

Identifying Financial Frauds on Darkweb

by

Krishna Tushar Dharaiya

A Thesis Presented in Partial Fulfillment  
of the Requirements for the Degree  
Master of Science

Approved April 2018 by the  
Graduate Supervisory Committee:

Paulo Shakarian, Chair  
Adam Doupe, Member  
Yan Shoshitaishvili, Member

ARIZONA STATE UNIVERSITY

May 2018

## ABSTRACT

Data breaches have been on a rise and financial sector is among the top targeted. It can take a few months and upto a few years to identify the occurrence of a data breach. A major motivation behind data breaches is financial gain, hence most of the data ends up being on sale on the darkweb websites. It is important to identify sale of such stolen information on a timely and relevant manner. In this research, we present a system for timely identification of sale of stolen data on darkweb websites. We frame identifying sale of stolen data as a multi-label classification problem and leverage several machine learning approaches based on the thread content (textual) and social network analysis of the user communication seen on darkweb websites. The system generates alerts about trends based on popularity amongst the users of such websites. We evaluate our system using the K-fold cross validation as well as manual evaluation of blind (unseen) data. The method of combining social network and textual features outperforms baseline method i.e only using textual features, by 15 to 20 % improved precision. The alerts provide a good insight and we illustrate our findings by cases studies of the results.

*Dedicated to my parents for their undying faith in me*

## ACKNOWLEDGMENTS

I would like to take this opportunity to express my deep gratitude to my advisor, Dr. Paulo Shakarian for his unwavering support and guidance for my Master's thesis research. I shall forever be thankful for the mentorship and invaluable lessons I learned working under him.

I am eternally grateful to Jana Shakarian for her valuable inputs and continuous support.

I would like to thank Dr. Adam Doupe and Dr. Yan Shoshitaishvili for being on my committee and the support for my dissertation.

I thank Ph.D student Eric Nunes for reviewing the document and for his insightful suggestions.

I would also like to thank all my fellow CySIS lab members for their support and motivation.

This research is supported by the Office of Naval Research(ONR) Neptune program, the ASU Global Security Initiative(GSI) and the Intelligence Advanced Research Projects Activity (IARPA) via the Air Force Research Laboratory (AFRL). The U.S Government is authorized to reproduce and distribute reprints for Governmental purposes not withstanding any copyright annotation there on. Disclaimer: The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements either expressed or implied, of ONR, IARPA, AFRL or the U.S. Government.

# TABLE OF CONTENTS

	Page
LIST OF TABLES.....	vi
LIST OF FIGURES.....	vii
CHAPTER	
1 INTRODUCTION.....	1
2 BACKGROUND.....	3
3 APPROACH AND SYSTEM OVERVIEW.....	4
3.1 Automatic Tagging.....	4
3.1.1 Contextual.....	5
3.1.2 Social Network.....	6
3.2 Alerts/Trending Threads.....	7
4 DATASET.....	9
4.1 Automatic Data Labeling.....	9
5 EXPERIMENTS.....	12
5.1 Automatic Tagging.....	12
5.1.1 Tagging By Contextual Features: Experiment Results.....	12
5.1.2 Tagging By Social Network Features: Experiment Results.....	18
5.2 Trending Threads/Alerts.....	20
6 RELATED WORK.....	28
7 CONCLUSION.....	30

	Page
REFERENCES.....	31
APPENDIX	
A The List Of Regular Expressions Used For Cleaning Data.....	33
B List Of Regular Expressions Used For Labeling Each Of The Categories.....	35

## LIST OF TABLES

Table	Page
4.1 Statistics Of The Data.....	9
4.2 Regular Expression Used For Every Category.....	10
4.3 Statistics Of The Training Data.....	11
5.1 Top Threads For Categories Based On LR Model Scores.....	16
5.2 Precision In Finding Top-100 Threads.....	17
5.3 Precision In Finding Top-100 Threads.....	18

## LIST OF FIGURES

Figure	Page
3.1 System Overview.....	5
5.1 Performance Of The Classification Models.....	14
5.2 Number Of Threads Classified As A Category.....	15
5.3 Posts By A User Having High SN <sub>FULZ</sub> Score.....	19
5.4 Trending Threads Over Time.....	21
5.5 No Of Users Replying Over Time To Thread “ <i>H*Acked Carded Western Union Moneygram &amp; Bank Transfers (Only 15% Of Total Amount)</i> ”.....	24
5.6 No Of Users Participating In Trending Threads Relating To Different Fraud Categories Over Time.....	25
5.7 No Of Users Replying Over Time To Thread “ <i>Western Union In 15 Minutes Bank Transfers In 2 Hours</i> ”.....	26



## CHAPTER 1

### INTRODUCTION

The financial and retail sectors are among the top targeted sectors when it comes to data breaches [8]. The number of incidents of data breaches have been growing exponentially over the past few years, with personal data and credentials totaling in billions. The number of incidents reported in 2016 for data breaches relating to the financial industry sector and retail industry were 998 and 321 respectively. The sheer numbers associated with these frauds should be indicative enough of the intensity of the situation. Often, it can take a few months and in some cases up to years to actually detect that a data breach has occurred [22]. Generally, the motive behind such breaches/attacks is lucrative cash for the criminals. The data obtained from such attacks are often then sold in black markets like deep and dark web. Hence tracking such activities on the darkweb can potentially lead to timely detection of data breaches and provide insight for better protection against such attacks.

We see a large variety of products and services on sale on darkweb and deepnet websites from stolen credit cards to people selling credentials, account and personal information such as names, address, social security numbers, credit card numbers, medical records from retailer accounts, insurance companies, banks, payment services accounts etc. Going forward in this thesis we refer to all such information collectively as credentials, accounts, and personally-identifying (CAP) information, a super-set of personally identifying information (PII).

In this thesis, we identify and classify threads in Darkweb/Deepnet forum that sell sensitive information like CAP and PII into categories based on the source of information being sold using context analysis into 4 categories listed below:

1. Bank account fraud (Information obtained from compromised bank accounts)
2. Fulz fraud (Information obtained from physical or virtual stolen credit/debit cards)
3. Payment Services fraud (Information obtained from compromised payment services accounts like paypal accounts/venmo accounts)
4. Retail fraud (Information obtained from retail accounts like Amazon, Ebay etc)

The specific contributions of this thesis are listed below:

- We assemble a filtered dataset of labeled Darkweb/Deepnet forum threads that sell stolen information categorized into four categories (Bank account, Fulz, Payment services, Retail fraud) based on keywords/regular expression searches.
- We frame identifying financial threats as a multi-label classification problem and leverage several machine learning approaches. We build features based on the thread content (textual) and social network analysis of the users communicating in the thread.
- Experimental results achieving an accuracy of more than 80% in identifying the financial threats and comparison between several machine learning approaches.
- An alert system for trending financial threats on the darkweb by monitoring thread activity – identifying anomalous activity showing increased interest of users in a financial threat category.

The rest of the thesis is organized as follows, section 3 and 4 describe the approach and dataset. Section 5 gives details regarding the results of the various experiments and finally related work is discussed in section 6.

## CHAPTER 2

### BACKGROUND

‘Deepnet’ refers to the part of the internet that is undiscoverable by standard search engines. ‘Darkweb’, colloquially, refers to a distinct network supporting cryptographically hidden sites [1]. Darkweb can be considered as a part of Deepnet that is not accessible by standard browsers. Criminal activities in cyberspace are increasingly facilitated by such burgeoning black markets and forums in both the tools (e.g., exploit kits) and the information (credit/debit card information etc) [2]. Markets for hacking tools, hacking services, and the fruits of hacking are gaining widespread attention as more attacks and attack mechanisms are linked in one way or another to such markets [2]. Carding shops are another important part of the global hacker and cybercriminal community. Carding shops help facilitate cyber carding crimes as they provide a supply chain for carders who wish to sell stolen cards [2].

