



HAL
open science

Amélioration de l'identification du type des objets connectés par classification supervisée

Nesrine Ammar, Ludovic Noirie, Sébastien Tixeul

► To cite this version:

Nesrine Ammar, Ludovic Noirie, Sébastien Tixeul. Amélioration de l'identification du type des objets connectés par classification supervisée. CORES2019 - Rencontres Francophones sur la Conception de Protocoles, l'Évaluation de Performance et l'Expérimentation des Réseaux de Communication, Jun 2019, Jardins de Saint Benoît, France. hal-02126555

HAL Id: hal-02126555

<https://hal.science/hal-02126555v1>

Submitted on 11 May 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Amélioration de l'identification du type des objets connectés par classification supervisée

Nesrine Ammar^{1,2}, Ludovic Noirie¹ et Sébastien Tixeuil^{2 †}

¹Nokia Bell Labs, Nokia Paris Saclay, route de Villejust, 91620 Nozay, France

²Sorbonne Université, Laboratoire d'Informatique de Paris 6, LIP6, F-75005, Paris, France

L'identification des objets connectés est primordiale pour aider les utilisateurs à mieux les gérer et assurer la sécurité de leurs réseaux domestiques. Nous traitons ce problème en utilisant une classification supervisée basée sur les attributs du trafic réseau des objets ainsi que les données textuelles transportées par les couches applicatives. Nos résultats montrent une précision de notre solution égale à 99% en moyenne.

Mots-clefs : Internet des objets, identification des objets connectés, trafic réseau, classification supervisée.

1 Introduction

Avec la croissance de l'internet des objets (*Internet of Things*, IoT) [ARC18], les gens achètent de plus en plus des objets connectés afin de bénéficier de nouveaux services leur permettant d'être informés et de contrôler leur maison à distance n'importe où et n'importe quand. Pour aider à la gestion de tous ces objets connectés dans un réseau domestique, nous avons besoin d'identifier leurs types, ce qui constitue un défi majeur à cause du nombre d'objets connectés et de leur hétérogénéité.

Dans notre travail précédent [ANT18], cette identification était basée sur les informations partagées par différents protocoles tel que SSDP, MDNS et DHCP. Nous obtenions une signature unique pour chaque objet testé. Malheureusement, certains objets ne partagent aucune information. De plus, lorsque le nombre d'objets augmente, le nombre de signatures augmente aussi, ce qui rend la gestion des signatures difficile en cas d'évolution et de mise à jour. Pour ces raisons, nous proposons d'automatiser et d'améliorer l'identification du type des objets connectés en utilisant les techniques d'apprentissage supervisé. Les contributions de cet article concernent la définition, l'implémentation et l'évaluation d'une telle solution.

Certains travaux ont déjà utilisé des techniques d'apprentissage automatique pour identifier les types des objets connectés. Les auteurs de [MMH⁺17] basent l'apprentissage sur les caractéristiques des flux uniquement. Ils obtiennent une précision de 95% pour 17 objets mais seulement autour de 50% pour 10 autres objets. Notre solution présentée dans cet article présente de bien meilleurs résultats. Les auteurs de [MBS⁺17] et [SBZD18] utilisent les attributs des sessions TCP pour classer les objets connectés. Le nombre d'objets est très limité pour entraîner leurs modèles, et cette technique exige la capture de trafic pendant une longue durée, incompatible avec une détection rapide du type des objets connectés. Notre solution n'a besoin que des premiers paquets émis et reçus par l'objet connecté dont on veut identifier le type, et nous l'avons évaluée sur un plus grand nombre d'objets.

2 Méthodologie et sélection des attributs sélectionnés

Notre solution extrait les données pertinentes du trafic émis et reçu par un nouvel objet connecté pour constituer un vecteur d'attributs qui nourrit un système de classification pour identifier le type de cet objet.

Attributs sélectionnés :

Nous analysons les premiers paquets émis et reçus par chaque objet. Un nouvel objet est détecté lorsqu'une nouvelle adresse MAC est observée. Deux sortes d'attributs sont extraits du trafic de l'objet :

[†]Ce travail a été partiellement réalisé dans le cadre du LINC'S (*Laboratory of Information, Networking and Communication Science*, <http://www.lincs.fr/>).

