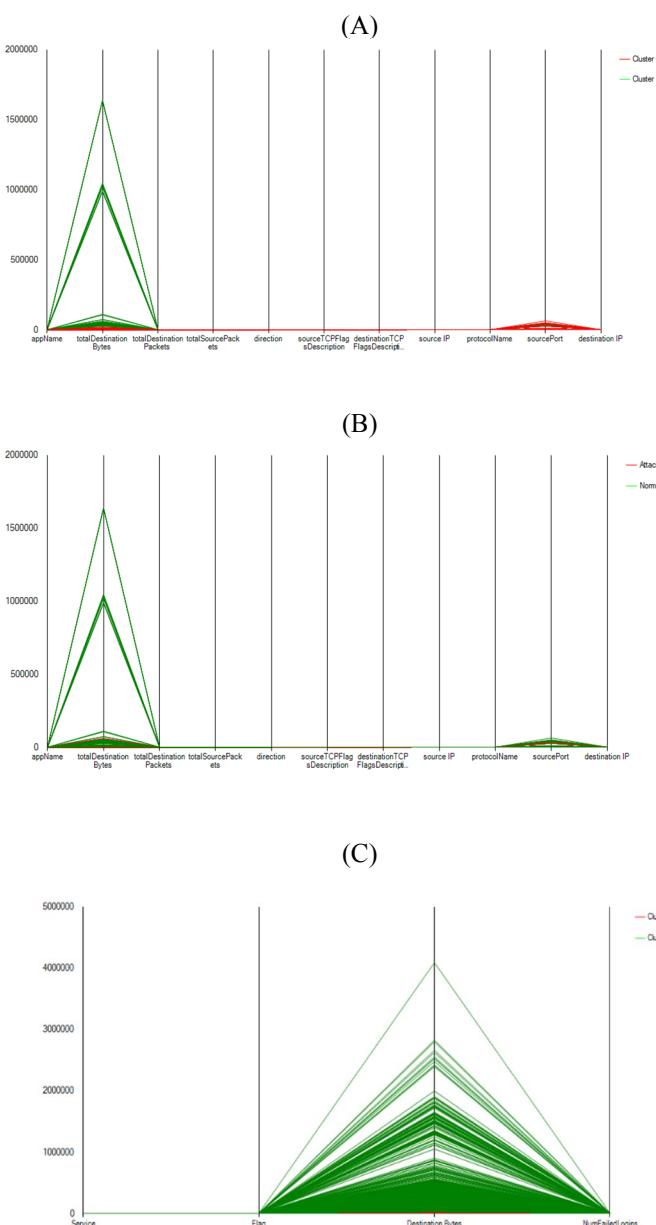


Fig. 9. Visualization of the parallel coordinate

## ACKNOWLEDGMENT

This research was supported by the Department of Computer Systems, Faculty of Computer Science, Sriwijaya University and the COMNETS research laboratory.

## REFERENCES

- [1] D. Stiawan, M. Y. Idris, A. H. Abdullah, F. Aljaber, and R. Budiarso, "Cyber-attack penetration test and vulnerability analysis," International Journal of Online Engineering, vol. 13, pp. 125-132, 2017.
- [2] A. Joshi, M. Wazid, and R. H. Goudar, "An Efficient Cryptographic Scheme for Text Message Protection Against Brute Force and Cryptanalytic Attacks," Procedia Computer Science, vol. 48, pp. 360-366, 2015.
- [3] K. Kaynar, "A taxonomy for attack graph generation and usage in network security," Journal of Information Security and Applications, vol. 29, pp. 27-56, 8// 2016.
- [4] S. Anandita, Y. Rosmansyah, B. Dabarsyah, and J. U. Choi, "Implementation of dendritic cell algorithm as an anomaly detection method for port scanning attack," in 2015 International Conference on Information Technology Systems and Innovation (ICITSI), 2015, pp. 1-6.
- [5] J. Vykopal, "A Flow-Level Taxonomy and Prevalence of Brute Force Attacks" in Advances in Computing and Communications. vol. 191, A. Abraham, J. Lloret Mauri, J. F. Buford, J. Suzuki, and S. M. Thampi, Eds., ed: Springer Berlin Heidelberg, 2011, pp. 666-675.
- [6] [S. Mukherjee and N. Sharma, "Intrusion Detection using Naive Bayes Classifier with Feature Reduction," Procedia Technology, vol. 4, pp. 119-128, 2012.
- [7] M. Kumagai, Y. Musashi, D. A. L. Romana, K. Takemori, S. Kubota, and K. Sugitani, "SSH Dictionary Attack and DNS Reverse Resolution Traffic in Campus Network," in 2010 Third International Conference on Intelligent Networks and Intelligent Systems, 2010, pp. 645-648.
- [8] W. Yassin, N. I. Udzir, Z. Muda, and M. N. Sulaiman, "Anomaly-based intrusion detection through k-means clustering and naives bayes classification," in Proc. 4th Int. Conf. Comput. Informatics, ICOCI, 2013, pp. 298-303.
- [9] H. Choi, H. Lee, and H. Kim, "Fast detection and visualization of network attacks on parallel coordinates," Computers & Security, vol. 28, pp. 276-288, 7// 2009.
- [10] A. Riad, I. Elhenawy, A. Hassan, and N. Awadallah, "Visualize network anomaly detection by using k-means clustering algorithm," International Journal of Computer Networks & Communications, vol. 5, p. 195, 2013.
- [11] R. Zuech, T. M. Khoshgoftaar, N. Seliya, M. M. Najafabadi, and C. Kemp, "A New Intrusion Detection Benchmarking System," in FLAIRS Conference, 2015, pp. 252-256