



UNIVERSITÀ
DEGLI STUDI
FIRENZE

Scuola di Scienze Matematiche, Fisiche e Naturali
Corso di Laurea Magistrale in Informatica

Multivariate Analysis & Statistical Learning

INTRUSION DETECTION SYSTEMS AND MACHINE LEARNING

FRANCESCO TERROSI

6326113

Anno Accademico 2018-2019

INDICE

1	Introduction	3
2	Intrusion Detection and Prevention Systems	5
	2.0.1 Machine Learning and IDS combined	5
3	Feature Selection	7
	3.1 Principal Component Analysis	7
4	Classification Approaches	9
5	IDS Implementation	11
6	Results & Conclusions	17

INTRODUCTION

Nowadays, we are seeing how much reliance is put into Artificial Intelligence, especially when statistical techniques come in.

The combination of AI and MASL (known as Machine Learning in computer science), gives us a set of very powerful tools in order to accomplish the hardest objectives. Most of the times, these techniques are used for financial business but what happens when we try to apply Machine Learning to Computer Security?

In this paper we will see how powerful can be statistical models when applied to computer security, how they can be applied and in what field and what are the techniques to be used and the steps to build a resilient and secure system.

The focus will be on Intrusion Detection and Prevention System (roughly, you can think to an Anti-Virus as an IPS) based on Anomaly Detection techniques.

Anomaly Detection relies on defining a model for the *normal* behaviour of a user into a system and comparing the observed behaviour of actual users (or programs) with the established model, in order to detect attacks and intrusion. for sake of understanding, you can think to a malware (e.g. a virus) as a "bad user" trying to gain control of your computer system.

In this paper we will go through these topics:

- A general overview of Intrusion Detection and Prevention Systems
- An overview of the techniques used for *Feature Selection*
- An overview of the techniques used for *Classification*
- A practical implementation

INTRUSION DETECTION AND PREVENTION SYSTEMS

Intrusion Detection (and Prevention) Systems (from now on, IDS and/or IPS) are particular HW/SW systems developed to detect security threats in a computer system.

Historically, they were developed using signature-based techniques, that is: you define an attack pattern and, if it recognizes that exact pattern in the system's activity, it raises an alarm.

Since Machine Learning is one of the increasing trends, gaining more and more popularity, we began to think if, instead of defining the "attack pattern" (or signature) for a particular malicious activity, we could define the "good behaviour" for users and program to detect malicious behaviour and here is where anomaly-based IDS comes in.

2.0.1 *Machine Learning and IDS combined*

Let's get more into the interesting things now. The first question is: How do I define the "normal" behaviour for users and programs? We need to have a dataset representing these behaviours. We have two alternatives:

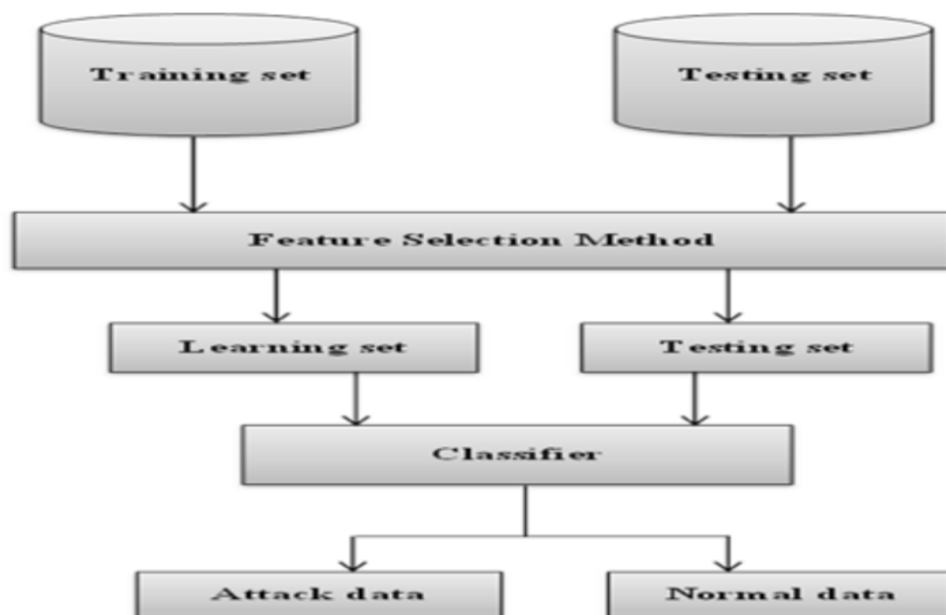
- Stateful Protocol Analysis
 - We buy these datasets from authorized vendors. These data are usually very accurate and very little work is needed on them. The bad thing? The price is usually very high
- Collect data by ourselves
 - This approach is usually less expensive but requires much more work on the data (such as *feature selection*). For the purpose of this course we will have a better look at this approach, since it allows us to use some more statistical techniques

Data are collected by the IDS using a *sensor* (an HW or SW module). These data are sent to an *analyzer* that compares them with the ones collected in the training phase, in order to detect intrusions. However, the very big problem with these kind of systems is that we want a very small number of *false positives* and *false negatives*, otherwise the whole system is useless!

