



NTNU – Trondheim
Norwegian University of
Science and Technology

Detecting DNS tunneling using machine learning

Terje Kristoffer Skow

Submission date: December 2015
Responsible professor: Than Van Do, ITEM
Supervisor: Hai Ngyuen, Telenor Research

Norwegian University of Science and Technology
Department of Telematics

Abstract

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

This is the second paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

And after the second paragraph follows the third paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

After this fourth paragraph, we start a new paragraph sequence. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of

the original language. There is no need for special content, but the length of words should match the language.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Preface

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Contents

List of Figures	ix
List of Tables	xi
List of Algorithms	xiii
List of Acronyms	xv
1 DNS	1
1.1 Introduction	1
1.2 Structure	1
1.3 How it works	3
2 DNS Tunneling	5
3 DNS Tunneling Detection	7
3.1 Traffic analysis	7
4 Machine Learning	9
4.1 The Basics	9
4.2 Anomaly Detection	9
4.3 Scikit-learn	10
4.3.1 Models	10
5 Results	11
5.1 Scaled data vs unscaled data	12
5.2 Different axes	13
5.3 Different parameters	14
6 Conclusion	15
References	17
Appendices	

List of Figures

1.1	Example of name spaces of a root with MIL, EDU and ARPA as immediate subdomains. Each leaf is a domain [Moc87].	2
4.1	Example of anomalies in a 2D dataset[CBK09].	10
5.1	Difference between scaled and unscaled data	12
5.2	Uplink vs duration of session.	13
5.3	Average speed in each direction	14

List of Tables

1.1	Example of Resource Record (RR) for telenor.no	3
-----	--	---

List of Algorithms

5.1	Recall and precision definitions	12
5.2	F_1 measure	12

List of Acronyms

DNS Domain Name System.

DPI Deep Packet Inspection.

MCD Minimum Covariance Determinant.

OCSVM One-Class SVM.

RR Resource Record.

SVM Support Vector Machine.

VPN Virtual Private Network.

Chapter 1

DNS

1.1 Introduction

Domain Name System (DNS) is an important part for the internet. It is a system of distributed databases which contains the information about all the domains. In the mid and late 1980s did the previous system, `HOST.TXT`, encounter problems [MD88] which lead to the creation and standardizing called DNS. Since that has the DNS system been updated and configured many times. It needed to be able maintain a fast response time as the database grew larger, this was solved by using a hierarchical set up. This means that each server only has a limited information and sends the request to a new server until it reaches the correct server. It started with one root server, which has expanded to 13 today. The each layer of the hierarchy is called a zone, and it delegates the responsibility for underlying zones delimited by the `dot` in the request name Figure ??.

1.2 Structure

The data in the databases are called RR and contains the information about what the server do with the request. It has the following fields [Moc83]:

- NAME – the owner name of the record.
- TYPE – what type of record this is, name-to-IP (A) or IP-to-name (PTR).
- CLASS – define the class of the record, usually `IN` for internet. It is not widely use and not important for this paper.
- TTL – an integer which says how long the record should be cached by the server receiving the response.
- RDLLENGTH – Specifies the length of the payload in number of octets. One octet is one octet of bits which corresponds to one character

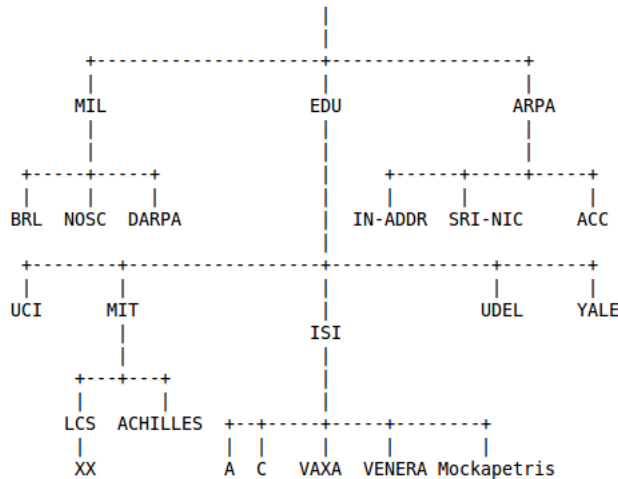


Figure 1.1: Example of name spaces of a root with MIL, EDU and ARPA as immediate subdomains. Each leaf is a domain [Moc87].

