

Gaining deep knowledge of Android malware families through dimensionality reduction techniques

Rafael Vega Vega^a, Héctor Quintián^{a,*}, José Luis Calvo-Rolle^a, Álvaro Herrero^b, Emilio Corchado^c

^a*University of A Coruña,
Department of Industrial Engineering,
Avda. 19 de febrero s/n, 15405, Ferrol, A Coruña, Spain*

^b*Universidad de Burgos,
Department of Civil Engineering,
University of Burgos, Burgos, Spain*

^c*Departamento de Informática y Automática,
Universidad de Salamanca,
Plaza de la Merced, Salamanca, Spain*

How to cite: Rafael Vega Vega, Héctor Quintián, José Luis Calvo-Rolle, Álvaro Herrero, Emilio Corchado, Gaining deep knowledge of Android malware families through dimensionality reduction techniques, Logic Journal of the IGPL, Volume 27, Issue 2, April 2019, Pages 160–176, <https://doi.org/10.1093/jigpal/jzy030>

Abstract

This research proposes the analysis and subsequent characterization of Android malware families, by means of low dimensional visualizations using dimensional reduction techniques. The well-known Malgenome dataset, coming from the Android Malware Genome Project, has been thoroughly analysed through six dimensionality reduction techniques: Principal Component Analysis, Maximum Likelihood Hebbian Learning, Cooperative Maximum Likelihood Hebbian Learning, Curvilinear Component Analysis, Isomap and Self Organizing Map. Results obtained enable a clear visual analysis of the structure of this high-dimensionality dataset, letting us gain deep knowledge about the nature of such Android malware families. Interesting conclusions are obtained from the real-life dataset under analysis.

Keywords: Android malware, malware families, dimensionality reduction, artificial neural networks

1. Introduction

Since the first smartphones came onto the market in the late 1990s, sales on that sector have increased constantly to the present-day. Among all the available operating systems, Google's Android has been, and increasingly is, the most

*Corresponding author

Email addresses: rafael.alejandrov.vega.vega@udc.es (Rafael Vega Vega), hector.quintian@udc.es (Héctor Quintián), jlcalvo@udc.es (José Luis Calvo-Rolle), acosio@ubu.es (Álvaro Herrero), escorchado@usal.es (Emilio Corchado)

popular mobile platform [1]. The number of Android units sold in Q1 2017 worldwide raised to 379.98 million out of 432.79 million units, that is a share of 87.79%. It is not only the number of devices but also the number of apps; those available at Google Play (Android’s official store) constantly increase, up to more than 3.4 million that are available nowadays [2]. With regard to the security issue, the number of malicious Android apps has greatly risen in the last four years; from the half million of them that were identified in 2013 to the nearly 3.5 million in 2017 [3]. Furthermore, it has been forecast that increase in malware for Android devices is expected to continue [3, 4]. This operating system is an appealing target for bad-intentioned apps, mainly because of its open mentality, in contrast to iOS or some other operating systems.

Smartphone security and privacy still are nowadays major concerns although great efforts have been devoted over past years [5]. In order to address these issues, it is required to understand the malware and its nature. Otherwise, it will not be possible to practically develop an effective solution [6]. According to this idea of gaining deeper knowledge about malware nature, present study is focused on the analysis of Android malware families. To do so, Malgenome (a real-life publicly-available) dataset [7] has been analyzed by means of several Dimensionality Reduction Techniques (DRTs). From the samples contained in such dataset, several alarming statistics were found [6], that motivate further research on Android malware. That is the case of the 36.7% of the collected samples that leverage root-level exploits to fully compromise the security of the whole system or the fact that more than 90% of the samples turn the compromised phones into a botnet controlled through network or short messages.

To characterize malware families, present study proposes a comprehensive comparison of many DRTs, that are able to visualize a high-dimensionality dataset (further described in section 2), to gain deep knowledge of Android malware families. Each individual from the Malgenome dataset (a malware app) encodes the subset of selected features by using a binary representation (details on section 3). These individuals are grouped by families and then visualized trying to identify patterns that exist across dimensional boundaries in the high dimensional dataset by changing the spatial coordinates of malware family data. The main goal is to obtain an intuitive visualization of the malware families to draw conclusions about the structure of the dataset and to characterize malware families subsequently.

Neural networks have been applied to a wide variety of fields in recent decades [8],[9],[10],[11],[12]; additionally, neural DRTs have been previously applied to massive security datasets, such as those generated by network traffic [13], [14], SQL code [15], [16], honeynets [17], and HTTP traffic [18]. In present paper, such methods are applied to a new problem, related to the characterization and knowing of malware families. On the other hand, several different techniques have been used to differentiate between legitimate and malicious Android apps, such

as machine learning [19],[20],[21], knowledge discovery [22], and weighted similarity matching of logs [23], among others as well as hybridization approaches [24]. Although some visualization techniques have been applied to the detection of malware in general terms [25], few dimensionality-reduction proposals for Android malware detection are available at present time. In [26] Pythagoras tree fractal is used to visualize the malware data, being all apps scattered, as leaves in the tree. Authors of [27] proposed graphs for deciding about malware by depicting lists of malicious methods, needless permissions and malicious strings. In [28], visualization obtained from biclustering on permission information is described. Behavior-related dendrograms are generated out of malware traces in [29], comprising nodes related to the package name of the application, the Android components that has called the API call and the names of functions and methods invoked by the application. Unlike previous work, Android malware families (instead of malware apps) are visualized by DRTs in present paper. Up to the authors knowledge, this is the first time that dimensionality-reduction models are applied to visualize Android malware.

The rest of this paper is organized as follows: the applied neural methods are described in section 2, the setup of experiments for the Android Malware Genome dataset is described in section 3, together with the results obtained and the conclusions of the study that are stated in section 4.

2. Dimensionality Reduction Techniques

This work proposes the application of several DRTs for the visualization of Android malware data. Visualization techniques are considered a viable approach to information seeking, as humans are able to recognize different features and to detect anomalies by means of visual inspection [30]. The underlying operational assumption of the proposed approach is mainly grounded in the ability to render the high-dimensional traffic data in a consistent yet low-dimensional representation [17], [18], [25]. In most cases, security visualization tools have to deal with massive datasets with a high dimensionality, to obtain a low-dimensional space for presentation [13], [15], [17], [18], [31], [32].

