

A Machine Learning Approach for IoT Device Identification Based on Packet Features Analysis

Rodrigo Caldas

Master in Electrical and Computer Engineering
Faculty of Engineering, University of Porto
Email: up201708987@fe.up.pt

Tiago Ribeiro

Master in Electrical and Computer Engineering
Faculty of Engineering, University of Porto
Email: up201708988@fe.up.pt

Abstract—IoT device is considered a device that is intended to perform a specific task such as a sense sleep monitor, a security camera or a smart speaker. We are thus surrounded by these devices which makes it easy for any unknown device to connect to the network without being invited.

Therefore, it is essential to know which devices are connected to our network to protect us from any kind of privacy invasion or malware attack.

Our work consists on a machine learning algorithm to identify different IoT Camera devices and brands by analyzing packet characteristics and features present in traffic captures made available by the University of New South Wales (UNSW) [1]. The results obtained during the experimental phase were quite positive, in the order of 96%, which shows that the algorithm can predict with great precision to which type of device the packets belong based on their extracted characteristics.

I. INTRODUCTION

The increased availability and diversity of IoT devices has revolutionized the technological world allowing everyone to take advantage of these new smart friends everywhere[2].

In order to make the network more secure it is essential to know what kind of devices are connected[3].

Let's imagine the case of a company that is used to using a specific supplier of security cameras for the visual control of different office floors.

It is expected that in network traffic report it is only possible to capture the traffic belonging to these known cameras and not to other unknown cameras. If traffic from other types of cameras is captured, we are faced with a vulnerability that can have serious consequences such as invasion of privacy or some type of malware attack[4].

This is why recognition of the type of devices is also important to block access to the network of devices considered vulnerable/unknown[5][6][7].

IoT device recognition can also be used by an observer to discover vulnerable devices by performing passive network traffic analysis. In fact, in a wireless IoT network, it is easier for an observer to capture network traffic to perform additional analysis in order to recognize the type of connected devices, even if the traffic is encrypted because encryption alone does not guarantee adequate protection[8].

The network traffic of IoT devices follows a stable pattern and therefore the generated network traffic is very predictable and well suited for machine learning techniques.

Given this, with this project we implemented a method to recognize IoT devices by analyzing some packet features on the network traffic generated.

To do this, the raw network traffic is processed to extract flows described by characteristics of the packets sent and received by each device so that it is possible to define a pattern corresponding to each device.

II. FEATURES EXTRACTION

From the network traffic captures provided by the University of New South Wales[1], a Python language script was developed which through the MAC address of each chosen device stores in a text file the selected features for each packet found in the traffic of the IoT device[9].

Later these text files are used for the construction of the datasets to be used in the machine training.

The vector of features used to distinguish packet patterns consists of: IPLength, IPHeaderLength, TTL, Protocol, SourcePort, DestPort, SequenceNumber, AckNumber, WindowSize, TCPHeaderLength, TCPLength, TCPStream, TCPUrgentPointer, IPFlags, IPID, IPchecksum, TCPflags and TCPChecksum.

III. EXPERIMENTAL WORK

The work elaborated consists of three main studies where from some intelligent devices conclusions were made about the easiness of standardization of certain traffic and about the correct identification or not of certain devices.

We used a binary classification where 0 corresponds to a certain characteristic and 1 to another (e.g. 0 if the packet belong to a non camera and 1 if the packet belongs to a camera).

The first study consists of a basic identification of devices based on the packet features analysis.

In which the dataset consists of four IoT cameras and four other IoT devices. For this case the test dataset is a 30% split of the training dataset.

The second study is similar to the first one, i.e. also identifies the type of device considering the patterns found in the characteristics of the packets, but this time the test data set is composed of packet features from different IoT devices

(cameras and not cameras) that are not part of the training dataset.

The third and last study done, has as main objective to distinguish different camera brands: to distinguish if a packet belongs to a camera of a specific brand (in this case the chosen brand was Samsung) or not.

One way to corroborate the results obtained was to perform at least three executions in order to make them more consistent[10][11][12].

Another way to corroborate the results obtained was through the confusion matrix[13].

The confusion matrix is a table that allows the visualization of the performance of a classification algorithm[14].

We used TensorFlow which is an open-source library developed by Google for machine learning for training and testing all the studies done[10][11][12].

IV. EXPERIMENTAL RESULTS

A. First Experimental Study

From a training dataset consisting of 8 intelligent devices: 4 different IoT cameras (1 Samsung SmartCam, 1 Insteon Camera, 1 Dropcam and 1 TP-Link Day Night Cloud camera) and 4 IoT devices of different categories (1 WeMo Switch, 1 Smart Portable Speaker, 1 Amazon Echo and 1 Belkin WeMo motion sensor), the test dataset has been elaborated which is a 30% split of the training set.

Remembering that we are dealing with a binary classification where 0 means that a packet belongs to a device different from a camera and 1 represents that a packet comes from a camera. Table 1 shows the results obtained for this experiment where 98% accuracy was obtained.

We can see that in general the machine learning algorithm had a very good performance, being able to distinguish almost perfectly packets from a camera and from a non camera.

In this case, it wrongly predicted 3 packets from cameras when in reality they were from not cameras and also wrongly predicted 28 from not cameras when in reality they were from cameras.

On the other hand, it correctly predicted 1004 packets from non cameras and 1426 from cameras.

TABLE I
CONFUSION MATRIX OF THE FIRST STUDY.

Actual Class	Predicted Class	
	0	1
0	1004	3
1	28	1426

B. Second Experimental Study

From the training data set used in the previous study (Section A), a test dataset containing packets belonging to different devices (cameras and not cameras) that were not included in the training set was used this time to check if there really was a viable pattern in the registered characteristics of

the packets belonging to these different types of devices.

