

The rest of this section describes the structure of the whitelist as well as the features used in the SVM classifier.

3.1. Structure of the whitelist

The whitelist in our approach is an XML file that contains a user's login pages' URLs and a set of keywords (see Figure II). The key words are composed of the domain names of the page's URL and a set of terms from the Document Object Model (DOM) tree for the site. More precisely:

- The content of the "title" tag, ex: <title> text </title>.
- The content of the "meta keywords" tag, ex: <meta name="keywords" content="text"/>.
- The content of the "meta description" tag, ex: <meta name="description" content="text"/>.

These three tags provide a precise description of the webpage content. Stop words and punctuation symbols are omitted from the parsed field and the remaining words are concatenated and stored in the whitelist with the URL of the page. Figure II gives an example of the whitelist content.

```
<WhiteList>
<Website>
<url> http://www.facebook.com/ </url>
<Keywords>facebook helps connect friends
share posts people life...</Keywords>
</Website>
....
<Website>
<url> http://www.papal.com/ </url>
<Keywords>paypal send money payments
credit debit email...</Keywords>
</Website>
</WhiteList>
```

Figure 2: The structure of the whitelist.

The keywords are used to calculate similarity between a visited webpage and the pages of the whitelist. We used a "bag of words" model [20] for the construction of the keywords' frequency vector of each page and a "cosine" distance for similarity calculation. The similarity between a visited page " P_v " and a whitelist page " P_l " is calculated as follows:

Where: $f_{x,t}$ is the frequency of the term " t " in the set " x ".

Two pages are similar if their cosine similarity is close to 1.

The whitelist contains the login pages of the top 10 most attacked sites (by means of phishing) [21], as well as the user's typically visited pages. This avoids any user intervention, which is often a source of error, and facilitates installation and use.

3.2. The Classification Features

If the visited web page has no similarity with the pages of the whitelist, a feature vector is constructed for the page and it is subsequently processed by our SVM classifier. We use eight features to represent a page $P = \langle F_1, F_2, F_3, F_4, F_5, F_6, F_7, F_8 \rangle$. Some features are constructed according to the URL of the webpage, and others from its content.

- *Feature 1(F_1): URL with IP address*

For cost minimization reasons, many phishing pages use an IP address instead a domain name, contrary to a legitimate site which is accessed most commonly through a hostname. A link that does not contain a domain name and requests sensitive information to users is probably a phishing site. The feature F_1 is a binary feature: if the URL of a webpage contains an irresolvable IP address then $F_1 = 1$, otherwise $F_1 = 0$.

- *Feature 2(F_2): special characters in the URL*

For this feature we tested the existence of the "@" character in the URL. The presence of this character in a URL forces characters before it to be treated as user login credentials to URL's intended site. Phishers often use this character to mislead users: for example, if a user encounters the URL <http://paypal.com@www.phishpaypal.com>, he may be fooled into thinking the site is a legitimate Paypal page. Our F_2 feature is also binary, set to 1 if "@" is present in the URL and 0 otherwise.

- *Feature 3(F_3): presence of a Secure Sockets Layer (SSL) certificate*

The majority of financial and commercial institutions have an SSL certificate for their websites, which is not typically the case for phishing sites. The binary feature F_3 is set to 1 if an SSL certificate is present and 0 otherwise.

- *Feature 4(F_4): whether the identity of the webpage conforms to its URL*

The webpage identity is the most frequent base domain in the page's hyperlinks. For example, the identity of the page <http://www.facebook.com/> is "facebook.com"; the majority of the page's hyperlinks should contain the base domain. In general, legitimate web pages have an identity that matches their URL domain despite the existence of links that point to a foreign domain. Phishing pages, on the other hand, behave different. A phishing page that imitates a legitimate page usually retains the same hyperlinks, thus leading to a mismatch between the page's identity and its URL base domain. The binary feature F_4 is 1 if the webpage identity matches its URL domain and 0 otherwise.

- *Feature 5(F_5): search engine*

The principle of this feature is to apply the algorithm TF-IDF (Term Frequency-Inverse Documents Frequency) on a webpage to extract a document signature (the most relevant words). This feature is used by [16] and [17]. The signature is used as a keyword in a search engine query. If the page's URL is among the first results, the page is defined as legitimate, if not it is classified as a phishing page. We use the same principle but, instead of applying the TF-IDF technique, we use the keywords extraction method used when generating the whitelist (see section III-A). The keywords of the search engine query are composed of the page identity and the four most frequent words among those resulting from the extraction phase. We used the "matacrawler1" search engine which aggregates and returns

the relevant results from three search engines (Google, Yahoo and Bing). The binary F5 feature takes the value 1 if the URL of the page is among the top 20 search results and 0 otherwise.