This problem of identifying patterns that exist across dimensional boundaries in high dimensional datasets can be solved by changing the spatial coordinates of data. Projection methods project high-dimensional data points onto a lower dimensional space in order to identify "interesting" directions in terms of any specific index or projection. Having identified the most interesting projections, the data are then projected onto a lower dimensional subspace plotted in two or three dimensions, which makes it possible to examine the structure with the naked eye [30].

2.1. Principal Component Analysis

Principal Component Analysis (PCA) is a well-known statistical model, introduced in [33], that describes the variation in a set of multivariate data in terms of a set of uncorrelated variables each, of which is a linear combination of the original variables. From a geometrical point of view, this goal mainly consists of a rotation of the axes of the original coordinate system to a new set of orthogonal axes that are ordered in terms of the amount of variance of the original data they account for.

PCA can be performed by means of neural models such as those described in [34] or [35]. It should be noted that even if we are able to characterize the data with a few variables, it does not follow that an interpretation will ensue.

2.2. Maximum Likelihood Hebbian Learning

Maximum Likelihood Hebbian Learning [30] which is based on Exploration Projection Pursuit (EPP). The statistical method of EPP [30], [36], [37] was designed for solving the complex problem of identifying structure in high dimensional data by projecting it onto a lower dimensional subspace in which its structure is searched for by eye. To that end, an “index” must be defined to measure the varying degrees of interest associated with each projection. Subsequently, the data is transformed by maximizing the index and the associated interest. From a statistical point of view the most interesting directions are those that are as non-Gaussian as possible.

2.3. Cooperative Maximum Likelihood Hebbian Learning

The Cooperative MLHL (CMLHL) model [38] extends the MLHL model, by adding lateral connections between neurons in the output layer of the model. Considering an N-dimensional input vector (x), and an M-dimensional output vector (y), with W_{ij} being the weight (linking input neuron j to output neuron i), then CMLHL can be expressed as defined in equations 1-4.

1. Feed-forward step:

$$y_i = \sum_{j=1}^N W_{ij}x_j, \forall i \quad (1)$$

2. Lateral activation passing:

$$y_i(t+1) = [y_i(t) + \tau(b - Ay)]^+ \quad (2)$$

3. Feedback step:

$$e_j = x_j - \sum_{i=1}^M W_{ij}y_j, \forall j \quad (3)$$

4. Weight change

$$\Delta W_{ij} = \eta \cdot y_i \cdot \text{sign}(e_j)|e_j|^{p-1} \quad (4)$$

Where: η is the learning rate, τ is the “strength” of the lateral connections, b the bias parameter, p a parameter related to the energy function and A a symmetric matrix used to modify the response to the data. The effect of this matrix is based on the relation between the distances separating the output neurons.

2.4. ISOMAP Algorithm

ISOMAP nonlinear DRT [39] attempts to preserve pairwise geodesic (or curvilinear) distance between data points. Geodesic distance is the distance between two points measured over the manifold. ISOMAP defines the geodesic distance as the sum of edge weights along the shortest path between two nodes. The doubly-centered geodesic distance matrix K in ISOMAP is of the form given by equation 5.

$$K = \frac{1}{2}HD^2H \quad (5)$$

Where $D^2 = D_{ij}^2$ means the element wise square of the geodesic distance matrix $D = [D_{ij}]$, and H is the centring matrix, given by equation 6.

$$H = I_n - \frac{1}{N}e_N e_N^T \quad (6)$$

In which $e_N = [1\dots 1]^T \in R^N$

The top N eigenvectors of the geodesic distance matrix represent the coordinates in the new n -dimensional Euclidean space.

2.5. Curvilinear Component Analysis Algorithm

Curvilinear Component Analysis (CCA) [40], [41] is a non-linear projection method that preserves distance relationships in both input and output spaces. CCA is a useful method for redundant and non-linear data structure representation and can be used in dimensionality reduction. CCA is useful with highly non-linear data, where PCA or any other linear method fails to give suitable information.

CCA brings some improvements to other methods like Sammon’s Mapping [42], although when unfolding a nonlinear structure, Sammon’s Mapping cannot reproduce all distances. One way to get round this problem consists in favoring local topology: CCA tries to reproduce short distances firstly, long distances being secondary. Formally, this reasoning led to the following error function (without normalization) defined in equation 7.

$$E_{CCA} = \sum_{i,j=1}^N (d_{i,j}^n - d_{i,j}^p)^2 F_\lambda(d_{i,j}^p) \quad (7)$$

In comparison with E_{Sammon} , E_{CCA} has an additional weighting function F depending on d_{ij}^p and on parameter λ . The F factor is a decreasing function of its argument, so it is used to favour local topology preservation. For example, F could be a step function of $(\lambda - d)$.

2.6. Self Organizing Maps

Among the great variety of tools for multidimensional data visualization, several of the most widely used are those belonging to the family of the topology preserving maps [43],[44],[45],[46],[47],[48]. Probably the best known among these algorithms is the Self-Organizing Map (SOM) [43], [45], [49], [50]. It is based on a type of unsupervised learning called competitive learning; an adaptive process in which the units in a neural network gradually become sensitive to different input categories or sets of samples in a specific domain of the input space. The main feature of the SOM algorithm is its topology preservation. When not only the winning unit, but also its neighbors on the lattice are allowed to learn, neighboring units gradually specialize to represent similar inputs, and the representations become ordered on the map lattice.

An input vector (x) is presented to the network and the node of the network in which the weights (W_i) are closest (in terms of Euclidean distance) to x , is chosen:

$$c = \operatorname{argmin}(\|x - W_i\|) \quad (8)$$

The weights of the winning node and the nodes close to it are then updated to move closer to the input vector. There is also a learning rate parameter that usually decreases as the training process progresses. The weight update rule for inputs is defined as follows:

$$\Delta W_i = \eta h_{ci}(x - W_i), \quad \forall i \in N \quad (9)$$

Where, W_i is the weight vector associated with neuron i , x is the input vector, and h is the neighborhood function.

3. Experiments & Results

As previously mentioned, several different DRTs (see Section 2) have been applied to analyze Android malware. Present section introduces the analyzed dataset as well as the main obtained results.

3.1. Malgenome Dataset

The Malgenome dataset [6], coming from the Android Malware Genome Project [7], has been analysed in present study. It is the first large collection of Android malware (1,260 samples) that was split in malware families (49 different ones). It covered the majority of existing Android malware, collected from the beginning of the project in August 2010.

Data related to many different apps from a variety of Android app repositories were accumulated over more than one year. Additionally, malware apps were thoroughly characterized based on their detailed behavior breakdown, including the installation, activation, and payloads.