This is the main reason why we will put some emphasis to *feature selection* and *feature classification*

Roughly, the process used to build this kind of system is depicted in Figure 1:

- 1) First of all we need a *Train Set*, to tell the system what are the accepted behaviours
- 2) Now that we have only the relevant components in our train set, we can classify data. Empirically, it has been observed that *Decision Tree* provides the higher accuracy for anomaly detection
- 3) At this point we are able to test our system with a *Test Set*, to check the accuracy of our model of correctly detecting intrusions



FEATURE SELECTION

Feature selection is, no doubt, the critical part when we need to prepare the dataset for the training phase.

If one collects data for the training on its own (i.e. you don't buy data from authorized vendors), typically they will have a lot of redundant attributes, a lot of useless attributes (for intrusion detection) and what he/she wants is to extract from all the observed data the most representative attributes for attacks and behaviours.

In this section we will briefly describe the techniques that can be used and the pros and cons of each approach.

- *Principal Component Analysis* One of the most used techniques for determining uncorrelated attributes and the one that has been used in this implementation
- *Correlation Based Feature Selection Method* works under the hypothesis that: "Good feature subsets contain features highly correlated with the class, yet uncorrelated with each other"
- *Information Gain Based Feature Selection*: roughly, we combine entropy calculation and Information Gain to select the attributes that give us the more information about the class we want to predict
- *Minimum Redundancy Maximum Relevance* is an approach based on the assumption that the most redundant features are the less relevant ones.

3.1 PRINCIPAL COMPONENT ANALYSIS

Principal Component Analysis is a technique used for feature selection, which involves algebra techniques. The idea is to calculate the variance-

covariance matrix and to calculate its eigenvectors and eigenvalues, in order to select the biggest ones (i.e. the ones that explain most of the variance in the dataset).

$$V = \begin{pmatrix} \sigma_1 & \sigma_{12} & \sigma_{13} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_2 & \sigma_{23} & \dots & \sigma_{2n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \sigma_{n3} & \dots & \sigma_n \end{pmatrix}$$

Given this matrix, calculate $\det(V - \lambda I) = 0$ to obtain all the eigenvectors and eigenvalues, and choose the "best" ones that is:

$$\max |x|_x \text{ eigenvalue} \in V$$

To use this techniques, it is *very* important to scale data: through standardization (mean=0, variance=1) or normalization (make values shrink [0,1]). This is because PCA produces a number n of eigenvalues and eigenvector, representing for each feature the impact it has on variance. If we have values with very different scales, we could give much more importance to integer values than to, for example, a rate values.

CLASSIFICATION APPROACHES

A classifier is a tool that tries to categorize unseen patterns in suitable classes.

An accurate (or a "good") train-set is crucial to build a model that provides the best accuracy to Anomaly Detection.

The classifier uses the train-set to build the model that will be used to classify users (and programs) behaviour and to take decision about whether the observed actions are legitimate or not. Empirical studies will help in choosing the most accurate model (in terms of false-positives and false-negatives ratio)

Here there are the most used approached towards classification in Anomaly Detection:

- *Decision Trees* are certainly one of the most used and accurate tools for classification. It builds a model from a train-set that takes decision based on the comparison between the observed features and the nodes of the tree, starting from the root until reaching a decision (leaf) node
- *Nave Bayes* simply uses a form of Bayesian Network to infer intrusions
- *Support Vector Machine* is a technique in which we try to minimize the risk (but we sacrifice some accuracy) in order to optimize speed and scalability
- *Neural Network* is a "connectionist" approach that tries to simulate the human's brain with artificial (digital) neurons. The pros of this approach is, no doubt, its accuracy when the network is trained

for a long period, consisting in multiple generations of NN. The cons is that if an error is made in the first generations, it will be propagated to the next ones

- *K-Nearest Neighbours* is the simplest approach (but surprisingly one of the most accurate in the field of Intrusion Detection!) among all these Machine Learning techniques. The idea is to classify "patterns" instead of features

Collecting some results from empirical studies, here it is a comparison table from all these classification approaches combined with the feature selections technique. As can be seen, the *Decision Tree* classification technique is the one with the *overall* best accuracy.