- RDATA – the payload of the record. The format and length varies depending on the TYPE and CLASS of the RR.

DNS was first implemented with around 15 different RR TYPE, which has now increased to over 30 [Far13] as a result of the development of the internet. The most notable values for TYPE are:

- A – the payload will contain the ipv4 address of the hostname requested. This is the most used TYPE
- AAAA – contains the ipv6 address of the hostname.
- CNAME – canonical name, respond with the correct alias of the hostname.
- MX – the mail exchange for the domain
- TXT – a text response with large payload, can be used in many ways and are an important type in DNS tunneling.
- PTR – pointer record. Used to map a hostname to an IP-address, commonly known as a reverse lookup.
- NS – authoritative name server for the domain

The 'A' type RR for telenor.no at the name server will look something like this:

Field	Value
NAME	telenor.no
TYPE	A
CLASS	IN
TTL	300
RDLENGTH	15
RDATA	153.110.156.145

Table 1.1: Example of RR for telenor.no

1.3 How it works

To explain of DNS works is an example the easiest way

When a request goes through DNS it starts in the root zone, where it sent down the hierarchy to the `.no` zone.

Normally a DNS server in an enterprise does not send requests directly to the internet, but use an internal DNS server instead. If you are the owner of the authoritative server for a domain, you can control the responses. This is what a DNS tunnel exploits, which will be explained more in the next section.

Chapter 2

DNS Tunneling

DNS tunneling was first used by people who exploited that DNS was not monitored in network you had to pay to use, e.g. hotels and cafés. It was used as an Virtual Private Network (VPN) tunnel. In later years it has been discovered that in enterprises the DNS are not monitored as much as other traffic on the network. People has therefore figured out that it is a good way to ex filtrate data in secure networks. DNS could also be used for a "command and control" attack, where commands are sent over DNS.

The way DNS works it that if you control the authoritative DNS server for a domain you can easily send commands.

With the increase of smartphones it has been discovered that DNS tunneling could again be used as the it started, to use the network without having to pay for it. Carriers can not start charging for regular queries since just regular use of a the internet produces a lot of DNS traffic. Which an user would not see and it would be hard for the carrier to explain for an user what he has been charged for.

Chapter 3

DNS Tunneling Detection

There has been done some research in detecting DNS tunneling over the years, but as it is still a problem no one has found a solution that is cost efficient. The best way for detecting tunnels is still Deep Packet Inspection (DPI) which slows down the DNS requests as the amount of requests increase. DPI looks into each request and response for payload information which can indicate a DNS tunnel. For instance if requests maximizes the size of the labels and the overall name it should be looked at [Far13], this since tunnels would try to minimize the number of packages and maximize speed. Looking at the hostname should also be an indication since regular DNS names is dictionary words or have some meaning, while an encoded name would be meaningless. Traffic analysis is the other main alternative to detecting tunnels. Looking at volume, frequency and other attributes of DNS traffic could give indication of a tunnel. Earlier research has covered different techniques, looking at the volume of DNS traffic from a IP address or the volume of DNS traffic to a specific domain [Far13]. The overarching way of detecting tunnels with traffic analysis is looking for anomalies and stand out cases.

3.1 Traffic analysis

Data that is tunnelled through DNS is normally limited to 512 bytes per request, which leads to clients to send and receive lots of requests and responses. If the server should have the possibility to send data to the client will the client have to constantly send requests to get the data as a response from the server. All this leads to lots of DNS traffic which is not similar to normal use.

Chapter 4

Machine Learning

Machine learning is a way of using statistics to solve problems either by learning from a data set how what the output should be for an input or by figure out different patterns in the data. This is called supervised and unsupervised learning respectively. It is widely used in spam filtering and search engines.