Collected malware was split in families, that were obtained by “carefully examining the related security announcements, threat reports, and blog contents from existing mobile antivirus companies and active researchers as exhaustively as possible and diligently requesting malware samples from them or actively crawling from existing official and alternative Android Markets” [6]. The defined families are: *ADRD*, *AnserverBot*, *Asroot*, *BaseBridge*, *BeanBot*, *Bg-Serv*, *CoinPirate*, *Crusewin*, *DogWars*, *DroidCoupon*, *DroidDeluxe*, *DroidDream*, *DroidDreamLight*, *DroidKungFu1*, *DroidKungFu2*, *DroidKungFu3*, *DroidKungFu4*, *DroidKungFuSapp*, *DoidKungFuUpdate*, *Endofday*, *FakeNetflix*, *FakePlayer*, *GamblerSMS*, *Geinimi*, *GGTracker*, *GingerMaster*, *GoldDream*, *Gone60*, *GPSSMSSpy*, *HippoSMS*, *Jifake*, *jSMShider*, *Kmin*, *Lovetrapp*, *NickyBot*, *Nickyspy*, *Pjapps*, *Plankton*, *RogueLemon*, *RogueSPPush*, *SMSReplicator*, *SndApps*, *Spitmo*, *TapSnake*, *Walkinwat*, *YZHC*, *zHash*, *Zitmo*, and *Zsone*. Samples of 14 of the malware families were obtained from the official Android market, while samples of 44 of the families came from unofficial markets. As some families are present in both markets (official and unofficial), the final dataset to be analysed consists of 49 samples (one for each family) and each sample is described by 26 different features derived from a study of each one of the apps. The features are divided into six categories, as can be seen in Table 4.1.

Tabla 1: Features describing each one of the malware families in the Malgenome dataset.

Category #1: Installation	Category #3: Privilege escalation
1 Repackaging	14 exploit
2 Update	15 RATC/zimperlich
3 Drive-by download	16 ginger break
4 Standalone	17 asroot
Category #2: Activation	18 encrypted
5 BOOT	Category #4: Remote control
6 SMS	19 NET
7 NET	20 SMS
8 CALL	Category #5: Financial charges
9 USB	21 phone call
10 PKG	22 SMS
11 BATT	23 block SMS
12 SYS	Category #6: Personal information stealing
13 MAIN	24 SMS
	25 phone number
	26 user account

The features describing each family take the values of 0 (if that feature is not present in that family) or 1 (if the feature is present).

3.2. Results

For comparison purposes, some different projection models have been applied, whose results are shown below.

PCA Projection

Fig. 1 shows the principal component projection (components 1 and 2), obtained by applying PCA to the previously described data.

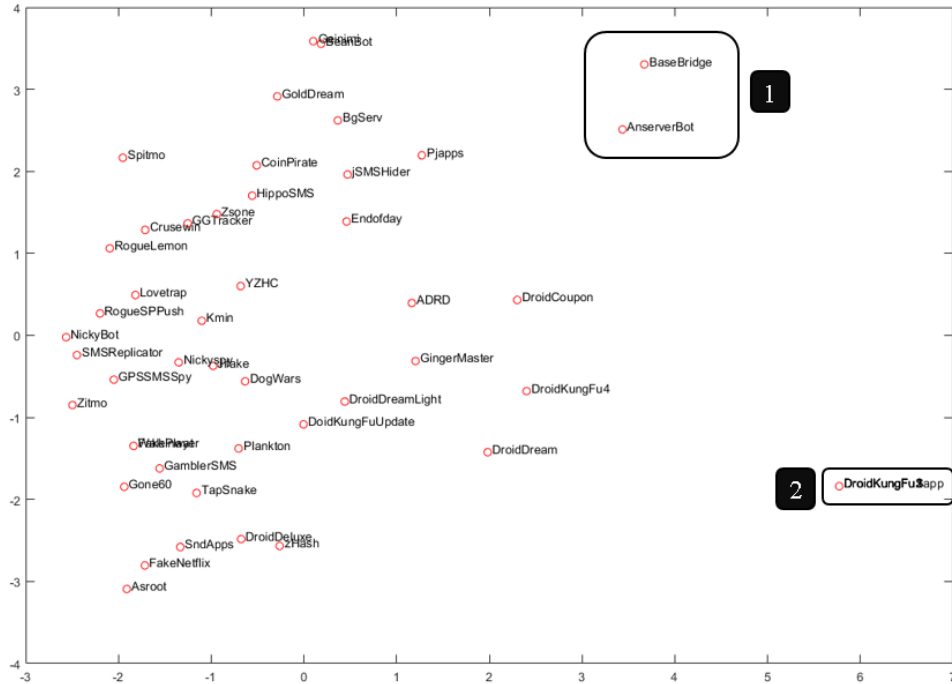


Figura 1: PCA projection of Malgenome families.

In Fig 1 it can be seen that most of the malware families are grouped in a main group (left side of the figure) while just a few families can be identified away from this cluster (groups 1 and 2). Group 1 gathers two families (*BaseBridge* and *AnserverBot*), that are the only two families in the dataset that combine repackaging and update installation. Group 2 gathers four families (*DroidKungFu1*, *DroidKungFu2*, *DroidKungFu3* and *DroidKungFuSapp*) that are the only ones in the dataset presenting the encrypted privilege escalation.

Additionally, this first projection let us identify that some families are projected at the very same place. By getting back to the data we have realized that these families take the very same values for all the features. This is the case of *Walkinwat* and *FakePlayer* on the one hand and for *DroidKungFu1*, *DroidKungFu2*, *DroidKungFu3* and *DroidKungFuSapp* on the other hand. It means that, by taking into account the features in the analysed dataset, it will not be possible to distinguish *Walkinwat* from *FakePlayer* malware or any of the 4 mentioned variants of *DroidKungFu* malware.

MLHL Projection

Fig. 2 shows the MLHL projection of the analyzed data (two main components). MLHL projection shows the structure of the data in a way that a kind of ordering can be seen in the dataset. However, as it is more clearly shown in the CMLHL projection (Fig. 3), MLHL is not further described.



Figura 2: MLHL projection of Malgenome families.

The parameter values of the MLHL model for the projections shown in Fig. 2 are: Number of output dimensions: 3. Number of iterations: 100, learning rate: 0.2872, p : 0.4852.

CMLHL Projection

When applying CMLHL to the analysed dataset, the projection (two main components) shown in Fig. 3 has been obtained. As expected, CMLHL obtained a sparser projection than MLHL and PCA, revealing the structure of the dataset in a clearer way.