As opposed to the Darknet/Deepnet market places which have a more structured form, the forum websites are places for discussions where threads vary from discussing the new phone release to selling items like zero day exploits, malwares and sensitive information from various sources [5]. As compared to the marketplaces, the number of forums is much larger and correspondingly generates much more data.

## CHAPTER 3

### APPROACH AND SYSTEM OVERVIEW

In order to address the challenges of identifying CAP and PII information on the deepweb and darkweb at scale, one must consider certain desired characteristics. The main challenges deal with identifying information relating to CAP and PII that is both timely and relevant. Hence, we view the below aspects as desirable in such a system.

- **Automatic Tagging:** Given that the amount of data generated every day in these forums is large and the discussions have a wide variety, it becomes important to identify relevant threats and prioritize the top threats at a given time.

The system analyzes a given thread and tags it as belonging to one or more of the four fraud categories described in section 1 or none.

- **Trending Threads/Alerts:** Trends/Anomaly in the darkweb threads can be indicative of new items being out for sale or imminent attacks. Hence our proposed system uses temporal features of the threads to generate trending threads/alerts to indicate unusual behavior in the activity of threads that could indicate such events.

The following block diagram gives an overview of the system. The data received from the system described in [5] is used to create labeled data set automatically using a set of regular expressions given in Appendix A. This labeled data is then used to train machine learning models for automatic labelling.

The tags from this stage are then used along with the temporal features to generate alerts.

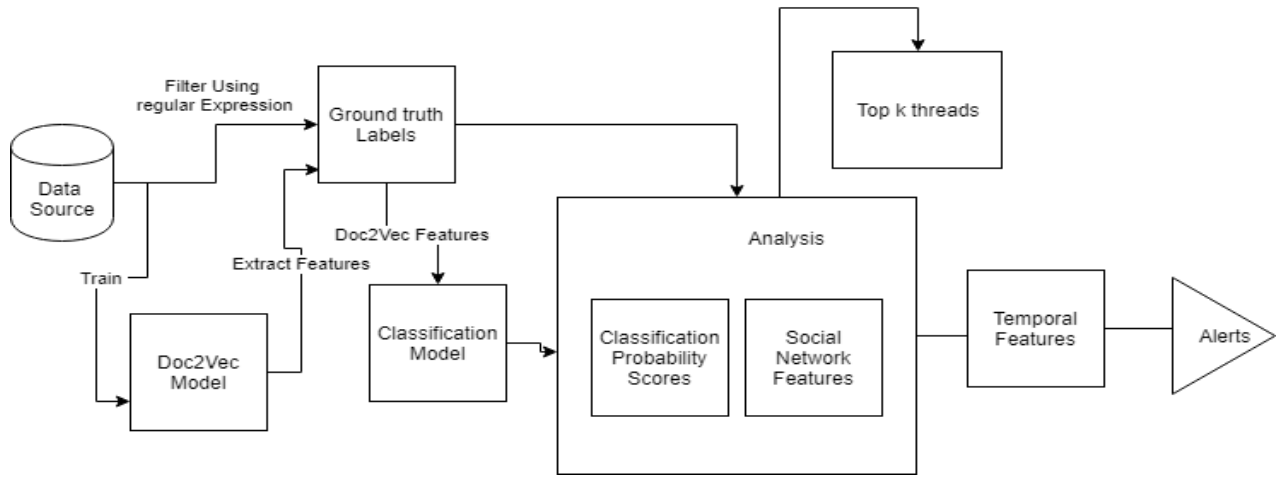


Figure 3.1 System Overview

Since discussions on darkweb data range from pornography, selling armed goods, gold to hacking of computers, vulnerabilities and even gift cards. To restrict our scope, we do an initial filter to remove threads that are related to pornography, drugs and guns/arms. The regular expressions used for data cleaning can be found in Appendix A.

### 3.1 AUTOMATIC TAGGING

For the automatic tagging of the threads we use two feature sets described as following

- Contextual: Features deriving from the actual text of the post

We use the text of the post to understand the intention of the user. Naturally certain words and phrases can be an obvious indication that a particular text is talking about a fraud. The document embedding (Doc2vec) helps understanding the semantics of the dataset and these vectors i.e, essentially a representation of words and phrases of the text, are then used by the classification model.

- Social Network: Features deriving from the user that actually posted the text.

Past studies have demonstrated that users having similar interests tend to participate in the similar threads, hence a large number of users that have past history of interest in financial fraud replying in a thread can directly mean that the thread itself is about financial fraud.

### 3.1.1 CONTEXTUAL

There are many methods to extract features from contextual data. Some of the popular ones include bag-of-words, TF-IDF (term frequency-inverse document frequency), word frequencies but the drawback of such methods is they do not capture semantics and co-occurrences in different documents.

To derive contextual features which also capture the semantics we use Doc2Vec Modelling. Doc2vec is an unsupervised algorithm implemented in the python genism library [9] [10] [11]. It is built on top of the word2vec model and uses word embeddings to create vectors of text. The Doc2Vec model is trained using 10,000 posts which includes all the labelled posts for all categories and the remaining posts are arbitrarily picked from our dataset. We use dimensionality as 100 for the model. We then infer vectors from this doc2vec model for our training data to train our supervised classification models.

### 3.1.2 SOCIAL NETWORK

In this section, we describe the construction of social network data from the training data set we constructed in section 3.1.

Various researches in the past have found social communities among users having similar interests [13] [14]. Based on the same intuition, if a thread has many users that are

interested in a category then it increases the chances of that thread of being of the same category.

We take the unique users from our data set and assign them scores based on the number of posts by that user in the training dataset and assign zero for any users which have no posts in the training dataset. We will refer to these Social Network scores as the  $SN_{category}$  score. For example, if user1 has 2 posts in the training set for payment frauds and 6 for bank account frauds then user1 gets  $SN_{payment}$  2 and  $SN_{bank}$  6. We then calculate the total  $SN_{category}$  for each thread as a sum of  $SN_{category}$  of all the users participating in that thread.

### 3.2 ALERTS/TRENDING THREADS

In many intrusion detection and user behavior analysis systems, user activity plays a very important role [15] [16]. User Activity can range from simple independent activity like number of webpages visited per unit time to a sequence of activities like tasks performed for logging into a system in a unit of time. Research work in the past has been successfully able to learn the patterns of activity and provide alerts when encountering outlying or unusual behaviors [15] [16].

The problem statement for such algorithms can be simplified as what defines outliers and detecting outliers that do not fit an ongoing pattern. Similarly, for detecting trends on social network websites like twitter, patterns in interests expressed by users in particular topics is used as a metric. We can view threads in a darkweb as an entity and try to point out threads that are trending or have outlying patterns of activity. We consider two components for monitoring activity of a thread: number of users and number of posts

As mentioned in section 2 each post has a date associated with it. We use this to create temporal data for each thread. We divide the data into bins of week and then maintain the count of number of users for each week of the year for each thread. Using the number of replies/posts as a metric can be misleading because many a times a large number of posts are observed in a thread in short period of time but only a couple of users are involved in the discussion, hence it does not have a larger influence. We create the following features for each thread for thread behavior analysis:

- Standard Deviation in Number of Users
- Moving average in Number of Users

Here standard deviation is given by  $\sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N-1}}$  where n= number of points,  $x_i$  is each point from 1 to n and  $\bar{x}$  is the mean of the data set. N is the sample size.



