

Utah State University

DigitalCommons@USU

All Graduate Theses and Dissertations

Graduate Studies

5-2017

Detecting Malicious Campaigns in Crowdsourcing Platforms

Hongkyu Choi
Utah State University

Follow this and additional works at: <https://digitalcommons.usu.edu/etd>



Part of the [Computer Sciences Commons](#)

Recommended Citation

Choi, Hongkyu, "Detecting Malicious Campaigns in Crowdsourcing Platforms" (2017). *All Graduate Theses and Dissertations*. 6504.

<https://digitalcommons.usu.edu/etd/6504>

This Thesis is brought to you for free and open access by the Graduate Studies at DigitalCommons@USU. It has been accepted for inclusion in All Graduate Theses and Dissertations by an authorized administrator of DigitalCommons@USU. For more information, please contact digitalcommons@usu.edu.



DETECTING MALICIOUS CAMPAIGNS IN CROWDSOURCING PLATFORMS

by

Hongkyu Choi

A thesis submitted in partial fulfillment
of the requirements for the degree

of

MASTER OF SCIENCE

in

Computer Science

Approved:

Kyumin Lee, Ph.D.
Major Professor

Curtis Dyreson, Ph.D.
Committee Member

Haitao Wang, Ph.D.
Committee Member

Mark R. McLellan, Ph.D.
Vice President for Research and
Dean of the School of Graduate Studies

UTAH STATE UNIVERSITY
Logan, Utah

2017

Copyright © Hongkyu Choi 2017

All Rights Reserved

ABSTRACT

Detecting Malicious Campaigns in Crowdsourcing Platforms

by

Hongkyu Choi, Master of Science

Utah State University, 2017

Major Professor: Kyumin Lee, Ph.D.

Department: Computer Science

Crowdsourcing systems enable new opportunities for requesters with limited funds to accomplish various tasks using human computation. However, the power of human computation is abused by malicious requesters who create malicious campaigns to manipulate information in web systems such as social networking sites, online review sites, and search engines. To mitigate the impact and reach of these malicious campaigns to targeted sites, we propose and evaluate a machine learning based classification approach for detecting malicious campaigns in crowdsourcing platforms as a first line of defense, and build a malicious campaign blacklist service for targeted site providers, researchers and users.

Specifically, we (i) conduct a comprehensive analysis to understand the characteristics of malicious campaigns and legitimate campaigns in crowdsourcing platforms, (ii) propose various features to distinguish between malicious campaigns and legitimate campaigns, (iii) evaluate a classification approach against baselines, and (iv) build a malicious campaign blacklist service. Our experimental results show that our proposed approaches effectively detect malicious campaigns with low false negative and false positive rates.

(37 pages)

PUBLIC ABSTRACT

Detecting Malicious Campaigns in Crowdsourcing Platforms

Hongkyu Choi

Crowdsourcing sites such as Mechanical Turk and Crowdflower provide a marketplace where requesters create tasks and recruit workers, who may perform certain tasks in order to get financial compensation. Anyone in the world can be a requester and/or a worker as long as he/she has the Internet connection. Crowdsourcing creates a new way to solve various tasks by using “human computation power”. However, crowdsourcing has been misused by malicious requesters and unethical workers for account generation, search engine optimization, content and link generation, ad posting and spam mailing, and social network linking. It creates new threats to the Web system. The consequences of the malicious tasks are receiving spam emails and spam messages from online social networks, polluted social links by fake link farming, and irrelevant search results because of manipulated web page link structure. Eventually, these have degraded the information trustworthiness on the Web. To solve this problem, we build a predictor that detects whether campaign/task in crowdsourcing sites is malicious or not so that malicious campaigns/tasks can be removed by crowdsourcing site providers as soon as they are created. In particular, we (i) analyze characteristics of malicious and legitimate campaigns; (ii) extract commonly available features from four crowdsourcing sites; and (iii) build predictors and evaluate their performance. Our experimental results show that our predictors are more effective and robust compared with several baselines. In the end, we design and build a malicious campaign blacklist service, which provides users with various information.

CONTENTS

	Page
ABSTRACT	iii
PUBLIC ABSTRACT	iv
LIST OF TABLES	vi
LIST OF FIGURES	vii
CHAPTER	
1 INTRODUCTION	1
2 RELATED WORK	4
3 DATASETS	6
4 ANALYZING MALICIOUS CAMPAIGNS AND LEGITIMATE CAMPAIGNS	9
4.1 Tasks and market sizes	9
4.2 Hourly wages	10
4.3 Requesters and workers	11
4.4 Clustering malicious campaigns	12
4.5 Real-world impact of malicious campaigns	15
5 FEATURES	17
6 EXPERIMENTS	21
6.1 Detecting malicious campaigns	21
6.2 Robustness of our proposed approach	23
7 BUILDING A MALICIOUS CAMPAIGN BLACKLIST SERVICE	25
8 CONCLUSION	27
REFERENCES	28

LIST OF TABLES

Table	Page
3.1 Datasets	8
4.1 Statistics for tasks and market sizes of malicious campaigns and legitimate campaigns.	9
4.2 Goals and median values of malicious campaign clusters.	14
4.3 The five most targeted sites by malicious campaigns and their corresponding median values.	14
5.1 Pearson correlation coefficient results for 8 features excluding text features.	17
6.1 Classification results.	22

LIST OF FIGURES

Figure	Page
3.1 A list of campaigns	6
3.2 A campaign description	7
4.1 Box plots for hourly wages in legitimate campaigns and malicious campaigns.	10
4.2 The number of campaigns created by requesters in each country and the total cost of the campaigns.	11
4.3 The number of tasks performed by workers in each country and their total earnings.	12
4.4 Sum of squared error obtained by increasing the number of cluster	13
4.5 Actual task description targeting on Facebook Like.	15
4.6 Consequence of manipulating the number of Likes in Facebook.	16
5.1 CDF of hourly wage by legitimate and malicious campaigns.	18
5.2 CDF of the number of tasks by legitimate and malicious campaigns.	19
5.3 CDF of estimated time to complete by legitimate and malicious campaigns.	19
5.4 CDF of the number of URLs by legitimate and malicious campaigns.	20
5.5 CDF of the number of words in title by legitimate and malicious campaigns.	20
5.6 CDF of the number words in task instruction by legitimate and malicious campaigns.	20
6.1 A macro-scale view between the 2nd and 12th weeks.	23
6.2 A micro-scale view between the 3rd and 12th weeks.	24
7.1 A malicious campaign blacklist service.	26

CHAPTER 1

INTRODUCTION

Crowdsourcing platforms such as Mechanical Turk (MTurk) and Crowdflower provide a marketplace where requesters recruit workers and request the completion of various tasks. Since anyone in the world can be a worker, the labor fees are relatively low, and workers are available at virtually all hours of the day. Due to these benefits, requesters have used crowdsourcing platforms for various tasks such as labeling datasets [1, 2], searching a boat from satellite images to find a lost person [3], proofreading a document [4], and adding missing data [5].

However, some requesters abuse crowdsourcing platforms by creating malicious campaigns to manipulate search engines, write fake reviews, and create accounts for additional attacks [6–10]. Using crowdsourced manipulation, malicious requesters and workers can potentially earn hundreds of millions of dollars [11, 12]. As a result, crowdsourced manipulation threatens the foundation of the free and open web ecosystem by reducing the quality of online social media, degrading trust in search engines, manipulating political opinion, and ultimately compromising the security and trustworthiness of cyberspace [13–15].

Prior research [14, 15] has identified the threat of malicious campaigns by quantifying their prevalence in several crowdsourcing platforms. Specifically, a large collection of loosely-moderated crowdsourcing platforms serves as launching pads for these malicious campaigns. Unfortunately, there is a significant gap in (i) our understanding of how to detect malicious campaigns at the source (i.e., crowdsourcing platforms), which would mitigate their impact and reach before they influence targeted sites, and (ii) building a malicious campaign blacklist service which provides various functions (e.g., searching malicious campaigns by a keyword search and grouping relevant malicious campaigns by various categories/goals) to targeted service providers, researchers and users.

Hence, in this thesis we aim to automatically predict and detect malicious campaigns in crowdsourcing platforms, and build a malicious campaign blacklist service by answering following research questions: What kind of malicious campaigns exist in crowdsourcing platforms? Can we find distinguishing patterns/features between malicious campaigns and legitimate campaigns? Can we develop a statistical model that automatically detects malicious campaigns? Can we build a malicious campaign blacklist service providing various functions?