The parameter values of the CMLHL model for the projections shown in Fig. 3 are; Number of output dimensions: 3. Number of iterations: 100, learning rate: 0.0406, p : 1.92, τ : 0.44056.

In Fig 3 it is easy to visually identify at least two main groups of data, labelled as 1 and 2. It has been checked that families in each one of these groups are similar in a certain way; group 1 gathers all the families with dangerous SMS activity, as “SMS activation” and “SMS financial charges” are present in all the families in this group. On the other hand, none of the families in this group present any

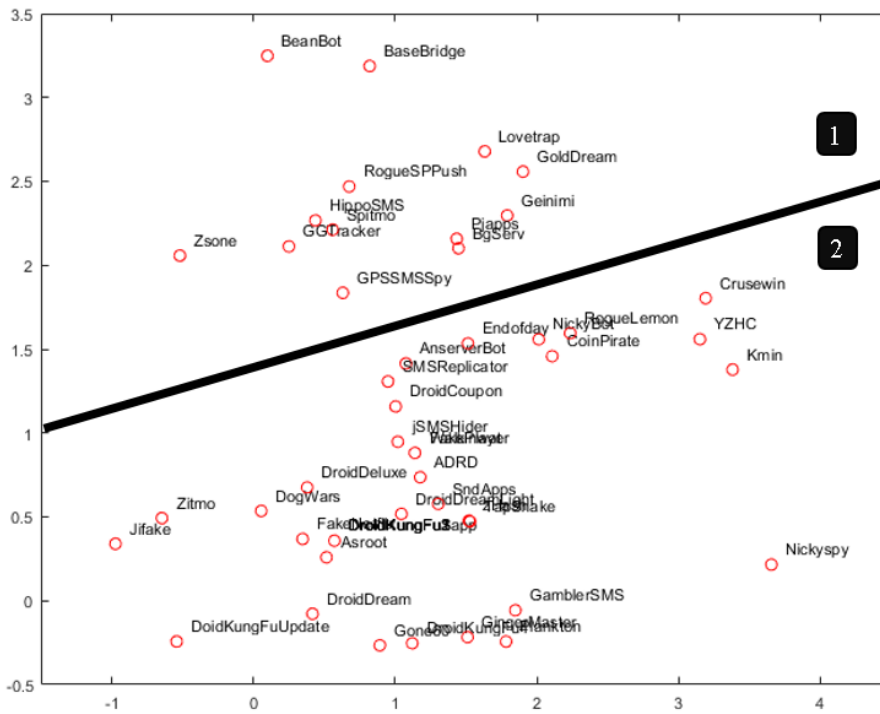


Figure 3: CMLHL projection of Malgenome families.

of the following features: USB or “PKG activation”, and “user-account personal information stealing”. This group is also characterized by the almost complete absence of privilege escalation, as only one of those features (RATC/Zimperlich) is present in only one of the families (*BaseBridge*). Regarding group 2, none of the families in Group 2 present the feature “phone-call financial charges”. From a deeper analysis of such groups, some subgroups can be distinguished and are identified in Fig. 4. Additionally, the families located in each one of these groups are listed in Table 2.

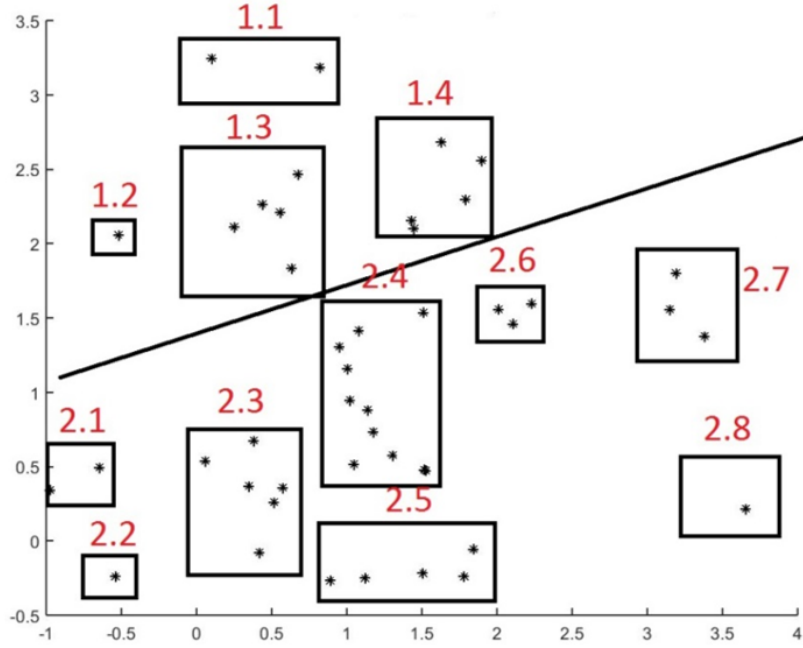


Figure 4: CMLHL projection of Malgenome families with identified subgroups.

Tabla 2: Families allocation to subgroups defined in CMLHL projection.

Group	Subgroup	Families
1	1.1	<i>BaseBridge, BeanBot</i>
	1.2	<i>Zsone</i>
	1.3	<i>GGTracker, GPSSMSSpy, HippoSMS, RogueSPPush, Spitmo</i>
	1.4	<i>BgServ, Geinimi, GoldDream, Lovetrap, Pjapps</i>
2	2.1	<i>Jifake, Zitmo</i>
	2.2	<i>DroidKungFuUpdate</i>
	2.3	<i>Asroot, DogWars, DroidDeluxe, DroidDream, DroidKungFu1, DroidKungFu2, DroidKungFu3, DroidKungFuSapp, FakeNetflix</i>
	2.4	<i>ADDRD, AnserverBot, DroidCoupon, DroidDreamLight, Endofday, FakePlayer, jSMShider, SMSReplicator, SndApps, TapSnake, Walkinwat, zHash</i>
	2.5	<i>DroidKungFu4, GamblerSMS, GingerMaster, Gone60, Plankton</i>
	2.6	<i>CoinPirate, NickyBot, RogueLemon</i>
	2.7	<i>Crusewin, Kmin, YZHC</i>
	2.8	<i>Nickyspy</i>

All the variants of *DroidKungFu* malware are located in the bottom-left side of the projection (groups 2.2, 2.3, and 2.5). *Jifake* and *Zitmo* are gathered in the same subgroup (2.1) as they are the only two families in group 2 presenting the drive-by download installation feature

ISOMAP Projection

In Fig. 5 it is shown the projections obtained by ISOMAP algorithm where each

sample is labelled with the name of the family it belongs.

