

Cracking Classifiers for Evasion: A Case Study on the Google’s Phishing Pages Filter

Bin Liang, Miaoqiang Su, Wei You, Wenchang Shi, Gang Yang

Renmin University of China

{liangb, sumiaoqiang, youwei, wenchang, yanggang}@ruc.edu.cn

ABSTRACT

Various classifiers based on the machine learning techniques have been widely used in security applications. Meanwhile, they also became an attack target of adversaries. Many existing studies have paid much attention to the evasion attacks on the online classifiers and discussed defensive methods. However, the security of the classifiers deployed in the client environment has not got the attention it deserves. Besides, earlier studies have only concentrated on the experimental classifiers developed for research purposes only. The security of widely-used commercial classifiers still remains unclear. In this paper, we use the *Google’s phishing pages filter* (GPPF), a classifier deployed in the *Chrome* browser and with over one billion users, as a case to investigate the security challenges for the client-side classifiers. A new attack methodology targeted to client-side classifiers, called *classifiers cracking*, is presented. According to the methodology, we successfully crack the classification model of GPPF and extract sufficient knowledge from it for performing effective evasion attacks, including the classification algorithm, scoring rules and features, etc. Most importantly, we completely reverse engineer 84.8% scoring rules, covering most of high-weighted rules. Based on the cracked information, we perform two kinds of evasion attacks to GPPF, using 100 real phishing pages as the target of evaluation. The experiments show that all the phishing pages (100%) can be easily manipulated to bypass the detection of GPPF. Our study demonstrates that the existing client-side classifiers are very vulnerable to classifiers cracking attacks.

Categories and Subject Descriptors

D.4.6 [Security and Protection]: Invasive software; I.2 [Artificial Intelligence]: Learning

General Terms

Security

Keywords

Phishing Detection, Machine Learning, Classifiers, Cracking, Collision Attacks, Evasion Attacks

1. INTRODUCTION

Machine learning techniques have been commonly adopted in security applications. Various classifiers were trained for detecting malicious web pages [23], spam [53], phishing [56], malware [47], etc. Not surprisingly, the classifiers themselves also became an attack target of adversaries. The adversary can attempt to fool classifiers by purposely modifying their behaviors. For example, a spammer can manipulate the spam mails to evade spam filters by inserting some *good* words indicative of legitimate

mails or misspelling *bad* words indicative of spam mails [41]. This requires the classifier can resist potential attacks.

Many existing studies have paid attention to the security of classifiers. According to the taxonomy of attacks against classifiers proposed in [8][9][32], the influences of attacks on the classifier are categorized into two types: *causative* attacks interfere training process with control over training data to downgrade the performance of the classifier, and *exploratory* attacks exploit knowledge of the trained classifier to cause misclassifications but do not affect training.

In causative attacks, the adversary has the opportunity to inject (poison) specially crafted samples during the collection of training samples and cause the learner to misclassify security violations (false negatives), such as [16][17][18][19]. For example, a poisoning attack method against support vector machines (SVM) is presented in [17]. It was demonstrated that the SVM’s classification accuracy can be largely impacted by feeding malicious training data. Fortunately, in practice, the adversary doesn’t always have an opportunity to effectively control over training data. In fact, the training process of most classifiers, especially the ones deployed in commercial products, is not open to the public. The adversary needs to fight with trained classifiers. For example, according to described in [56], the Google’s phishing pages classifier is developed in an offline training process. The training dataset consists of millions of samples from various domains. In this case, it is very difficult, if not impossible, for an adversary to craft enough amounts of malicious inputs to effectively poison the training process.

On the other hand, exploratory attacks attempt to learn enough knowledge about the trained classifiers and find a way to evade the classification. Some existing studies about evasion attacks made the unrealistic assumption that the adversary has perfect knowledge of classification model [26]. In practice, the adversary often needs to send some probes (e.g., membership queries) to the classifier and observe its response to deduce desirable knowledge [23], perform an adversarial learning to get sufficient knowledge about the target classifier to construct evasion attacks [42], or reconstruct an imitation of the target classifier based on the available public information (e.g., training data) to gain key knowledge [51]. In theory, the success of evasion attacks heavily depends on the amount of knowledge possessed by the adversary. Especially, the knowledge of features contributes most to the success of the attacks as discussed in [51]. Accordingly, some mitigation techniques have been proposed to against evasion attacks by reducing the leakage of exploitable knowledge as much as possible [8][11] or making the learning method more robust to evasion [15][36].

However, existing studies often overlooked an important fact that some classifiers are deployed in the client environment fully

controlled by users (*client-side classifier* for short) rather than in a remote server. For example, the classifiers for filtering spam emails and phishing pages are often embedded in the email clients or web browsers respectively. In the scenario, the classifiers will face more serious security challenges. Instead of collecting the information via indirectly observations, the adversaries can freely and directly analyze the implementation and configuration of the classifiers to evade them. Consequently, it should be investigated carefully that how an adversary can learn the exploitable knowledge from a classifier deployed in the user clients and how effective the knowledge are exploited in launching an evasion attack. Additionally, the existing studies generally focused on the experimental classifiers developed for research purposes only. The security of widely-used classifiers in commercial products still remains unclear. From a practical point of view, evaluating the security of commercial classifiers is more significant for protecting end users from evasion attacks.

To this end, in this study, we investigate the security challenges for the client-side classifiers via a case study on the *Google’s phishing pages filter* (GPPF), a very widely-used classifier for automatically detecting unknown phishing pages. The classifier is completely integrated within the *Chrome* browser and invoked for every web page visited by users to determine whether it is phishing or not. Due to the popularity of *Chrome*, there are over one billion users using GPPF to against potential phishing attacks [1]. It is probably the most widely-used classifier as we know. If the adversary can easily evade it, countless users will be exposed to out-of-control phishing attacks.

In this paper, we demonstrate a practical and effective attack methodology, named *classifiers cracking*, in which various reverse engineering techniques are leveraged to extract sufficient knowledge from the client-side classifier for launching evasion attacks. Specifically, by performing some static and dynamic analysis on the implementation of *Chromium* (the development version of *Chrome*), we successfully extract the classification model of GPPF from it, mainly involving the classification algorithm, 2,130 scoring rules and corresponding weights, and 1,009 hashed features composing the scoring rules. With the help of some public datasets (e.g., large corpora), we then launch a collision attack to the hashed features and decrypt 815 (80.8%) of them only within a dozen of hours. As a result, we can completely reverse engineer 1807 (84.8%) scoring rules, covering most of high-weighted rules. Additionally, 196 (9.2%) scoring rules are partially cracked and can also be exploited to compromise the classification. There are only 127 (6.0%) rules surviving from the collision attack.

Based on the extracted information, we design two kinds of evasion attacks, i.e., *good features insertion* and *bad features elimination*. The basic idea behind these attacks is to add or remove some features with remarkable contributions to GPPF scoring into or from the target phishing pages to reduce their phishing scores, making the computed scores lower than the positive threshold defined by GPPF. We evaluate the effectiveness of the attacks with the 100 latest real phishing pages collected from PhishTank [2], a famous phishing URLs tracking site. The results show that we can easily manipulate all the phishing pages under the direction of the cracked knowledge, to make them successfully evade the detection of GPPF in the latest version of *Chrome*.