1. **Les attributs de flux de trafic** – Un flux est une succession de paquets entre une source et une destination caractérisée par le 4-tuplet (IP source, IP destination, port source, port destination). Les attributs considérés sont des données statistiques sur ces flux (valeurs numériques) ou des données de types de flux (valeurs binaires pour « oui/non » sur chaque type) :
 - (a) Tailles des paquets (moyenne, maximum et minimum) ;
 - (b) Temps moyen des inter-arrivées des paquets d'un flux ;
 - (c) Taille du flux mesurée en nombre de paquets ;
 - (d) Protocoles utilisés (HTTP, HTTPS, SSDP, mDNS, TFTP, DHCP, DNS, NTP, BOOTP, TCP, UDP).
2. **Les attributs textuels** – Certaines données peuvent être extraites des messages en clair au niveau applicatif. Ces données textuelles contiennent des mots clés caractérisant l'objet connecté. Tout d'abord, les données sont nettoyées en éliminant les symboles et les chiffres isolés, afin d'obtenir un ensemble de mots uniques. Ensuite, on génère un vecteur binaire à partir de l'ensemble des mots uniques extraits dont la valeur est définie à 1 si le mot existe dans la description de l'objet, 0 sinon. Les données utilisées sont, lorsqu'elles sont disponibles, celles utilisées dans [ANT18] :
 - (a) Nom du fabricant à partir de l'adresse MAC ;
 - (b) Nom de l'objet à partir de DHCP ;
 - (c) Noms du fabricant, du modèle, de l'objet et du type à partir de la description XML partagée par le protocole de découverte UPnP ;
 - (d) Nom local et services offerts par l'objet à partir des enregistrements du protocole mDNS ;
 - (e) Système d'exploitation, modèle et type à partir de l'agent utilisateur inclus dans la requête HTTP.

Méthodologie de classification et algorithmes de classification :

Les vecteurs d'attributs obtenus nous permettent d'utiliser une classification supervisée. Notre objectif est d'obtenir un mappage entre le vecteur d'attributs généré et le type d'objet correspondant, en utilisant un algorithme de classification. Afin de choisir le meilleur algorithme pour ce problème, nous testons les 5 algorithmes connus suivants : arbre de décision (DT), machines à vecteurs de support (SVM), classification naïve bayésienne (NB), forêts aléatoires (RF) et K plus proches voisins (K-NN).

Ces algorithmes sont entraînés sur un jeu de données d'apprentissage puis testés sur un second jeu de données de tests. Ils donnent une probabilité d'appartenance à chacune des classes. Un objet est considéré de classe inconnue si plusieurs classes ont des probabilités élevées proches ou si toutes les classes ont des probabilités inférieures à un seuil fixé à 50%.

3 Implémentation de la solution

Nous avons collecté le trafic des objets connectés à une passerelle résidentielle afin de construire une partie de nos données d'apprentissage et de test. Pour cela nous avons implémenté un prototype dans un environnement laboratoire afin de collecter le trafic des objets et tester notre solution.

Configuration matérielle du réseau domestique expérimental :

Nous utilisons une passerelle résidentielle constituée d'un routeur Wi-Fi connecté à un commutateur Ethernet via une connexion filaire. Nous connectons un ordinateur d'écoute à un port d'écoute du commutateur afin de regarder tout le trafic qui passe par le commutateur. Pour capturer le trafic des objets sans fil, nous connectons au commutateur Ethernet un point d'accès Wi-Fi transparent. La capture de trafic dure une minute après la détection du premier paquet concernant un nouvel objet connecté afin de s'assurer d'avoir un ensemble suffisant de paquets émis ou reçus par cet objet en début de connexion.

Implémentation logicielle de notre solution :

Notre solution est implémentée dans l'ordinateur d'écoute, mais elle pourrait l'être dans la passerelle résidentielle pour une version industrielle. Elle comprend deux modules, représentés sur la figure 1.

Le premier module de capture est implémenté en *node.js*. Il est composé d'un sous-module de capture de trafic utilisant *Win10Pcap*, un sous-module gérant les objets connectés et identifiant les nouvelles connexions, et un sous-module servant les requêtes de l'assistant de classification (HTTP/REST).