To answer these questions, we make the following contributions:

- First, we collect a large number of campaigns from popular crowdsourcing platforms: MTurk, Microworkers, Rapidworkers, and Shorttask¹. Then, we cluster malicious campaigns to understand what types of malicious campaigns exist in crowdsourcing platforms.
- Second, we analyze characteristics of malicious campaigns and legitimate campaigns in terms of their market sizes and hourly wages. Then, we propose and evaluate various features for distinguishing between malicious campaigns and legitimate campaigns, and we visualize each feature to concretely illustrate the differing properties for malicious and legitimate campaigns.
- Third, we develop a predictive model, and evaluate its performance against baselines in terms of accuracy, false positive rate and false negative rate. To our knowledge, this is the first study to focus on detecting malicious campaigns in multiple crowdsourcing platforms.
- Finally, we build a malicious campaign blacklist service as a repository and retrieval service.

The rest of this thesis is organized as follows. After reviewing the related work in the next chapter, we describe the dataset and ground truth on campaign types. Then we introduce the characteristics of malicious and legitimate campaigns. After that, we present

¹MTurk, Microworkers, Rapidworkers and Shorttask represent www.mturk.com, microworkers.com, rapidworkers.com, and shorttask.com, respectively.

our proposed features for building malicious campaign classifiers. The following chapter reports the experimental results on real-world datasets with performance and robustness of the model. In the next chapter, we propose a malicious campaign blacklist service. Finally, the last chapter concludes the thesis.

CHAPTER 2

RELATED WORK

Since the emergence of crowdsourcing platforms (e.g., MTurk and Crowdflower), researchers have studied how to use crowd wisdom. Wang et al. [16] hired workers to identify fake accounts in Facebook and Renren. Workers have identified improper tasks in a Japanese crowdsourcing site [17] and proofread documents in near-real time [4]. Other researchers were interested in analyzing the demographics of workers [18] and quantifying the evolution of campaigns/tasks in MTurk [19]. Ge et al. analyzed a supply-driven crowdsourcing marketplace regarding key features that distinguish “super sellers” from regular participants [20].

Another research topic is to measure the quality of workers and outcomes (and determining how to control that quality). Due to the openness of these crowdsourcing systems, anyone can be a worker. Consequently, workers might be lazy or dishonest and seek money by quickly completing tasks with low quality answers. Venetis and Garcia-Molina [21] proposed three scoring methods such as gold standard, plurality answer agreement, and Task Work Time to filter low quality answers. A machine learning technique was applied to detect low quality answerers [22]. Soberón et al. [23] showed that adding open-ended questions (i.e., explanation-based techniques) into tasks was useful for identifying low quality answers.

With the rising popularity of crowdsourcing systems, malicious campaigns and tasks have been created by some requesters. To understand the problems, Motoyama et al. [14] introduced possible web service abuse in Freelancer.com. Wang et al. [15] analyzed two Chinese crowdsourcing platforms and found that up to 90% of campaigns are malicious campaigns. Lee et al. [13] found that social networking sites and search engines were mainly targeted by malicious campaigns. Researchers began analyzing crowdsourced manipulation

and the characteristics of workers in targeted sites such as Facebook and Twitter. Fayazi et al. [24] proposed a reviewer-reviewer graph clustering approach based on a Markov Random Field to identify workers that posted fake reviews on Amazon. Song et al. [25] proposed a crowdsourced manipulation detection method to detect target objects such as a post, page, and URL on Twitter.

In contrast to this previous research, we collected a large number of campaigns from four crowdsourcing platforms, analyzed characteristics of malicious campaigns and legitimate campaigns, developed predictive models to automatically identify malicious campaigns, and built a malicious campaign blacklist service. Our research will complement existing research base.

CHAPTER 3

DATASETS

In a crowdsourcing platform, there are two types of users – (i) a requester and (ii) a worker. A requester is a user who creates a campaign with detailed instructions for one or more tasks. Each task is then performed by one worker. If the requester is satisfied with the worker’s outcome, the requester will approve it, and compensation (i.e., money) will be passed to the worker by the crowdsourcing platform.

To collect a dataset, we developed a crawler for four popular crowdsourcing platforms: Amazon Mechanical Turk (MTurk), Microworkers, Rapidworkers, and Shorttask. The crawler collected campaign listings and detailed campaign descriptions. We ran the crawler for 3 months between November 2014 and January 2015, and it collected 23,220

Job name	Payment
Firma: Sign up + Screenshot ●	\$0.41
iOS App Testing (Chip N Dales): Download + Install +... ●	\$0.47
iOS App Testing (Indian Chief): Download + Install +... ●	\$0.45
iOS App Testing (AIT): Download + Install + Screenshot + Bonus	\$0.56
Properties: Sign up + Video Post	\$0.51
TTV-Make 10 Q&A about "Chengdu Tour" (USA-Western)	\$0.60
iOS App Testing (Wild West): Download + Install + Screenshot... ●	\$1.59
TTV-Write a Post: Self-Motivation (MentalNotes)	\$0.80
Android App Testing (VSM): Download + Test + Screenshot	\$1.01
Julia Jackson (Cambria Image): Image Search + Obtain Info +... ●	\$0.18
iOS App Testing: Download + Install + Honest Feedback	\$1.25
iOS App Testing (More Chill): Download + Install + Screenshot ●	\$1.75
Julia Jackson (Arcanum): Image Search + Bookmark + Screenshot ●	\$0.30
Quora: Answer + Screenshot	\$0.50

Fig. 3.1: A list of campaigns

Facebook Like: Post

Work done: **145**/³⁶⁰

You will earn **\$0.30**

Task takes less than **3** min to finish

Job ID: 2effbe9e1c58


Employer: [I Help Promote](#)

[add to Exclude List](#)

[add to Include List](#)

Tasks will be rated within **1** day

You can accept this job if you are from any of these countries:

 Australia, Canada, New Zealand, United Kingdom, United States

Facebook → Facebook Like (no Friends)

? What is expected from Workers?

1. Go to <https://www.facebook.com/OfficialJoeyoung/videos/1416646708374476/>
2. Like the post

! Required proof that task was finished?

1. Your Facebook display name
 2. URL to your Facebook profile
- Make sure you have set your profile to Public View in order for your task to be verified.

Fig. 3.2: A campaign description

campaigns consisting of 3,356,153 tasks¹. Figure 3.1 shows a list of campaigns and Figure 3.2 represents a sample campaign description in Microworkers. The campaign description contains the number of available tasks, compensation for each task, estimated time to complete a task, and task instructions that describe what a worker is supposed to do.

As we mentioned earlier, our ultimate goal is to understand characteristics of malicious campaigns and legitimate campaigns, and build a predictive model to automatically predict a campaign's class to either a malicious campaign or a legitimate campaign. We define a malicious campaign as one that requires workers to manipulate information in targeted sites such as social media sites and search engines. For example, a malicious campaign might require workers to post fake reviews on Amazon [6], artificially create backlinks to boost a

¹A campaign contains multiple tasks, and each task is assigned to one worker.

Table 3.1: Datasets

Types of Campaigns	Number of Campaigns
Malicious Campaigns	5,010
Legitimate Campaigns	18,210
Total	23,220

specific website’s search engine ranking, or “Like” a specific Facebook page (as shown in Figure 3.2).

Using this definition, two annotators manually labeled each campaign in our dataset as a malicious campaign or a legitimate campaign, based on the campaign description. When the two annotators disagreed about a particular campaign’s label, a third annotator labeled the campaign. The annotators made the same labeling decision on 23,079 out of 23,220 campaigns, achieving 99.4% agreement.

Table 3.1 shows the number of malicious and legitimate campaigns in our collected dataset. 5,010 of the 23,220 campaigns (21.6%) were malicious campaigns, and each campaign contained 145 tasks on average. Overall, the malicious campaigns contained about 800K tasks. Thus, it is important to understand and analyze how malicious campaigns are different from legitimate campaigns so that we can build models to automatically detect malicious campaigns and develop a malicious campaign blacklist service.

CHAPTER 4
ANALYZING MALICIOUS CAMPAIGNS AND LEGITIMATE CAMPAIGNS

Now we turn our attention to analyzing characteristics of malicious campaigns and legitimate campaigns in the crowdsourcing platforms.

4.1 Tasks and market sizes

First, we calculated the number of tasks associated with malicious campaigns and legitimate campaigns, and then, we measured the market sizes for malicious and legitimate campaigns. To estimate the market size for a collection of malicious and legitimate campaigns, we used the following equation:

$$MarketSize(C) = \sum_{i=1}^n r(i) * count(i) \quad (4.1)$$

where C is a set of malicious or legitimate campaigns $\{c_1, c_2, c_3 \dots c_n\}$ in crowdsourcing platforms, n is the number of malicious or legitimate campaigns, $r(i)$ is the reward (i.e., compensation) per task for campaign i , and $count(i)$ is the number of tasks associated with campaign i .