The parameter values of the ISOMAP model for the projection shown in Fig. 5

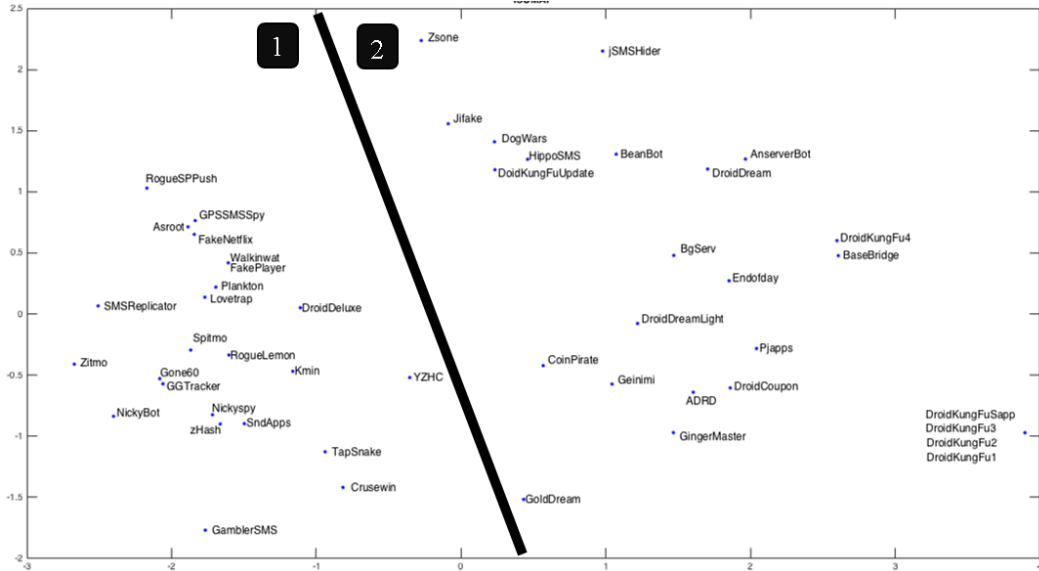


Figura 5: ISOMAP projection of Malgenome families.

are: number of neighbours: 10.

ISOMAP clearly visualize the internal dataset structure, with a main division into 2 groups (1 and 2 in Fig. 5). For a deeper analysis, these two main groups are split in different subgroups as shown in Fig. 6.

From groups in Fig 5, it can be highlighted that group 1 contains Malgenome families that present “standalone” but not “repackaging” installation features (see Table 3). However, in case of group 2, none of its samples present “standalone installation” feature and all of them present the “repackaging installation” feature.

As shown in Fig. 6, group 1 is divided in 2 subgroups (G1a and G1b), where G1a gathers families that do not present the “BOOT activation” feature as opposed to G1b, where its samples present this “BOOT activation” feature. Similarly, G2 is clearly divided into 2 subgroups (G2a and G2b), with analogous characteristics to samples in G1a and G2a respectively (“BOOT activation” for G2a and the opposite for G2b). Finally, G2a presents samples with dangerous activity of “NET remote control” (group G2a.1), and samples without such dangerous feature (group G2a.2).

Table 3 shows the Malgenome families contained in each one of the identified groups, and the features characterizing all the families in that group.

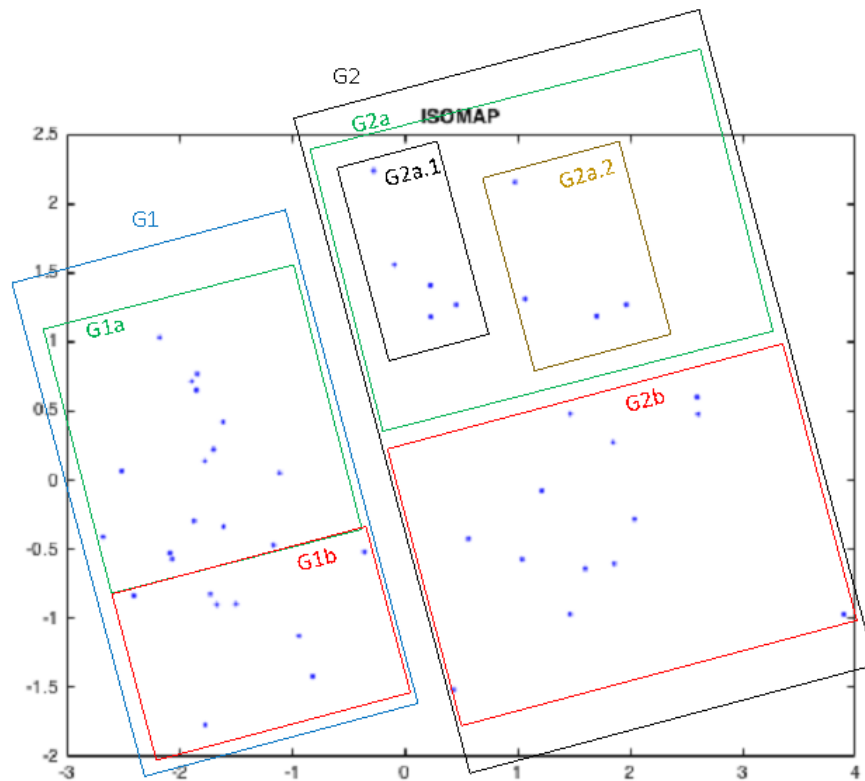


Figura 6: ISOMAP projections of Malgenome families with identified subgroups.

Tabla 3: Families allocation to subgroups defined in ISOMAP projection.