Table 2 shows the results obtained for this experiment where an accuracy of 95% was obtained.

We can see that in general the machine has performed well, being able to effectively distinguish the characteristics of the packets of a camera from a non camera.

In this case, we are dealing with a much more voluminous set of tests compared to the previous study, where 9845 packets from cameras were wrongly predicted when in reality they were from non cameras and also wrongly predicted 15005 packets from not cameras when in reality they were from cameras.

On the opposite, it correctly predicted 125260 from non cameras and 314342 from cameras.

TABLE II
CONFUSION MATRIX OF THE SECOND STUDY.

Actual Class	Predicted Class	
	0	1
0	125260	9845
1	15005	314342

C. Third Experimental Study

The third case study aims to identify different camera brands, i.e. to distinguish whether a packet belongs to a camera of a specific brand (in this case the specific brand chosen was Samsung) or not.

The binary classification assigned a 0 if a packet belongs to a camera that is not Samsung and a 1 if a packet comes from a camera that is Samsung.

With this it will be possible to verify if there is a stable pattern that is able to identify, characterize and distinguish cameras based on their packet features (Chapter II).

In table 3 it is possible to observe the experimental results that demonstrate a good performance in the distinction with 94% accuracy.

Thus it was possible to wrongly predict 3523 packets from cameras that were not Samsung and surprisingly 0 packets from Samsung cameras were wrongly predicted.

On the other hand, were correctly predicted 27339 packets from cameras that were not Samsung and 23856 packets from Samsung cameras.

The high performance in packet identification of Samsung cameras is mainly due to strong standards in the characteristics of the packets assessed that have unique aspects compared with non-Samsung cameras, for example, Samsung cameras have for the same features much longer information compared with cameras of other brands.

TABLE III
CONFUSION MATRIX OF THE THIRD STUDY.

Actual Class	Predicted Class	
	0	1
0	27339	3523
1	0	23856

V. POSSIBLE PRACTICAL APPLICATION

A possible practical application of this developed project would be to integrate this detection module in the firewall of the command center of the surveillance camera. The objective of this firewall is to detect if the packets arriving at the command center are really from the surveillance system's equipment, that is, to detect that they have not been adulterated, as well as to detect if there is any kind of cloning of the equipment. Being the last point of defense of the surveillance center, since it requires computational effort, it should be used in parallel with the control and verification of IP and MAC addresses of the equipment because these two addresses are easy to counteract with their spoofing. This is why it is increasingly important to integrate such systems in these centers to control and secure traffic.

VI. CONCLUSION

All the work done shows that, in this case, these intelligent devices (IoT cameras) can be precisely identified based on the patterns of the features of the packets generated by the devices. Our method can classify camera devices very precisely, including brand recognition, which makes our method feasible and usable at company level in the future. We believe our method can be used to accurately identify the connection of a camera device to a company's computer network, thereby minimizing possible violations of company policies by providing a higher level of security. As future work, this project can be adapted to other types of intelligent devices, as well as different types of features can be used to identify traffic patterns on the devices, such as at a deeper level the content of the packets received and sent by the devices.

REFERENCES

- [1] Iot security unsw sydney. <https://iotanalytics.unsw.edu.au/>.
- [2] Tai-hoon Kim, Carlos Ramos, and Sabah Mohammed. Smart city and iot, 2017.
- [3] R. Mahmoud, T. Yousuf, F. Aloul, and I. Zuolkernan. Internet of things (iot) security: Current status, challenges and prospective measures. In *2015 10th International Conference for Internet Technology and Secured Transactions (ICITST)*, pages 336–341, 2015.
- [4] I. Andrea, C. Chrysostomou, and G. Hadjichristofi. Internet of things: Security vulnerabilities and challenges. In *2015 IEEE Symposium on Computers and Communication (ISCC)*, pages 180–187, 2015.
- [5] Yair Meidan, Michael Bohadana, Asaf Shabtai, Juan David Guarnizo, Martín Ochoa, Nils Ole Tippenhauer, and Yuval Elovici. Profiliot: A machine learning approach for iot device identification based on network traffic analysis. In *Proceedings of the Symposium on Applied Computing*, SAC '17, page 506–509, New York, NY, USA, 2017. Association for Computing Machinery.
- [6] A. Aksoy and M. H. Gunes. Automated iot device identification using network traffic. In *ICC 2019 - 2019 IEEE International Conference on Communications (ICC)*, pages 1–7, 2019.
- [7] Ola Salman, Imad H. Elhadj, Ali Chehab, and Ayman Kayssi. A machine learning based framework for iot device identification and abnormal traffic detection. *Transactions on Emerging Telecommunications Technologies*, n/a(n/a):e3743. e3743 ETT-19-0273.R1.
- [8] Noah Apthorpe, Dillon Reisman, and Nick Feamster. A smart home is no castle: Privacy vulnerabilities of encrypted iot traffic, 2017.
- [9] Gerald Combs et al. Wireshark-network protocol analyzer. *Version 0.99*, 5, 2008.
- [10] Tensorflow learn. <https://www.tensorflow.org/learn>.
- [11] Tensorflow tutorials. <https://www.tensorflow.org/tutorials>.
- [12] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 265–283, Savannah, GA, November 2016. USENIX Association.
- [13] Sofia Visa, Brian Ramsay, Anca Ralescu, and Esther Knaap. Confusion matrix-based feature selection. volume 710, pages 120–127, 01 2011.
- [14] Confusion matrix. https://en.wikipedia.org/wiki/Confusion_matrix.