As shown in Table 4.1, the malicious tasks for the four crowdsourcing platforms amounted to 45% of the entire market size, while the number of malicious tasks only represented 24%

Table 4.1: Statistics for tasks and market sizes of malicious campaigns and legitimate campaigns.

	Number of Tasks	Market Sizes
Malicious	798,796	\$148,911
Legitimate	2,557,357	\$179,696
Total	3,356,153	\$328,607

of the total campaign tasks. This analytical result reveals that the reward per malicious task is much higher than the reward per legitimate task.

4.2 Hourly wages

Next, we concretely evaluate if the hourly reward for malicious campaigns is actually higher than the hourly reward for legitimate campaigns. A campaign's description contains estimated time to complete (ETC) information and a reward per task. We calculated the hourly wage for each campaign using $\frac{\text{reward per task} * 60}{ETC}$ because ETC's unit of measurement is a minute. Figure 4.1 shows box plots for hourly wages in legitimate and malicious campaigns. The median hourly wage in malicious campaigns (\$2.48) is larger than the median hourly wage in legitimate campaigns (\$1.88). One explanation for this result is that malicious requesters provide higher rewards to workers so that they can attract these workers, who may have ethical concerns about these malicious campaigns/tasks.

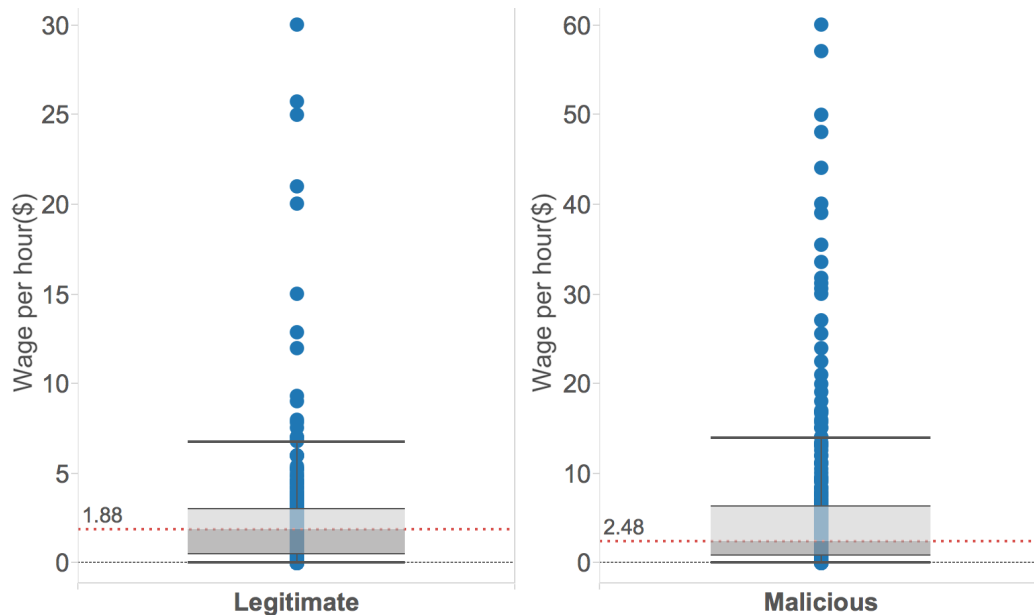


Fig. 4.1: Box plots for hourly wages in legitimate campaigns and malicious campaigns.

4.3 Requesters and workers

Next, we analyze the characteristics of requesters and workers. In particular, we are interested in answering the following research questions: what is the geographic distribution of requesters? How much do requesters typically spend on campaigns? What is the geographic distribution of workers, and how many tasks do workers typically complete? Fortunately, the Microworkers platform publicly discloses location information of requesters and workers as well as compensation data.

We analyzed our Microworkers dataset, which consists of 3,971 campaigns that were created by 518 distinct requesters. The top 5 originating countries for requesters in our dataset is US (51.7% of all requesters), UK (6.2%), Canada (4.2%), Bangladesh (3.9%), and India (3.3%). We can clearly see that many requesters originate from English-speaking countries. This observation also suggests that most campaigns target Western websites (where the majority of the user base is also English-speaking). As shown in Figure 4.2, US, UK and Canada requesters created 3,472 campaigns (87.4% of all campaigns), and they spent \$116,684 (85.6% of the market size). In the figure, the size of a circle represents the

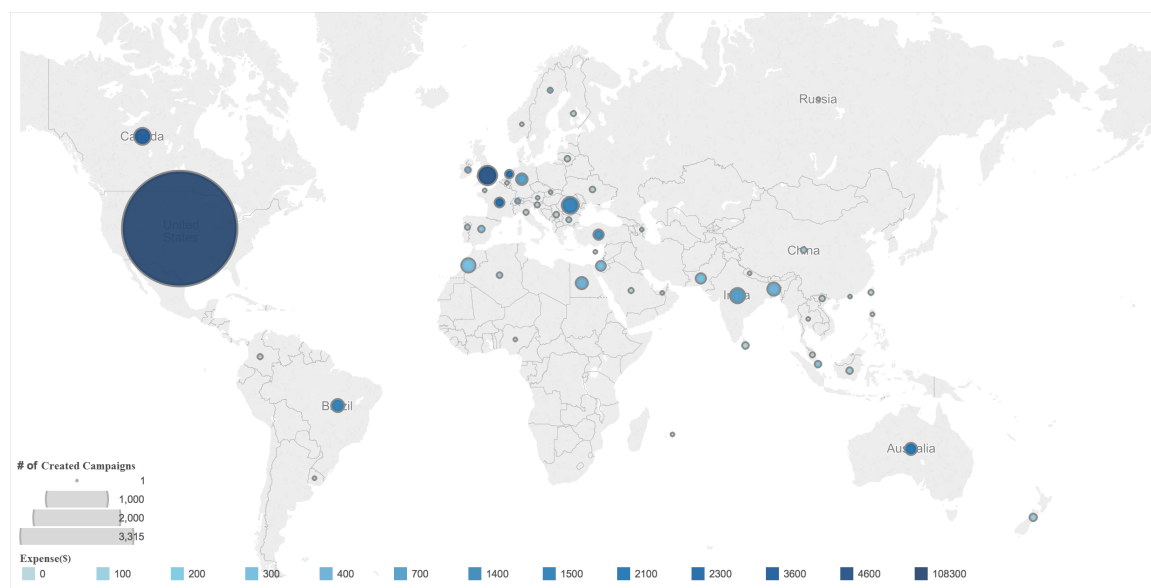


Fig. 4.2: The number of campaigns created by requesters in each country and the total cost of the campaigns.

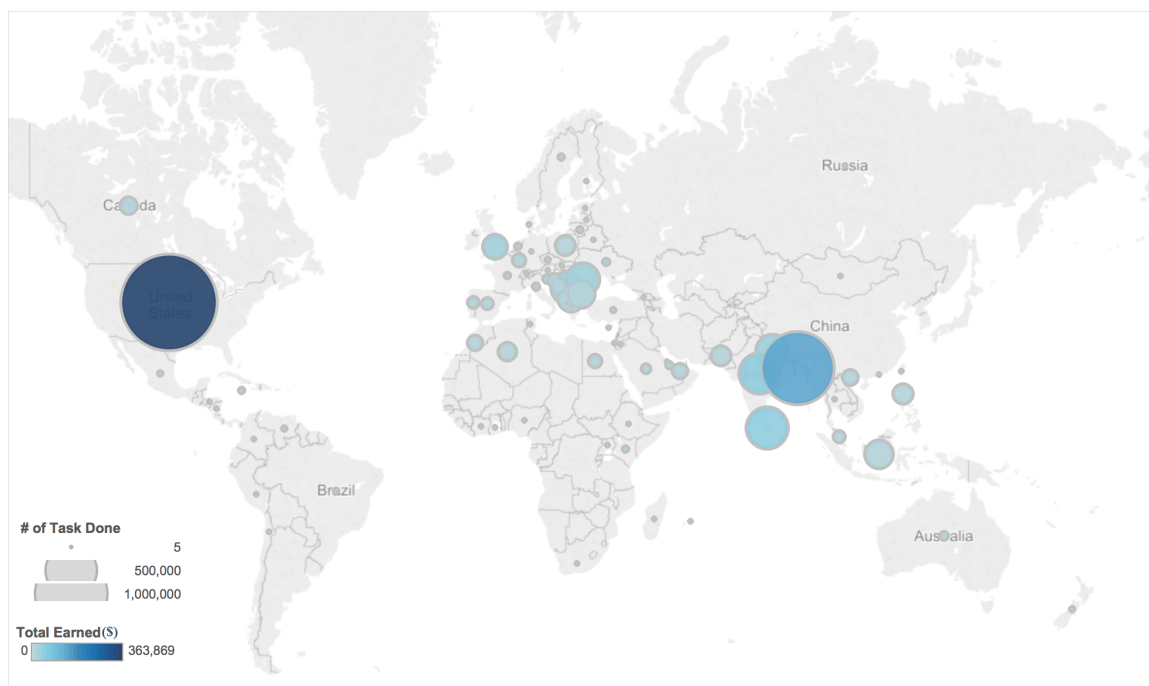


Fig. 4.3: The number of tasks performed by workers in each country and their total earnings.

number of campaigns created in a given country, and the color of the circle represents the amount spent by requesters in that country.