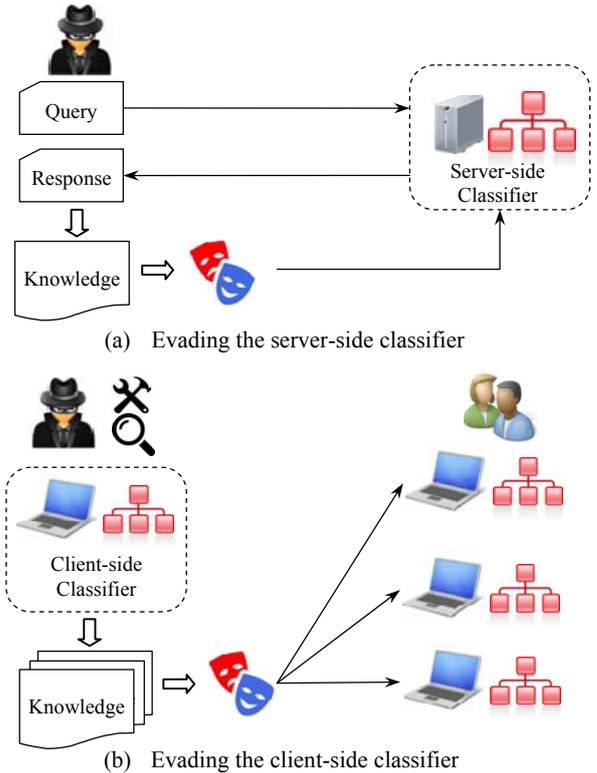


Figure 1. Threats to classifiers.

We also analyze the weakness of existing defense techniques when applying them to client-side classifiers, and introduce a potential defense strategy from the viewpoint of against cracking. The main idea is to employ the machine learning method to construct the client-side classifier such that its features are not prone to be reverse engineered inherently.

This paper makes the following main contributions.

- We propose a new attack methodology, classifiers cracking, aiming at the client-side classifiers. The adversary can follow it to readily acquire exploitable knowledge from the target classifier to launch effective evasion attacks.
- We successfully crack and evade the GPPF, a commercial classifier with over one billion users. It is demonstrated that the existing client-side classifiers are indeed vulnerable to classifiers cracking attacks.
- We discuss a potential defense strategy against classifiers cracking. We believe that it can be employed to enhance the security of various client-side classifiers, not only GPPF.

2. BACKGROUND

2.1 Threat Model

As shown in Figure 1(a), how to classify an instance in a server-side classifier is often a black-box to the adversary. The adversary can only send some queries and analyze responses to learn the information about it. In many cases, this way is already enough to launch an evasion attack. The adversary can construct a malformed instance to fool the classifier based on the information learned in advance.

However, when a classifier is deployed in the client-side computer, the situation may become worse. As shown in Figure 1(b), for a client-side classifier, its operations are performed in a white-box. The adversary can leverage almost all kinds of analysis techniques, such as debugging, disassembling, code analysis, dynamic taint tracking, etc., to thoroughly analyze the target classifier. As a result, the adversary has an opportunity to get more comprehensive knowledge about the classifier to develop more sophisticated evasion attacks. The malformed instance can be applicable for all the users using the classifier. Besides, if the adversary gets perfect knowledge about the classifier, she can even reengineer a new classifier for commercial purposes. In this study, it is assumed that all the implementation and configuration of the client-side classifier are available for the adversary. The adversary can figure out the type of classification model, the classification algorithm and the feature extraction method by leveraging various techniques. Considering the advancement of modern analysis techniques, this assumption is reasonable.

Some client-side classifiers, have introduced some defense techniques to prevent the adversary from learning crucial information. For example, GPPF employs the cryptography technique to protect the classification model. Unfortunately, it is proved to be ineffective to against classifier cracking (discussed in Section 3 and 4).

2.2 Phishing and GPPF

According to the latest report [3] of Anti-Phishing Working Group (APWG), phishing attacks remain widespread: the number of unique phishing reports submitted to APWG during Q4 of 2014 was 197,252, and there is an increase of 18 percent from the 163,333 received in Q3. To minimize the impact of phishing attacks, a variety of methods have been proposed to detect phishing pages, involving machine learning [39][52][56] or other techniques [24] [25] [27] [31] [33] [46] [57] [58].

Modern web browsers also provide detection tools to assist end users against phishing attacks. *Safe Browsing*, a service offered by *Chrome*, is aiming at providing not only blacklists of malicious URLs but also a trained classifier (GPPF) which automatically detects phishing pages as a countermeasure to the phishing problem [4]. In *Chrome*, *Safe Browsing* serves as a guard when a request comes, and the request URL will be checked before the content is allowed to begin loading. The URL is checked against two blacklists: malware and phishing. If the URL is matched with the blacklists, *Chrome* will block the request and jump to a warning page as shown in Figure 2. More importantly, for the URL that is not present in the blacklists, *Chrome* will further invoke GPPF to determine whether it is legitimate or phishing. In practice, the phishing blacklist needs to be updated constantly and users will be vulnerable to newly created phishing websites. GPPF acts as an indispensable role in protecting end users from unknown phishing pages.

In fact, GPPF is the local version of a Google’s internal classifier. Google developed and trained a scalable machine learning classifier in its servers to detect phishing websites and use it to maintain Google’s phishing blacklist automatically [56]. Training the classifier is a constant offline process. The training process uses a sample of roughly ten million URLs analyzed over the past three months as the training dataset. The number of URLs from a single domain is also limit to 150 per week to prevent a single domain from having too much contribution to the classification

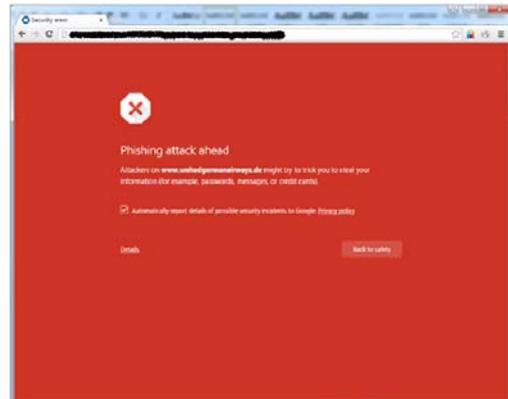


Figure 2. Phishing warning page.

model. Consequently, the adversaries don’t have an opportunity to alter the training dataset enough to make the trained classifier misclassify phishing pages as legitimate. However, to provide the real-time detection of unknown phishing pages, the trained classifier is also implemented as a part of *Safe Browsing*, i.e., GPPF. As an internal component of the *Chrome* browser, GPPF is completely deployed and running in the user environment. This actually allows the adversary to freely analyze its implementation and configurations to construct more sophisticated phishing attacks.

According to the report of StatCounter [5], from Aug 2014 to Aug 2015, *Chrome* shares an average of 48.6% market and is the most popular web browser. In May 2015, Google announced that *Chrome* has over one billion active users [1]. This means over one billion users’ web surfing are protected by GPPF. Note that if a phishing page can fool GPPF, it will have more chances to keep away from the Google’s phishing blacklist. Furthermore, the phishing blacklist provided by Google is also employed in *Firefox* and *Safari* browsers, as well as by Internet Service Providers (ISPs) [6]. We have reason to believe that the security breach of GPPF will potentially impact many more people besides just the users of *Chrome*.

3. CRACKING GPPF

There is very limited public information about the design and implementation of GPPF. We choose to directly analyze the development version of the *Chrome* browser, *Chromium*, to crack GPPF. The cracking includes two main steps: (1) extracting the classification model of GPPF from *Chromium*; and (2) decrypting the hashed features of the model. **It needs to be mentioned that some sensitive details of the cracking are intentionally omitted to prevent them from being used for malicious purposes.**

3.1 Extracting the Classification Model

3.1.1 Classification Algorithm

The multi-process architecture that *Chrome/Chromium* adopts helps it be more robust. According to a very brief description in [4], we can know that *Browser* process will periodically fetch an updated model from Google’s server and send it to every *Render* process via an IPC channel. This allows the classification to be done in the *Render* process, which will score the request page to tell whether it is phishing or not.

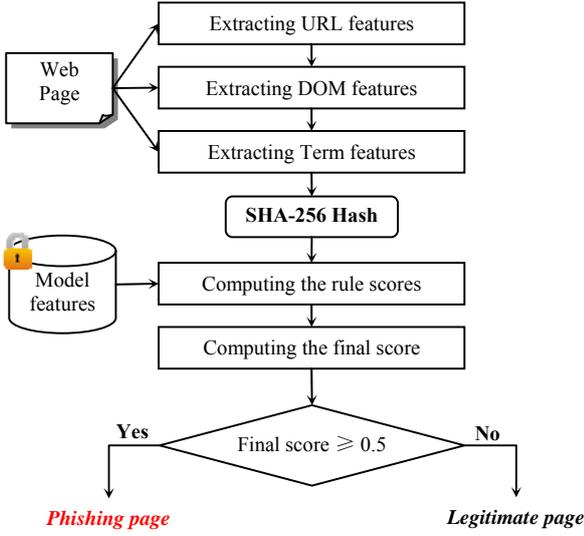


Figure 3. Classification workflow.

