# ELG5166 Cloud Analytics
# "Real-Time Detection of DNS Exfiltration and Tunnelling from Enterprise Networks"

## Assignment 1

**Instructor:** Dr. Paula Branco

**Student's Name:** Ali Amin El-Sayed Mahmoud El-Sherif          **ID:** 300327246

**Article's Authors:** Jawad Ahmed, Hassan Habibi Gharakheili, Qasim Raza, Craig Russelly and Vijay Sivaraman

# Summary:

Exploitation of channels through DNS channels by cyber attackers has been a headache and an increased liability to enterprises as they pass through their firewalls without inspection. In other words, hidden channels are created allowing breaches to occur through exfiltration and tunneling. This has been avoided through two steps, firstly, development, tuning and training of "isolation forest (iForest)" algorithm to detect anomalies in DNS and through using benign dataset of two enterprise networks, also, the algorithm was tuned using scikit-learn with its APIs. Secondly, work on 10 Gbps of traffic -one week period- from two enterprises network edges and inject more than one million malicious DNS queries generated by "DNS exfiltration toolkit". Three evaluation techniques were applied, one through applying cross validation and testing the accuracy of the trained model of benign instances. Secondly, assessment of the detection rate for malicious DNS queries generated by DNS exfiltration toolkit. Lastly, evaluation of 10 Gbps live traffic streams from the two enterprises. Anomaly detection in both enterprises showed accuracy of 97.99~98.44% on detecting benign instances as "normal" during cross-validation and testing. On the other hand, the model showed false alarms on the same instances of 1.56~2.01% as "anomalous". Generally, the whole project showed promising accuracies of almost 98% in detecting the artificially generated malicious DNS queries at a rate of 1250 query per second (each takes 800 μ sec).

# Critical Review:

Research Goal: To develop, tune and train machine learning model to detect malicious DNS queries' names in real time produced by exfiltration and tunneling in data passing through a DNS. This can be done through several questions, firstly, which machine learning model should be effective in detecting the attacks and how to develop and tune it, second of all, how to get or generate benign and malicious data, then how to test the results of this model, and finally, if the detection rate is sufficient considering the overall data rate of the connection.

Clarity: The paper information was clearly illustrated and the specific short scientific terms with quite good English language was such an interest. The paper is relatively understandable to machine learning learners with no or modest knowledge in networks with few searches for specific terms.

Related Works: A variety of other papers were demonstrated clearly by showing the used approaches like Kolmogorov complexity. This section also discussed how the classification method is not sufficient for live data feed. The in-text citation was clear, correct and well organized in terms.

Methods: Isolation Forest algorithm (iForest) was used since its efficacious for detecting the anomalies in high-dimensional datasets with the low memory and time delay compared to other algorithms. DNS queries' anomalies were detected through identifying attributes for the FQDN of queries. In training, benign data of four days was used. In tuning, 3 parameters were used "n_estimators", "max_samples" and contamination rate. Methods were clearly illustrated, highly related and detailed.

Results and Claims: It was mentioned that iForest is better than one-class SVM and Replicator Neural Network that give higher false alarms. Applying the input "others" resulted an accuracy of 70.57~78.43% in detecting normal queries. Claiming that the system works in real time and can detect malicious DNS queries with good accuracy.

**Support of Results and Claims:** The system did work in real time with detection rate of 1250 DNS queries per second which is higher than the DNS queries rate in both enterprises. The system detected the malicious DNS queries with benign domains as inputs with around 98% accuracy.

**Missing Claims and Results:** The iForest was not proven to be better than one-class SVM and Replicator Neural Network regarding the false alarms. The input "others" size was not mentioned so the false detection is not clear to the reader.

**Discussion:** The research has a good structure with straight forward delivery. It is mostly understandable for machine learning readers with no proper knowledge in the article's context. The research was tangible enough to be considered a strong one, yet, there was no limitations to be mentioned. After applying few modifications, the defined problem can be generalized on many levels as there are a lot of applications for it in real world.

**Future Work:** The authors did not mention future work. Generally, collecting the states for hosts which generate anomalies queries can be a priority to end the mitigate malicious DNS exfiltration and tunneling for good. Advanced deep learning techniques can be used for better system outputs.