技术研究

2012年第04期

doi 10.3969/j.issn.1671-1122.2012.04.009

网络钓鱼防御技术研究

黄华军,王耀钧,姜丽清

(中南林业科技大学计算机与信息工程学院,湖南长沙 410004)

摘 要:网络钓鱼正步入快速发展期,对电子商务健康发展已造成严重威胁,网络钓鱼防御技术也日益重要。文章深入剖析网络钓鱼基本概念、攻击步骤、类型及演化进程。对网络钓鱼防御研究领域已有研究成果进行了分类和总结。重点对钓鱼网站检测技术的基本原理、特点进行阐述,还详细分析各类钓鱼网站检测技术的典型应用。探讨现有网络钓鱼防御技术性能的问题,指出网络钓鱼防御技术未来发展方向。

关键词:网络钓鱼;网络钓鱼防御;网络钓鱼检测;身份窃取

中图分类号:TP393.08 文献标识码:A 文章编号:1671-1122(2012)04-0030-06

Countermeasures Technology of Phishing

HUANG Hua-jun, WANG Yao-jun, JIANG Li-qing

(College of Computer and Information Engineering, Central South University of Forestry & Technology, Changsha Hunan 410004, China)

Abstract: In this paper, a thorough overview of a phishing attack and its countermeasure techniques, which is called anti-phishing, is presented. Firstly, technologies used by phishers and the definition, classification and future works of deceptive phishing attacks are discussed. Following with the existing anti-phishing techniques in literatures and research-stage technologies are shown, and a thorough analysis that shortcomings of countermeasures is given.

Key words: phishing; anti-phishing; phishing detection; identity theft

0 引言

网络钓鱼(phishing)是基于社会工程学的一种攻击手段^[1]。它通过垃圾邮件、即时聊天工具、手机短信或网页虚假广告发送声称来自于银行或其他知名机构的欺骗性信息,意图引诱用户登录看起来极其真实的假冒网站,给出敏感信息(如用户名、口令、账号 ID、ATM PIN 码、信用卡)的一种攻击方式。目前,网络钓鱼防御主要集中在钓鱼网站检测、垃圾邮件过滤、钓鱼网站追踪、网络钓鱼行为分析、终止钓鱼网站域名解析等多个研究领域,但以钓鱼网站检测最活跃。钓鱼网站检测是在浏览器中安装检测插件,当用户浏览可疑钓鱼网页时,浏览器的插件将提醒用户当前网页为钓鱼网页。插件的检测算法采用 URL 检测技术、基于启发式检测技术、基于视觉相似检测技术等判别钓鱼网页。

1 网络钓鱼概念与分类

网络钓鱼攻击可分为以下3类攻击方法[2]:

- 1) 欺骗攻击:发送欺骗的电子邮件,欺骗用户登陆钓鱼网站的方法来骗取用户的机密信息;
- 2) 恶意程序攻击:利用恶意程序,如键盘记录程序(Keylogger)或截屏程序(Screenlogger),来盗窃用户的机密信息;
- 3) 基于域名攻击:通过修改主机名,误导用户登陆到钓鱼网站的方法来获取用户的机密信息。

随着网络钓鱼防御技术不断的发展,如垃圾邮件过滤、URL 黑名单过滤、钓鱼网站检测、停止域名解析等,钓鱼者不断改进网络钓鱼攻击方式。网络钓鱼已从最初仿冒经典网站向其它形式发展,如内容感知网络钓鱼、分布式网络钓鱼、网络钓鱼黑色产业链、钓鱼网站软件工具套件,移动环境网络钓鱼、快速改变域名网络钓鱼以对抗现有的防御技术;同时网络钓鱼传输途径

收稿时间 2012-02-05

基金项目 湖南省自然科学基金资助项目[10JJ4043,10JJ5062];湖南省教育厅资助项目[08B091];湖南省科技重大专项项目[2010J05]湖南省科技计划重点项目[2010NK2003];湖南省科技计划项目[2010TZ4012];长沙市科技局科技计划项目[K1005180-61];湖南省公安厅科学研究项目(湘公装[2008]14号)