We collect the execution traces of a phishing page and a legitimate page by monitoring the Render process of them in *Chromium* using a debug tool *gdb*. With a differential analysis of the traces, we find the GPPF’s scoring function *ComputeScore()*, which is a method of the *Scorer* class located in the file *scorer.cc*. Combining a dynamically backward tracking of the execution path started from *ComputeScore()* and a static analysis on the source code, we conclude the workflow of the classification. As shown in Figure 3, the classifier first extracts three kinds of page features from the current web page in order, i.e., *URL*, *DOM* and *Term* features. Second, the collected page features are hashed with the SHA-256 algorithm and sent to the function *ComputeRuleScore()* to compute the rule score for every scoring rule, along with the hashed model features. Third, *ComputeScore()* combines all the rule scores to generate a final score for the current page. Finally, the score is compared with a predefined threshold (fixed in 0.5). If the score is smaller than the threshold, the page will be regarded as legitimate; otherwise, it will be identified as a phishing page and be blocked.

Based on the analysis of the scoring process, we find the GPPF is a *logistic regression* classifier, which uses the following two expressions to compute the phishing score for the target page. GPPF computes the total score for the page in log odds using the expression (1), and uses the normalization expression (2) to transform the score in log odds to the final score.

$$Logodds = W_1 + \sum_{i=2}^{2130} \left(W_i \prod_{j=1}^{n_i} V_{i,j} \right) \quad (1)$$

$$score = \frac{e^{Logodds}}{1 + e^{Logodds}} \quad (2)$$

According to expression (1), the computing of the log odds of the page involves 2,130 scoring rules. Every rule has a weight, namely $W_1 \sim W_{2130}$. Except for the first rule, every rule consists of from one to four (i.e., n_i for the i th rule) model features. Before computing the rule scores, the page features are first mapped to string forms, which will be hashed and compared with the model features. For every rule, the classifier creates a set of feature values (i.e., $V_{i,1} \sim V_{i,n_i}$) for all matched model features. For Boolean feature, *True* is converted to 1.0 and *False* is converted

Table 1. Scoring rule examples

Rule	Features #	Hashed Feature	Weight
R_{1494}	2	32ffbec120ed857f57f3d7bb37 e6652955b21da7a7efd81d9a9 aa2865173eb35	-1.26907706
		ec92914c7db4483437c84975 8c45cf8bbc6dd0148cdb2f72b ec0a728e8c91a7d	
R_{2050}	1	760e98536a709d0fcb9b717eb 542cc5af77bbabf60a501dbf7d f81a111d1e807	2.5238471

to 0.0. The continuous features are scaled to be between 0.0 and 1.0. If a model features of the rule is absent from the target page, its feature value will be set to 0.0. The score of the rule will be computed by combining the product of all the feature values and its weight. Finally, the log odds of the page will be produced by summing up all the rules scores.

To crack the classification model, we need to recover the weight and model features for every scoring rule. The rule weight can be collected by debugging the Render process. We set a breakpoint in *ComputeRuleScore()*, in which an extractor written in *gdb* script is invoked to read the weight information from the rule objects in the memory and save them in a file. In a similar way, we also get the number of model features for every rule. However, the model features are not stored in plaintext; instead, they are hashed with the SHA-256 algorithm and hidden in some complex data structures. With carefully tracking of the scoring process, we locate their addresses and designed a *gdb* script to extract them from rule objects. Take two extracted rules as examples. As listed in Table 1, the rule R_{1494} is a negative rule with two features. This kind of rule is used to identify the good property indicative of legitimate pages. On the contrary, the rule R_{2050} is a positive rule, including only one feature. Some of the model features are present in different scoring rules. After eliminating duplicates, in total, we collect 1,009 individual hashed model features. The decryption of them will be described in Section 3.2.

3.1.2 Model Features

To decrypt hashed model features, we should first get clear about their semantics and how the page features are mapped to them. When computing the score, three kinds of page features will be mapped to corresponding model features in different ways.

URL features. In practice, the phishers often obfuscate their URLs to hide suspicious addresses or confuse victims into believing they come from a trusted party. Based on the observation, in GPPF, some characteristics of the URL are employed to identify phishing pages. By analyzing the implementation of the classifier, we recover all seven kinds of properties of the URL being extracted as the page URL features, as shown in Table 2. The page URL features will be converted to string forms, which will be hashed and compared with the encrypted model URL features during computing scores.

The page URL features can be categorized into two groups. For one of the first group of features (the first three in Table 2), if it is present in the URL, a hashed predefined string will be taken as its corresponding model feature. Take the first page URL feature as an example. If the hostname part of the URL is a numeric IP address, the string “*UrlHostIsIpAddress*” is hashed with the SHA-256 algorithm to act as the model feature. For the second group of

Table 2. URL features

No.	Page URL Features	Model URL Features
1	The hostname is an IP address?	$UrlHostIsIpAddress$
2	The number of other host components is greater than one?	$UrlNumOtherHostTokens>1$
3	The number of other host components is greater than three?	$UrlNumOtherHostTokens>3$
4	Top level domain	$UrlTld=*$
5	The first host component below top level domain	$UrlDomain=*$
6	Other host components	$UrlOtherHostToken=*$
7	Path token in URL	$UrlPathToken=*$

features (the last four in Table 2), a string in equation form will be generated by concatenating a predefined string and the concrete URL property. For example, for the fifth page URL feature, if the URL is *www.phishing.com*, the string “ $UrlDomain=phishing$ ” will be hashed as the model feature. In scoring rules, all the URL features will be assigned a Boolean feature value, i.e., 1.0 if it is present in the page or 0.0 if it is absent.

The predefined strings used to generate the model feature (shown in the third column of Table 2) can be inferred from the implementation of the classifier. However, we cannot directly recover the complete plaintexts from the hashed model URL features in equation forms. In GPPF, there are hundreds of model features about the URL in equation modes. Based on their semantics discussed above, we design a collision attack to decrypt them as far as possible (described in Section 3.2).

DOM features. GPPF also uses some features about the Document Object Model (DOM) elements of the page to tell whether or not it is phishing. As shown in Table 3, we recover all 12 kinds of DOM features being employed by GPPF. In a similar way to the URL features, these page DOM feature will also be converted to string forms.

As listed in Table 3, the first seven page DOM features are used to identify the structure property of the page, e.g., to determine whether the page has some kinds of DOM elements or not. These features directly correspond to seven predefined strings respectively, which will be hashed and compared with the model features. For example, if the page has the `<form>` element, the string “ $PageHasForms$ ” will be hashed to act as the corresponding model feature. The eighth page DOM feature records all external domains that the page links to, which will be mapped to a string in equation mode for every individual external domain. In scoring rules, all the above DOM features will be assigned a Boolean feature value. The last four page DOM features indicate the fraction of some kinds of DOM elements in all elements. They correspond to four predefined strings. In scoring rules, the values of matched features are set to the fraction value scaling between 0.0 and 1.0.

For the DOM features, related predefined strings can be directly recovered and are shown in the third column of Table 3. For the eighth page DOM feature, there are many related hashed model features in equation forms to identify different external domains.