To investigate characteristics of workers, we randomly collected 3,261 workers' profiles. Based on this data, the top 5 originating countries for workers is Bangladesh (33.7%), US (11.9%), India (9.8%), Nepal (5.6%), and Sri Lanka (4.5%). Unlike the requesters, most of the workers are from developing countries. Figure 4.3 shows the number of tasks performed by workers in each country as well as the earnings for workers in each country. Interestingly, US workers earned \$364,067 (44% of the market size) by performing 1,633,344 tasks (33.4% of the total), and Bangladesh workers earned \$128,099 (15.6%) by performing 927,894 tasks (19%). This suggests that US workers receive significantly higher earnings than Bangladesh workers (\$0.22/task vs. \$0.14/task).

4.4 Clustering malicious campaigns

To investigate characteristics of malicious campaigns associated with specific goals, we

clustered the campaigns based on their goals and targeted sites. From 5,010 malicious campaigns, we extracted a title from each campaign and tokenized it by unigram. Then, we removed stop words and measured term frequency-inverse document frequency (TF-IDF). Now, each campaign is represented as a vector based on TF-IDF. Given a list of vectors, we used a k -means clustering algorithm to cluster the vectors (i.e., campaigns). To obtain the optimal number of clusters, we experimented with k values in the range of 2 through 10, and we measured Sum of Squared Error (SSE) for each value. Figure 4.4 shows how SSE was changed as we increased the k value from 2 to 10. When k was incremented from 7 to 8, the curve flattened, which means 7 is the optimal number of clusters.

After clustering the malicious campaigns, we investigated every cluster and found objectives for the campaigns in each cluster as shown in Table 4.2. The median values for time, reward, and hourly wage are presented in the table. The goals for most of the campaigns were to manipulate content on social networking sites (e.g., Google Plus, Twitter, Yahoo Answer and Facebook) and manipulate search engine results by searching a specific keyword and clicking a certain web page link. “Download and install a new application”

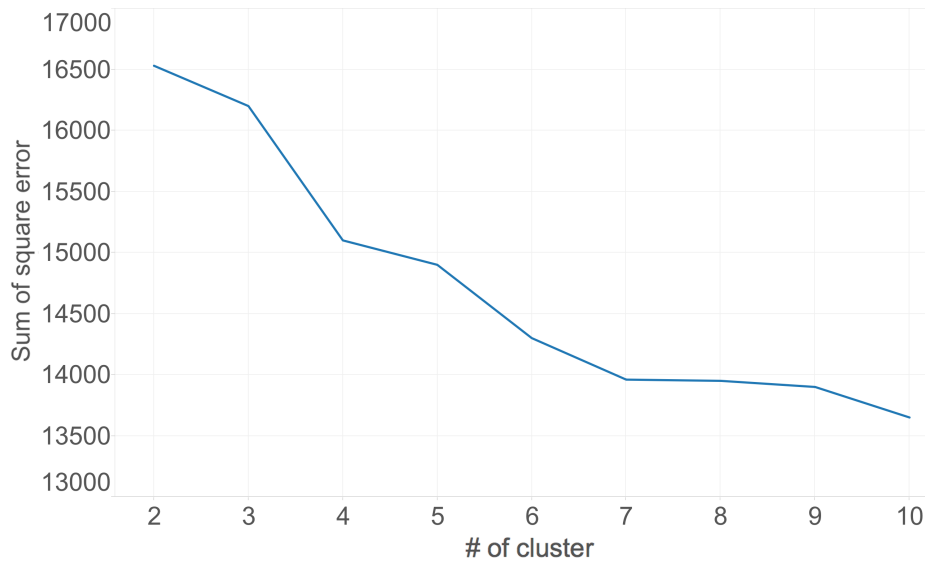


Fig. 4.4: Sum of squared error obtained by increasing the number of cluster

campaigns provided workers with larger compensation per hour than the other campaigns.

To identify the sites that were targeted the most by malicious campaigns, we extracted a list of the top 500 companies from Alexa, and we searched for each of those company names (and company hostnames) in malicious campaign descriptions. Table 4.3 shows the five most targeted sites. Nine hundred and two (18%) malicious campaigns targeted Google. Social networking sites such as Twitter, Instagram, Facebook, and Youtube were also targeted frequently.

Table 4.2: Goals and median values of malicious campaign clusters.

Campaign Goal	Malic. C.	%	ETC(min)	Reward	\$/hour
Social network associated (Review, Link, Share, Retweet and Like)	2,987	60	33.5	\$0.45	\$2.13
Search and click	863	17	8	\$0.22	\$2.65
Search and visit	654	13	5	\$0.21	\$3.85
Add a comment	197	4	8	\$0.31	\$2.91
Register in a forum and post a message	168	3	12	\$0.35	\$1.02
Create a new pin at Pinterest	96	2	3	\$0.11	\$2.20
Download and install a new application	45	1	12.5	\$0.81	\$4.50
Average			12	\$0.35	\$2.75

Table 4.3: The five most targeted sites by malicious campaigns and their corresponding median values.

	Malic. C.	%	Reward	ETC(min)	\$/hour
Google	902	18	\$0.21	6.0	2.3
Twitter	600	12	\$0.16	7.5	1.9
Instagram	210	4	\$0.13	6.5	1.0
Facebook	154	3	\$0.35	7.0	3.0
Youtube	153	3	\$0.20	9.5	2.2

4.5 Real-world impact of malicious campaigns

Thus far, we've identified important characteristics of malicious campaigns. Now, we need to determine if malicious campaigns have any real-world impact on targeted sites and if existing security algorithms/systems can detect manipulations in the targeted sites. To investigate these issues, we tracked 29 malicious campaigns targeting Facebook in which workers manipulated Facebook Likes. We collected daily snapshots of the malicious campaigns from crowdsourcing platforms and daily snapshots of the targeted Facebook pages. The 29 malicious campaigns consisted of 8,268 tasks, each task required adding one fake Like. Out of 8,268 fake Like, 7,160 of the Likes were successfully attributed to the target pages when we checked those pages later, which means only 1,108 (13.4%) of the fake Likes were deleted by the Facebook security team.

Figures 4.5 and 4.6 show examples of the malicious campaigns manipulating the number of Facebook Likes. Figure 4.5 shows the campaign description containing a total number of tasks, the number of available tasks, and task instructions for workers. Eight hundred of the campaign's tasks were completed within 4 days although it had 7 days duration.

Figure 4.6 shows the number of completed tasks reported to a requester and the number of fake Likes completed by workers for the targeted Facebook page. We can clearly observe

Facebook Like: [Redacted]

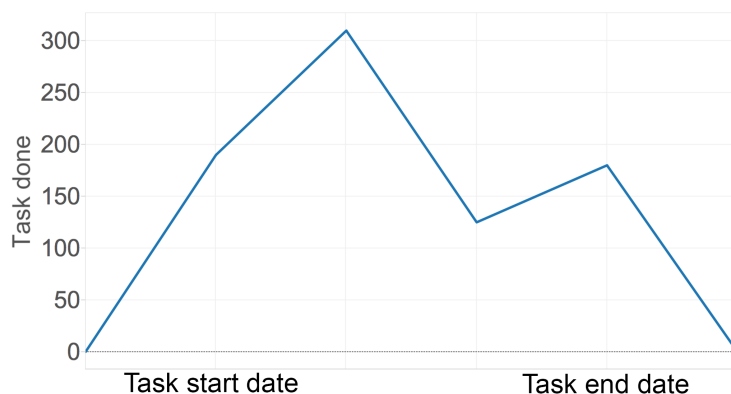
Work done: **174/800** Employer: [Redacted]
 You will earn **\$0.10** Tasks will be rated within **7** days
 Task takes less than **1** min to finish
 Job ID: df9c8e85b85f

Facebook → Facebook Like (no Friends)

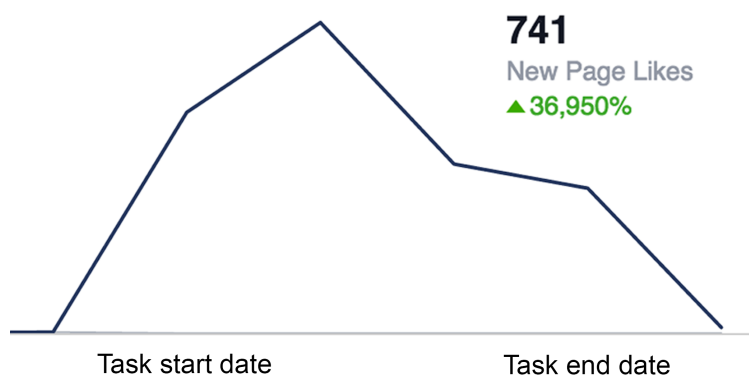
? What is expected from Workers?