On the other hand we got the best accuracy combining K-NN with Information Gain

Feature Selection technique	Classifier	Accuracy (in %)
CFS	NB	82.66
	SVM	76.61
	DT	98.99
	NN	83.8
	k-NN	97.65
IGR	NB	90.29
	SVM	94.39
	DT	97.83
	NN	97.7
	k-NN	99.07
PCA	NB	89.91
	SVM	96.78
	DT	98.95
	NN	97.5
	k-NN	98.87
Minimum redundancy maximum relevance feature selection	NB	87.56
	SVM	88.93
	DT	98.78
	NN	94.6
	k-NN	98.05

IDS IMPLEMENTATION

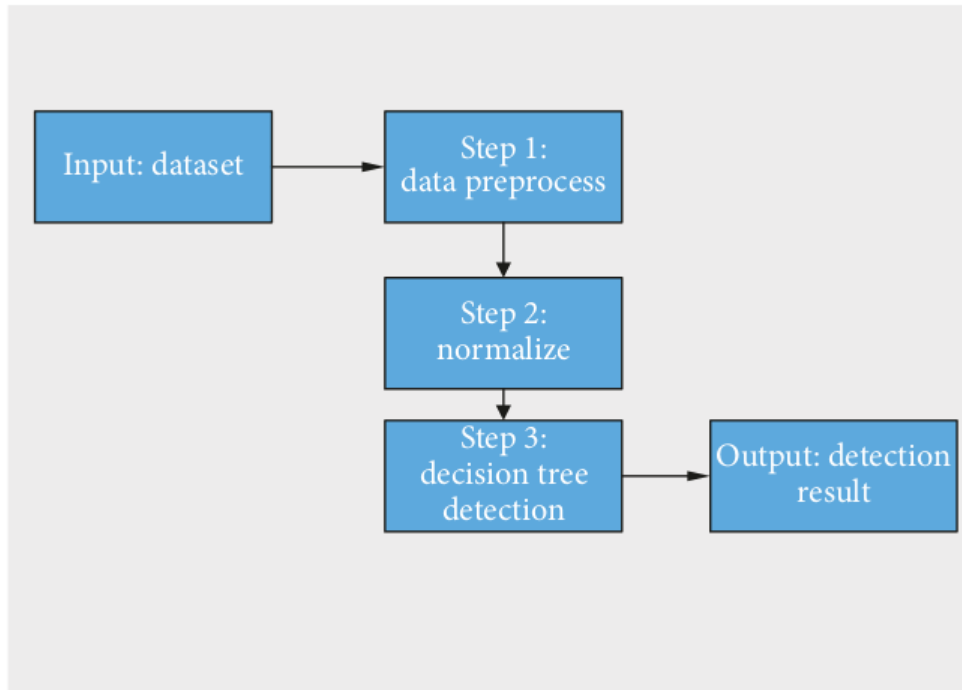
Let's get serious now! An overview of IDS and how statistical techniques can be used with it has been proposed the next step is how to make all these things working together and what are the problems when you go from theory to practice!

In this example, it has been used the KDDCUP99 dataset, a well known set of network data collected by a military company

The steps one must follow in order to build a good system, both in intrusion detection efficiency and computational efficiency are the followings:

- 1) **Data preprocess** - In this step the main issue is that most of the times data are numbers and strings mixed together. A numerical representation for all the strings in the dataset is mandatory for the next steps
- 2) **Data standardization/normalization** - At this point, if we observe the dataset D , it can be seen that the range of numbers in D is not uniform. What must be done here is to run an algorithm that produces an output consisting of a new dataset D' where the values are standardized (or normalized, according to your needs)
- 3) **Decision Tree method** Last step. The classification tool (the Decision Tree in this case) must be fed with the new dataset D' to build the data structure that will be then used to detect intrusions and malicious behaviours

A visual model for these steps is given in the next figure.



So, the first step is to preprocess data in order to work with them. This step is nothing but technical and requires really little skills

```

values = []
index = []
for i in range(0, len(data[0])):
    !isinstance(data[0][i], int) and !isinstance(data[0][i], float):
        index.append(i)

for i in range(0, len(data)):
    for j in range(0, len(index)):
        if data[i][index[j]] not in values:
            values.append(data[i][j])

for i in range(0, len(data))
    for j in range(0, len(values))
        for k in range(0, len(index)):
            if data[i][index[k]] == values[j]:
                data[i][index[k]] = j
                continue
  
```

At this point, there is the need to standardize the data in order to be able to do the Principal Component Analysis. This is a very important step because, if wrongly implemented, it can really blast down the accuracy of predictions.

Standardizing data mean that the new dataset must have all the distribu-

tions associated to the features with mean = 0 and variance = 1.

Luckily, there is a python library that implements a lot of these features for us: the Pandas library.