Table 3. DOM features

No.	Page DOM Features	Model DOM Features
1	Page has <code><form></code> element?	$PageHasForms$
2	Page has <code><input type=text></code> element?	$PageHasTextInput$
3	Page has <code><input type=password></code> element?	$PageHasPswdInputs$
4	Page has <code><input type=radio></code> element?	$PageHasRadioInputs$
5	Page has <code><input type=checkbox></code> element?	$PageHasCheckInputs$
6	The number of <code><script></code> elements in the page is greater than 1?	$PageNumScriptTags>1$
7	The number of <code><script></code> elements in the page is greater than 6?	$PageNumScriptTags>6$
8	Token feature containing each external domain that is linked to	$PageLinkDomain=*$
9	Fraction of form elements whose action points to an external domain	$PageActionOtherDomainFreq$
10	Fraction of links in the page which point to an external domain	$PageExternalLinksFreq$
11	Fraction of page links that use https	$PageSecureLinksFreq$
12	Fraction of images whose src points to an external domain	$PageImgOtherDomainFreq$

A collision attack is performed to recover their plaintexts (described in Section 3.2).

Term features. In GPPF, the terms appearing in the page are taken as a kind of feature. A term feature can be a single word or a compound of multiple words (at most five).

When fetching the page terms, the page text is first converted to a list of words in lowercase. In practice, using every word of the page text to construct features will greatly overburden the learning process. Instead, GPPF only makes features of the words contained in a predefined set. A fast hash algorithm, *Murmurhash3*, is employed to implement a word filter. GPPF maintains a list of candidate words, which are hashed with the *Murmurhash3* algorithm. It was generated by collecting the words with the highest term frequency-inverse document frequency (TF-IDF) values [50] from a large dataset.

GPPF uses an array *previous_words* to construct the page term features, which can store at most five continuous candidate words of the page text. The array is initially empty. The first word is fetched and removed from the page word list. Its *Murmurhash3* value is computed to determine whether it is contained in the candidate list or not. If it is a candidate, the word will be added in the first element of *previous_words*. GPPF then checks the subsequent word in the list and adds it to the array in sequence if it is also a candidate word. For every time a word is added, all words currently contained in the array (at most five) are connected and combined with a predefined prefix (“ $PageTerm=$ ”) to construct a phrase. It will be hashed with SHA-256 algorithm and compared with the hashed model term features. For example, if three continuous words (“*abc*”, “*def*”, and “*ghi*”) have been added in the array, the generated corresponding phrases will be “ $PageTerm=abc$ ”, “ $PageTerm=abc def$ ”, and “ $PageTerm=abc def ghi$ ”. In scoring rules, the values of a term feature will be set to 1.0 if there is a matched phrase; otherwise to 0.0. When

Table 4. Model features needed to be decrypted

Category	Model Features	Total	Decrypted	%
URL-related	<i>UrlTld=*</i>	563	69	75.7%
	<i>UrlDomain=*</i>		21	
	<i>UrlOtherHostToken=*</i>		28	
	<i>UrlPathToken=*</i>		201	
	<i>PageLinkDomain=*</i>		107	
term-related	<i>PageTerm=*</i>	432	375	86.8%
Sum		995	801	80.5%

encountering a non-candidate word or the array is full, GPPF will clear the array, fetch the next word and repeat the above steps until the list is empty.

In GPPF, there are 432 hashed model term features. Every one corresponds to a phrase that may consist of one to five words. We also use a collision attack to recover their plaintexts.

3.2 Collision Attacks

As discussed in Section 3.1, besides 14 features being directly recovered in the model extraction, there still are 995 hashed model features needed to be decrypted. As shown in Table 4, they can be divided into two categories: *URL-related* and *term-related*. According to their semantics, we design different collision attacks to decrypt them.

3.2.1 Decrypting URL-related Features

In total, there are 563 hashed URL-related features. So far, it is impossible to directly construct a collision for a given SHA-256 hash value. Instead, we collect four datasets related to URLs to perform targeted brute force attacks to find potential collisions as much as possible. To prevent the adversary from reproducing the attacks, the sources of the datasets are not presented in this paper.

- 1) We use a dataset with about 8,000 top level domain names to decrypt *UrlTld* features. We select the name from the set one by one and add the prefix "*UrlTld=*" to generate a test case. By hashing it with SHA-256 and comparing the hash value with all URL-related features, we successfully recover 69 *UrlTld* features with a desktop computer in about five minutes.
- 2) We collect over 30,000 URLs of history phishing pages, and use the different elements of the URLs (e.g., hostname) to generate test cases for other four kinds of URL-related features. In a similar way as above, 171 features are successfully decrypted in about four minutes, including 20 *UrlDomain*, 27 *UrlOtherHostToken*, 17 *UrlPathToken*, and 107 *PageLinkDomain* features.
- 3) With the URLs of legitimate pages in thousands of top sites, we get 3 *UrlDomain* features and 34 *PageLinkDomain* features in less than one minute.
- 4) A very large URL database with over 2,000,000 records is leveraged to construct test cases. The decryption process takes about 20 minutes. As a result, we get 46 *UrlTld* features, 21 *UrlDomain*, 28 *UrlOtherHostToken*, 201 *UrlPathToken* and 107 *PageLinkDomain* features.

After removing duplicates, as listed in Table 4, we eventually recover a total of 426 (75.7%) URL-related features, including 69

Table 5. Decrypting the term features with seven corpora

Language	Decrypted	Time
English	201	1.7 hours
French	6	2.3 hours
German	51	3.2 hours
Spanish	5	1.1 hours
Dutch	1	6 minutes
Chinese	27	20 minutes
Japanese	1	5 minutes
Sum	292	8.8 hours

UrlTld, 21 *UrlDomain*, 28 *UrlOtherHostToken*, 201 *UrlPathToken* and 107 *PageLinkDomain* features.

3.2.2 Decrypting Term-related Features

GPPF employs 432 hashed term features to detect phishing pages based on the page text. In practice, the text of a phishing page can be written in various languages. To this end, we collect some full-text corpora for seven popular natural languages (English, French, German, Spanish, Dutch, Chinese and Japanese) to perform collision attacks. The basic steps are as follows.

- According to the semantics of the term feature, we build a candidate word filter based on the implementation of the Murmurhash3 algorithm in *Chromium*. With it, we extract all possible word sequences consisting of one to five continuous candidate words from these corpora respectively.
- For every word sequence, adding the prefix "*PageTerm=*" to generate a test case.
- Hashing every test case with SHA-256 and comparing the hash value with all term features to find potential collisions.

Via the above steps, we successfully recover 292 (67.9%) term features in various languages in about 8.8 hours. The result is detailed in Table 5.

To further improve the cracking result about term features, we also perform blind brute force attacks. We construct an alphabet consisting of letters in western languages. With the alphabet, all possible combinations of no more than eight letters are produced. After filtering, they are used as candidate words to generate test cases to find collisions. Surprisingly, in about 16 hours, we recover 281 term features only using a part of test cases. In a similar way, we also quickly recover 40 term features based on a set of Chinese, Japanese and Korean (CJK) ideographs. The related results are detailed in Table 6 and Table 7 respectively.

After combining all above attacks results and removing duplicates, we eventually recover a total of 375 (86.8%) term features.

Table 6. Decrypting the term features with an alphabet

Term Size	Candidate Words	Decrypted	Time
1-word	1-letter to 8-letter	186	1 minute
2-word	1-letter to 8-letter	76	8.6 hours
3-word	1-letter to 6-letter	15	14 minutes
4-word	1-letter to 4-letter	4	7.4 hours
Sum		281	16.25 hours

Table 7. Decrypting the term features with CJK ideographs

Term Size	Candidate Words	Decrypted	Time
1-word	1-ideograph to 3-ideograph	31	1 minute
2-word	1-ideograph to 3-ideograph	7	< 1 minute
3-word	1-ideograph to 3-ideograph	2	< 1 minute
4-word	1-ideograph to 3-ideograph	0	2 minute
Sum		40	5 minutes

3.3 Result Analysis

As shown in Table 4, we successfully decrypt 801 (80.5%) model features with collision attacks. Together with 14 features being recovered in the model extraction, we eventually get the complete plaintexts of a total of 815 (80.8%) model features.