4.1 The Basics

The basics for machine learning is to use to the computer to create a model from which the computer is able to predict the output or category of an input based on the values of the input. Machine learning most often consists of two phases, training and testing. The training phase is where the algorithm learns the model and works out how to categorize the data, and testing is where a new dataset is used to see how accurate the algorithm is. In this report supervised and unsupervised learning is used. Supervised learning means the data has a label of which it is meant to be categorized as, and unsupervised use unlabeled data with the assumption that the majority of data is considered normal.

4.2 Anomaly Detection

Anomaly detection, also called outlier detection, is a problem very suited for machine learning. It is a way of identify observations or data which doesn't fit an expected pattern. These observations will be referred to as anomalies or outliers in this report. Illustrated in Figure 4.1 is an example of anomalies in a two-dimensional data set. N_1 and N_2 is the normal areas where most of the observations are, o_1 and o_2 are anomalies and O_3 is an area with multiple anomalies.

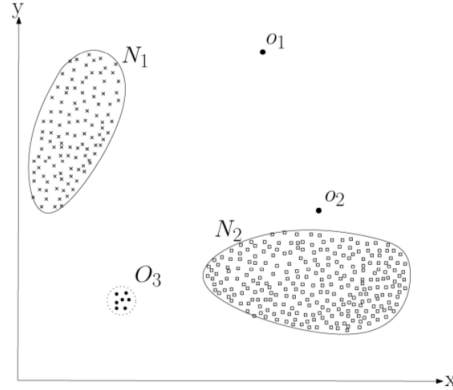


Figure 4.1: Example of anomalies in a 2D dataset[CBK09].

4.3 Scikit-learn

Scikit-learn was chosen as the machine learning library in this project. It is a Python library and I am most comfortable programming in Python. Alternatively could either Weka or R be used. Weka is a complete program with GUI, and has an API making it possible to integrate in a Java program. R is a special programming language which is specifically designed for statistical analysis. Scikit-learn depends on `numpy` and `scipy` to have the data in correct arrays and to do statistical analysis, and to create graphs is it necessary to use `matplotlib`. These do not follow in the regular installation of scikit-learn, and will have to be install on its own.

4.3.1 Models

Support Vector Machines (SVMs) is general terms for models that use supervised learning to analyse and recognize patterns in the data given as a training set. The problem must be binary, which means each point must belong to one of two categories. One-Class SVM (OCSVM) is a unsupervised SVM model, which is used used to solve outlier detection problems. In this project that was needed since the data was not labeled. Elliptical envelope was also used to see the difference between machine learning models. It is a covariance model, and is used to calculate a Minimum Covariance Determinant (MCD). MCD is a function that tries to find a proportion value of the correct observations which is then used to weight the observation to give a better representation [PVG⁺11].

Chapter 5

Results

With scikit-learn was a program created in python that reads through the dataset and puts the DNS calls in to an ndarray, which is a N-dimensional array. This is to easier use to right values for the machine learning model. The program also cleans up the dataset by removing some of the fields, which for this project seemed unnecessary. I created multiple scripts for taking out different sections of the dataset, to see if there would be any difference. The code for the parser program is in A

The main program read through a `csv` file, the dataset, and put each line as an array in a ndarray. N-dimensional array is a numpy class which is beneficial to use with scikit-learn. Further the ndarray was preprocessed to scale everything to a similar level. This level is set by calculating the mean and standard deviation of the ndarray.

The different classifiers are initiated with parameters that is changed to see what best fit the data. As the data is unlabeled, it is impossible to know how many of the observations are a part of a DNS tunnel if there is any. The classifiers must not have any inputs, but it will help making them more precise. In this project the contamination level, which is the level of data which is viewed as incorrect, was first calculated as the percentage of observations which was over average as a base. The value has been change up and down, but kept in the same area.

To be able to show the solutions and really understand what the machine learning model did, all of the experiments used two dimensions. The dimensions change between *downlink*, *uplink*, *duration*, *downlink/duration* and *uplink/duration*. The values was not marked with any type, but based on network theory is duration set to the of length of conversation in seconds. The uplink and downlink is the number of bytes transferred up to the server or downloaded from the server, respectively.