## CHAPTER 4

### DATASET

For this thesis, we use data collected by a commercial version of the system described in [5], below is a table showing the statistics of the entire dataset.

Forums	132
Users	195,499
Posts	177,0975
Threads	205,441

Table 4.1 Statistics of the Dataset

The data consists of forums that are in the English language. The conversations in such forums can be viewed as the question-answers on sites like StackOverflow, Quora, Yahoo Answers etc. where there are multiple replies by various users to a question posted, similarly in a Darknet forums we have topics instead of questions and rest follows the same pattern. In this thesis, we will refer to the unit of one topic and all replies(posts) to it as a thread. The following section describes our process of automatically generating labeled dataset.

#### 4.1 AUTOMATIC DATA LABELING

One of the major challenges for identification of categories is the construction of training data as in any other supervised learning approach. Given the size of the dataset, it becomes difficult to manually identify enough ground truth for training the model in a reasonable amount of time. In this approach, we use rules built on top of a set of regular expressions to identify posts related to each category. We then manually narrow down from the derived data set to create ground truth for each category in a much faster way. The following table

shows some of the regular expressions used for each category. The full list of regular expressions used appears in Appendix B.

<b>Category</b>	<b>Regular expression</b>
Bank account frauds	Bank account, routing number, bank login, Credentials + bank, bank drop, name, address, ssn Password, username, email address, address, bin, account number
Fulz Frauds (Cards)	Cvv, credit card, fulz, debit card, cv, name, address, ssn, email address, address, phone number, mastercard, visacard
Payment Services Frauds	Paypal, login, credentials, venmo, name, address, ssn Password, username, email address, address, account, webmoney, phone number, contact me
Retail Frauds	Target.com, credentials, account, amazon.com, name, address, ssn Password, username, email address, address, phone number, ebay, Netflix, hulu, bestbuy, nordsdrom

Table 4.2 Regular Expression used for every Category

Each of these keywords present in the table are turned into regular expressions to identify various forms of the word and identify patterns, for example, for the keyword phone number we have regular expression to identify the word phone and number along with +1-123-123-1234 i.e a phone number mentioned. These are then used in combination to form multiple rules that are then used to identify the posts. We built the following training dataset after manually selecting posts from the dataset derived using regular expressions.

<b>Category</b>	<b>No. of Identified Posts</b>	<b>No. of unique Users</b>
Bank account frauds	378	290
Fulz Frauds (Cards)	567	269
Payment Services Frauds	317	260
Retail Frauds	240	121

Table 4.3 Statistics of the Training Data

## CHAPTER 5

### EXPERIMENTS

In this section, we evaluate the results for the experiments conducted for classification and identification of top-k threads and Alert systems. We discuss results based on various evaluation methods and also illustrate findings using case studies of some results

#### 5.1 AUTOMATIC TAGGING

The intensity of impact of the stolen data varies depending on its source. Considering the large amount of data generated on darkweb websites on daily basis, we need a robust automatic system that provides relevant information. In this section we discuss the automatic tagging results.

##### 5.1.1 TAGGING BY CONTEXUAL FEATURES: EXPERIMENT RESULTS

In order to do a robust evaluation, we do both a 10-fold cross validation as well as manual evaluation of blind(unseen) data. The advantage of k-fold cross validation is that every data point gets tested. It gives a good idea about the model's performance if the training data is a good representation of the entire dataset. In cases where dataset is constantly evolving and training set is much less, manual evaluation of the performance of the model on unseen data can provide insight and help improve the model. We now discuss the results from both evaluation methods.

##### 10-FOLD CROSS-VALIDATION RESULTS

We train 4 machine learning models for each category mentioned using the doc2vec word embedding features discussed in section 3.2.1. The supervised methods include the well-

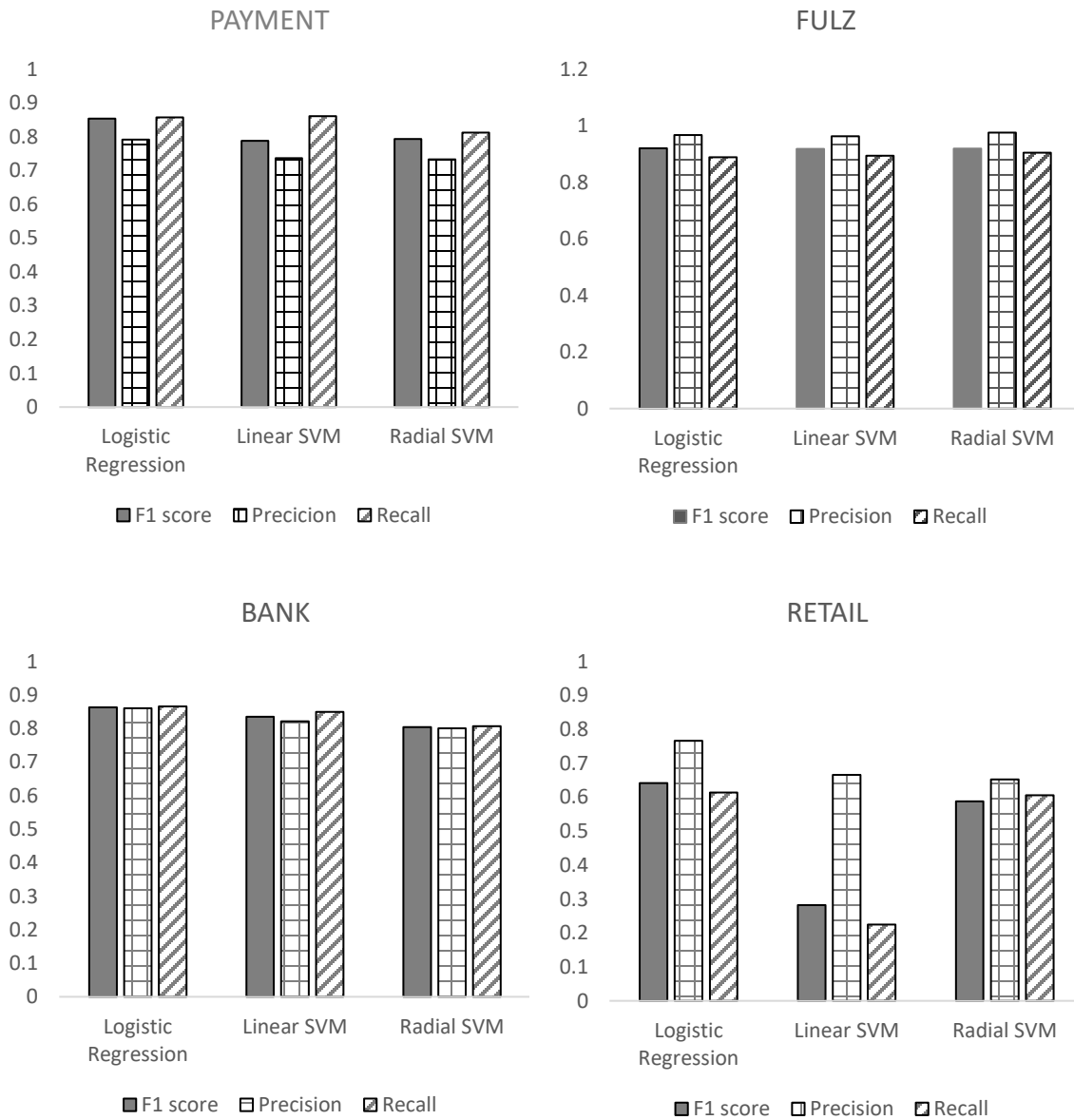
known classification techniques of radial support vector machine, linear Support Vector Machine and logistic Regression. As a first step of evaluation we use cross validation. Cross-validation is a technique to evaluate predictive models by partitioning the original sample into a training set to train the model, and a test set to evaluate it. In k-fold cross-validation, the original sample is randomly partitioned into k equal size subsamples. Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k-1 subsamples are used as training data. The cross-validation process is then repeated k times (the folds), with each of the k subsamples used exactly once as the validation data [12]. We use 10-fold cross validation. The results of the cross validation are then quantified using 3 metrics, Precision, Recall and F1 score. Precision is the fraction of relevant instances among the retrieved instances, while recall is the fraction of relevant instances that have been retrieved over the total amount of relevant instances [21]. They can be defined as

$$\text{Precision} = \frac{\text{True Positives}}{\text{True positives} + \text{false positives}}$$

$$\text{Recall} = \frac{\text{True Positives}}{\text{True positives} + \text{false negatives}}$$

Here For classification tasks, the terms true positives, true negatives, false positives, and false negatives compare the results of the classifier under test with trusted external judgments. The terms positive and negative refer to the classifier's prediction (sometimes known as the expectation), and the terms true and false refer to whether that prediction corresponds to the external judgment (sometimes known as the observation) [21].