After applying the decryption result to 2,130 extracted scoring rules, we can completely reverse engineer 1807 (84.8%) rules, namely every feature of them is decrypted. Besides, there are also 196 (9.2%) rules we cannot completely crack, but at least one of their features is decrypted. Only 127(6.0%) rules remain confidential, no one of their features is cracked.

According to their weights, GPPF’s scoring rules can be categorized into two types: *positive rules* and *negative rules*. As their names suggest, the former are assigned with a positive weight and can cause a rise in the phishing score for the page, while the latter are just the opposite. Naturally, the top-weighted positive or negative rules will make remarkable contributions to tell whether a page is phishing. After analyzing top 100 most weighted positive rules, we learn that 66 of them are completely reverse engineered, and 20 are partially cracked. For the top 100 most weighted negative rules, 77 of them are completely reverse engineered, and 21 are partially cracked. In other words, given the cracking result, the adversary has a great chance to disguise a phishing page as a legitimate one by targetedly manipulating its content.

4. EVASION ATTACKS

In this section, we perform some evasion experiments to demonstrate the effectiveness of the classifiers cracking via exploiting the recovered knowledge.

For a specific phishing page, we can infer adding or removing what features can reduce its phishing score based on the cracking result presented in Section 3. If a feature can provide negative contributions to the phishing scoring for a page, we call it as a *good* feature from the adversary’s point of view. On the contrary, if a feature only has positive contributions, we call it a *bad* feature. Correspondingly, we design two kinds of evasion attacks, *good features insertion* and *bad features elimination*. The basic idea behind them is to add or remove appropriate good or bad features into or from a phishing page to make its phishing score lower than the threshold, resulting in a misclassification. The latest 100 real phishing pages are collected from PhishTank as the attack dataset. We will try to use the two evasion attacks to manipulate them to evade the detection of GPPF. To minimize the potential side-effects, we will use pseudonyms when referring to specific good features or bad features in the following part of this section.

Table 8. The required number of Good features

Feature	MIN	MAX	Average
URL	1	10	2.5
DOM	1	6	2.2
Term	1	17	3.7

4.1 Good Features Insertion

Given a phishing page, there may be many features can be leveraged to reduce its phishing score. By utilizing plenty of negative rules having been completely reverse engineered; we can adopt a very primitive but effective way to choose desirable good features. In fact, we can sort all negative rules only with one recovered feature by their weights, and directly use the features of top-weighted rules as good feature candidates for all target pages. More surprisingly, for many phishing pages in the dataset, we can easily convert them to legitimate pages only by inserting just one such good feature. Moreover, as detailed in Table 8, we find that only using one kind of good feature can also be effective. For example, we can reduce the scores of all test pages lower than 0.50 by inserting at most six good DOM features into the page. On average, 2.2 good DOM features are required.

It should be noted that a sophisticated adversary can carefully introduce the good features to preserve the utility of phishing pages. For example, to prevent the inserted terms from attracting the attention, their color can be set to background color.

After introducing above good features, the manipulated test pages are deployed in our Web server. We then use the latest version of *Chrome* (45.0.2454) to visit them one by one to check whether they can successfully evade the detection of GPPF. We find all the dressed-up pages (100%) are regarded as legitimate pages and display properly in the browser. For example, there is a phishing page to imitate the login page of Chase Bank. When browsing it, the *Chrome* can successfully block it as a phishing page and jump to the warning page as shown in Figure 2. In fact, the page is given a very high phishing score 0.9986. However, after inserting six good term features $T_1 \sim T_6$ into its text, the score is reduced to only 0.2784. As a result, the dressed-up page can be normally visited with *Chrome* as shown in Figure 4.



Figure 4. The dressed-up phishing page can evade GPPF.

4.2 Bad Features Elimination

Compared with the good features insertion, selecting proper bad features from a given phishing page to perform an effective evasion attack is not a trivial task. The number of available bad

features is limited for a given page. Additionally, some features can be referred by multiple scoring rules. A feature may be not only present in a positive rule but also in a negative rule. Directly removing the features in positive rules may also result in some negative rules losing their efficacy.

To this end, we design a search-based method to automatically select proper bad features for a given page. Specifically, we implement a script to compute the contribution of one feature or a feature set to the final score, by removing it or them from the page and re-computing the score. For a given page, we apply the script to all its recovered features to search a feature or a feature set whose contribution is enough to the exploitation. Namely, after removing the feature(s), the score of the page will be lower than the threshold, allowing it to be classified as a legitimate page.

With the method, we successfully find proper bad features for every test page respectively. By eliminating corresponding bad features from the pages respectively, all test pages (100%) can evade the classification and normally display in *Chrome*. Take the phishing page shown in Figure 4 as an example. We find four bad term features $BT1 \sim BT4$ for it and eliminate them with some obfuscation techniques, such as changing a word from singular to plural form. As a result, we succeed in reducing its score from 0.9986 to 0.4591 and dressing it up as a legitimate page. In the experiment, we find that removing at most five bad features is enough to make the page evade the classifier. On average, 3.1 bad DOM features are required.

5. MITIGATION

Google developers have discussed the potential adversarial attacks that GPPF might encounter [56]. They believe that possible attacks on GPPF are either limited or expensive. From their point of view, the adversary who tries to evade GPPF by disguising the phishing page as a legitimate one cannot preserve its utility and visual similarity at the same time. However, thanks to the cracking results, we can purposely introduce some easy-to-hide good features to evade GPPF with a very low cost. For example, we are able to make the newly added term features invisible by setting their color to background color of the target page.

In practical applications, GPPF is proved to be a very valuable tool against phishing attacks under the non-adversarial environment. Tens of thousands phishing sites are detected by Safe Browsing per week [6]. To this end, the developers may want to improve its robust as well as change the architecture as little as possible. A natural and direct idea is to select the features difficult to being recovered by brute force attacks. For example, the developers can just select the comparative long phrases, 5-word phrases or even longer, as the term features. This would result in a combinatorial explosion when the adversary performs a blind brute force attack for cracking. The computation of enumerating and hashing all possible compounds of five words is unacceptable. Unfortunately, this idea is not effective enough when the adversary is aware of the feature extraction method. In fact, the adversary can still reverse engineer sufficient features by collecting appropriate page-related data as test cases to perform a collision attacks. The adversary can take the data as a web page to extract the possible word sequences according to the feature extraction method, and hash the sequences to check whether they are a term feature. For a concrete feature extraction method, the amount of sequences is actually limited regardless of how long the sequence is. Given appropriate test cases, the adversaries have

a fair chance to find sufficient collision instances. As presented in Section 3.2.2, we recover 292 (67.9%) term features only using seven full-text corpora in 8.8 hours. These features are already enough for evasion attacks.

Based on above discussions, we can learn that the most effective defense way is to essentially increase the complexity of reverse engineering the classification model, especially the semantics of features. To effectively solve this problem, we propose a potential defense strategy: employing deep learning method to construct naturally robust client-side classifiers.

Deep learning, as a novel powerful machine learning method, has been widely applied in fields of multimedia, natural language procedure, data mining, etc. Based on the basic thoughts of deep learning, especially the layer-wise learning and fine tuning, some powerful deep neural networks such as LeNet [38], stacked auto-encoder (SAE) [44] and deep belief net (DBN) [10][30], have been proposed to detect high level features and produce complicated decision functions [10].

In deep learning method, gradually from the low layers to the high layers, deep neural networks (DNN) could realize efficiently feature extractions. Generally, the features imported into the low layers are the raw data describing the basic original properties of the problem instances. After the effect of multiple layers abstraction, the features extracted from high layers possess complicate semantic information, which is hard to be comprehended for researchers [54]. In other words, no one could exactly explain the mapping between the raw data and the high level features.

