

CITY UNIVERSITY OF HONG KONG

香港城市大學

ArrogateText: a Content-based Approach to Shanzhai

Website Detection

ArrogateText:基於內容分析的山寨網站偵測方法

Submitted to

Department of Computer Science

電腦科學系

in Partial Fulfillment of the Requirements

for the Degree of Master of Philosophy

哲學碩士學位

by

LU Zhi

蘆 智

December 2010

二零一零年十二月

Abstract

Shanzhai website attack is a new kind of security threat that is harmful to the revenue and reputation of legitimate organizations. Despite this problem seems similar to the Phishing attack, which has been attracting a great many researchers for decades and many Anti-Phishing tools have been developed, Shanzhai website problem has its distinct features that make the current Anti-Phishing solutions infeasible to this problem.

Shanzhai websites usually use semantically similar content with the real websites, which is attacked by the Shanzhai websites, to deceive the visitors. In this thesis, a content-based Shanzhai website detection method, called ArrogateText, is proposed. The method first collects the candidate web pages from searching results of a given brand, in which the most strongly associated group of results are selected by extended TermCut. A matrix of these web pages' semantic relationships is constructed by analyzing the given web pages contents, where both the text semantic similarity and sentence structure similarity are considered. Finally, these web pages are ranked by multiplying the matrix and the ranking scores vector of the importance of the web pages. For a given website claiming a brand, if we find it is above certain position in the ranking list, we can determine it as the Shanzhai website. Otherwise, it should be treated as a legitimate website.

We design three experiments to evaluate the proposed method ArrogateText. As an early attempt to fight against the Shanzhai website problem, the proposed method achieves fairly good detection accuracy while remaining a relatively low false positive rate.

Table of Contents

Abstract	I
Acknowledgement	II
Table of Contents	III
List of Figures	V
List of Tables	VI
Chapter 1 Introduction	1
1.1 Background and motivation	1
1.2 Contributions	3
1.3 Thesis organization	4
Chapter 2 Related work	5
2.1 User Education for Anti-Phishing	6
2.2 User Interface design for Anti-Phishing	7
2.3 Techniques for Anti-Phishing	10
3 Detecting Shanzhai website	17
3.1 Shanzhai website definition	17
3.2 System overview	18
3.3 Adopting TermCut for Shanzhai websites detection	20
3.4 ArrogateText Component	22
3.4.1 Sentence structural similarity	23
3.4.2 Text semantic similarity	28
3.5 Fusion of features	30
3.6 Ranking	31
3.7 Heuristic for reducing the false positive rate	32

Chapter 4 Experiments.....	34
4.1 Experiment 1 - Evaluation of Extended TermCut	35
4.2 Experiment 2 - Fixing optimal position in ranking list.....	37
4.3 Experiment 3 - ArrogateText performance evaluation	40
Chapter 5 Conclusion and future work	45
5.1 Conclusion	45
5.2 Future work.....	46
References.....	48

List of Figures

Figure 1. Average F1-Scores at each position	39
Figure 2. Comparison of ArrogateText and ArrogateText + Heuristic	42
Figure 3. ROC curve of the three methods	43

List of Tables

Table 1. Examples of sentences extracted from Shanzhai websites	24
Table 2. Brands used for finding optimal position	38