1. Go to [https://www.facebook.com/\[Redacted\]](https://www.facebook.com/[Redacted])
2. Like the page

Fig. 4.5: Actual task description targeting on Facebook Like.



(a) Trend of task completion after task posting.



(b) Changes in the number of Likes in Facebook page.

Fig. 4.6: Consequence of manipulating the number of Likes in Facebook.

that the middle and right figures show similar temporal patterns. Out of 800 Likes, 741 Likes remained on the Facebook page, which means Facebook only labeled 59 (7%) Likes as “fake” Likes.

This example and the previous analysis (for 29 fake liking campaigns) show that malicious campaigns have a real-world impact on targeted sites, and current security systems are unable to detect most of the manipulated content. The previous work [11] also confirmed that Twitter’s safety team only detected 24.6% of fake followers. These results motivated us to investigate an automated approach for detecting malicious campaigns in crowdsourcing platforms using predictive models.

CHAPTER 5
FEATURES

In this chapter, we describe proposed features for building malicious campaign classifiers. To build a universal classifier which can be applied to any crowdsourcing platform regardless what information is available, we proposed and extracted commonly available features across the four crowdsourcing platforms. Our proposed features are reward, number of tasks, estimated time to complete (ETC), hourly wage, number of URLs in task instruction, $\frac{\text{Number of URLs in task instruction}}{\text{Number of words in task instruction}}$, number of words in a task title, number of words in task instruction, and text features extracted from task title and task instruction.

To avoid the overfitting problem by removing features that are too similar [26,27], we measured the Pearson correlation coefficient of each pair of the first 8 features excluding text features (we conducted another feature selection for the text features). Table 5.1 presents

Table 5.1: Pearson correlation coefficient results for 8 features excluding text features.

	reward	tasks	ETC	hourly wage	URLs	URLs / words	words in title	words in instruction
reward	1							
tasks	-.021*	1						
ETC	.316*	-.021*	1					
hourly wage	.084*	-.005	-.061*	1				
URLs	.005	-.017*	-.107*	.059*	1			
URLs / words	-.088*	-.001	-.137*	.127*	.366*	1		
words in title	.298*	-.003	.263*	-.043*	-.049*	-.225*	1	
words in instruction	.408*	-.020*	.049*	.010	.088*	-.180*	.334*	1

*Correlation is significant at the 0.01 level (2-tailed).

the correlation coefficient of each pair. The number of word in task instruction and reward had the highest correlation among the pairs, achieving 0.408. Since 0.408 represented no significant correlation, we kept all of the first 8 features.

From task title and task instruction, we extracted text features as follows: (i) first, we removed stopwords from the title and task instruction, and then, we applied stemming to them; (ii) second, we extracted unigrams, bigrams, and trigrams from the text; (iii) third, we measured χ^2 values for the extracted unigram, bigram, and trigram features [28]; (iv) finally, we only used text features with positive χ^2 values. Through this process, we used thousands of text features.

Next, Figure 5.1 shows cumulative distribution functions (CDFs) on hourly wage for malicious campaigns and legitimate campaigns. Interestingly, requesters for 80% of the legitimate campaigns paid less than one dollar to each worker in terms of hourly wage, while requesters for 10% of the malicious campaigns paid the same hourly wage to workers. This suggests that performing malicious campaigns was more profitable, which is consistent with our previous results.

Malicious campaigns also contain a larger number of tasks than most legitimate campaigns, and malicious campaigns have shorter ETC than legitimate campaigns. Task instructions in malicious campaigns contain more URLs than legitimate campaigns, which

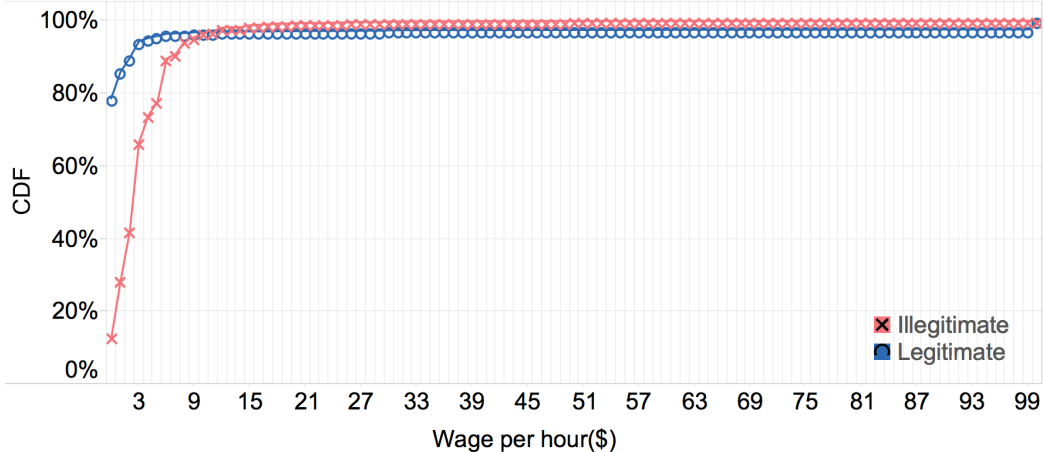


Fig. 5.1: CDF of hourly wage by legitimate and malicious campaigns.

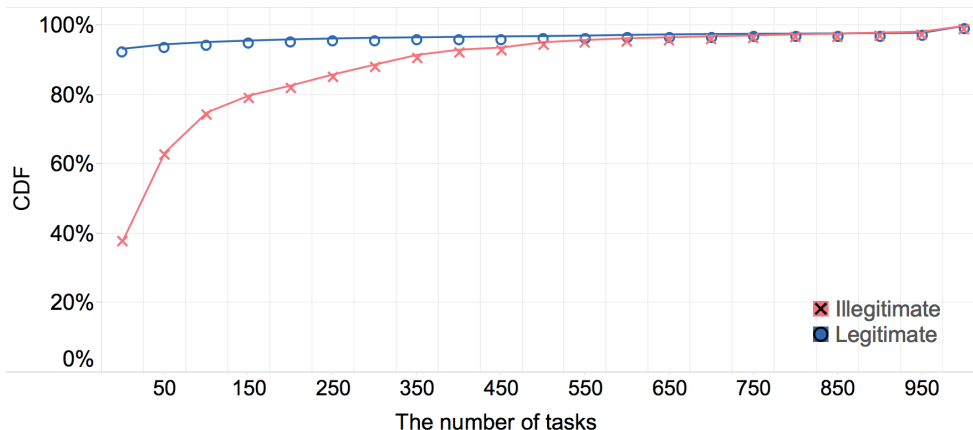


Fig. 5.2: CDF of the number of tasks by legitimate and malicious campaigns.

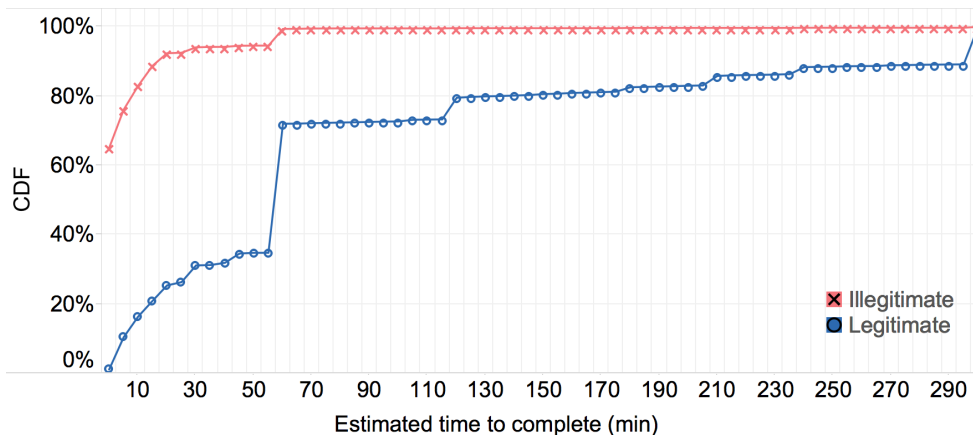


Fig. 5.3: CDF of estimated time to complete by legitimate and malicious campaigns.

suggests that malicious campaigns require workers to access external websites (potentially targeted sites) more often.