Since it was no known DNS tunnels in the dataset, it is impossible to use precision, recall or the F_1 measure as results. The precision and recall measurements are defined

in algorithm 5.1. The F_1 measurement was defined van Rijsbergen [?] as a way of combining precision and recall. The formula is in algorithm 5.2. As seen from the definitions without the knowledge of tunnels in the dataset will there not be a way of measuring the correctness of the machine learning algorithms.

Algorithm 5.1 Recall and precision definitions

$$recall = \frac{\text{Number of items of a category identified}}{\text{Number of items in the category in the dataset}}$$

$$precision = \frac{\text{Number of items of a category identified}}{\text{Number of items classified to the category}}$$

Algorithm 5.2 F_1 measure

$$F_1(R, P) = \frac{2RP}{R + P}$$

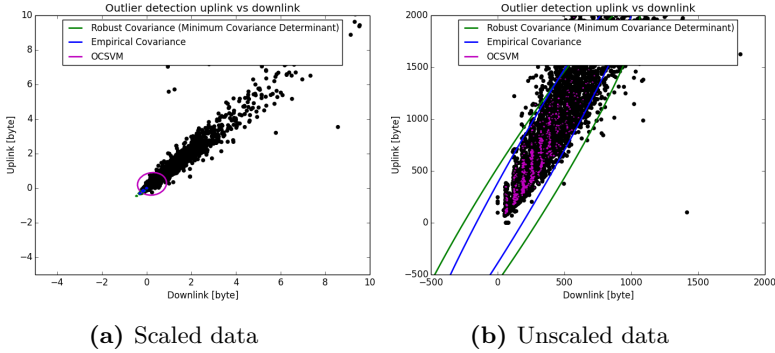


Figure 5.1: Difference between scaled and unscaled data

5.1 Scaled data vs unscaled data

As seen in Figure 5.1 the scaled and unscaled graphs bears the same characteristics. Most of the observations has almost the same value for uplink and downlink. The results were better for OCSVM with scaled data, while the covariance calculations worked better with the unscaled. The unscaled data spreads so far out that it is hard to see all the data points, while the scaled data is more compact giving a better overview. The ellipses on the figures is the learned decision of the classifier model. Meaning that inside the ellipses is the area where an observation is considered normal. Since the OCSVM is the main focus, scaled data was used for the rest of the experiments.

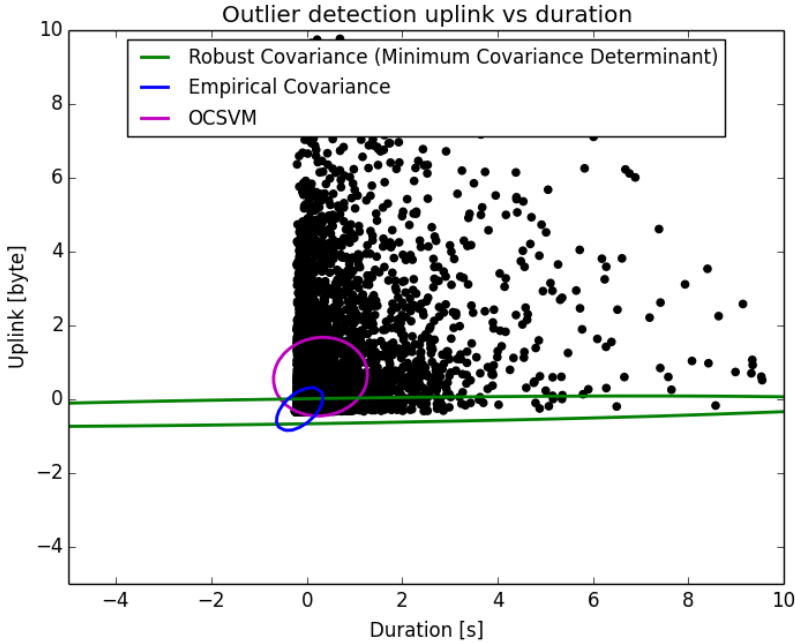


Figure 5.2: Uplink vs duration of session.