For an instance, as shown in Figure 5, the deep convolutional neural network (DCNN) proposed by Hinton [37] could classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into 1,000 different classes. The DCNN has eight learning layers, i.e., five convolutional and three fully-connected. The network adopts several effective strategies, such as *ReLU nonlinearity*, *max pooling*, *dropout* and *local response normalization*, to increase its classifier ability. This enhances DCNN complexity. In detail, the input data are RGB images with 224*224 pixels. The last fully-connected layers output the high level features with 4,096 dimensions. As a result, every dimension in the high level features is relative with all dimensions of the inputting image. The weight connections among layers, partially embedded with non-linear activation functions, are greatly complicated and huge (almost millions). Consequently, even though researchers have grasped the meaning of one specific dimension in the highest layer, they could not find a feasible method that corresponds to modify partial dimensions of the raw data, and meanwhile affect only the specific dimension. This characteristic of deep learning is especially fit for building a robust classifier to against classifiers cracking.

Similarly, we can build a phishing pages classifier based on the deep learning method, to make it hard to be reverse engineered by analyzing its implementation and configurations. As shown in Figure 6, the classifier has n layers, i.e., $n-1$ layers for feature extraction, and one for classification. The classifier has been trained off-line and deployed in the client-side. The raw html represented as a large sparse feature vector, e.g., the features of bag-of-words model [35], is imported into the first layer L_1 originally. The feature vector is transformed by several functions, e.g. *long short term memory* (LSTM) [28], *max pooling*, and *sigmoid activation*, in different layers to produce high level

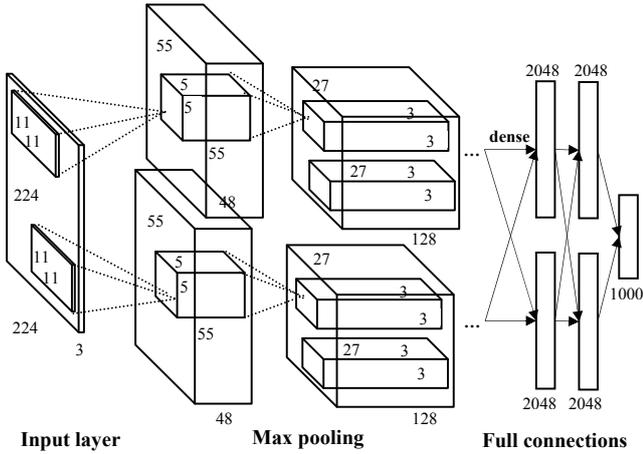


Figure 5. The deep architecture adopted in DCNN for ImageNet classification.

features in high layers. The features, output as high abstractive properties in layer L_{n-1} , are exported to the classification algorithm (as L_n). The final classification algorithm can be non-linear or linear. The original features are eventually transformed to the final high level features with $n-2$ times of complex many-to-many mapping. As a result, the complex relationship among them heavily increases the difficulties of inferring the original features from the final ones.

When the classification model is cracked by the adversary, even though the deep architecture is completely recovered, the complexity of feature mappings can still effectively ensure the robust of the classification. In fact, the feature extraction in classifier can be regarded as a black-box in applications, which performs complex non-linear transforms with mutual affection of multiple layers weight connections. The amount of weight connections are millions, as exhibited in [55]. This makes it impossible for the adversary to figure out the exploitable relationships among the input and output features. We have reasons to believe that the deep learning method inherently has the energy to against classifiers cracking.

6. DISSCUSSION

In this study, we present an attack methodology, classifiers cracking, aiming at client-side classifiers and successfully demonstrate its effectiveness with a widely-used classifier, GPPF. In theory, the methodology is generic and applicable to other client-side classifiers. However, when applying the methodology to a specific classifier, we need to develop a specially designed crack techniques according to its implementation. In fact, there are many classifiers equipped with different classification algorithms, such as [7][40]. To further demonstrate the security challenges brought by classifiers cracking, in the future, we will pay attention to some other types of classifiers and investigate their security from the point of cracking. These classifiers may take security into consideration to different extents and be deployed in different ways. More reverse engineer techniques may need to be employed to crack them.

As described in Section 3 and 4, we eventually completely reverse engineer 84.8% scoring rules of the GPPF classification model, which is proved to be sufficient for launching effective

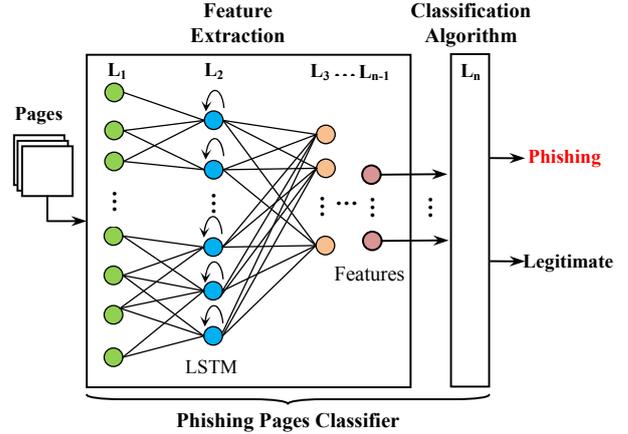


Figure 6. DNN model for phishing pages classification.

evasion attacks. However, in fact, we can get better cracking results by introducing more appropriate corpora. For example, using a comprehensive database of history phishing pages can decrypt more term features. Sometimes, the adversary may want to get perfect knowledge about a classifier for some special purposes, such as stealing its techniques to reengineer a new classifier. Besides, it needs to be emphasized that some seemingly unrelated dataset, e.g., a corpus, can also be leveraged to compromise the security of client-side classifiers. The developers should collect as much as possible dataset, especially publicly available, to evaluate the robust of their classifier before releasing it.

We have got sufficient knowledge about the GPPF classification model by cracking it. This allows us to easily find exploitable good and bad features for a given page. In this study, it is not necessary to design a sophisticated algorithm to more effectively and efficiently find exploitable features. However, if the adversary has only limited knowledge about the target classifier, she can develop a powerful algorithm to discover exploitable features. Furthermore, in theory, combining the good features insertion and bad features elimination can produce better performance. It is also helpful for the adversary to attack a classifier. To this end, developers should prevent the information of their classifier from being inferred by the adversaries as far as possible.

7. RELATED WORK

Many existing studies have paid much attention to the security of classifiers, and the arm race between adversaries and defenders will never end.

Attacks. The attacks can be categorized into two types by their influences: causative attacks and exploratory attacks.

In causative attacks, the adversary has the chances to affect the training process by contaminating training data (e.g. injecting many specially crafted samples). This kind of attack has been used to degrade the performance of a lot of learning-based applications, such as biometric authentication [16][18], spam filtering [45], and network intrusion detection [34][49]. In [16], a method is proposed to mislead an adaptive biometric system to perform self-update by submitting a proper sequence of spoofed biometric traits to the sensor and cause a misclassification

eventually. A further work [18] reveals that poisoning attacks can be used to compromise face templates in a more general case. Another study [45] succeeds in exploiting machine learning to compromise a spam filter by manipulating the filter’s training data. They proposed two kinds of poisoning attacks by inserting different sets of words into attack emails: *dictionary attacks* inject words indicative of legitimate emails to increase misclassifications, and *focused attack* tries to introduce words to have the filter block one specific kind of emails (e.g. emails from business rivals). Besides, as discussed in [20][34][49], the intrusion detection systems may also be vulnerable to causative attacks. The adversary can inject carefully crafted malicious traffic samples into training dataset and finally force the classifier to learn a wrong model of the normal traffic.

In exploratory attacks, the adversary tries to figure out as much knowledge (e.g. type of classifier, features, and threshold) of the classifiers as possible to effectively evade them. Exploratory attacks have been applied to various security applications. Lowd and Meek conduct an attack that minimizes a cost function [42]. They further propose attacks against statistical spam filters that add the words indicative of non-spam emails to spam emails [41]. The same strategy is employed in [45]. In [43], a simple but effective attack methodology called *reverse mimicry* is designed to evade structural PDF malware detection systems. The main idea is injecting malicious content into a legitimate PDF while introducing minimum differences within its structure. The related experiments show that some very popular classification algorithms (e.g. SVMs and neural networks) can also be evaded with this method. A recent work [51] uses PDF_{RATE}, an online learning-based system for detection of PDF malware, as a case to investigate the effectiveness of evasion attacks. The study reconstructs a similar classifier through training one of the publicly available datasets by a few deduced features, and then evaded PDF_{RATE} by insertion of dummy content into PDF files. Additionally, in [17], a simple algorithm is proposed for evasion of classifiers with differentiable discriminant functions. The study empirically demonstrated that very popular classification algorithms, e.g., SVMs and neural networks, can still be evaded with high probability even if the adversary can only learn limited knowledge.