Finally, malicious campaigns have shorter titles and task instructions than legitimate campaigns. This observation might indicate that some of the legitimate campaigns are more complicated to perform and require longer ETC. Overall, the CDFs illustrate distinct differences between malicious campaigns and legitimate campaigns.

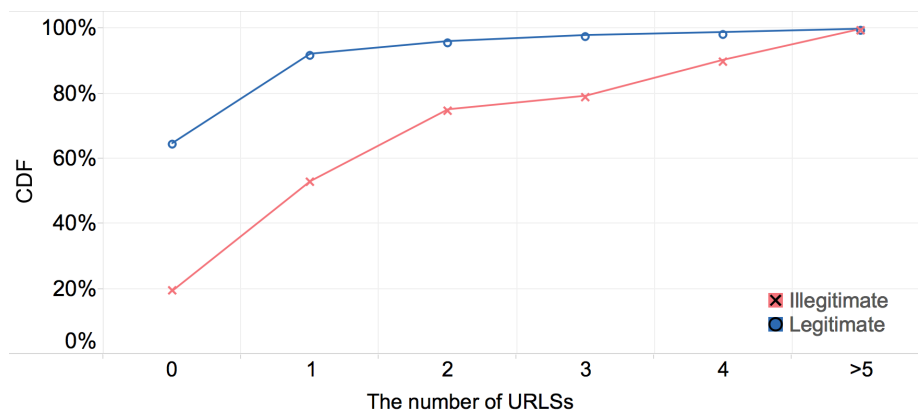


Fig. 5.4: CDF of the number of URLs by legitimate and malicious campaigns.

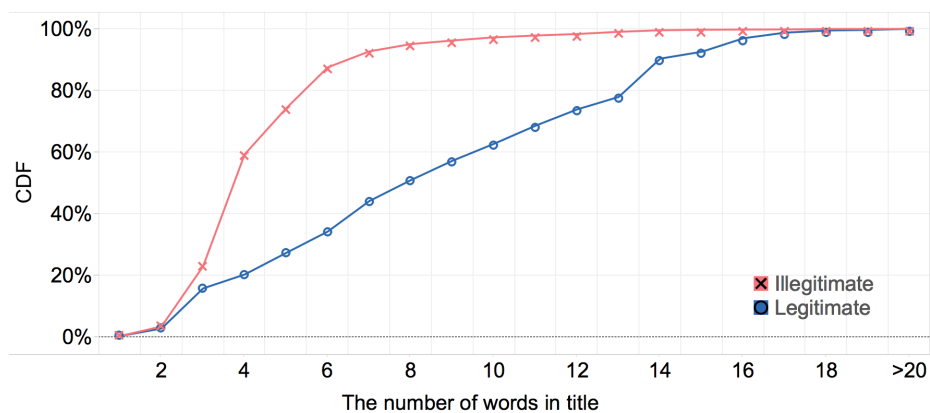


Fig. 5.5: CDF of the number of words in title by legitimate and malicious campaigns.

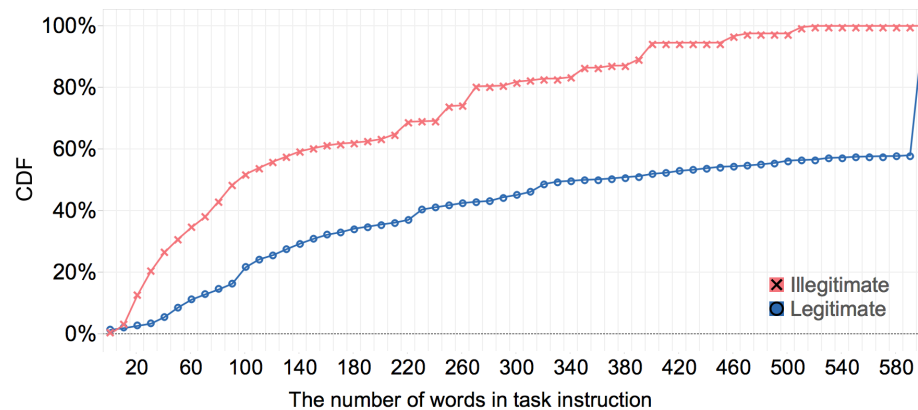


Fig. 5.6: CDF of the number words in task instruction by legitimate and malicious campaigns.

CHAPTER 6

EXPERIMENTS

In the previous chapter, we observed that malicious campaigns and legitimate campaigns have different characteristics. In this chapter, we build classifiers to detect malicious campaigns by exploiting these differences. First, we build malicious campaign classifiers and compare their performance with baselines. Then, we further evaluate robustness of our classification approach.

6.1 Detecting malicious campaigns

As we mentioned in chapter 3, we collected campaign descriptions for 3 months between November 2014 and January 2015. The dataset consists of 18,210 legitimate campaigns and 5,010 malicious campaigns. We built and tested statistical models with 10-fold cross validation. We compared the performance of three classification algorithms: Naive Bayes, J48, and Support Vector Machine (SVM).

We compared our statistical models/classifiers with following baselines: (i) *majority selection approach* which always predicts a campaign’s class as the majority instances’ class (i.e., a legitimate campaign in the dataset); (ii) *URL-based filtering* approach which classifies a campaigns as a malicious campaign if its description contains at least one URL whose host name is one of top K sites; and (iii) *principal component analysis* (PCA) approach, an unsupervised machine learning technique, inspired from the previous work [29]. In PCA approach, we projected campaigns (using the same features with our classifiers) onto the normal and residual subspaces to classify malicious and legitimate campaigns. The space spanned by top principal components is the normal subspace and the remaining space is known as residual subspace. From our dataset, we achieved 85% variance from the top 35 principal components out of 1,835 components. We computed L2 norm and set the squared

Table 6.1: Classification results.

Approach	Accuracy	FPR	FNR
Majority Selection	78.4%	1	0
URL-based filtering@100	72.4%	0.708	0.157
URL-based filtering@500	72.3%	0.688	0.164
URL-based filtering@1000	71.9%	0.635	0.183
PCA - 12% threshold	85.2%	0.999	0.031
our Naive Bayes	89.0%	0.044	0.147
our J48	99.1%	0.023	0.058
our SVM	99.2%	0.019	0.055

prediction error as the threshold value to find the malicious campaigns. We changed the threshold value from 1% to 50% by 1% increment each time to get the best classification result. Campaigns, whose L2 norm was greater than the threshold value, were classified as malicious campaigns.

To evaluate the performance of classifiers, we used following evaluation metrics: accuracy, false positive rate (FPR) and false negative rate (FNR). FPR means malicious campaigns were misclassified as legitimate campaigns while FNR means legitimate campaigns were misclassified as malicious campaigns.

In experiments, we ran majority selection approach, URL-based filtering approach at top 100, 500 and 1000 sites, PCA approach, and our three classification approaches (Naive Bayes, J48 and SVM). Table 6.1 shows experimental results of the baselines and our classification approaches. Majority selection approach achieved 78.4% accuracy, 1 FPR and 0 FNR, URL-based filtering@100 achieved 72.4% accuracy, 0.708 FPR and 0.157 FNR, and PCA approach with 12% threshold (only reporting the best result) achieved 85.2% accuracy, 0.999 FPR and 0.031 FNR. Overall, our SVM-based classifier significantly outperformed the other approaches, achieving 99.2% accuracy, 0.019 FPR and 0.055 FNR, and balancing between low FPR and low FNR.

6.2 Robustness of our proposed approach

In the previous experiment, we learned SVM classifier achieved the best prediction results for detecting malicious campaigns. Now, we analyze (i) how much training data we need to achieve a high prediction rate and (ii) whether a predictive model (i.e., a classifier) would remain robust over time.

To investigate these issues, we split the dataset chronologically based on weeks (i.e., the 3 month data was split into 12 weeks). Then, we trained a SVM classifier using the first week of data, and we used the classifier to test the data for each of the next weeks. Next, we added the following week’s data (e.g., the second week of data) into the training set and tested the data for each of the next weeks. Incrementally, we added each week’s data to the training set until the training set included data for the first 11 weeks.