Le second module est l'assistant de classification d'objets connectés. Il est développé en *python* et il interagit avec le module de capture afin de recevoir les traces des premiers paquets émis par un nouvel

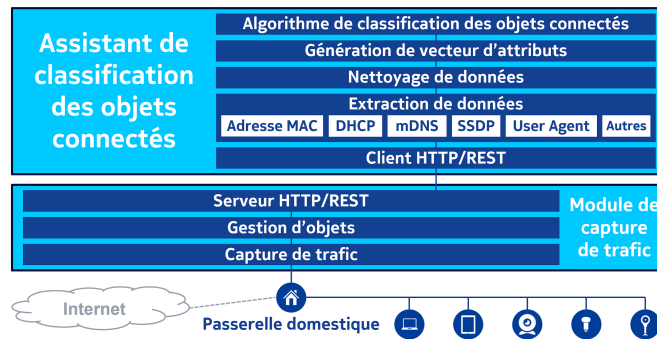


FIGURE 1: Architecture de l'assistant d'identification des objets connectés dans le réseau domestique.

objet connecté. Les différents sous-modules au-dessus permettent d'analyser les captures afin d'extraire les attributs définis dans la sous-section 2 qui sont ensuite nettoyés pour générer un vecteur d'attributs. Les vecteurs d'attributs sont utilisés pour l'apprentissage des classificateurs et l'évaluation de notre solution.

4 Évaluation de la solution

Données récoltées et utilisées pour l'évaluation :

Nous avons collecté le trafic d'objets connectés de notre laboratoire en utilisant notre implémentation de la section 3. Ces objets sont variés et représentatifs d'un réseau domestique, avec connexion Wi-Fi et Ethernet. Ces données forment une matrice creuse dont la taille de vecteur d'attributs est 83 et le nombre d'instances total pour tous les objets utilisés est 12000. Nous avons complété ces données avec des données de 27 objets venant des travaux de [MMH⁺ 17], pour obtenir un nombre total d'objets connectés de 33.

Avant de démarrer l'apprentissage, nous séparons nos données en deux ensembles. Le premier correspond à 70% des données de tous les objets et est utilisé pour l'apprentissage afin d'assurer un mappage entre les vecteurs d'attributs et l'objet. Les 30% restantes sont utilisées pour tester les modèles sur des instances non existantes dans l'ensemble d'apprentissage. Cette répartition est communément utilisée dans la littérature.

Évaluation des différents algorithmes de classification :

Les algorithmes de classification listés dans la section 2 possèdent des hyper-paramètres. Nous utilisons la méthode de validation croisée afin de fixer les meilleurs valeurs de chacun d'entre eux. Une fois les algorithmes entraînés et validés, nous évaluons leurs performances sur l'ensemble de données de tests, les 30% restantes qui n'ont pas servi à l'apprentissage, en se basant sur la métrique précision de classification. Les résultats montrent que l'algorithme par arbre de décision (DT) a une meilleure performance que les autres avec une précision de 99% en moyenne. SVM retourne la précision la plus faible 88%. Enfin, NB, K-NN et RF montre une précision moyenne égale à 98%, 94% et 94% respectivement. Pour la suite, nous considérons donc l'algorithme DT.

Intérêt de combiner les données textuelles et les données de flux :

Les résultats précédents ont été obtenus avec l'ensemble des attributs sélectionnés dans la section 2. Si nous retirons les attributs textuels et ne conservons que les attributs de flux de trafic, nous obtenons une précision égale à 72% seulement : les attributs textuels sont donc nécessaires pour une bonne classification. Si nous retirons les attributs de trafic et ne conservons que les attributs textuels comme ceux mentionnés dans [ANT18], seuls les objets ayant ce genre d'attributs peuvent être identifiés, soit 93% des objets.

Analyse détaillée pour la classification par arbre de décision (DT) avec tous les attributs :

Nous utilisons la matrice de confusion pour l'évaluation du mappage entre les attributs et les modèles d'apprentissage obtenus pour chaque classe par l'algorithme DT, sur l'ensemble de données de tests. Elle indique le nombre de fois où des instances d'une classe A sont classées en tant que classe B , par exemple le nombre de fois que le classificateur a confondu le type caméra D-Link avec un capteur D-Link. Cela nous permet de calculer notamment les critères de mesure des performances *précision*

$$= \frac{\text{vrais positifs}}{\text{vrais positifs} + \text{faux positifs}} \text{ et } \text{rappel} = \frac{\text{vrais positifs}}{\text{vrais positifs} + \text{faux négatifs}} \text{ pour chaque type d'objet.}$$