5.2 Different axes

Next up was looking which values would be best for the axes to represent the data. In Figure 5.1 the fields of the dataset used were `uplink` and `downlink`. This were the first thoughts of finding irregularities, as mentioned in chapter 3 this is a volume analysis comparing total volume from a user in one DNS session. Figure ?? shows graphs where other fields of the data set were used as input for the model. In Figure 5.2 the size of the uploaded data were compared to the time the session were used, this shows that there are a number of users in the same area which the model learn as the normal area, which in scale is around 0. The result were similar when changing uplink with downlink. There are users uploading and downloading lots of data in short time, this lead to the idea of introducing a new field to the dataset. Combining the discoveries from Figure 5.1a and Figure 5.2, by seeing how much data where uploaded or downloaded during a session. This is depicted in Figure 5.3. In this graph it is possible to see that the observation is even more clustered. This seemed like a good representation.

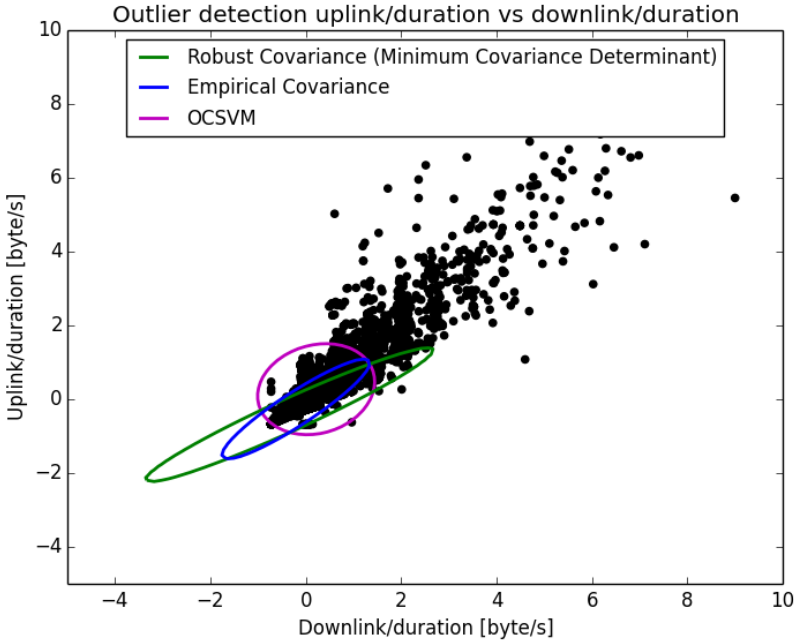


Figure 5.3: Average speed in each direction

5.3 Different parameters

The fields of the dataset is, as mentioned earlier, not the only variable when finding the best way for a model to fit to the dataset. The machine learning model has parameters that is given when instantiated. For OCSVM it is possible to set what kind of kernel it should use and the number of training errors among others. `Kernel` is the function the model should use. The number of training errors, `nu` in scikit-learn, is the upper bound fraction of outliers in the training set, and the lower bound of training examples used as support vectors. For this project the kernels used where `rbf` and `linear` where `rbf` is the one used in the graphs so far. The `nu` was first set to 0.15 which was calculated as the fraction of observations above average, this has been changed to see how it affect the learned decision. In do you see how changing `nu` is affecting the model. It is clear to see that in this dataset there are a low fraction of observations that does not fit the model. Using the linear kernel the model tries to find a linear line it can draw where all observations on one side is correct observations and on the other side is outliers. This is seen in . It was worth a try, but it seems that `rbf` if the best kernel.

Chapter 6

Conclusion

References

- [CBK09] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3):15, 2009.
- [Far13] Greg Farnham. Detecting dns tunneling. *InfoSec Reading Room*, 2013.
- [MD88] P. Mockapetris and K. J. Dunlap. Development of the domain name system. In *Symposium Proceedings on Communications Architectures and Protocols*, SIGCOMM '88, pages 123–133, New York, NY, USA, 1988. ACM.
- [Moc83] Paul V Mockapetris. Domain names: Implementation specification. 1983.
- [Moc87] Paul V Mockapetris. Domain names-concepts and facilities. 1987.
- [PVG⁺11] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

Appendix

Parser program

