

Description of Kyoto University Benchmark Data

Jungsuk SONG¹, Hiroki Takakura², and Yasuo Okabe³

¹ National Institute of Information and Communications Technology (NICT), Japan
song@nict.go.jp

² Information Technology Center (ITC), Nagoya University
takakura@itc.nagoya-u.ac.jp

³ Academic Center for Computing and Media Studies (ACCMS), Kyoto University
okabe@i.kyoto-u.ac.jp

Our benchmark data consist of the following 24 statistical features; 14 conventional features and 10 additional features. Among them, the first 14 features were extracted based on KDD Cup 99 data set, which is a very popular and widely used performance evaluation data in intrusion detection research field[1]. Among 41 original features of KDD Cup 99 data set, we have extracted only 14 significant and essential features from the raw traffic data obtained by honeypot systems[2] that are deployed in Kyoto University. Addition to those 14 features, we have also extracted additional 10 features which may enable us to investigate more effectively what happens on our networks. Of course, they also can be utilized for training and testing our data with 14 convention features. Note that the order of the below features is exactly the same to that of the actual data.

=====14 conventional features=====

1. Duration: the length (number of seconds) of the connection
2. Service: the connection's service type, e.g., http, telnet, etc
3. Source bytes: the number of data bytes sent by the source IP address
4. Destination bytes: the number of data bytes sent by the destination IP address
5. Count: the number of connections whose source IP address and destination IP address are the same to those of the current connection in the past two seconds
6. Same_srv_rate: % of connections to the same service in Count feature
7. Error_rate: % of connections that have "SYN" errors in Count feature
8. Srv_error_rate: % of connections that have "SYN" errors in Srv_count(the number of connections whose service type is the same to that of the current connection in the past two seconds) feature
9. Dst_host_count: among the past 100 connections whose destination IP address is the same to that of the current connection, the number of connections whose source IP address is also the same to that of the current connection
10. Dst_host_srv_count: among the past 100 connections whose destination IP address is the same to that of the current connection, the number of connections whose service type is also the same to that of the current connection

11. Dst_host_same_src_port_rate: % of connections whose source port is the same to that of the current connection in Dst_host_count feature
12. Dst_host_serror_rate: % of connections that have “SYN” errors in Dst_host_count feature
13. Dst_host_srv_serror_rate: % of connections that “SYN” errors in Dst_host_srv_count feature
14. Flag: the state of the connection at the time the summary was written (which is usually when the connection terminated). The different states are summarized in the below section.

=====10 additional features=====

1. IDS_detection: reflects whether IDS(Intrusion Detection System) triggered an alert for the connection; ‘0’ means any alerts were not triggered, and an arabic numeral(except ‘0’) means the different kinds of the alerts. Parenthesis indicates the number of the same alert observed during the connection. We used Symantec IDS[3] to extract this feature.
2. Malware_detection: indicates whether malware, also known as malicious software, was observed in the connection; ‘0’ means no malware was observed, and a string indicates the corresponding malware observed at the connection. We used ‘clamav’ software to detect malwares. Parenthesis indicates the number of the same malware observed during the connection.
3. Ashula_detection: means whether shellcodes and exploit codes were used in the connection by using the dedicated software[4]; ‘0’ means no shellcodes and exploit codes were observed, and an arabic numeral(except ‘0’) means the different kinds of the shellcodes or exploit codes. Parenthesis indicates the number of the same shellcode or exploit code observed during the connection.
4. Label: indicates whether the session was attack or not; ‘1’ means the session was normal, ‘-1’ means known attack was observed in the session, and ‘-2’ means unknown attack was observed in the session.
5. Source_IP_Address: indicates the source IP address used in the session. Due to the security concerns, the original IP address on IPv4 was properly sanitized to one of the Unique Local IPv6 Unicast Addresses (private IP addresses)[5]. Also, the same private IP addresses are only valid in the same month: if two private IP addresses are the same within the same month, it means their IP addresses on IPv4 were also the same, but if two private IP addresses are the same within the different month, their IP addresses on IPv4 are also different.
6. Source_Port_Number: indicates the source port number used in the session.
7. Destination_IP_Address: indicates the source IP address used in the session. Due to the security concerns, the original IP address on IPv4 was properly sanitized to one of the Unique Local IPv6 Unicast Addresses (private IP address)[5]. Also, the same private IP addresses are only valid

in the same month: if two private IP addresses are the same within the same month, it means their IP addresses on IPv4 were also the same, but if two private IP addresses are the same within the different month, their IP addresses on IPv4 are also different.

8. Destination_Port_Number: indicates the destination port number used in the session.
9. Start_Time: indicates when the session was started.
10. Duration: indicates how long the session was being established.

Connection State Summaries

- S0: Connection attempt seen, no reply.
- S1: Connection established, not terminated.
- SF: Normal establishment and termination.
- REJ: Connection attempt rejected.
- S2: Connection established and close attempt by originator seen (but no reply from responder).
- S3: Connection established and close attempt by responder seen (but no reply from originator).
- RSTO: Connection established, originator aborted (sent a RST).
- RSTR: Established, responder aborted.
- RSTOS0: Originator sent a SYN followed by a RST, we never saw a SYN ACK from the responder.
- RSTRH: Responder sent a SYN ACK followed by a RST, we never saw a SYN from the (purported) originator.
- SH: Originator sent a SYN followed by a FIN, we never saw a SYN ACK from the responder (hence the connection was “half” open).
- SHR: Responder sent a SYN ACK followed by a FIN, we never saw a SYN from the originator.
- OTH: No SYN seen, just midstream traffic (a “partial connection” that was not later closed).

References

1. The third international knowledge discovery and data mining tools competition dataset KDD99-Cup <http://kdd.ics.uc.i.edu/databases/kddcup99/kddcup99.html>, 1999.
2. Jungsuk Song, Hiroki Takakura and Yasuo Okabe, “Cooperation of Intelligent Honey-pots to Detect Unknown Malicious Codes”, WOMBAT Workshop on Information Security Threat Data Exchange (WISTDE 2008), The IEEE CS Press, Amsterdam, Netherlands, 21-22 April 2008.
3. Symantec Network Security 7100 Series.
4. <http://www.secure-ware.com/contents/product/ashula.html>
5. RFC4193:<http://www.ietf.org/rfc/rfc4193.txt>