F1 score is a measure of a test's accuracy. The F1 score is the harmonic average of the precision and recall, where an  $F_1$  score reaches its best value at 1 (perfect precision and recall) and worst at 0 [20]. We calculate these results only on positive labeled dataset. The following figure shows the performance comparison of the 3 supervised learning approaches for each category.



### Figure 5.1 Performance of the Classification Models

Since Logistic Regression gives the best performance in all categories except Fulz where radial SVM gives slightly better results, we conduct the remaining experiments using the logistic Regression(LR) Models for all categories.

### EVALUATION ON BLIND DATA

We now use the classification scores for all posts in our test data set and try to find top 100 threads for each category. The following figure shows the distribution of the categories when classified by the logistic regression (LR) model. Here test data is the entire data set, statistics of which are mentioned in Table 1. For this we take the average classification score of each post in a thread and find the top 100 threads having the highest average classification scores. We set a threshold and only consider posts that have a classification score more than 0.90 while calculating this average.

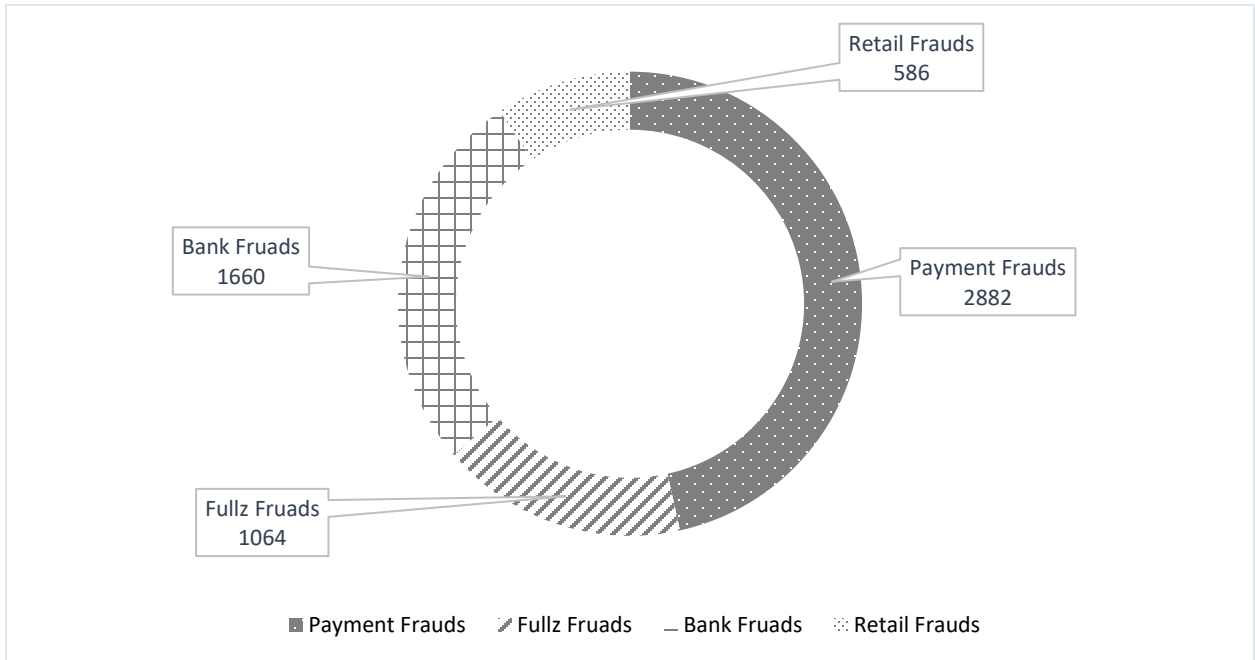


Figure 5.2 Number of threads classified as a Category

The next table shows few thread names from the top-100 threads for all categories

<p><b>Payment Frauds</b></p> <ol style="list-style-type: none"> <li>1. selling, cvv2, fulz, dumps+track 1\2, paypal+contact me : icq : 710530177</li> <li>2. need paypal account sending 1k usd</li> <li>3. i need a paypal account with balance</li> <li>4. sell cvv good-dumps-wu transfer-bank login-acc paypl-shipping</li> </ol>	<p><b>Fulz Fruads</b></p> <ol style="list-style-type: none"> <li>1. selling, cvv2, fulz, dumps+track 1\2, paypal+contact me : icq : 710530177</li> <li>2. selling 100% valid cards ( us - uk - ca -eu- spain - fulz+dob+ssn+pass</li> <li>3. selling, cvv2, fulz, dumps+track 1\2+acc, paypal+contact: icq : 710530177</li> </ol>
<p><b>Bank account frauds</b></p> <ol style="list-style-type: none"> <li>1. free dumps track2</li> <li>2. i'm seller cvv + dumps + track1&amp;2 paypal + do wu transfer + bank login + atm plastic</li> <li>3. free full info ssn dob mmn dl daily - csu.su - carding forum - (carding, card fraud, carding forum, carders forum, carders board, darknet,</li> </ol>	<p><b>Retail Frauds</b></p> <ol style="list-style-type: none"> <li>1. cc, gift cards, wu, paypal and banks transfers</li> <li>2. cloned credit cards, gift cards, wu, bank, paypal transfers</li> </ol>



blackhat, darknet forum, crdclub, shadow brokers, darknet markets, credit card fraud, fraud, atm fraud, credit card scams, atm skimmer, dumps shop, credit cards cvv, credit cards cvv2, dumps, dumps with pin, cvv2, buy dumps, buy credit cards, dumps with pin)	
--	--

Table 5.1 Top Threads for Categories based on LR model scores

To measure the performance of our LR models in finding top-100 threads, we manually evaluate the threads and calculate a precision score for each category. For the evaluation of each financial fraud, we go through every post of the thread and see if the thread contains one or more posts that claim selling stolen data for example, they show samples of the data available or provide descriptions and pay rates etc, from the correct source i.e if we are evaluating bank fraud results the post should be selling information relating to banks.

The following table shows the precision.

Fraud Categories	Precision
Bank	0.42
Fulz	0.50
Payment	0.40
Retail	0.27

Table 5.2 Precision in finding top-100 threads

## DISCUSSION

It can be observed from the results that there is some difference in the accuracy when checked with unseen data. Some of the things to be noted here would be that the results

on blind data are not simple text document classification results as with the cross validation. We then decide if a thread, which is a series of text documents, should be assigned a tag based on the individual scores of each document. Generally, only 10-20% of the entire threads have the actual content that describes the fraud or products that are sold, rest of the threads are conversational replies. Hence the variance in the results.