作者简介 黄华军(1978-), 男, 湖南, 博士, 副教授, 主要研究方向:信息隐藏、网络钓鱼防御; 王耀钧(1976-), 男, 湖南, 硕士研究生, 主要研究方向: 网络与信息安全、网络钓鱼防御; 姜丽清(1986-), 女, 浙江, 硕士研究生, 主要研究方向: 网络与信息安全、网络钓鱼防御。

技术研究

已从单一发送大量垃圾邮件,向发送垃圾邮件、即时聊天信息、 手机短信或网页虚假广告等多种途径发展,使得网络钓鱼危 害更大,防御更难。

2 网络钓鱼防御技术

网络钓鱼防御是网络钓鱼的对抗技术^[3-20]。由于网络钓鱼具有钓鱼网站的来源和危害具有跨越国界、制作成本低、存在周期短、危害大等特点,使得网络钓鱼防御已成为一个令全球关注并感到棘手的问题。网络钓鱼防御不是一个国家、一个部门或者一种技术能解决的问题,而是需要多方参与和多种防御技术联合起来共同防御,才能取得成效。

为了有效地防御网络钓鱼,先要了解网络钓鱼泛滥的深层次原因,网络钓鱼防御的参与方以及分类。

网络钓鱼泛滥的原因主要有:

- 1)制作成本低廉,巨大经济利益,犯罪事实无法追踪:不论是盗用真实网站还是新开发钓鱼网站,以及网上发布钓鱼信息的实施成本低,但只要能诱使部分网络用户上当,所获得的经济回报将是巨大的。钓鱼者诱使受害者在钓鱼网站提交用户ID和密码后,通过盗用受害者的身份登录到真实网站,并从中获取经济利益。对此,不论是用户还是经济机构都无法追踪并确定钓鱼者的身份。
- 2) 静态、单向用户名/口令认证体制:网络用户需要向网站服务器提交身份信息,认证用户的身份,却无法认证网络服务器的真实身份。这导致网络钓鱼攻击存在的根本原因。
- 3) 网站身份易被复制、窃取,网站服务器身份真实性难以保证:钓鱼者能从目标网站下载源文件,进行修改。受害者通常是处于弱势地位,无法区分钓鱼网站与真实网站的身份,使得网络钓鱼技术简单,但十分有效。

不同的网络钓鱼防御参与方对于网络钓鱼防御的动机、 付出精力是不一样的。网络钓鱼防御的参与方包括如下:

直接受害者,用户、组织、金融机构及网络商家是网络钓鱼的直接受害者;基础设施提供者,Internet 服务提供商、电子邮件服务商、浏览器、域名注册主管机构、域名注册机构等是基础设施提供者;受益的保护者,垃圾邮件过滤软件、反病毒安全软件等网络钓鱼防御服务提供商是网络钓鱼防御受益的保护者;公众保护者,国家管理机构、法律部门、国家信息安全主管部门、学术团队等是网络钓鱼防御的公众保护者。

从网络钓鱼防御的参与方来分,网络钓鱼防御可分为服务器端防御、用户端防御和第三方防御。服务器端防御是指网站服务器端借助其他技术,如数字水印、数字指纹、动态安全皮肤,双重认证协议等,向用户证明网站身份的真实性。用户端防御是指在用户浏览器安装插件,检测到钓鱼网页后提示用户、或保护用户敏感信息输入等。第三方防御包括钓鱼

垃圾邮件过滤、第三方认证机制、浏览器提供商的 URL 黑名单过滤机制、安全软件厂商防御机制、公众保护机构、用户防范网络钓鱼培训等。

网络钓鱼防御技术研究取得了一定的进展,其中又以用户端的网络钓鱼检测研究最活跃,研究成果丰富,下面将重点介绍用户端的网络钓鱼检测技术的研究进展。

3 网络钓鱼检测技术

网络钓鱼检测是在客户端浏览器安装插件,检测用户当前浏览的网页是否是钓鱼网页,从而提示用户被欺骗的威胁。插件中检测算法有以下3类方法:基于 URL 检测、基于启发式检测和基于视觉相似检测。