- *Feature 6(F6): Nil anchors*

A nil anchor is an anchor that does not point anywhere, e.g.: ``, ``. Some phishing pages that imitate a legitimate page replace links to external pages with nil anchors. A high percentage of nil anchors are a likely sign of phishing. The (non-binary) feature F6 is calculated as follows:

$$F6 = \frac{LN}{LT} \text{ If } LT > 0; F6 = 0 \text{ If } LT = 0. \quad (2)$$

Where: *LN* is the number of nil anchors and *LT* is the total number of anchors.

- *Feature 7(F7): frequency of links*

Some phishing pages use images instead of html code to imitate a legitimate website's appearance, thus reducing the number of links pointing to other pages. Feature F7 models the frequency of links pointing to pages compared to links pointing to images or scripts and is calculated as follows:

$$F7 = \frac{LP}{LT} \text{ If } LT > 0; F7 = 0 \text{ If } LT = 0. \quad (3)$$

Where: *LP* is the number of links to pages and *LT* is the total number of links.

- *Feature 8(F8): action complies with the page identity*

A login page requests access information from users with a form that contains "input" fields, as represented in the following example:

```
<form method="post" action="action.jsp">
<input name="login" id="username" />
<input name="passwd" id="passwd" />
<button type="submit" > Connection </button>
</form>
```

The information entered in the input fields are processed by the function whose URL is specified in the action field. Usually, phishing web pages claim a legitimate page identity but the action field contains a different URL compared to this identity. The trinary F8 feature models this behaviour as presented in the Algorithm 1.

4. Evaluation

To evaluate our approach, we will test the validity of the whitelist followed by the performance of the SVM classifier.

4.1. The whitelist

To evaluate the performance of the whitelist we used 400 pages, of which 200 were legitimate and 200 were phishing pages. The legitimate pages are composed of the following:

- 10 login pages of top targeted websites [21].
- 50 login pages of the most visited websites according to Alexa2.
- 140 pages from Yahoo 3.

All of the 200 phishing pages were collected from PhishTank4.

As mentioned earlier, the whitelist contains only the information about login pages from the top 10 websites targeted by phishers[21].

Table 1: Evaluation Results of the Whitelist.

Threshold	WebPages (WP)	Phish detected	Legitim detected
≥ 0.7	Legitimate WP (200)	02	10
	Phishing WP (200)	36	0
≥ 0.8	Legitimate WP (200)	01	10
	Phishing WP (200)	34	0
≥ 0.9	Legitimate WP (200)	00	10
	Phishing WP (200)	30	0

To choose an adequate threshold we have tested our whitelist with three threshold values: 0.7, 0.8 and 0.9. Table 1 summarizes the test results. The results show that the higher the threshold similarity is, the lower whitelist wrong decisions are. With a threshold of 0.7 the whitelist detects more phishing pages (36 pages against, 34 and 30 for 0.8 and 0.9 respectively), however there is also an increase in incorrect decisions (2 legitimate pages were classified as phishing pages, versus 1 and 0 for 0.8 and 0.9 respectively). As we seek to avoid any incorrect classifications in the whitelist level, we have adopted a value of 0.9 for our threshold similarity. With this value, the whitelist detected about 5% of all legitimate pages and 15% of phishing pages. Despite this low percentage, we note the absence of incorrect classifications (i.e., the rate of false positives and false negatives is 0%). Moreover, we note that pages with low similarity are not processed at the whitelist level.

Lastly, whitelist effectiveness will increase with growing use as a user's frequented pages are automatically added to the whitelist. This will reduce the risk of a user falling victim to a phishing page imitating one of his or her frequented pages.

4.2. The SVM classifier

If a page has a low similarity with the pages of the whitelist, this page is transformed into a feature vector to be classified. We used an SVM classifier [22], a binary classifier well adapted to our case, as we have only two classes (phish or legitimate).

We used a database of 850 pages, of which 400 are the ones used to evaluate the whitelist, 200 are legitimate pages from Yahoo Random, and 250 are phishing pages collected from PhishTank. We trained our classifier against 400 pages (200 legitimate and 200 phishing) and tested the classifier again the remaining 450 pages (200 legitimate and 250 phishing). Before transforming the testing dataset into feature vectors, we applied the whitelist as a filter. The whitelist filtered 41 phishing pages and 0 legitimate pages. The remaining 409 pages are transformed into feature vectors and passed to the SVM classifier.