Figures 6.1 and 6.2 show experimental results for macro-scale and micro-scale views of our approach¹. In particular, Figure 6.1 shows experimental results of the 2nd week to the 12th week in a macro-scale view. When we used *the first week* data as a training set

¹We did not show FPR and FNR lines because of the limited space.

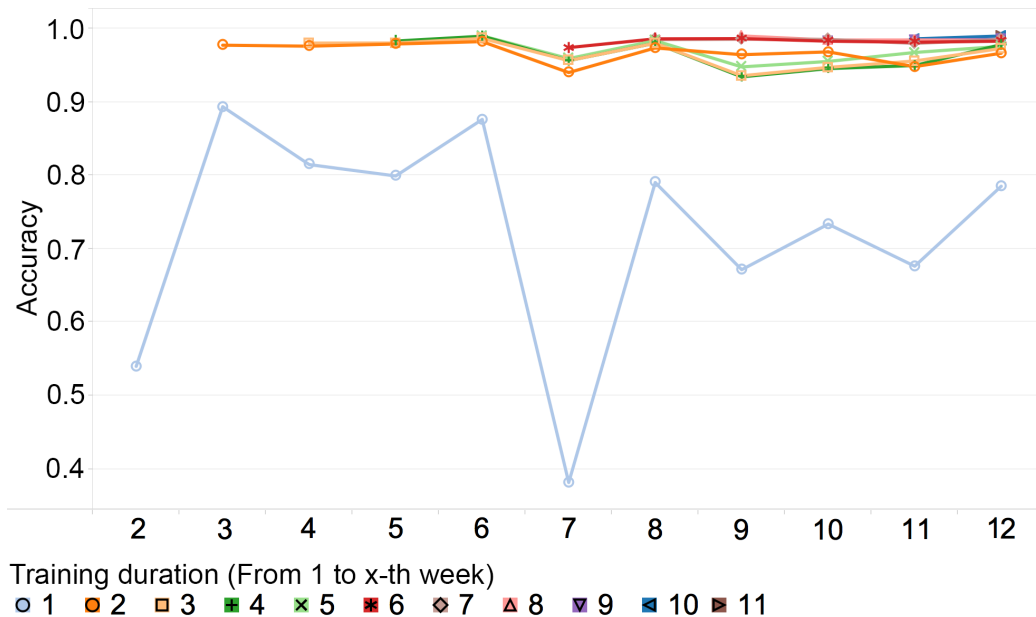


Fig. 6.1: A macro-scale view between the 2nd and 12th weeks.

and applied a classifier to each of the following weeks, the classifier achieved low accuracy. However, when we added one more week of data to a training set (i.e., the training set contained the first and second week of data), the classifier achieved significantly high accuracy. Note that 7th testing week’s classification results were slightly lower than earlier testing weeks because there were very small number of legitimate campaigns posted in the 7th week (e.g., 62% malicious campaigns and 38% legitimate campaigns in the 7th testing week vs. 11% malicious campaigns and 89% legitimate campaigns in the 6th testing week). We conjecture that the 7th week is a week containing Christmas and New Year holidays, so very less number of legitimate campaigns were created while almost same number of malicious campaigns was created compared with the 6th week.

Figure 6.2 shows experimental results in a micro-scale view by removing the first week training result (i.e., the classifier that was only trained with a single week of data). Based on this figure, we clearly observe that a SVM classifier based on data for the first 2 weeks achieved high accuracy even though the performance was up and down over time. Overall, the lowest accuracy, highest FPR and highest FNR among all the cases were 93.4%, 0.19, 0.03, respectively. Based on these experimental results, we conclude that two weeks of data is enough to build an effective predictive model for identifying malicious campaigns. We also conclude that our proposed classification approach consistently and robustly worked well over time.

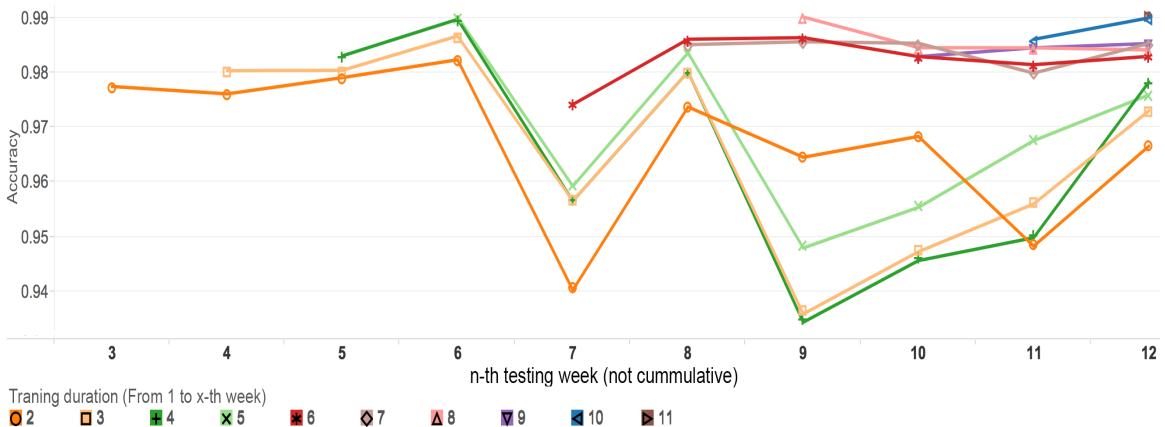


Fig. 6.2: A micro-scale view between the 3rd and 12th weeks.

CHAPTER 7

BUILDING A MALICIOUS CAMPAIGN BLACKLIST SERVICE

So far we learned that our proposed classification approach can detect malicious campaigns accurately and effectively. In practice, crowdsourcing platform providers can identify malicious campaigns in near-real time by using our approach when malicious requesters create malicious campaigns. As we learned from the previous study [15], some crowdsourcing platform providers may not be willing to remove these malicious campaigns because they may lose their revenue like commission from the malicious requesters if they filter/remove these malicious campaigns. In this case, unfiltered malicious campaigns would affect targeted services and targeted service providers, and eventually online users would get manipulated information. To protect these victims in targeted services side, we designed and built a web service-based malicious campaign blacklist which provides users (e.g., administrators of targeted sites or researchers) with various information such as campaign descriptions and statistical time-line charts.

The core component of the malicious campaign blacklist service is the malicious campaign classifier that we built and tested in the previous chapter. A current version of our blacklist service¹ supports a keyword search, group retrieved relevant malicious campaigns by various categories (so that users can further access a specific category), and visualizes a statistical time-line chart showing how many relevant malicious campaigns have been generated/detected over time.

Figure 7.1 shows a snapshot of our blacklist service, given a search keyword “google”. In this example, the blacklist service retrieved relevant malicious campaigns and showed the most relevant results in the first page. Simultaneously, it grouped the malicious campaigns and showed names of clusters/categories in the left pane with the number of subcategories

¹Our malicious campaign blacklist service will be released in public soon.

for better navigation. It also showed a trend of related malicious campaigns.

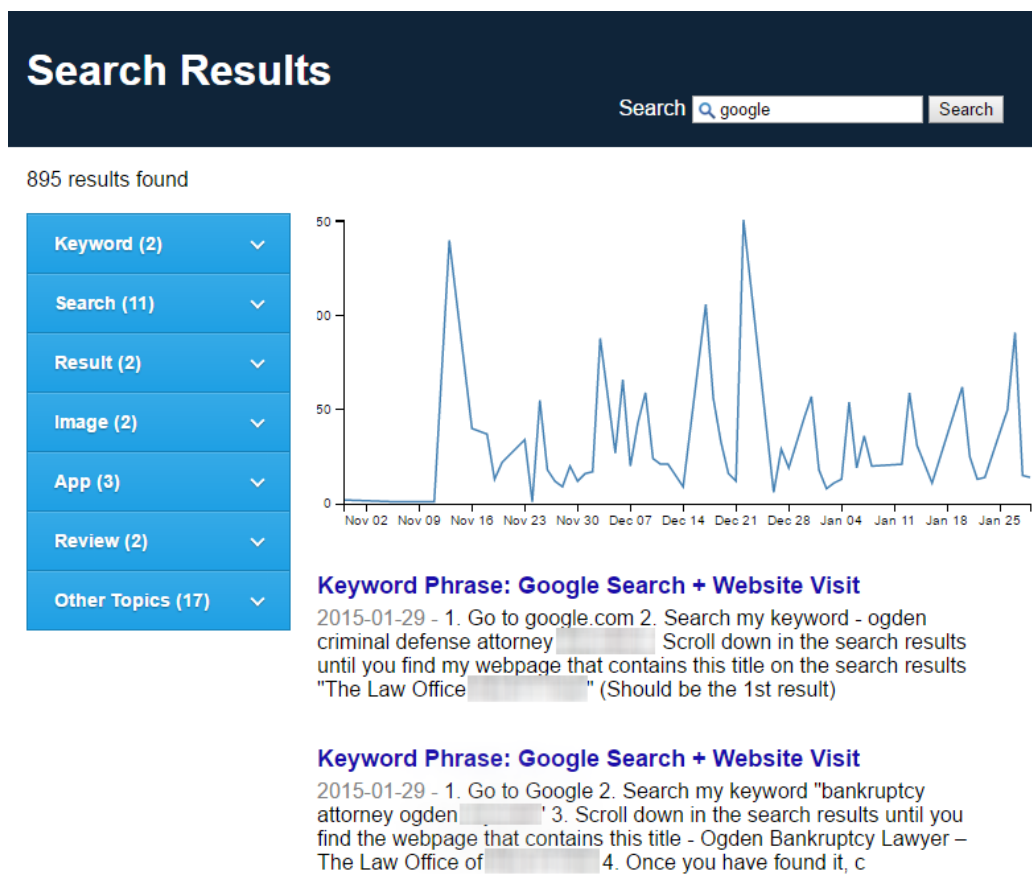


Fig. 7.1: A malicious campaign blacklist service.