The very important drawback here is that we can't standardize the train-set and the test-set separately (or the accuracy of prediction will be mined): the train-set must be standardized first and the same criteria must be used to standardize the test-set as well. In other words: the scaler must be "fit" using the train-set and then use the *same* scaler to apply the transformation to both sets

```
import pandas as pd
from sklearn.preprocessing import StandardScaler

test_set = pd.read_csv(test_set_name)
train_set = pd.read_csv(train_set_name)

# Separating out the features
test_set1 = test_set.loc[:, features].values
train_set1 = train_set.loc[:, features].values

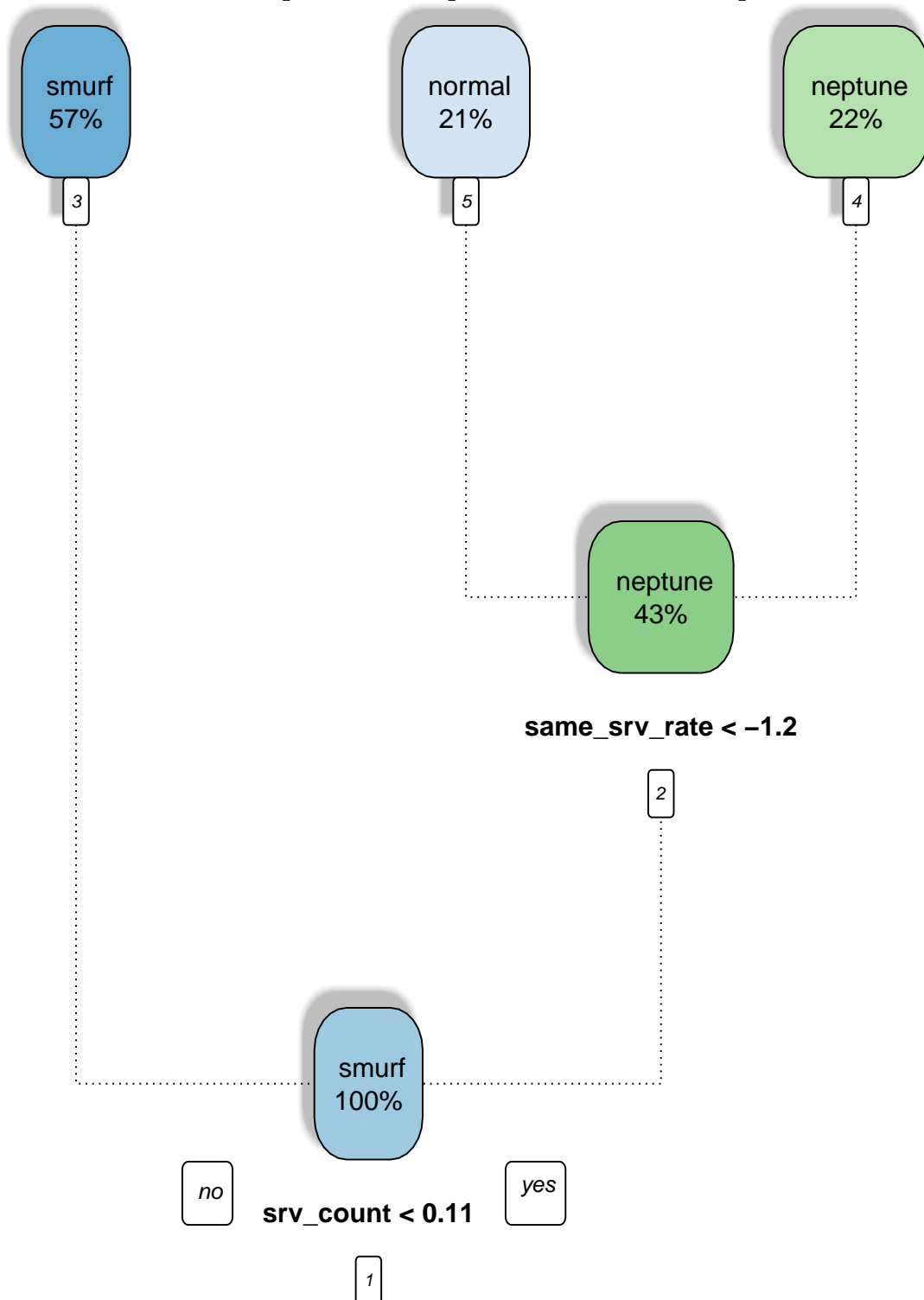
# Separating out the target
y = test_set.loc[:, ['target']].values
y1 = train_set.loc[:, ['target']].values

scaler = StandardScaler()

scaler.fit(train_set1)

# Standardizing the features
test_set1 = scaler.transform(test_set1)
train_set1 = scaler.transform(train_set1)
```


Small and simple tree example to be read bottom-up



RESULTS & CONCLUSIONS

Now that that all the technical things are done we can run our algorithm to see its accuracy.