After going through the threads, we discover that many of the misclassified threads are very similar to the frauds in the identified categories or act as an auxiliary. It leads to discovering new aspects that are not present in our training set.

We now illustrate some examples of our findings. Some of the misclassified threads were trying to sell passport, visas, driver's license, green card etc. The data sold in such postings is also essentially CAP information. An excerpt of a post advertising to sell such information is given below.

```
“buy registered passportsvisasid carddriver license green cardsshello we are the best
producers of high quality counterfeit banknotesgetting a fake and a real(genuine) passport
id or drivinglicense or any other document is simple. we can make you both real andfake
documents.however the real documents are more expensive than the fake becauseit takes
time skill and contacts to get it done. note that the fakeis going to be 100% unique and in
very good quality. the difference isbased on the registration of the numbers. the real
document will beregistered with the country^s database so you can use..”
```

Many of the misclassified threads were socks list, proxy servers are used to access fake account details and are often sold with bank or payment services account details.

Consider this excerpt from a misclassified post

```
“have fun:http://pastebin.com/xqzhy1mthacked payment gateway details      ^authc^
=> ^4b4hvwdjrg^      ^gatew^ => ^8mk2wwfye9^snippet:      // process the
transaction      $result=$this-
```

```
>do_curl("https://platformpay.com/client/5ppmrc/ccprocess.php"$fields) ^authc^  
=> ^4b4hvwdjrg^ “
```

it contains payment gateway details that can be further exploited.

### 5.1.2 TAGGING BY SOCIAL NETWORK FEATURES: EXPERIMENT RESULTS

We now combine the contextual and social network features to find the top-k threads. Since we have the social network i.e user information for only those forums that are represented in the ground truth, we limit our dataset to only those 60 forums for this experiment.

#### EVALUATION

As in section 3.2.1 we consider the average classification score of the thread and the total  $SN_{category}$  score of the thread as the sum of  $SN_{category}$  score of each post in the thread. We then find the top-100 threads for each fraud category and calculate the accuracy by manually analyzing the results. In order to have a comparison, we find the top-100 threads using only the contextual features on the same 60 forums dataset used for conducting the social network experiment.

Table 5.3 shows the precision scores of the experiment.

Fraud Category	Precision: Social Network + Contextual	Precision : Contextual
Bank	0.72	0.49
Fulz	0.85	0.72
Payment	0.60	0.42
Retail	0.42	0.25

Table 5.3 Precision in finding top-100 threads

As compared to the accuracy when we only consider classification scores, we see that the combination of the two features performs about 15-25% better for classifying all the 4 fraud categories.

We can extend this by maintaining the user information over all forums in the future. Such an approach can also help in identifying and tracking key actors that contribute to such frauds and can also be helpful in generating alerts.

## DISCUSSION

In most social network websites users generally have some set of interests and will only interact in topics relating to those interests. We use this intuition to better the classification of financial frauds on darkweb. If a thread has many users who having a past history of showing interest in a particular type of financial fraud, the chances of that thread also being about the said financial fraud is high. As we see in our results that when we combine classification scores with these social network features we see an increased precision on the blind dataset. Most of the users that have a high SN score (more than 5) are vendors for financial fraud services hence a large SN score indicates that a lot of these vendors are posting in the thread. To give a better idea of this, let us consider one user in our database which has a  $SN_{FULZ}$  of 6. The following are few of the posts by this user in darkweb

darkode hacked selling quality fullz cvv dumps+pin carding service selling quality fullz cvv dumps+p
hello all buy credit elite carding forum seller cvv good: us.uk.eu.au.ca
carder forums 2016 sell cvv fullz 2016, transfer wu, bank transfer,

Figure 5.3 Posts by a user having high  $SN_{FULZ}$  score

We see that all these posts by the user are advertisements about selling fulz products.

## 5.2 TRENDING THREADS/ALERTS

In this section we describe our approach for the trending threads/alert systems and elaborate on our findings using specific examples.

## APPROACH

Hackers use darkweb forums for communication as well as for exchange and sale of goods and services. Increased participation in a conversation in a thread may imply important or interesting information being given out on the thread. If a thread is being used for selling data/services, increased activity may indicate discounts being offered, new products on sale. With respect to financial frauds, such activities can be indicative of new data leaks. Detection/Monitoring of such out of norm activity in a thread can help in identifying trending threads as well as important actors involved in such frauds. Such trends can be useful in giving an insight on imminent threats.

As described in section 4.2, we define the activity levels of threads in units of weeks and then identify the trending threads based on the interest shown in a thread. We use the classification scores from the earlier experiments and generate trending threads for financial frauds for each week over the time period of January 2016 to August 2017. These topics are the most active or have exponential growth in the number of users and can also be viewed as alerts indicating threads that are potentially making large amount of sales. In this time of about 18 months the system identifies 190 threads in total, that were trending for the financial frauds in darkweb. Figure 6.4 shows the number of trending threads or alerts over time. Here the graph plots all threads irrespective of the fraud category it belongs to it, hence we can say it shows the threads trending in the financial fraud sector as a whole.

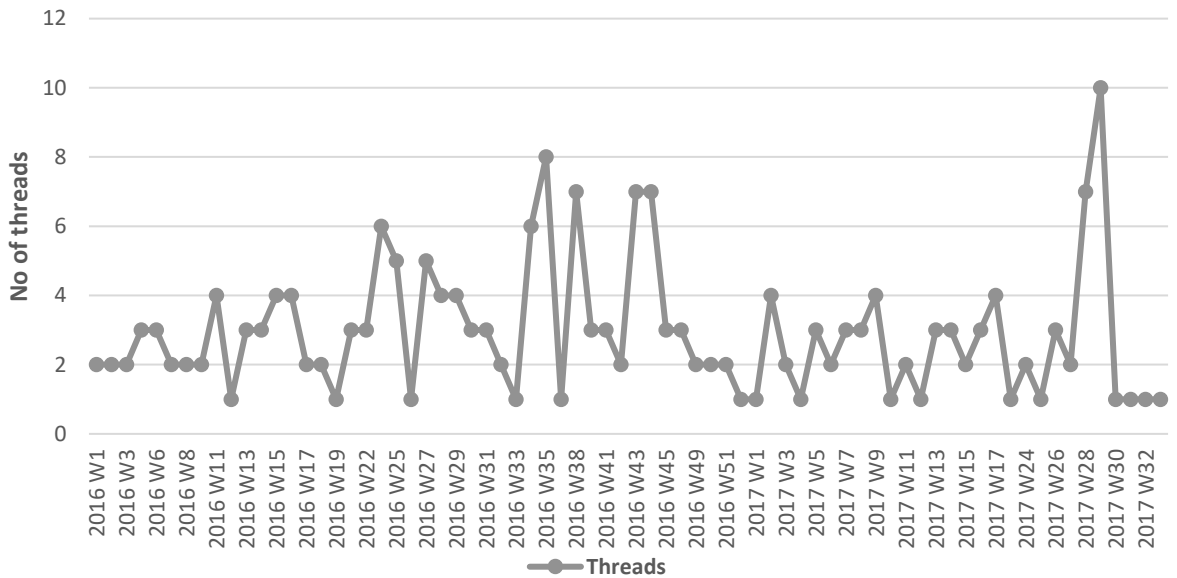


Figure 6.4 Trending Threads over time

After analyzing the results further, we observe that financial fraud threads are among the top threads in terms of number of users replying to the thread. In fact in our entire dataset, we find that the thread having the highest number of users replying in a week happens to be “ h\*acked carded western union moneygram & bank transfers (only 15% of total amount) ” having 873 users posting to the thread in the 29<sup>th</sup> week of 2017, which is discussion about basically selling bank account information and western union account information that can be used to make money. This shows that there is a great deal of interest from users on darkweb in financial frauds and they are among the popular topics. The following is an excerpt of the post by the user(the vendor) that is offering the services.