3.1 基干URL检测技术

3.1.1 URL黑名单检测技术

基于 URL 检测是指利用 URL 地址,判断当前的网站是否为钓鱼网站。最初的方法是利用黑名单中存储被确认的钓鱼网站 URL 地址,当浏览器浏览时,提醒用户当前网站为钓鱼网站。Microsoft IE、Google Safe Browser、Netcraft Tool Bar、eBay Tool Bar、McAfee SiteAdvisor 等知名 IT 企业采用黑名单防御钓鱼网站。为验证黑名单性能,Ludl 收集了三周内 10000条钓鱼网站的 URL 地址,测试微软 IE 浏览器和谷歌安全浏览器的检测性能。他们发现 Google 能够识别 90% 的钓鱼网站URL 地址。Sheng 等人使用 191 新的钓鱼网站 URL,测试 8种网络钓鱼防御工具中黑名单更新速度。他们发现小于 20%的网络钓鱼防御工具能在短时间内识别钓鱼网站。尽管 URL 黑名单检测技术简单、检测率精确,但存在无法检测不在黑名单内的钓鱼网站,且确认黑名单需要人工验证,费时耗力的问题。

3.1.2 基于机器学习检测技术

基于机器学习的 URL 检测技术是直接利用 URL 检测钓鱼网站,主要流程如下:选择钓鱼网站 URL 特征向量、生成训练数据、训练构建分类器模型、应用分类器分类 URL。此类检测特征选取和分类器构建是关键。

1) Garera 算法 [11]。Garera 等人分析钓鱼网站 URL 结构,详细介绍特征集合选取过程,利用回归滤波器(Logistic Regression Filter)分类 URL。对钓鱼网站 URL 结构进行分析,得出 4 种类型的 URL 结构 表 1 列出了 4 种钓鱼网站 URL 结构。特征集合由页面特征、域名特征、类型特征和单词特征共 18 个特征构成。页面特征借助谷歌搜索引擎,选取 URL 页面排名、域名页面排名、页面排名存在爬行数据库、页面存在索引数据库、两个页面质量评价共 6 个特征;域名特征选择域名在白名单表 1 个特征;类型特征选择类型 I、II、和 III,3 个特征;单词特征选择 secure、account、webscr、login、ebayiapi、

2012年第04期

signin、banking 和 confirm, 8 个单词。

表1 钓鱼网站URL结构

Type	Examples
1	http://210.80.154.30/~test3/.signin.ebay.com/ebayisapidllsignin.html
	http://0xd3.0xe9.0x27.0x91:8080/.www.paypal.com/uk/login.html
Ш	http://21photo.cn/https://cgi3.ca.ebay.com/eBayISAPI.dIISignIn.php
	http://2-mad.com/hsbc.co.uk/index.html
Ш	http://www.volksbank.de.custsupportref1007.dllconf.info/r1/vm/
	http://sparkasse.de.redirector.webservices.aktuell.lasord.info
IV	http://www.wamuweb.com/IdentityManagement/
	http://mujweb.cz/Cestovani/iom3/SignIn.html?r=7785

为得到较好的分类效果,选择回归滤波器训练 18 个特征的系数,得到每个特征在分类钓鱼网站 URL 的相对权重的让步比(Odds Ratio),记为 $e^{coefficient}$ 。让步比定义为 URL 为钓鱼网站 URL 的概率与 URL 为合法 URL 的概率比。让步比大于1 的较适合判断钓鱼 URL。表 2 给出了18 个特征的让步比。