Group	Subgroup	Families and features
G1	G1a	RogueSPPush, GPSSMSSpy, Asroot, FakeNetflix, Walkinwat, FakePlayer, Plankton, SMSReplicator, Lovetrap, DroidDeluxe, Spitmo, Zitmo, Spitmo, RogueLemon, Gone60, GGTracker, Kmin Present features: 4 (standalone installation) and 5 (BOOT activation) Not-present feature: 1(repackaging installation)
	G1b	NickyBot, Nickyspy, zHash, SndApps, YZHC, TapSnake, Crusewin, GamblerSMS Present features: 4 (standalone installation) Not-present feature:1(repackaging installation) and 5 (BOOT activation)

Group	Subgroup	Families and features
G2	G2a.1	Zsone, Jifake, DogWars, HippoSMS, DoidKungFuUpdate Present features: 1 (repackaging installation) and 5 (BOOT activation) Not-present feature: 4 (standalone installation) and 19 (NET remote control)
	G2a.2	jSMShider, BeanBot, AnserverBot, DroidDream Present features: 1 (repackaging installation), 5 (BOOT activation) and 19 (NET remote control) Not-present feature: 4 (standalone installation)
	G2b	DroidKungFu4, BaseBridge, BgServ, Endofday, DroipDreamLight, Pjapps, DroidCoupon, CoinPirate, Geinimi, ADRD, GingerMaster, DroidKungFuSapp, DroidKungFu3, DroidKungFu2, DroidKungFu1 Present features: 1 (repackaging installation) and 19 (NET remote control) Not-present feature: 4 (standalone installation), and 5 (BOOT activation)

CCA Projection

Fig. 7 presents the projection obtained by CCA of Malgenome families, where it can be seen that a clear internal structure of the dataset can not be identified, and malware families can not be clearly gathered in groups, as it happened in previous results.

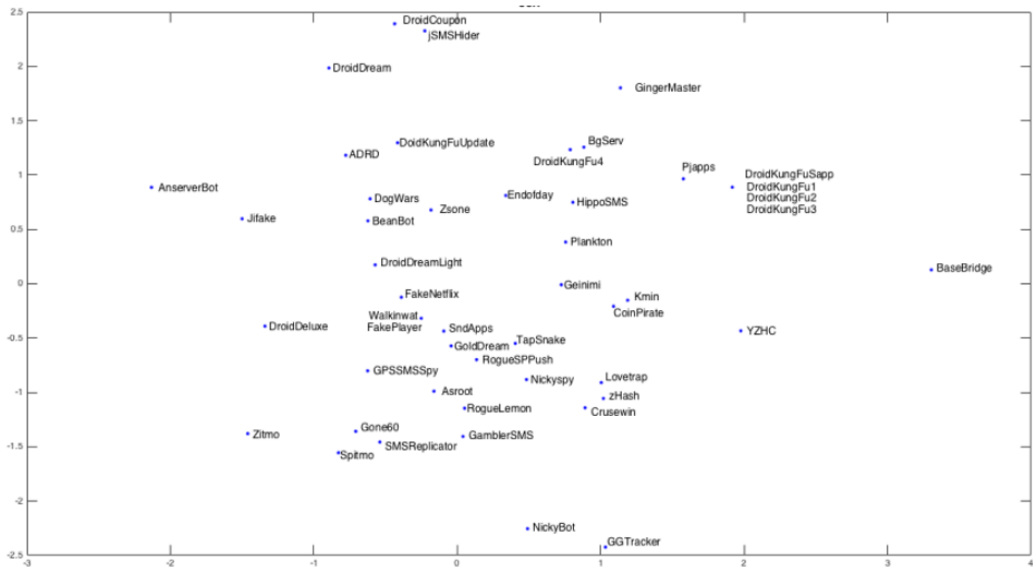


Figura 7: CCA projection of Malgenome families.

Different combinations of ν values were tested for the parameters of the CCA model. The best projection obtained is the one shown in Fig. 7, that was generated with 1,000 epochs, $\alpha=0.5$ and $\lambda=1.5152$.

SOM results

Finally, SOM has been also applied to the Malgenome dataset and the obtained U-matrix is shown in Fig. 8. Each one of the neurons in the map has been labelled with the names of the malware families to which the neuron responds. From this figure, and according to the inter-neuron distances, neurons in the map could be

easily split in two main groups (G1 and G2). At the same time, G1 could be divided in two subgroups (G1a and G1b), and G1a could also be divided into three subgroups as neuron distances are high between them (blue color means high distance in Fig. 8). In the case of G2, authors believe that it can not be divided into subgroups, as neuron distances within G2 are quite small, so it can not be said that families in this group (G2) are very different.

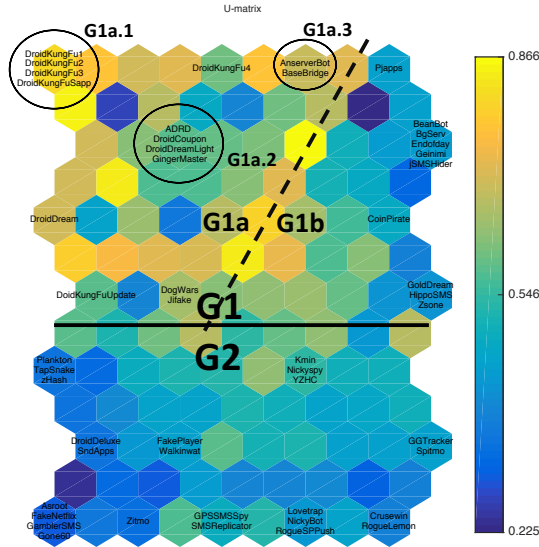


Figura 8: SOM U-matrix for Malgenome families with identified groups.

The parameter values of the SOM model for the mapping shown in Fig. 8 are; map size: $[7, 5]$, lattice: hexagonal, neighbourhood function: Gaussian. On the other hand, some metrics about the obtained mapping are: quantization error = 1.1, and topographic error = 0.0. In general terms, it can be said that the “repackaging installation” and “standalone installation” features are the only ones that let distributing samples in groups G1 and G2. In the case of subgroups G1a and G1b, it is the presence of the “SMS financial charges” feature what characterize malware families in each one of them.

Table 4 shows the Malgenome families contained in each one of the identified groups by the SOM network, and the features characterizing all the families in that group. General information (present and not-present features) of a group (i.e. group 1) is applicable to all malware families contained in its subgroups (i.e. subgroups G1a1, G1b, etc.).

Tabla 4: Families allocation to subgroups defined in SOM u-matrix.