last updated: 20th april 2017my official email for orders - mayor@tutamail.com{ { hacked & carded - western union moneygram & bank transfers instantly worldwide for just 15% of total amount. } }hello guys i specialize in carding and transferring funds from hacked bank accounts & credit card data to you by western union money gram or bank transfers worldwide. please note this is illegal and fraud com please if you are new to this then contact me for advice before you proceed on ordering anything from me. prices are in usd

@ 15%: (i can also send in usdgbpetc as it will be automatically converted to your country currency after transfer).world wide western union & moneygram prices:you receive: €/£/\$ 1200 you pay: €/£/\$ 180 [up to 2 separate transfers x2 mtcn codes - €600 each]you receive: €/£/\$ 2000 you pay: €/£/\$ 300[up to 3 separate transfers x3 mtcn codes - €660 each]you receive: €/£/\$ 3000 you pay: €/£/\$ 450 [up to 4 separate transfers x4 mtcn codes - €750 each]\*these are the maximum amounts i do per transfer.\*security question & answer is available and no i.d is necessary just request this option (not available for all countries)world wide bank transfer priceswe are offering transfers to your bank accounts for 20% (lesser for high amount transfers.you receive: €/£/\$ 1600 you pay: €/£/\$ 250you receive: €/£/\$ 2500 you pay: €/£/\$ 350you receive: €/£/\$ 5000 you pay: €/£/\$ 600i have endless suppliers of hacked bank accounts in most countries com transfers are instant if i don^t have hacked bank account in your country transfers will take up to 2 days. \*please scroll down to the bottom for more pictures.i can provide mtcn codes receipts etc. of completed transfers from other clients. just ask me

---

## special custom listing for transfers above \$5000.## us pre-pays worked for us during testing. can^t assure for 100%## the funds are semi-clean by channeling through a number of payment processors and bank accounts.## chargeback possibility is there but its very low. in our testing only 1 account got chargebacked that too after 3 weeks.## while placing your order you are requested to provide correct bank account details to avoid further hassles.## pm us in forum to check availability of transfers to countries other than usa & eu.quote:what information is needed to do your western union & moneygram transfer:first namelast namecountrycity/statecurrencyemail addressquote:what information is needed to do your bank transfer:bank namebank addresszip codeaccount holderaccount numberaccount typerouting numberswift numberbic and iban \*not all info is needed depends on what country your bank account is located in.faq^s how quick are transactions completed ? answer : total completion time is 30 minutes maximum. from payment to receiving mtcn collection codes. all western union transfers are sent via the western union money in minutes service. how do we guarantee complete anonymity ? answer : firstly you are advised to use fake names when collecting payment as id is not required with our western union transfers. instead all of our transfers have a security question attached. and you are provided with the answer. id is not required if the security question/answer feature is used. this security feature is used as default on all our western union transfers. (read more here <https://www.westernunion.com/us/en/s...curity-qa.html>)secondly and more importantly our total control of agent terminals (administration level control) gives us complete access to the western union network. we therefore have the ability to delete our transactions from the overall transaction list on the western union network. how do we actually do western union transfers what is our method?answer : we use kvm devices at various western union agents. once we place our kvm on a terminal we technically have full agent control. meaning we can transfer and also delete transactions on the western union network.we only target large western union agents and we always change locations. we never use cc^s fulls or software to do transfers.how successful are our western union transfers ?answer : our western union transfers have a 100% success rate. our service is gold rated by the mayor (admin) and the moderators on this forum.what are the rules ?answer : our rules are straight forward. do not waste our time and collect transfers within 48 hours of receiving the mtcn.

also remember to leave good feedback on our page when we have completed your transfer for you.why don't we send unlimited western union transfers to ourself ?answer : we do have cashiers we are currently using for per centage (%) based work and we do encourage repeat customers to go through the verification process and join our worldwide team of cashiers. however the amount of available transfers cannot be handled by our current team of cashiers in europe and the u.s. its simply too much to send. this surplus of transfers are instead sold here. we rather sell transfers for a secure value then to allow our controlled western union terminals to be unused or not used to their full potential.how do you pay?answer : we are now accepting all modes of payment for this service.how quickly can you get started?answer : simply contact the mayor and select our service. as you all know all first transactions are done via the mayor's escrow to insure we are all safe. once we complete a couple of transfers for you then we will have direct contact.terms & conditions:• we are not responsible for any account restriction make sure your accounts are fully functional.• we will need to access your account to confirm that funds have been reversed.• you are not allowed to open a dispute on darkgeo after a successful transaction.• you are not allowed to leave negative feedback if we deliver your order.• we may change any of the terms of service at any moment without prior notice.- escrow accepted -----  
-----: payment methods :-----• btc (bitcoin)• pm (perfectmoney)• wu (western union)• mg (money geam)any questions or would like to order please contact me in private message or email me (email on top of thread) escrow is accepted

The excerpt gives a good idea about the type of services being offered. The following graph shows the activity levels for this thread i.e the number of users replying in the thread. As we see the thread got unexpectedly large number of users replying in the first week itself. This is purely due to the popularity of the service. In the replies to this thread, people are asking for/buying services or people thanking the vendor for a successful buy. In short, the high point is showing large number of sales being made in that thread.



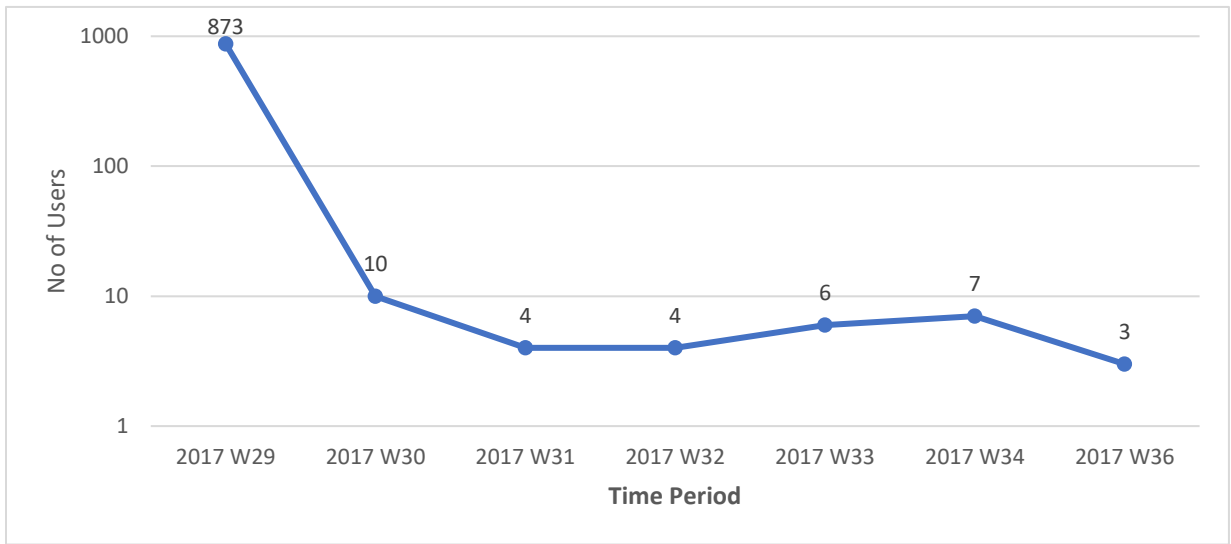


Figure 5.5 No of users replying over time to thread “h\*acked carded western union moneygram & bank transfers (only 15% of total amount)”

## DISCUSSION

In order to gain more insight about the popularity of a fraud category, we plot the number of users replying to a trending thread belonging to a fraud category over time.

