# Assignment 1

**Name:** Snehal Sudhir Bhole                                                    Student ID: 300318105

**Title:** Real-Time Detection of DNS Exfiltration and Tunneling from Enterprise Networks

**Authors**: Jawad Ahmedy, Hassan Habibi Gharakheili, Qasim Raza, Craig Russelly, and Vijay Sivaraman

**Summary:**   The paper above presents a notion to dynamically identify malicious Domain Name System (DNS) queries or anomalies in Enterprise networks using the Machine learning-based approach. The developed model was validated against high-rank primary domains for a large University and a mid-sized Government Research Institute. Authors used exfiltration tools to inject oddities in the DNS data stream from both organization. Initially, incoming and outgoing internet traffic, DNS query attributes and frequency of queries entering the enterprise host are analyzed. Dataset used is taken from Majestic which ranks sites by the number of subnets linking to that site. Data is trained using the iForest algorithm. The model is tuned by setting tree height, contamination rate and the number of trees whereas the model is evaluated on live 10Gpbs data from two organizations. The ML model was developed, tuned and trained to detect irregularity in DNS queries using a known dataset of non-malicious queries. Then they evaluated the efficacy of the scheme on a 10 Gbps traffic scheme from the borders of two organizations' networks. For both organizations, the mean anomaly score for normal and anomalous queries s 0.44 and 0.59 respectively. Approximately, 1250 DNS queries per second can be processed per second using a Virtual Machine with 4 cores of CPU, 6GB of Memory and 50GB of storage. For Known DNS exfiltration, the approach used here efficiently detects approximately 95.07% and 98.4% for research institute and university campuses respectively.

• **Research Goal**:  The main research goal is to maximize the detection of anomalous queries while reducing the rate of false alarms. The authors are trying to develop and train an ML network and detect queries from enterprise networks from the real-time stream of data.

• **Clarity:** The paper is written in a fluid and concise manner. It is written in an orderly manner explaining the required terminologies and strategies used. The content is easy to understand by a person from a non-technical background as well, but some basic ML knowledge is asset.

• **Related Works:**  Related work is discussed in the paper adequately, identifying the contributions of other authors. Issues with the approaches and benefits of the approach are discussed to some extent but not too detailed. The following approaches from past publications are discussed are Kolmogorov Complexity, Logistic Regression, Frequency Analysis and Anomaly-based Solution.

• **Methods:**  Authors have used isolation Forest(iForest) algorithms to detect abnormal queries in a huge dataset with minimal memory use and time complexities. Selection of a feature and split value from the available ranges of value for the chosen feature is randomly selected by the algorithm. The number of splittings required to separate an instance is the same as the path length from the root to a terminal node. The path length is averaged over such random trees, which conjointly produces smaller lengths for a particular instance which probably is an anomaly. Open-source ML package, scikit-learn and its APIs are used during development. Three tunning parameters used by the algorithm are the number of trees, height

limit of trees and contamination rate. For optimal tuning, the parameter sets the threshold value of the anomaly score to 0.54. 97% accuracy for non-malicious instances and more than 95% accuracy for malicious samples were obtained with 2 trees, tree height at 18 and contamination rate of 2%. The methods used are explained intelligently, following the technical flow of development.

• **Results and Claims:** Authors have used graphical and tabular representations to present their findings along with a precise explanation. Authors aim to enhance the security of enterprise networks and prevent stealing valuable and sensitive data over DNS channels. Authors claim that they have developed, tuned and validated machine learning algorithms to find oddities in DNS queries using a known dataset of the non-malicious domains as ground truth. They have evaluated the efficacy of the system on live 10 Gbps traffic streams from organizations after inserting exfiltration query. For Known DNS exfiltration, the approach used here efficiently detects approximately 95.07% and 98.4% for research institute and university campuses respectively.

• **Support of Results and Claims:** The authors have supported the claim with their observations and anomaly detection table. To support their results, they have validated the model performance on live 10 Gbps data from middle-size Government Research Organizations and a large university. The average anomaly score for instances classified as normal and anomalous is 0.44 and 0.59 respectively in both organizations.

• **Missing Claims and Results:** Precision, Recall or F-measure values for the trained model could also be included. Moreover, learning of model, comparing loss and accuracy of data with each epoch could be represented graphically to understand the model development.

• **Discussion:** Discussion is adequate and the accuracy and real-time performance of the model are detailed but limitations or the issues faced during development are not discussed. FQDN query attributes are used but the usage is not explained thoroughly. The contamination rate parameter used to tune the model is not explained. Its transition from 10% to 2% to increase accuracy is mentioned in the paper.

• **Future Work:** In future work the model developed could be tested against more data with actual infiltrated queries rather than using a query infiltration tool. Since iForest algorithm face issues dealing with multivariate time series information, biases and masking for deviant data, an algorithm that could deal with these efficiently can be used (like Neural Network, etc.).