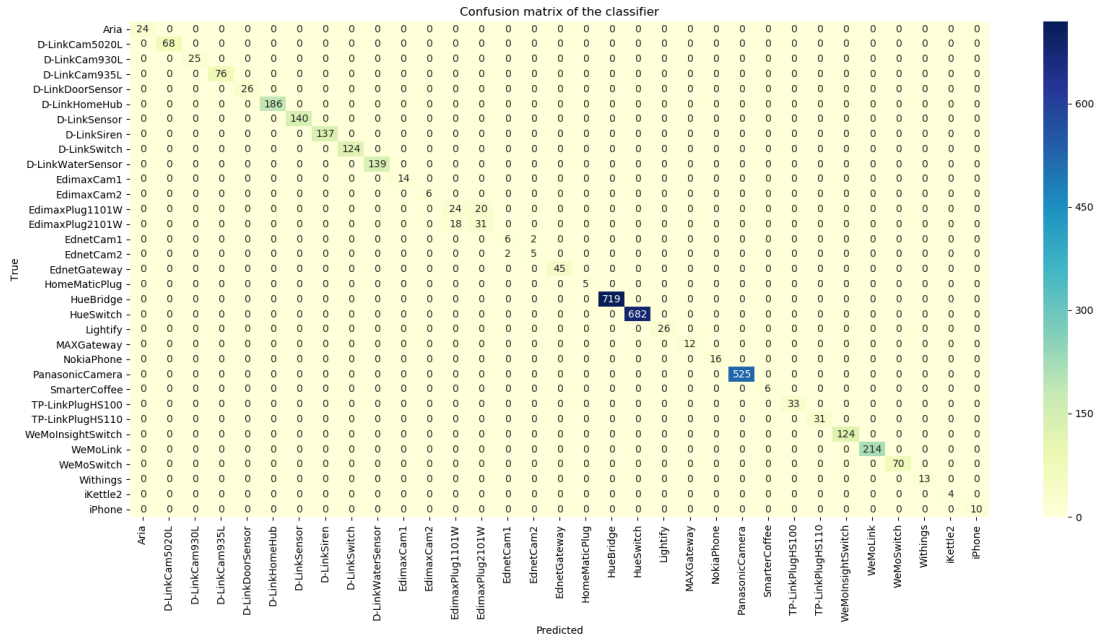


FIGURE 2: Matrice de confusion.

Cette matrice de confusion, représentée dans la figure 2, illustre l’excellente performance de DT avec l’ensemble des attributs utilisés. Nous remarquons uniquement des confusions sur le type d’objet pour des objets venant du même vendeur mais d’un modèle différent : entre EdimaxPlug1101W et Edimax-Plug2101W, et entre Ednetcam1 et Ednetcam2. Cela explique leurs faibles valeurs de rappel et précision autour de 50% pour ces objets.

5 Conclusion

Dans cet article, nous avons traité le problème de discrimination du type des objet connectés en utilisant des techniques de classification supervisée se basant sur leurs attributs de flux de trafic et leurs attributs textuels. Nos résultats montrent que l’algorithme de classification par arbre de décision a de meilleures performances que les autres algorithmes avec une précision égale à 99% en moyenne. Nos résultats obtenus sont basés sur la classification de 33 objets connectés. Nous souhaitons consolider notre approche en testant un plus grand nombre d’objets. Dans les travaux futurs, utilisant notre implémentation, nous allons donc augmenter notre jeu de données avec d’autres objets. Nous allons aussi investiguer l’identification du types d’objets autres que les objets IP, par exemple Bluetooth.

Références

- [ANT18] N. Ammar, L. Noirie, and S. Tixeuil. Identification du type des objets connectés par les informations des protocoles réseaux. In *Rencontres Francophones sur la Conception de Protocoles, l’Évaluation de Performance et l’Expérimentation des Réseaux de Communication*, Roscoff, France, May 2018.
- [ARC18] M. Ammar, G. Russello, and B. Crispo. Internet of things : A survey on the security of IoT frameworks. *Journal of Information Security and Applications*, 38 :8 – 27, 2018.
- [MBS⁺17] Y. Meidan, M. Bohadana, A. Shabtai, J. D. Guarnizo, M. Ochoa, and N. O. Tippenhauer and Y. 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*, pages 506–509, New York, NY, USA, 2017. ACM.
- [MMH⁺17] M. Miettinen, S. Marchal, I. Hafeez, N. Asokan, A. Sadeghi, and S. Tarkoma. IoT SENTINEL : Automated device-type identification for security enforcement in IoT. In *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)*, pages 2177–2184, June 2017.
- [SBZD18] M. R. Shahid, G. Blanc, Z. Zhang, and H. Debar. IoT devices recognition through network traffic analysis. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 5187–5192, Dec 2018.