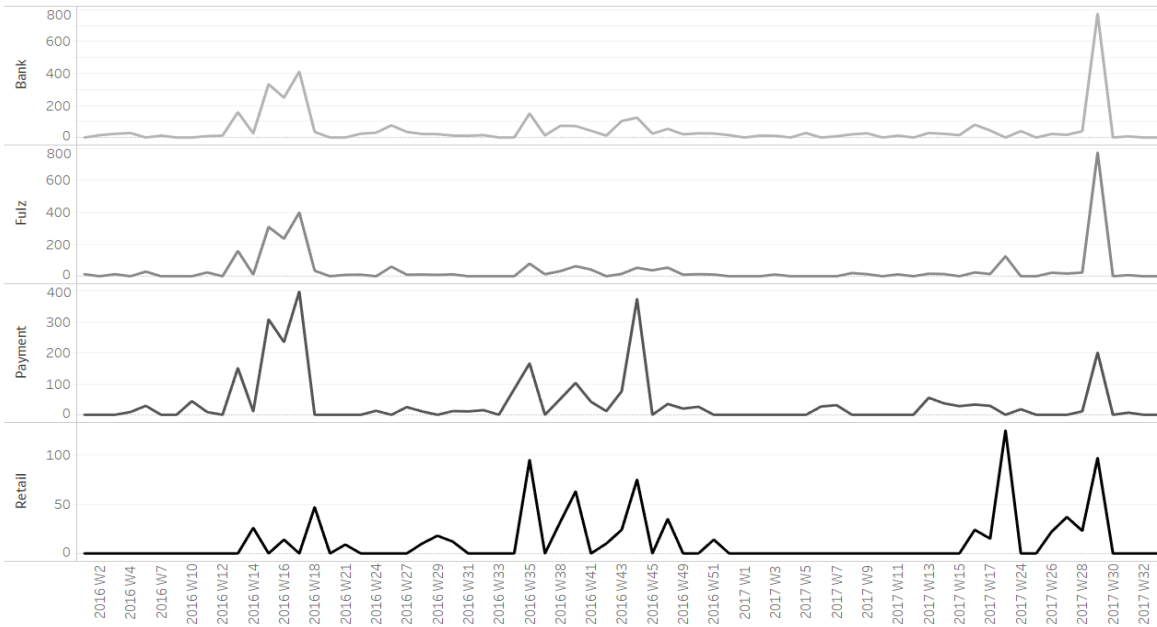


Figure 5.6 No of users participating in trending threads relating to different fraud categories over time

From the trending threads we observe that the users on darkweb show more interest in bank, fulz and payment frauds than in retail frauds. The average number of participants on a trending thread about bank and fulz in a week is about 50 users while that on a discussion about payment services is 40 users. The average number of participants in a week for retail is 11, much lower in comparison. We observe that threads about bank frauds trend more consistently than the other frauds.

In our results we also discover several threads that have been going on for years and have trended at several points in time. Such threads act as markets in themselves and alerts on activities on them indicate newer products available for sale or discounted products etc. Let us consider an example of one such thread that our system produced alerts on in 2017. This thread has a topic as “western union in 15 minutes bank transfers in 2 hours”. The thread started on August 2014, and has had replies on it sporadically since it started. Figure 5.7 shows the number of users replying to this thread over time divided in bins of week as earlier.

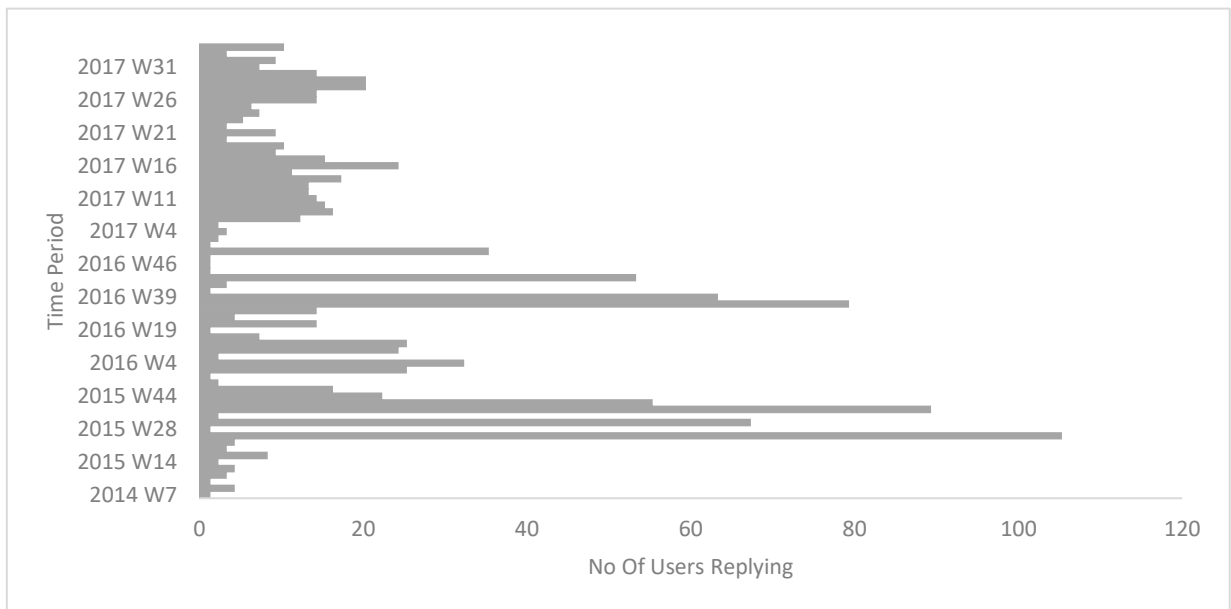


Figure 5.7 No of users replying over time to thread “western union in 15 minutes bank transfers in 2 hours”

The system produced warnings for this thread in 2016 in 35th, 38th, 44th, 48<sup>th</sup> week and in 16<sup>th</sup> week of 2017. Recurring alerts on the same threads over a large period of time can be indicative of such markets in a thread which can then be labeled.

## CHAPTER 6

### RELATED WORK

In this thesis, we leverage the context of the posts as well as social network information and the activity of the thread.

Data breach reports over the last couple of years have noted the rise in the number of data leaks relating to financial sector (Information leaked from banks, payment services etc.) [8] [18]. Authors of [8] report that 998 incidents of data breach were recorded for the financial sector alone in 2016 and 71% of the data compromised was credentials [8]. A large amount of this compromised data is then sold on darkweb.

Previous work analyzed forums and markets to understand the fundamental structure and working of the darkweb as well as detect threats that pose great risk to individuals, organizations, and governments [4] [5] [1] [2]. In [4] the authors studied 1899 threads on darkweb and through qualitative content analyses show that most products sold in these forums were some form of stolen data (84.3%). Most sellers offered dumps, referring to bank account or credit card data (44.7%), as well as CVV data from credit cards (34.9%). In papers [1] [2] [4], authors did not use any Automated/Artificial Intelligence methods for analysis, all the findings presented were based on qualitative analysis, hence the dataset was limited as opposed to the approach in this thesis. In paper [17] the authors use clustering technique to categorize the products sold in 17 marketplaces and found that the highest number of products were in clusters related to carding and PayPal.