G1	
Present feature: 1 (repackaging installation)	
Not-present feature: 4 (standalone installation)	
G1a	
Present feature: 22 (SMS financial charges)	
G1a1	DroidKungFu1, DroidKungFu2, DroidKungFu3, DroidKungFuSapp Present features: 11 (BATT activation), 14 (exploit privilege escalation), 18 (encrypted privilege escalation) Not-present feature: 2 (update installation), 24 (SMS personal information stealing)
	DroidCoupon, DroidDreamLight, GingerMaster, ADRD Not-present feature: 11 (BATT activation), 14 (exploit privilege escalation), 18 (encrypted privilege escalation), 2 (update insatalltion), 24 (SMS personal information stealing)
G1a2	AnserverBot, BAseBridge Present features: 11 (BATT activation), 2 (update installation), 24 (SMS personal information stealing) Not-present feature: 14 (exploit privilege escalation), 18 (encrypted privilege escalation)
G1a3	Present features: 11 (BATT activation), 2 (update installation), 24 (SMS personal information stealing) Not-present feature: 14 (exploit privilege escalation), 18 (encrypted privilege escalation)
Others samples	DogWars, DroidkungFuUpdate, DroidDream, DroidKungFu4, Jifake
G1b	
Not present features: 22 (SMS financial charges)	
BeanBot, BgServ, CoinPirate, Endofday, Geinimi, GoldDream, HippoSMS, jSMShider, Pjapps, Zsone	
G2	
Present features: 4 (standalone installation)	
Not-present feature: 1 (repackaging installation)	
Asroot, Crusewin, DroidDeluxe, FakeNetflix, FakePlayer, GamblerSMS, GGTracker, Gone60, GPSSMSSpy, Kmin, Lovetrap, NickyBot, Nickyspy, Plankton, RogueLemon, RogueSPPush, SMSReplicator, SndApps, Spitmo, TapSnake, Walkinwat, YZHC, zHash, Zitmo	

4. Conclusions

From the results shown in section 3, it can be concluded that dimensionality reduction techniques are an interesting proposal to visually analyse the structure of a high-dimensionality dataset in general terms. More specifically, when studying Android malware families, this kind of techniques let us gain deep knowledge about the nature of such app families. Thanks to the obtained projections, similarities and differences of the studied families are identified.

From the extensive set of applied DRTs, PCA, MLHL and CCA failed in generating an informative visualization of samples by reducing the dimensionality of them to 2D. On the other hand and generally speaking, it can be said that the DRTs that group malware families, are able to do that in a way consistent with the seminal characterization of Malgenome dataset [7]. More precisely, it is worth mentioning that installation features (repackaging and standalone) have been identified by ISOMAP and SOM as the most important ones for a general characterization of Android malware families (see section 3 for further details). It is an important result as repackaging is one of the most common techniques applied to hide malware (86 % of the malware apps in the original dataset were repackaged versions of different legitimate apps including paid apps, popular game apps, powerful utility apps, etc. [6]).

In a complementary way, CMLHL identified some activation (SMS, USB, and PKG) and financial charges (SMS, and phone call) as the most important ones for such a task. Knowledge generated by the application of DRTs could be applied to improve the detection rate of Android Malware at different stages (markets, devices, etc.) thanks to the characterization of the different families.

As a final conclusion, it can be said that the identification and characterization of Android malware is still an open challenge that requires great efforts to be devoted in coming years.

5. Future work

As future work it is planned to apply new DTRs models and compare them with other supervised algorithms such as decision trees in order to gain deep knowledge of the dataset. It will also be analysed other datasets related to cybersecurity applying the same approach followed in this research in order to generalize to other datasets the proposed method.

Referencias

- [1] A. Altaher, An improved android malware detection scheme based on an evolving hybrid neuro-fuzzy classifier (EHNFC) and permission-based features, *Neural Computing and Applications* 28 (2017) 4147–4157.
- [2] AppBrain, Android operating system statistics, ????
- [3] M. Proyect, Android malware genome project, ????
- [4] S. Arshad, M. A. Shah, A. Khan, M. Ahmed, Android malware detection & protection: a survey, *International Journal of Advanced Computer Science and Applications* 7 (2016) 463–475.
- [5] D. Atienza, Á. Herrero, E. Corchado, Neural analysis of http traffic for web attack detection, in: *Computational Intelligence in Security for Information Systems Conference*, Springer, pp. 201–212.
- [6] L. Sáiz, A. Pérez, Á. Herrero, E. Corchado, Analyzing key factors of human resources management, in: *International Conference on Intelligent Data Engineering and Automated Learning*, Springer, pp. 463–473.
- [7] I. J. Machón González, H. López García, J. L. Calvo Rolle, Neuro-robust controller for non-linear systems (controlador neurorobusto para sistemas no lineales), *Dyna* (2011).
- [8] L. Cen, C. S. Gates, L. Si, N. Li, A probabilistic discriminative model for android malware detection with decompiled source code, *IEEE Transactions on Dependable and Secure Computing* 12 (2014) 400–412.

- [9] H. Chang, D.-Y. Yeung, Y. Xiong, Super-resolution through neighbor embedding, in: Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004., volume 1, IEEE, pp. I–I.
- [10] N. Chen, B. Ribeiro, A. Vieira, A. Chen, Clustering and visualization of bankruptcy trajectory using self-organizing map, *Expert Systems with Applications* 40 (2013) 385–393.
- [11] G. Cirrincione, J. Héroult, V. Randazzo, The on-line curvilinear component analysis (oncca) for real-time data reduction, in: 2015 International Joint Conference on Neural Networks (IJCNN), IEEE, pp. 1–8.
- [12] E. Corchado, C. Fyfe, Connectionist techniques for the identification and suppression of interfering underlying factors, *International Journal of Pattern Recognition and Artificial Intelligence* 17 (2003) 1447–1466.
- [13] E. Corchado, Á. Herrero, J. M. Sáiz, Testing cab-ids through mutations: on the identification of network scans, in: International Conference on Knowledge-Based and Intelligent Information and Engineering Systems, Springer, pp. 433–441.
- [14] E. Corchado, Á. Herrero, Neural visualization of network traffic data for intrusion detection, *Applied Soft Computing* 11 (2011) 2042–2056.
- [15] E. Corchado, D. MacDonald, C. Fyfe, Maximum and minimum likelihood hebbian learning for exploratory projection pursuit, *Data Mining and Knowledge Discovery* 8 (2004) 203–225.
- [16] C. C. Turrado, M. d. C. M. López, F. S. Lasheras, B. A. R. Gómez, J. L. C. Rollé, F. J. d. C. Juez, Missing data imputation of solar radiation data under different atmospheric conditions, *Sensors* 14 (2014) 20382–20399.
- [17] P. Demartines, J. Héroult, Curvilinear component analysis: A self-organizing neural network for nonlinear mapping of data sets, *IEEE Transactions on neural networks* 8 (1997) 148–154.
- [18] J. J. Fuertes, M. Domínguez, P. Reguera, M. A. Prada, I. Díaz, A. A. Cuadrado, Visual dynamic model based on self-organizing maps for supervision and fault detection in industrial processes, *Engineering Applications of Artificial Intelligence* 23 (2010) 8–17.
- [19] C. Fyfe, A neural network for pca and beyond, *Neural Processing Letters* 6 (1997) 33–41.