Unfortunately, to our best knowledge, all of the existing studies don’t pay any special attention to the client-side classifiers. As demonstrated in this study, the client-side classifiers have a larger attack surface and hence larger number of potential attacks.

Defenses. Many countermeasures against evasion attacks have been proposed, such as using game theory [21][22] or probabilistic models [15][48] to predicted attack strategy to construct more robust classifiers, employing multiple classifier systems (MCSs) [12][13][14] to increase the difficulty of evasion, and optimizing feature selection [29][36] to make the features evenly distributed.

Game-theoretical approaches [21][22] model the interactions between the adversary and the classifier as a game. The adversary’s goal is to evade detection by minimally manipulating the attack instances, while the classifier is retrained to correctly classify them. However, the retraining procedure is very expensive in the situation where the classifier is cracked. The adversary always can construct an attack instance to evade the current classifier. Similarly, for approaches based on probabilistic models [15][48], the adversary can also easily construct a hard-to-predict attack instance based on cracked knowledge.

MCSs [12][13][14], as the name suggests, use multiple classifiers rather than only one to improve classifier’s robustness. The adversaries who want to effectively evade the classification have to fight with more than one classifier. Although MCSs actually increases the workload of classifiers cracking, it doesn’t improve the security of client-side classifiers fundamentally.

In [29], the method *weight evenness* via feature selection optimization is proposed. By appropriate feature selection, the weight of every feature is evenly distributed, thus the adversaries have to manipulate a larger number of features to evade detection. In [36], the features are reweighted inversely proportional to their corresponding importance, making it difficult for the adversary to exploit the features. However, given sufficient knowledge, the adversary can easily find enough exploitable features. Besides, in many cases, the adversary can hide the manipulation very deeply without attracting the attention. For example, a phisher can leverage various HTML techniques to make good features invisible.

These defense techniques are built on the assumption that the classification model is kept confidential to the adversary or can be updated timely. However, when the adversary learned sufficient knowledge by cracking classifiers, they can easily and quickly construct effective evasion attacks targeted to the defense techniques.

8. CONCLUSIONS

In this paper, we presented a new attack methodology, *classifier cracking*, for evading the client-side classifier. Our approach is different from existing attack methods is that various reverse engineering techniques are leveraged to directly extract desirable knowledge from client-side classifier for launching evasion attacks. Our study took GPPF, a learning-based filter for phishing pages deployed in *Chrome* as a case to study, which owns over one billion users. Employing various reverse engineering techniques, we successfully crack the GPPF model and completely recovered 84.8% encrypted scoring rules. Based on the information, we developed two kinds of evasion attacks: *good features insertion* and *bad features elimination*. The latest 100 real phishing pages collected from PhishTank were taken as the target of evaluation. The attack experiments showed that we can easily manipulate all the phishing pages (100%) to make them successfully evade the detection of GPPF in the latest version of *Chrome*. Additionally, a potential defense strategy was also discussed. We believe that the deep learning method can be employed to build client-side classifiers for essentially increasing the complexity of cracking.

Our research revealed an important fact that the client-side classifiers have a larger attack surface and hence larger number of potential attacks. In the future, we will further research potential defense techniques, especially based on the deep learning method, to develop more robust client-side classifier framework.

9. ACKNOWLEDGMENTS

The authors would like to thank the anonymous reviewers for their insightful comments. The work is supported by National Natural Science Foundation of China (NSFC) under grants 61170240, 91418206 and 61472429, and National Science and Technology Major Project of China under grant 2012ZX01039-004.

10. REFERENCES

- [1] Google has over a billion users of Android, Chrome, YouTube, and search. <http://www.theverge.com/2015/5/28/8676599/google-io-2015-vital-statistics>
- [2] PhishTank. <https://www.phishtank.com/>
- [3] Phishing Attack Trends Report of the 4th Quarter in 2014. http://docs.apwg.org/reports/apwg_trends_report_q4_2014.pdf
- [4] Design Documents of Safe Browsing. <http://www.Chromium.org/developers/design-documents/safebrowsing>
- [5] Market share of popular web browsers from Aug 2014 to Aug 2015. <http://gs.statcounter.com/#browser-ww-monthly-201408-201508>
- [6] Google's Safe Browsing service protects 1 billion Chrome, Firefox, and Safari users from malware and phishing. <http://thenextweb.com/google/2013/06/25/googles-safe-browsing-service-now-protects-1-billion-Chrome-firefox-and-safari-users-from-malware-and-phishing/>
- [7] I. Androustopoulos, G. Paliouras, V. Karkaletsis, G. Sakkis, C. D. Spyropoulos, and P. Stamatopoulos. Learning to filter spam e-mail: A comparison of a naive bayesian and a memory-based approach. In *Proceedings of the Workshop on Machine Learning and Textual Information Access*. 2000.
- [8] M. Barreno, B. Nelson, R. Sears, A. D. Joseph, and J. D. Tygar. Can machine learning be secure?. In *Proceedings of the 2006 ACM Symposium on Information, Computer and Communications Security*. ASIACCS'2006. ACM, 16–25.
- [9] M. Barreno, B. Nelson, A. Joseph, and J. Tygar. The security of machine learning. *Machine Learning*. 2010. Springer, 121–148.
- [10] Y. Bengio. Learning deep architectures for AI. 2009. *Foundations and trends® in Machine Learning*, 2(1), 1-127.
- [11] B. Biggio, G. Fumera, and F. Roli. Adversarial pattern classification using multiple classifiers and randomization. In *Proceedings of the 2008 Joint IAPR International Conference on Structural, Syntactic, and Statistical Pattern Recognition*. SPR'2008. Springer, 500-509.
- [12] B. Biggio, G. Fumera, and F. Roli. Multiple classifier systems for adversarial classification tasks. In *Proceedings of the 8th International Workshop on Multiple Classifier Systems*. MCS'2009. Springer, 132-141.
- [13] B. Biggio, G. Fumera, and F. Roli. Multiple classifier systems for robust classifier design in adversarial environments. 2010. *International Journal of Machine Learning and Cybernetics*, 1(1-4), 27-41.
- [14] B. Biggio, G. Fumera, and F. Roli. Multiple classifier systems under attack. In *Proceedings of the 9th International Workshop on Multiple Classifier Systems*. MCS'2010. Springer, 74-83.
- [15] B. Biggio, G. Fumera, and F. Roli. Design of robust classifiers for adversarial environments. In *Proceedings of the 2011 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. SMC'2011. IEEE, 977-982.
- [16] B. Biggio, G. Fumera, F. Roli, and L. Didaci. Poisoning adaptive biometric systems. In *Proceedings of the 2012 Joint IAPR International Conference on Structural, Syntactic, and Statistical Pattern Recognition*. SPR'2012. Springer, 417–425.
- [17] B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrđić, P. Laskov, G. Giacintoet , and F. Roli. Evasion attacks against machine learning at test time. *Machine Learning and Knowledge Discovery in Databases*. 2013. Springer, 387-402.
- [18] B. Biggio, L. Didaci, G. Fumera, and F. Roli. Poisoning attacks to compromise face templates. In *Proceedings of the 2013 International Conference on Biometrics Compendium*. 2013.
- [19] B. Biggio, I. Pillai, S. R. Bulò, D. Ariu, M. Pelillo, and F. Roli. Is data clustering in adversarial settings secure?. In *Proceedings of the 6th ACM Workshop on Artificial Intelligence and Security*. AISec'2013. ACM, 87–98.
- [20] B. Biggio, G. Fumera, and F. Roli. Security evaluation of pattern classifiers under attack. In *Proceedings of IEEE Transactions on Knowledge and Data Engineering*. TKDE'2014. IEEE, 26(4): 984–996.
- [21] M. Brückner and T. Scheffer. Stackelberg games for adversarial prediction problems. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery in Data Mining*. KDD'2011. ACM, 547-555.
- [22] M. Brückner, C. Kanzow, and T. Scheffer. Static prediction games for adversarial learning problems. 2012. *The Journal of Machine Learning Research*, 13(1), 2617-2654
- [23] D. Canali, M. Cova, G. Vigna and C. Kruegel. Prophiler: a fast filter for the large-scale detection of malicious web pages. In *Proceedings of the 20th International Conference on World Wide Web*. WWW '2011. ACM, 197-206.
- [24] Y. Cao, W. Han, and Y. Le. Anti-phishing based on automated individual white-list. In *Proceedings of the 4th ACM Workshop on Digital Identity Management*. DIM '2008. ACM, 51–60.
- [25] N. Chou, R. Ledesma, Y. Teraguchi, and J. C. Mitchell. Client-side defense against web-based identity theft. In *Proceedings of the Network and Distributed System Security Symposium*. NDSS'2004.
- [26] N. Dalvi, P. Domingos, Mausam, S. Sanghai, and D. Verma. Adversarial classification. In *Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD'2004. ACM, 99-108.
- [27] A. Y. Fu, W. Liu, and X. Deng. Detecting phishing web pages with visual similarity assessment based on earth mover's distance (EMD). In *Proceedings of IEEE Transaction on Dependable Secure Computing*. TDSC'2006. IEEE, 301–311.
- [28] F. A. Gers and J. Schmidhuber. LSTM recurrent networks learn simple context free and context sensitive languages. In *Proceedings of IEEE Transactions on Neural Networks*. 2001. IEEE, 12(6):1333–1340.
- [29] A. Globerson and S. Roweis. Nightmare at test time: robust learning by feature deletion. In *Proceedings of the 23rd International Conference on Machine Learning*. ICML'2006. ACM, 353-360.