We evaluate our classification model according to the percentage of correctly classified phishing pages or true positives rate (TP also called recall), legitimate sites wrongly classified as phishing or false positives rate (FP), precision

(P) that represent the degree to which pages identified as phishing sites are indeed malicious, and F-measure (FM), the harmonic mean between the precision and recall. These various metrics are calculated as follows:

$$TP(R) = \frac{P_P}{P_P + P_L} \tag{4}$$

$$FP = \frac{L_P}{L_P + L_L} \tag{5}$$

$$P = \frac{P_P}{P_P + L_P} \tag{6}$$

$$FM = 2 \times \frac{P \times R}{P + R} \tag{7}$$

Where P_P , P_L , L_P , L_L respectively represent: the number of correctly classified phishing web pages, the number of incorrectly classified phishing pages, the number of legitimate pages wrongly classified as phishing sites, and the number of correctly classified legitimate sites.

The following table summarizes our evaluation results on the aforementioned testing dataset.

Table 2: Evaluation results on the testing dataset (SVM classifier performances).

	TP	FP	P	FM
Values	98%	3.5%	96.6%	97.3%

We compared our classification model to those of CANTINA [16] and [17](the approach of M. He et al.).

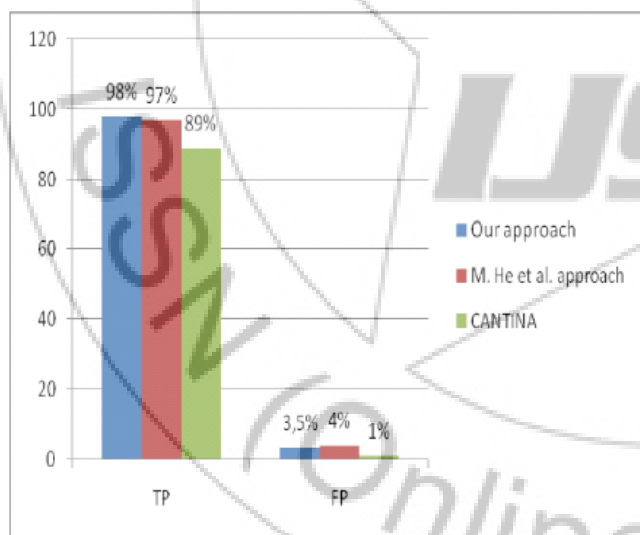


Figure 3: Comparison against previous works

Figure 3 shows a slight improvement of our model in terms of true positive rates: our model can detect 98% of phishing pages against 89% for [16] and 97% for [17]. If we count the 41 pages detected at the whitelist level (since the pages were correctly classified), results will be further improved as illustrated in Table 3.

Table 3: Evaluation results on the testing dataset (svm+whitelist).

	TP	FP	P	FM
Values	98.4%	3.5%	97.2%	97.7%

Our approach suffers from a high false positives rate (>3%); this high value is associated principally with the SVM classifier. As mentioned earlier, the application of the SVM classifier will decrease after the whitelist stabilization, as the whitelist is incrementally augmented by the user’s surfing history; at this point, most classification decisions are made at the whitelist level and the risk of a successful phishing attack against the user decreases considerably since any phishing pages that attempt to imitate a page frequented by the user will be detected by the whitelist. If a phishing page imitates a non whitelisted page, it is unlikely that the user would provide the site with sensitive information and, in the unlikely event that the user does provide such information; the page will likely be detected as a phishing attempt by our SVM classifier.

5. Conclusion

In this work we have described an anti-phishing solution that combines a personalized whitelist and an automated classification engine. Our combined approach benefits from the advantages of both techniques without suffering from the drawbacks of each method.

We maintain the accuracy of whitelist solutions and eliminate the difficulty of managing and updating large amounts of data by using a personalized whitelist. False positives traditionally present in these solutions are eliminated since, if a page does not belong to the whitelist, it is not classified as a phishing page but is instead treated by our SVM classifier. Moreover, the proposed whitelist is designed to be automatically updated without user intervention, significantly reducing configuration errors and improving usability.

Despite these advantages, our proposed approach does suffer from some shortcomings. For example, our approach is unable to detect whether legitimate websites are attached by a DNS spoof. We can solve this shortcoming by adding the IP addresses of each page to the whitelist, since the IP address of the majority of targeted websites is often stable [8]. Also the dependence of one feature of the classification model on a search engine can affect the ease of use and responsiveness of the tool in the case of the search engine dysfunction.