- [20] C. Fyfe, D. R. McGregor, R. Baddeley, Exploratory projection pursuit: an artificial neural network approach, Department of Computer Science, University of Strathclyde, 1994.
- [21] R. F. Garcia, J. L. C. Rolle, M. R. Gomez, A. D. Catoira, Expert condition monitoring on hydrostatic self-levitating bearings, *Expert Systems with Applications* 40 (2013) 2975–2984.
- [22] H. Fakhouri, L. Cherrat, M. Ezziyyani, Towards a new approach to improve the classification accuracy of the kohonen’s self-organizing map during learning process (????).
- [23] A. Herrero, U. Zurutuza, E. Corchado, A neural-visualization ids for honey-net data, *International Journal of Neural Systems* 22 (2012) 1250005.
- [24] S. Hou, P. D. Wentzell, Re-centered kurtosis as a projection pursuit index for multivariate data analysis, *Journal of Chemometrics* 28 (2014) 370–384.
- [25] J.-w. Jang, J. Yun, A. Mohaisen, J. Woo, H. K. Kim, Detecting and classifying method based on similarity matching of android malware behavior with profile, *SpringerPlus* 5 (2016) 273.
- [26] T. Kohonen, Essentials of the self-organizing map, *Neural networks* 37 (2013) 52–65.
- [27] T. Kohonen, The self-organizing map, *Neurocomputing* 21 (1998) 1–6.
- [28] Malwarebytes, Malwarebytes Labs Report: Cybercrime Tactics and Techniques Q3 2017, Technical Report, 2017.
- [29] E. Mohebi, A. Bagirov, Constrained self organizing maps for data clusters visualization, *Neural Processing Letters* 43 (2016) 849–869.
- [30] V. Moonsamy, J. Rong, S. Liu, Mining permission patterns for contrasting clean and malicious android applications, *Future Generation Computer Systems* 36 (2014) 122–132.
- [31] E. Oja, Principal components, minor components, and linear neural networks, *Neural networks* 5 (1992) 927–935.
- [32] M. Paliwal, U. A. Kumar, Neural networks and statistical techniques: A review of applications, *Expert systems with applications* 36 (2009) 2–17.
- [33] W. Park, K.-H. Lee, K.-S. Cho, W. Ryu, Analyzing and detecting method of android malware via disassembling and visualization, in: 2014 International Conference on Information and Communication Technology Convergence (ICTC), IEEE, pp. 817–818.

- [34] A. Paturi, M. Cherukuri, J. Donahue, S. Mukkamala, Mobile malware visual analytics and similarities of attack toolkits (malware gene analysis), in: 2013 International Conference on Collaboration Technologies and Systems (CTS), IEEE, pp. 149–154.
- [35] K. Pearson, Liii. on lines and planes of closest fit to systems of points in space, The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science 2 (1901) 559–572.
- [36] C. I. Pinzón, J. F. De Paz, A. Herrero, E. Corchado, J. Bajo, J. M. Corchado, idmas-sql: intrusion detection based on mas to detect and block sql injection through data mining, Information Sciences 231 (2013) 15–31.
- [37] C. Pinzón, J. F. De Paz, J. Bajo, Á. Herrero, E. Corchado, Aiida-sql: An adaptive intelligent intrusion detector agent for detecting sql injection attacks, in: 2010 10th International Conference on Hybrid Intelligent Systems, IEEE, pp. 73–78.
- [38] H. Quintián, E. Corchado, Beta hebbian learning as a new method for exploratory projection pursuit, International journal of neural systems 27 (2017) 1750024.
- [39] H. Quintián, E. Corchado, Beta scale invariant map, Engineering Applications of Artificial Intelligence 59 (2017) 218–235.
- [40] J. W. Sammon, A nonlinear mapping for data structure analysis, IEEE Transactions on computers 100 (1969) 401–409.
- [41] R. Sánchez, Á. Herrero, E. Corchado, Visualization and clustering for snmp intrusion detection, Cybernetics and Systems 44 (2013) 505–532.
- [42] B. Sanz, I. Santos, C. Laorden, X. Ugarte-Pedrero, J. Nieves, P. G. Bringas, G. Álvarez Marañón, Mama: manifest analysis for malware detection in android, Cybernetics and Systems 44 (2013) 469–488.
- [43] Statista, Smartphone sales by os worldwide 2009-2017, ????
- [44] O. Somarriba, U. Zurutuza, R. Uribeetxeberria, L. Delosières, S. Nadjm-Tehrani, Detection and visualization of android malware behavior, Journal of Electrical and Computer Engineering 2016 (2016).
- [45] SOPHOSLABS, SophosLabs 2018 Malware Forecast, Technical Report, 2017.
- [46] P. Teuff, M. Ferk, A. Fitzek, D. Hein, S. Kraxberger, C. Orthacker, Malware detection by applying knowledge discovery processes to application metadata on the android market (google play), Security and communication networks 9 (2016) 389–419.

- [47] M. Wagner, F. Fischer, R. Luh, A. Haberson, A. Rind, D. A. Keim, W. Aigner, A survey of visualization systems for malware analysis, in: Eurographics Conference on Visualization (EuroVis), pp. 105–125.
- [48] Y. Wu, T. K. Doyle, C. Fyfe, Multi-layer topology preserving mapping for k-means clustering, in: International Conference on Intelligent Data Engineering and Automated Learning, Springer, pp. 84–91.
- [49] Y. Zhou, X. Jiang, Dissecting android malware: Characterization and evolution, in: 2012 IEEE symposium on security and privacy, IEEE, pp. 95–109.
- [50] H.-J. Zhu, Z.-H. You, Z.-X. Zhu, W.-L. Shi, X. Chen, L. Cheng, Droiddet: effective and robust detection of android malware using static analysis along with rotation forest model, *Neurocomputing* 272 (2018) 638–646.

Acknowledgement: This is a pre-copied, author-produced version of an article accepted for publication in *Logic Journal of the IGPL* following peer review. The version of record: Rafael Vega Vega, Héctor Quintián, José Luís Calvo-Rolle, Álvaro Herrero, Emilio Corchado, Gaining deep knowledge of Android malware families through dimensionality reduction techniques, *Logic Journal of the IGPL*, Volume 27, Issue 2, April 2019, Pages 160–176, is available online at: <https://doi.org/10.1093/jigpal/jzy030>