CHAPTER 8

CONCLUSION

In this thesis, we analyzed characteristics of malicious campaigns and legitimate campaigns. The median hourly wage in malicious campaigns (\$2.48) was larger than the median hourly wage in legitimate campaigns (\$1.88), tempting workers to perform malicious campaigns in targeted sites such as social networking sites, online review sites, and search engines. To measure the real-world impact of malicious campaigns, we selected Facebook Liking campaigns and found that Facebook caught only 13% fake likes. This suggests that current defense systems in targeted sites are inadequate and potentially undetected malicious campaigns are deteriorating information quality and trust.

To overcome this problem, we proposed features which were distinguished between malicious campaigns and legitimate campaigns. Then, we built malicious campaign classifiers based on the features for mitigating the impact and reach of the malicious campaigns to targeted sites. Our classifiers outperformed the baselines – majority selection, URL-based filtering and PCA approaches –, achieving 99.2% accuracy, 0.019 FPR and 0.055 FNR.

By using the classifier, we built a malicious campaign blacklist service that provides a keyword search, retrieves relevant malicious campaign descriptions, groups these campaigns by various categories for better navigation, and shows a trend of relevant malicious campaigns. The blacklist service will help targeted service providers, researchers and users to understand what kind of malicious campaigns have been running under the targeted services and which information is manipulated, and potentially mitigate the impact of these malicious campaigns in the target sites as well.

REFERENCES

- [1] Z. Cheng, J. Caverlee, H. Barthwal, and V. Bachani, “Who is the barbecue king of texas?: A geo-spatial approach to finding local experts on twitter,” in *SIGIR*, 2014.
- [2] A. Sorokin and D. Forsyth, “Utility data annotation with Amazon Mechanical Turk,” 2008.
- [3] A. Doan, R. Ramakrishnan, and A. Y. Halevy, “Crowdsourcing systems on the world-wide web,” *Commun. ACM*, vol. 54, no. 4, pp. 86–96, Apr. 2011.
- [4] M. S. Bernstein, G. Little, R. C. Miller, B. Hartmann, M. S. Ackerman, D. R. Karger, D. Crowell, and K. Panovich, “Soylent: A word processor with a crowd inside,” in *UIST*, 2010.
- [5] M. J. Franklin, D. Kossmann, T. Kraska, S. Ramesh, and R. Xin, “Crowddb: Answering queries with crowdsourcing,” in *SIGMOD*, 2011.
- [6] C. Conner, “Amazon sues 1,114 fake reviewers on fiverr,” <http://www.forbes.com/sites/cherylsnappconner/2015/10/18/amazon-sues-1114-fake-reviewers-on-fiverr-com/>, October 2015.
- [7] E. De Cristofaro, A. Friedman, G. Jourjon, M. A. Kaafar, and M. Z. Shafiq, “Paying for likes?: Understanding facebook like fraud using honeypots,” in *IMC*, 2014.
- [8] N. Pham, “Vietnam admits deploying bloggers to support government,” <http://www.bbc.co.uk/news/world-asia-20982985>, January 2013.
- [9] G. Stringhini, M. Egele, C. Kruegel, and G. Vigna, “Poultry markets: On the underground economy of twitter followers,” in *Workshop on Online Social Networks*, 2012.

- [10] K. Thomas, D. McCoy, C. Grier, A. Kolcz, and V. Paxson, "Trafficking fraudulent accounts: The role of the underground market in twitter spam and abuse," in *USENIX Conference on Security*, 2013.
- [11] K. Lee, S. Webb, and H. Ge, "The dark side of micro-task marketplaces: Characterizing fiverr and automatically detecting crowdturfing," in *ICWSM*, 2014.
- [12] K. Thomas, D. Huang, D. Wang, E. Bursztein, C. Grier, T. J. Holt, C. Kruegel, D. McCoy, S. Savage, and G. Vigna, "Framing dependencies introduced by underground commoditization," in *Workshop on the Economics of Information Security*, 2015.
- [13] K. Lee, P. Tamilarasan, and J. Caverlee, "Crowdturfers, Campaigns, and Social Media: Tracking and Revealing Crowdsourced Manipulation of Social Media," in *ICWSM*, 2013.
- [14] M. Motoyama, D. McCoy, K. Levchenko, S. Savage, and G. M. Voelker, "Dirty jobs: The role of freelance labor in web service abuse," in *USENIX Conference on Security*, 2011.
- [15] G. Wang, C. Wilson, X. Zhao, Y. Zhu, M. Mohanlal, H. Zheng, and B. Y. Zhao, "Serf and turf: crowdturfing for fun and profit," in *WWW*, 2012.
- [16] G. Wang, M. Mohanlal, C. Wilson, X. Wang, M. J. Metzger, H. Zheng, and B. Y. Zhao, "Social turing tests: Crowdsourcing sybil detection," in *NDSS*, 2013.
- [17] Y. Baba, H. Kashima, K. Kinoshita, G. Yamaguchi, and Y. Akiyoshi, "Leveraging non-expert crowdsourcing workers for improper task detection in crowdsourcing marketplaces." *Expert Syst. Appl.*, vol. 41, no. 6, pp. 2678–2687, 2014.
- [18] J. Ross, L. Irani, M. S. Silberman, A. Zaldivar, and B. Tomlinson, "Who are the crowdworkers?: Shifting demographics in mechanical turk," in *CHI*, 2010.
- [19] D. E. Difallah, M. Catasta, G. Demartini, P. G. Ipeirotis, and P. Cudré-Mauroux, "The dynamics of micro-task crowdsourcing: The case of amazon mturk," in *WWW*, 2015.

- [20] H. Ge, J. Caverlee, and K. Lee, “Crowds, gigs, and super sellers: A measurement study of a supply-driven crowdsourcing marketplace,” in *ICWSM*, 2015.
- [21] P. Venetis and H. Garcia-Molina, “Quality control for comparison microtasks,” in *CrowdKDD workshop*, 2012.
- [22] H. Halpin and R. Blanco, “Machine-learning for spammer detection in crowd-sourcing,” in *Human Computation workshop in conjunction with AAAI*, 2012.
- [23] G. Soberón, L. Aroyo, C. Welty, O. Inel, H. Lin, and M. Overmeen, “Measuring crowd truth: Disagreement metrics combined with worker behavior filters,” in *CrowdSem 2013 Workshop*, 2013.
- [24] A. Fayazi, K. Lee, J. Caverlee, and A. Squicciarini, “Uncovering crowdsourced manipulation of online reviews,” in *SIGIR*, 2015.
- [25] J. Song, S. Lee, and J. Kim, “Crowdtarget: Target-based detection of crowdturfing in online social networks,” in *CCS*, 2015.
- [26] H.-H. Hsu and C.-W. Hsieh, “Feature selection via correlation coefficient clustering,” *Journal of Software*, vol. 5, no. 12, pp. 1371–1377, 2010.
- [27] A. Janecek, W. N. Gansterer, M. Demel, and G. Ecker, “On the relationship between feature selection and classification accuracy.”
- [28] Y. Yang and J. O. Pedersen, “A comparative study on feature selection in text categorization,” in *ICML*, 1997.
- [29] B. Viswanath, M. A. Bashir, M. Crovella, S. Guha, K. P. Gummadi, B. Krishnamurthy, and A. Mislove, “Towards detecting anomalous user behavior in online social networks,” in *USENIX Security*, 2014.