- [30] G. E. Hinton, S. Osindero, and Y. -W. Teh. A Fast Learning Algorithm for Deep Belief Nets. *Neural Computation*. 2006. MIT Press, 18(2):1527–1554.
- [31] T. Holz, C. Gorecki, K. Rieck, and F. C. Freiling. Measuring and detecting fast-flux service networks. In *Proceedings of the Network and Distributed System Security Symposium*. NDSS’2008.
- [32] L. Huang, A. D. Joseph, B. Nelson, B. Rubinstein, and J. D. Tygar. Adversarial machine learning. In *Proceedings of the 4th ACM Workshop on Artificial Intelligence and Security*. AISec’2011. ACM, 43–58.
- [33] M. Khonji, Y. Iraqi, and A. Jones. Phishing detection: a literature survey. *Communications Surveys & Tutorials*. 2013. IEEE, 15(4): 2091-2121.
- [34] P. Kloft and M. Laskov. A “poisoning” attack against online anomaly detection. In *Proceedings of Neural Information Processing Systems (NIPS) workshop on Machine Learning in Adversarial Environments for Computer Security*. 2007.
- [35] Y. Ko. A study of term weighting schemes using class information for text classification. In *Proceedings of the 35th International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR’2012. ACM, 1029-1030.
- [36] A. Kolcz and C. H. Teo. Feature weighting for improved classifier robustness. In *Proceedings of the 6th Conference on Email and Anti-Spam*. CEAS’2009.
- [37] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proceedings of Neural Information Processing Systems*. NIPS’2012. MIT Press, 1097–1105.
- [38] Y. LeCun, F. J. Huang, and L. Bottou. Learning methods for generic object recognition with invariance to pose and lighting. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. CVPR’2004. IEEE, 97-104.
- [39] P. Likarish, D. Dunbar, and T. E. Hansen. B-apt: Bayesian anti-phishing toolbar. In *Proceedings of IEEE International Conference on Communications*. ICC ’2008. IEEE, 1745-1749.
- [40] O. Linda, T. Vollmer, and M. Manic. Neural network based intrusion detection system for critical infrastructures. In *Proceedings of the 2009 International Joint Conference on Neural Networks*. IJCNN’2009. IEEE, 1827-1834.
- [41] D. Lowd and C. Meek. Good word attacks on statistical spam filters. In *Proceedings of the 2nd Conference on Email Anti-Spam*. 2005.
- [42] D. Lowd and C. Meek. Adversarial learning. In *Proceedings of the 11th ACM SIGKDD International Conference on Knowledge Discovery in Data Mining*. KDD’2005. ACM, 641-647.
- [43] D. Maiorca, I. Corona, and G. Giacinto. Looking at the bag is not enough to find the bomb: an evasion of structural methods for malicious pdf files detection. In *Proceedings of the 8th ACM SIGSAC symposium on Information, Computer and Communications Security*. ASIACCS’2013. ACM, 119-130.
- [44] J. Masci, U. Meier, and D. Cire. Stacked convolutional auto-encoders for hierarchical feature extraction. In *Proceedings of the 21st International Conference on Artificial Neural Networks*. ICANN’2011. Springer, 52–59.
- [45] B. Nelson, M. Barreno, F.J. Chi, A. D. Joseph, B. I. P. Rubinstein, U. Saini, C. Sutton, J. D. Tygar, and K. Xia. Exploiting machine learning to subvert your spam filter. 2008.
- [46] P. Prakash, M. Kumar, R. R. Kompella, and M. Gupta. Phishnet: predictive blacklisting to detect phishing attacks. In *Proceedings of the 29th Conference on Information Communications*. INFOCOM’2010. IEEE, 346–350.
- [47] M. A. Rajab, L. Ballard, N. Lutz, P. Mavrommatis, and N. Provos. CAMP: Content-agnostic malware protection. In *Proceedings of the Network and Distributed System Security Symposium*. NDSS’2013.
- [48] R. N. Rodrigues, L. L. Ling, and V. Govindaraju. Robustness of multimodal biometric fusion methods against spoof attacks. 2009. *Journal of Visual Languages & Computing*, 169–179.
- [49] B. I. P. Rubinstein, B. Nelson, L. Huang, A. D. Joseph, S. Lau, S. Rao, N. Taft, and J. D. Tygar. Stealthy poisoning attacks on PCA-based anomaly detectors. *ACM SIGMETRICS Performance Evaluation Review*. 2009. ACM, 37(2): 73-74.
- [50] G. Salton and M. J. McGill. Introduction to modern information retrieval. 1983.
- [51] N. Šrđić and P. Laskov. Practical evasion of a learning-based classifier: A case study. In *Proceedings of the 2014 IEEE Symposium on Security and Privacy*. SP’2014. IEEE, 197-211.
- [52] F. Toolan and J. Carthy. Phishing detection using classifier ensembles. *eCrime Researchers Summit*. eCRIME ’2009.
- [53] K. Tretyakov. Machine learning techniques in spam filtering. *Data Mining Problem-oriented Seminar*. 2004. MTAT, 60-79.
- [54] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th International Conference on Machine Learning*. ICML’2008. ACM, 1096–1103.
- [55] K. Wang, X. Wang, L. Lin, M. Wang, and W. Zuo. 3D human activity recognition with reconfigurable convolutional neural networks. In *Proceedings of the ACM International Conference on Multimedia*. MM ’2014. ACM, 97–106.
- [56] C. Whittaker, B. Ryner, and M. Nazif. Large-scale automatic classification of phishing pages. In *Proceedings of the Network and Distributed System Security Symposium*. NDSS’2010.
- [57] Y. Zhang, J. I. Hong, and L. F. Cranor. Cantina: a content-based approach to detecting phishing web sites. In *Proceedings of the 16th International Conference on World Wide Web*. WWW ’2007. ACM, 639–648.
- [58] H. Zhang, G. Liu, T. Chow, and W. Liu. Textual and visual content based anti-phishing: A bayesian approach. In *Proceedings of IEEE Transactions on Neural Networks*. 2011. IEEE, 1532–1546.