For empirical purposes, the algorithm have been run with different train-set and with different settings (e.g. complexity parameter), in order to see the differences between an analysis with and without principal components.

Before showing the results here there are some technical data:

The complexity parameter for the tree construction was set to 0.01, 0.0025 and 0.001 and train-sets were:

- 1) Standardized set, using all the features
- 2) Standardized set, using only the features *selected* by PCA (that is, indeed, a non-standard use of this technique)
- 3) The new set of features *produced* by PCA

Regarding test and train set:

- The train-set is the 10% of the Kddcup dataset (which is much much lower than the "standard" dimension, which usually is around 70/80%)
- The test-set is the *whole* Kddcup dataset, consisting of 42 features with almost 5 millions values
- When the PCA train-set was used, the test-set was transformed accordingly to the transformations done on the test-set (remember we had to "fit" the scaler with pandas?)

Results and some of the plots are shown in the next figures:

Complexity parameter	0.01	0.0025	0.001
StandardizedKddcup	0.9870385	0.9832191	0.9843356
KddcupPrincipalFeatures	0.9849946	0.9948412	0.9963878
PCA eigenvectors	0.9901552	0.9951242	0.9956888

If we observe data we can see how Principal Component Analysis raises the accuracy of our prediction. Moreover, the trees built using the selected components are much more balanced when we increase the complexity parameter (i.e. we build bigger trees), while using all the features we get an unbalanced tree.

Another thing to point out is that, with a c.p. of 0.001 we get more accuracy if we don't use the new features produced by PCA (but this could be a matter of the dimension of the train-set)

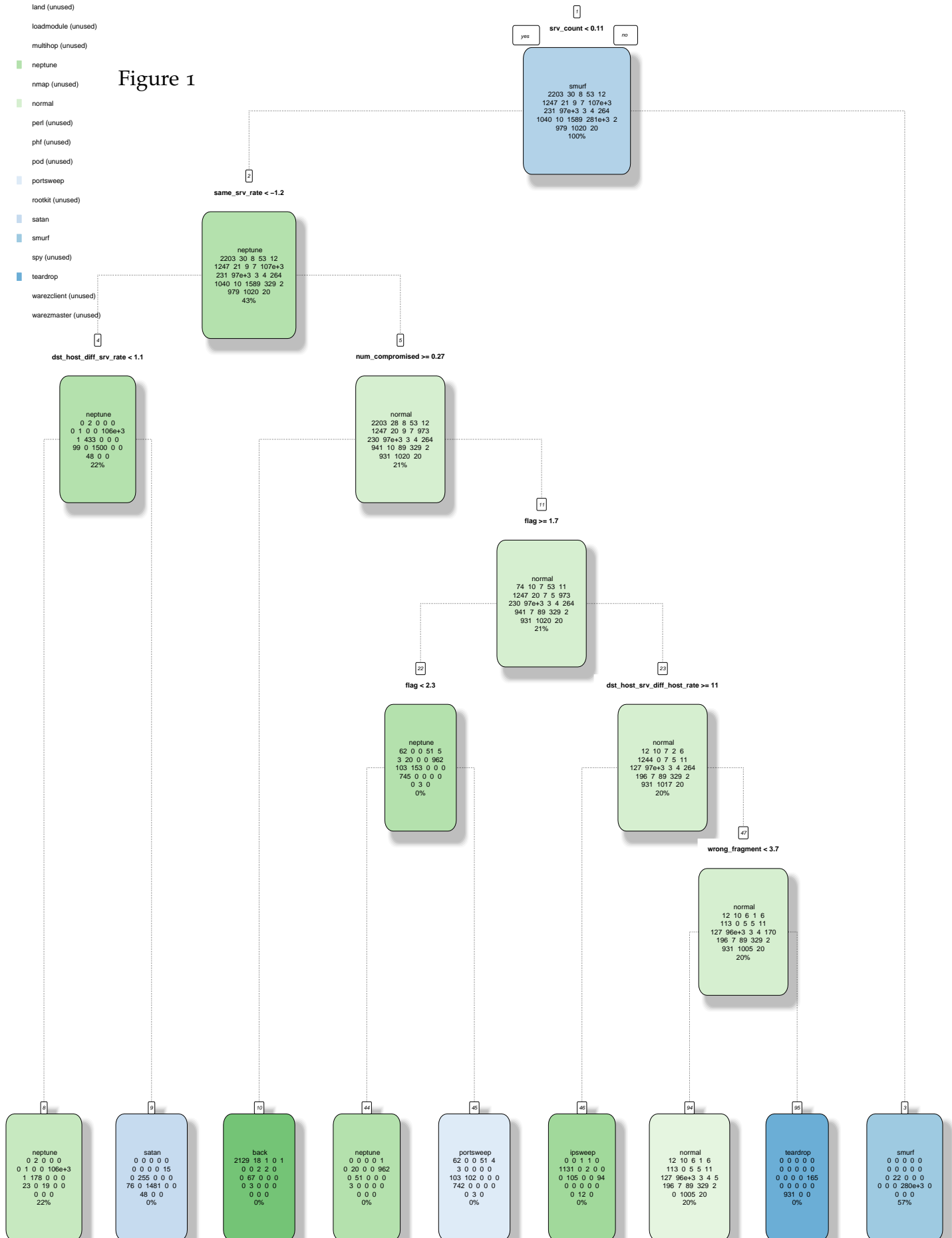
To conclude, it has been observed that PCA is indeed a useful technique to improve the overall accuracy of predictions in Decision Trees (and it also fasten execution time since we have less features!). For future works, considerations may be done regarding Random Forest (i.e. Distributed IDS) and how to combine it with PCA.