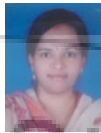
Lastly, our approach’s performance may be improved, classification features need to be tuned up to work better, while new relevant features may be discovered in the future to further differentiate the legitimate and phishing pages.

References

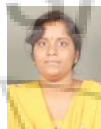
[1] Phishing Activity Trends Report, 1st Half 2011, <http://www.antiphishing.org/>

- [2] <http://www.social-engineer.org/>, query date: January 2012.
- [3] A. Emigh, "Phishing attacks: Information flow and chokepoints". In M. Jakobsson and S. Myers, editors, *Phishing and Countermeasures*, pages 31-64. Wiley, 2007.
- [4] A. Bergholz, J.H. Chang, G. Paass, F. Reichartz and S. Strobel, "Improved Phishing Detection using Model-Based Features". CEAS 2008.
- [5] Mozilla. Phishing protection, <http://www.mozilla.com/enUS/firefox/phishingprotection/>, query date: March 2011.
- [6] google safe browsing : <http://code.google.com/intl/fr/apis/safebrowsing/>, query date: March 2011.
- [7] <http://www.phishtank.com/> query date: March 2011.
- [8] Y. Cao, W. Han and Y. Le, "Anti-phishing Based on Automated Individual Whitelist", DIM'08, October 31, 2008, Fairfax, Virginia, USA.
- [9] V. P. Reddy, V. Radha and M. Jindal, "Client Side protection from Phishing attack", International journal of advanced engineering sciences and technologies Vol no. 3, Issue no. 1, 039 - 045, 2011.
- [10] Y. Wang, R. Agrawal and B. Choi, "Light Weight Anti-Phishing with User Whitelisting in a Web Browser", IEEE Region 5 Conference, April 2008.
- [11] V. I. Levenshtein, "Binary codes capable of correcting deletions, insertions, and reversals". Soviet Physics Doklady 10 (1966).
- [12] M. Khonji, A. Jonesy, Y. Iraqi, "A Brief Description of 47 Phishing Classification Features", http://khonji.org/upload/feature_desc. Accessed January 2012.
- [13] I. Fette, N. Sadeh, A. Tomasic, "Learning to Detect Phishing Emails". Proceeding of International World Wide Web Conference (WWW 2007), Banff, Alberta, Canada, May 2007.
- [14] O. Salem, A. Hossain, M. Kamala, "Awareness Program and AI based Tool to Reduce Risk of Phishing Attacks", CIT 2010, Bradford, UK, (2010).
- [15] S. Abu-Nimeh, D. Nappa, X. Wang, and S. Nair, "A Comparison of Machine Learning Techniques for Phishing Detection", APWG eCrime Researchers Summit, Pittsburgh, PA, USA. (2007).
- [16] Y. Zhang, J. Hong and L. Cranor, "CANTINA: A Content-Based Approach to Detecting Phishing Web Sites". Proceeding of International World Wide Web Conference (WWW 2007), Banff, Alberta, Canada, May 2007.
- [17] M. He, S.J. Horng, P. Fan, M. K. Khan, R.S. Run, J.L. Lai, R.J. Chen and A. Sutanto, "An efficient phishing webpage detector", Journal Expert Systems with Applications (12018-12027): Volume 38 Issue 10, September, 2011.
- [18] S. Garera, N. Provos, M. Chew, and A. D. Rubin, "A framework for detection and measurement of phishing attacks". In Proceedings of the WORM. (2007).
- [19] J. Ma, L.K. Saul, S. Savage, GM. Voelker, "Beyond Blacklists: Learning to Detect Malicious Web Sites from Suspicious URLs", KDD 09. Paris, France (2009).
- [20] G. Salton, "Mathematics and information retrieval". Journal of Documentation, 35(1), 1-29, (1979).
- [21] OpenDNS 2010 Report: Web Content Filtering and Phishing. <http://www.opendns.com/pdf/opendns-report-2010.pdf>.
- [22] C. C. Chang and C.-J. Lin, "LIBSVM: a library for support vector machines". ACM Transactions on Intelligent Systems and Technology, 1-27 (2011).

Author Profile



Komatla. Sasikala Obtained the M.sc degree in Computer Science from Govt. College for Women, Guntur. At present pursuing the M.Tech in Computer Science and Engineering (CSE) Department at Guntur Engineering College, Guntur.



P. Anitha Rani obtained the B.Tech Degree from Sri C.R. Reddy Engineering College and M. Tech (CSE) from JNTU, Hyderabad. She has 10 years of teaching experience and working in Computer Science and Engineering (CSE) Department at Guntur Engineering College, Guntur.