In [19] the authors analyze communication of malicious hackers on darkweb and twitter to generate alerts on current and imminent cyber-attacks. They used text mining and could generate 661 alerts in a time period of about 6 months with a precision of 84%. The system

generated alerts about several cyber-attacks like mirai as well as data leaks. This work focuses heavily on vulnerability exploitation and alerts give warning about terms rather than categories.

## CHAPTER 7

### CONCLUSION

In this thesis, we implemented a system for timely and relevant identification and categorization of financial frauds on the darkweb websites. We evaluated our model using 10-fold cross-validation as well as evaluated blind data. We augment the accuracy of the machine learning algorithm using social network features. To ensure timely detection of these frauds we used the popularity of the threads based on interest expressed by users and generated alerts on trending threads. After analyzing the results, we found that financial frauds are among the top most popular topics on the darkweb and that bank and payment frauds are the most popular among the financial frauds.

## REFERENCES

- [1] Daniel Moore, Thomas Rid "Cryptopolitik and the Darknet", Survival, 58 (1), 2016
- [2] Ablon L, Libicki, MC, Golay, AA "Markets for cybercrime tools and stolen data", Rand National Security Division, Santa Monica, California, 2014.
- [3] Jane Leclair "Protecting our future, volume 2: Educating a cybersecurity workforce", Albany, NY. Hudson Whitman/ Excelsior College Press, 2015
- [4] Holt, Thomas J., and Olga Smirnova. "Examining the Structure, Organization, and Processes of the International Market for Stolen Data." Washington DC: National Criminal Justice Reference Service, 2014.
- [5] Nunes, Eric, Ahmad Diab, Andrew Gunn, Ericsson Marin, Vineet Mishra, Vivin Paliath, John Robertson, Jana Shakarian, Amanda Thart, and Paulo Shakarian. "Darknet and deepnet mining for proactive cybersecurity threat intelligence.", Intelligence and Security Informatics (ISI), IEEE, 2016.
- [6] Xie, Yi, and Shun-Zheng Yu. "A large-scale hidden semi-Markov model for anomaly detection on user browsing behaviors." IEEE/ACM Transactions on Networking (TON) 17.1 (2009): 54-65.
- [7] Viswanath, Bimal, et al. "Towards Detecting Anomalous User Behavior in Online Social Networks." USENIX Security Symposium. 2014.
- [8] Verizon "Verizon Data Breach Investigations Report", <http://www.verizonenterprise.com/verizon-insights-lab/dbir/2017/> 2017.
- [9] Quoc Le and Tomas Mikolov. Distributed Representations of Sentences and Documents. <http://arxiv.org/pdf/1405.4053v2.pdf>
- [10] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient Estimation of Word Representations in Vector Space. In Proceedings of Workshop at ICLR, 2013.
- [11] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed Representations of Words and Phrases and their Compositionality. In Proceedings of NIPS, 2013.
- [12] Vanschoren, Joaquin. "OpenML." OpenML: Exploring Machine Learning Better, Together.N.p., n.d. Web. 2017.
- [13] L'huillier, Gastón, et al. "Thread-based social network analysis for virtual communities of interests in the dark web." ACM SIGKDD Workshop on Intelligence and Security Informatics. ACM, 2010.

- [14] Zhang, Yulei, et al. "Dark web forums portal: searching and analyzing jihadist forums." *Intelligence and Security Informatics*, 2009. IST'09. IEEE International Conference on. IEEE, 2009
- [15] Benevenuto, Fabrício, et al. "Characterizing user behavior in online social networks." *Proceedings of the 9th ACM SIGCOMM conference on Internet measurement conference*. ACM, 2009.
- [16] Viswanath, Bimal, et al. "Towards Detecting Anomalous User Behavior in Online Social Networks." *USENIX Security Symposium*. 2014.
- [17] E. Marin, A. Diab, P. Shakarian, Product Offerings in Malicious Hacker Markets, 2016 IEEE Conference on Intelligence and Security Informatics (ISI-16) (Sep. 2016)
- [18] Enterprise, Verizon. "Data breach investigations report." Report, Verizon Enterprise (2016).
- [19] Sapienza, Anna, et al. "Early Warnings of Cyber Threats in Online Discussions." *Data Mining Workshops (ICDMW), 2017 IEEE International Conference on*. IEEE, 2017.
- [20] "F1 Score." Wikipedia. Wikimedia Foundation, 10 Feb. 2018. Web. 17 Feb. 2018.
- [21] "Precision and Recall." Wikipedia. Wikimedia Foundation, 13 Feb. 2018. Web. 17 Feb. 2018.
- [22] "Breach Detection by the Numbers: Days, Weeks or Years?" Infocyte, [www.infocyte.com/blog/2016/7/26/how-many-days-does-it-take-to-discover-a-breach-the-answer-may-shock-you](http://www.infocyte.com/blog/2016/7/26/how-many-days-does-it-take-to-discover-a-breach-the-answer-may-shock-you).



## APPENDIX A

### THE LIST OF REGULAR EXPRESSIONS USED FOR CLEANING DATA

- porn
- sex
- drugs\*
- cocaine
- guns\*
- arms\*
- gold
- grams\*
- got\s\*milk
- marijuana

## APPENDIX B

### LIST OF REGULAR EXPRESSIONS USED FOR LABELING EACH OF THE CATEGORIES

CATEGORY	KEYWORD
BANK	<ul style="list-style-type: none"> <li>• login</li> <li>• credentials*</li> <li>• pas+wo*rd</li> <li>• pwd</li> <li>• user\s*name</li> <li>• ssn</li> <li>• name</li> <li>• ad+res+</li> <li>• street</li> <li>• city</li> <li>• state</li> <li>• dob</li> <li>• acc</li> <li>• license</li> <li>• bin</li> <li>• ac+ount</li> <li>• routing</li> <li>• bank</li> <li>• bank\s*drop</li> <li>• dumps*</li> <li>• balance</li> <li>•</li> </ul>
FULZ	<ul style="list-style-type: none"> <li>• login</li> <li>• credentials*</li> <li>• pas+wo*rd</li> <li>• pwd</li> <li>• user\s*name</li> <li>• ssn</li> <li>• name</li> <li>• ad+res+</li> <li>• street</li> <li>• city</li> <li>• state</li> <li>• dob</li> <li>• license</li> <li>• cc</li> <li>• cvv</li> <li>• card</li> <li>• debit</li> </ul>

	<ul style="list-style-type: none"> <li>• credit</li> <li>• visa</li> <li>• master</li> <li>• ful+z+</li> </ul>
RETAIL	<ul style="list-style-type: none"> <li>• login</li> <li>• credentials*</li> <li>• pas+wo*rd</li> <li>• pwd</li> <li>• user\s*name</li> <li>• ssn</li> <li>• acc</li> <li>• name</li> <li>• ad+res+</li> <li>• street</li> <li>• city</li> <li>• state</li> <li>• dob</li> <li>• license</li> <li>• walmart</li> <li>• target.com</li> <li>• costco</li> <li>• netflix</li> <li>• hulu</li> <li>• amazon\.*com</li> <li>• best\s*buy</li> <li>• nortsdorm</li> <li>• ebay</li> <li>• ecommerce</li> <li>• retail</li> </ul>
PAYMENT SERVICES	<ul style="list-style-type: none"> <li>• login</li> <li>• credentials*</li> <li>• pas+wo*rd</li> <li>• pwd</li> <li>• user\s*name</li> <li>• ssn</li> <li>• acc</li> <li>• name</li> <li>• ad+res+</li> <li>• street</li> <li>• city</li> </ul>

	<ul style="list-style-type: none"><li>• state</li><li>• dob</li><li>• license</li><li>• paypal</li><li>• account</li><li>• venmo</li><li>• paytm</li><li>• web\s*money</li></ul>
--	--