In the following pages there are some plots of the various Decision Tree in order to appreciate the differences between the ones that used PCA and the ones that didn't.

- Figure 1 - This is the tree built with the whole features-set, with a c.p. = 0.0025. As can be seen it is very unbalanced. The more complex is the model, the more unbalanced is the tree
- Figure 2 - This is the tree built using the new features produced by PCA. Compared to the previous one, it is much more balanced (and accurate!)
- Figure 3 - Just for some fun, this is the (very big) tree, built using PCA's features with a c.p. = 0.001
- Figure 4 - And here it is the DT built with the whole feature-set and a c.p. = 0.001. As can be seen, things exacerbate when we try to build complex models, while using PCA we get better data structures

guess_passwd (unused)
imap (unused)
ipsweep
land (unused)
loadmodule (unused)
multihop (unused)
neptune
nmap (unused)
normal
perl (unused)
pftl (unused)
pod (unused)
portsweep
rootkit (unused)
satan
smurf
spy (unused)
teardrop
warezclient (unused)
warezmaster (unused)

Figure 1



ftp_write (unused)
guess_passwd (unused)
imap (unused)
ipsweep
land (unused)
loadmodule (unused)
multihop (unused)
neptune
nmap (unused)
normal
perl (unused)
phf (unused)
pod (unused)
portsweep
rootkit (unused)
satan
smurf
spy (unused)
teardrop
warezclient (unused)
warezmaster (unused)

