

Clustering Malware's Network Behavior using Simple Sequential Features

Azqa Nadeem^{1*}, Carlos H. Gañán¹, and Sicco Verwer¹

¹Delft University of Technology, Delft, The Netherlands

*a.nadeem-2@tudelft.nl

ABSTRACT

Malware attacks are becoming more elusive and powerful. The sheer number of malware samples that are detected on a daily basis is exceeding our capacity to analyze them. In addition, developing malware variants is extremely cheap for attackers because of the availability of various obfuscation tools. These variants can be grouped in malware families based on the similarity of behavioral information retrieved from their system and network activities. Available research suggests that the network traffic generated by malware shows a different aspect of its behavior. This may also be considered as its core behavior since it captures the interaction with its developer. In order to capture such interactions, techniques using Deep Packet Inspection (DPI) are becoming increasingly more common due to their high detection rates. However, DPI has severe privacy implications as it involves inspecting payloads of the network traffic. DPI also does not work well in the presence of encrypted traffic.

We propose an exploratory study, the aim of which is to characterize and cluster malware behavior using high-level, non-privacy-invasive, sequential features extracted from its network activity. The key intuition behind the proposed solution is that if the underlying infrastructure of distinct malware samples is similar, the order in which they carry out their objectives should also be similar. We expect to see this similarity in high-level features that do not invade privacy. To this end, we develop a systematic framework using existing techniques that clusters similar network activities of malware samples using sequences of high-level features (i.e. packet sizes, interval between packets, source and destination port numbers) instead of statistical features.

The results of this research show that sequence clustering allows flexible and robust clusters, as opposed to using statistical features. A detailed cluster content analysis shows that the clusters capture attacking capabilities of malware, such as port scans. They also identify the reuse of the same C&C servers by multiple malware families. In addition, an analysis of distinct network activities exhibited by malware families helps identify unique behaviors associated to particular families. These can be used to extract behavioral signatures for Intrusion Detection Systems. Lastly, clusters also capture common behaviors exhibited by multiple malware families. These common behaviors may suggest author collaboration. The resulting insights obtained from the proposed framework can prove to be valuable threat intelligence information for malware and network analysts.