Figure 2

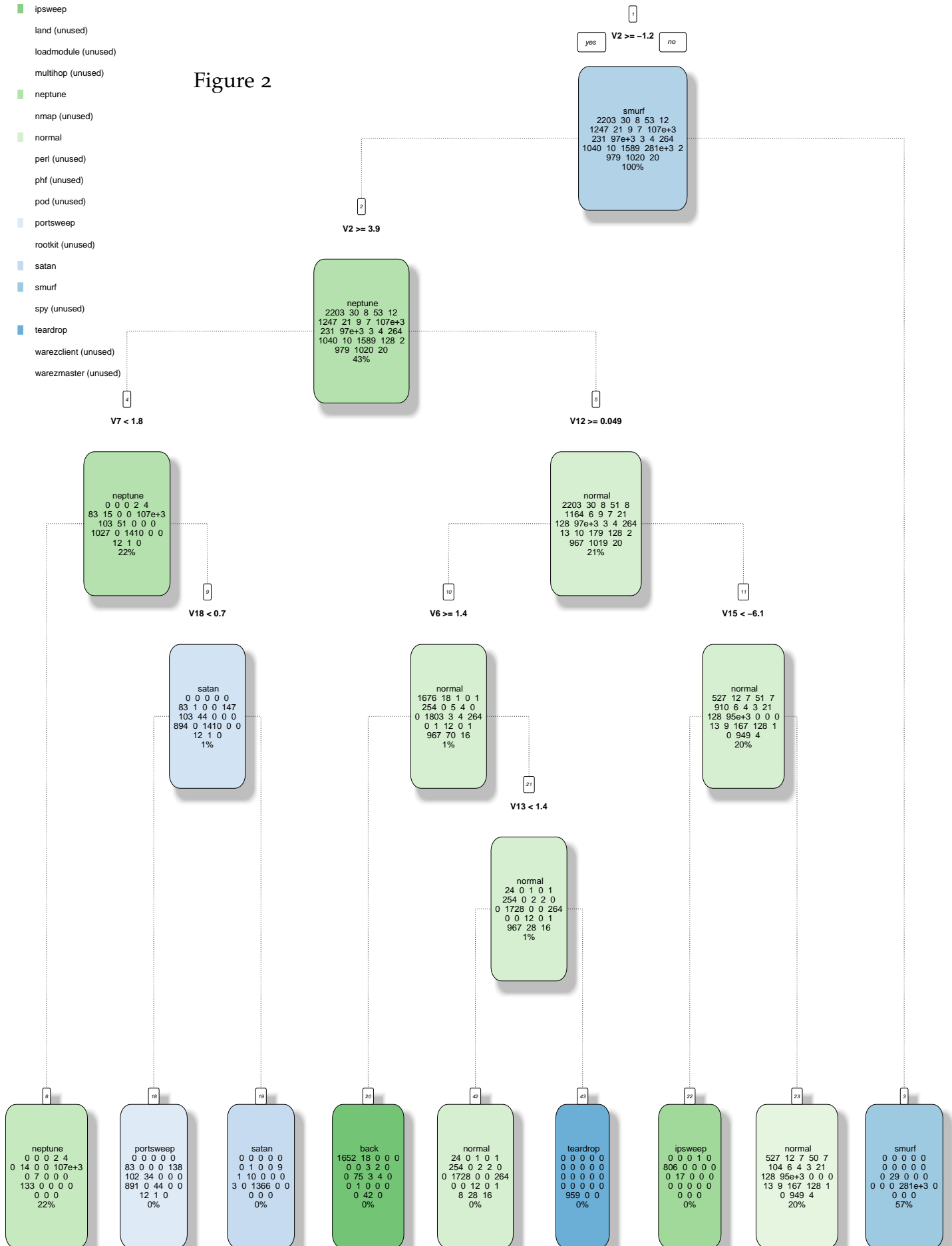
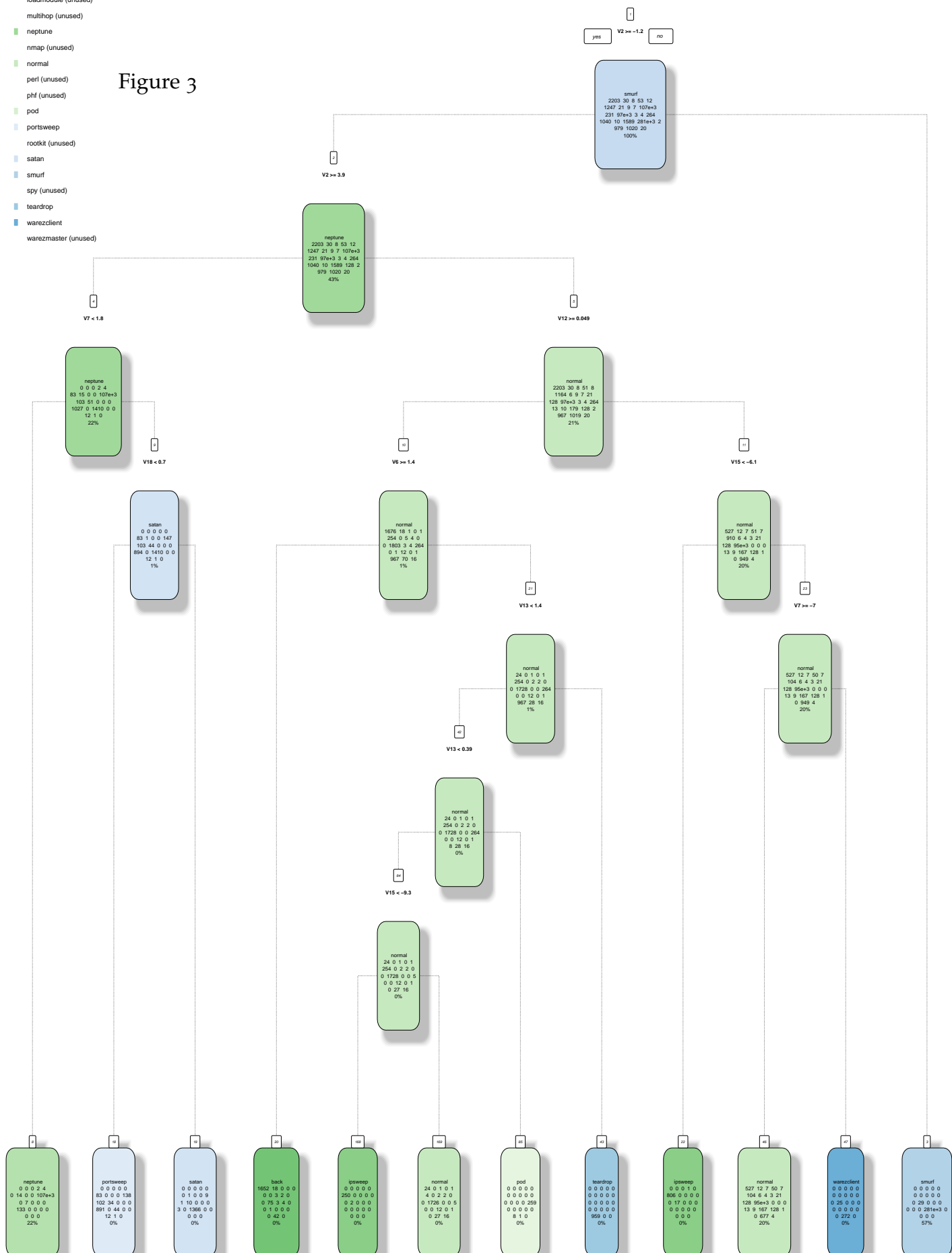
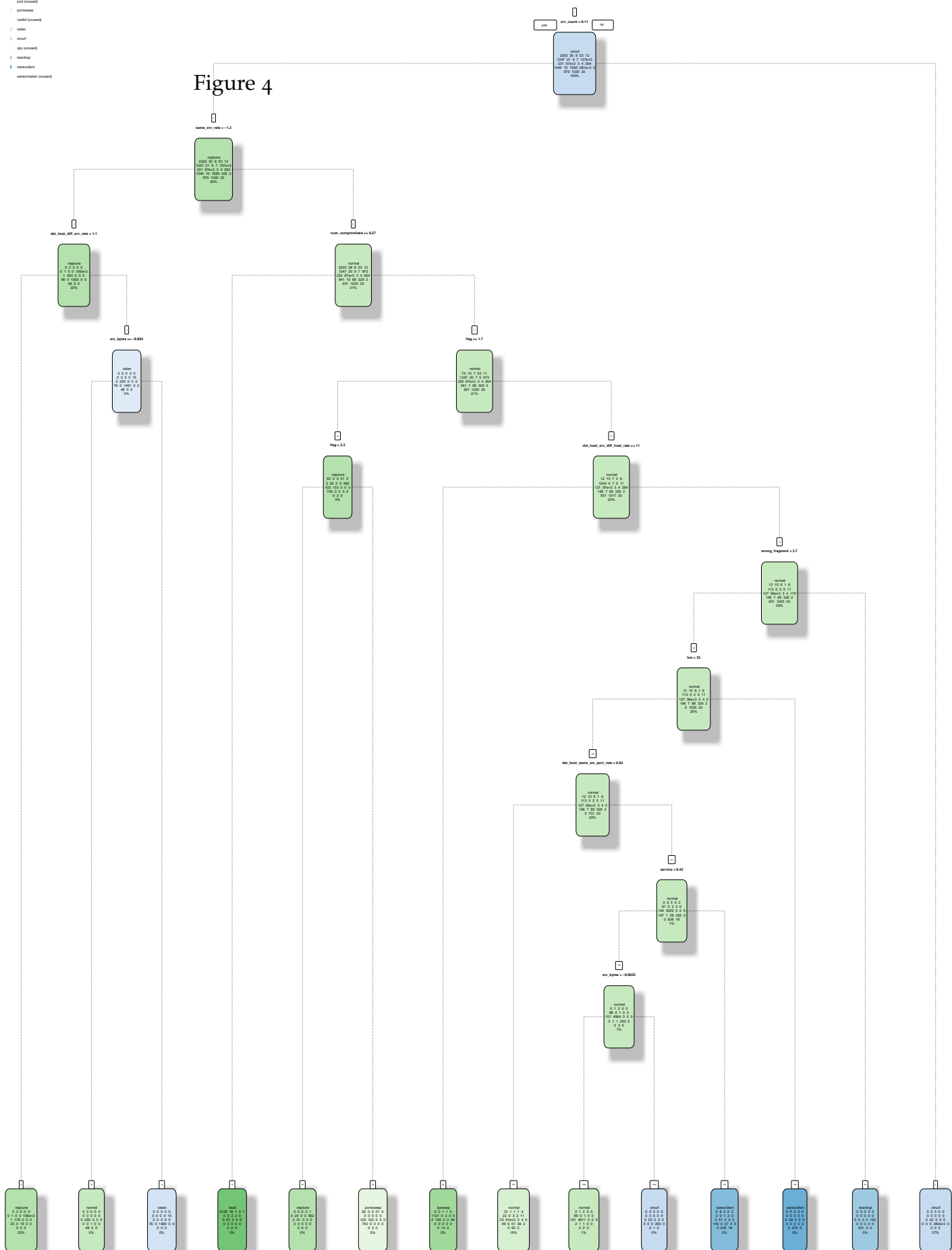


Figure 3





BIBLIOGRAFIA

- [1] Saroj Kr. Biswas - *Intrusion Detection Using Machine Learning: A Comparison Study*
- [2] Kai Peng, Victor C. M. Leung, Lixin Zheng, Shangguang Wang, Chao Huang, Tao Lin - *Intrusion Detection System Based on Decision Tree over Big Data in Fog Environment*