最后, Garera 等人测试了算法的性能:分类的精度为97.31%, 准确率为95.8%, 虚警率为1.2%。

表2 特征让步比

Feature	$e^{coefficient}$	Feature	$e^{^{coefficient}}$	Feature	$e^{coefficient}$
Quality Score II	0.141	Is Type I	597.815	Webscr	2.710
Host page rank	0.152	Is Type II	19.038	Login	6.416
URL page rank	0.283	Is Type III	1.259	Ebayisapi	8.722
Page rank presence	0.585	Is white Domain table	0.022	Signin	12.685
Quality score I	1.045	Secure	1.395	Banking	13.959
Page in Index	2.396	account	2.361	confirm	15.777

2) Ma 算法。与 Garera 等人采用 18 个特征作为分类网络 钓鱼 URL 不同的是 Ma 等人分析可疑 URL 的词汇和主机属性,采用词袋模型表示特征,获得了成千上万的特征 [12]。在词汇特征中,一方面考虑主机名长度、URL 长度、URL 中点号数等;另一方面,对于 URL 中主机和路径中每一个词汇符号,采用词袋模型建立一个二值特征。在主机特征中,考虑了 IP 地址属性、WHOIS 属性、域名属性和地理位置属性。

考虑到批量学习和在线学习性能要求,对于批量学习,Ma等人分析了朴素贝叶斯、支持向量机和回归滤波的分类性能;在线学习中,研究了感知器、随机梯度下降回归滤波、被动贪婪算法和秘密权证算法的分类性能。在批量分类中,取得了95-99%的精度,在线分类中,分类的精度99%。

与 Ma 相类似的算法还包括 Blum 等人提出的基于词汇特征的钓鱼 URL 的分类算法。McGrath 等人经过详细实验,分析了钓鱼者做法,剖析钓鱼 URL 结构和域名、钓鱼

网站域名注册信息、域名注册到使用的时间、钓鱼网站主机以及钓鱼网站生存周期。这些的实验分析对检测钓鱼 URL 起到前提依据。

3.2 基于启发式检测技术

基于启发式检测技术根据钓鱼网站存在的异常特征超出设定阈值和不合乎常规进行判断。

3.2.1 SpoofGuard软件[13]

SpoofGuard 是 Chou 等人提出浏览端网络钓鱼检测软件。

SpoofGuard 软件监视用户网络浏览行为,计算当前网页仿冒指数,当指数超出设定阈值时,警示用户当前威胁。SpoofGuard软件利用钓鱼网站7点共同属性检测钓鱼网站:

1) 网站标志,钓鱼网站的 logo 与真实网站的 logo 相同 2) 可疑的 URLs 地址,如 URL 地址过长、使用 IP 地址代替域名、采用不常用字符 @、可疑域名等;3) 用户输入,钓鱼网站要求用户输入帐号、密码、信用卡号、社会安全号等机密信息;4) 网站生存周期短,钓鱼网站生存周期通常为几天或者几小时;5) 复制,钓鱼网页复制真实网页,或修改很小部分;6) 异常,钓鱼网站存在愚蠢拼写、语法错误和自相矛盾;7) HTTPS,钓鱼网站不使用 HTTPS 来保护传送的信息。

SpoofGuard 测试下载的网页,评分机制评定检测结果。 评分机制用标准聚集函数表示:

$$TSS(page) = \sum_{i=1}^{n} w_i P_i + \sum_{i,j=1}^{n} w_{i,j} P_j P_j + \sum_{i,j,k=1}^{n} w_{i,j,k} P_i P_j P_k + ... (1)$$

其中,TSS(total spoof score) 是指仿冒总得分。 P_i 是每次测试 T_i 的取值,范围是 [0,1]。 P_i =1 表示网页是钓鱼网页, P_i =0 则相反。 w_i 是预先设定的权重,为了减少检测的虚警率,在很多情况下 w_i 都设为 0。

3.2.2 融合多特征的检测技术

1) CANTINA 系统。CANTINA 是 Zhang 等人基于网页内容,设计并实现的钓鱼网站检测系统。CANTINA 利用著名信息检索中 TF-IDF (Term Frequency/Inverse Document Frequency) 算法,分析网页中每一个术语的 TF-IDF 值,借助鲁棒超级链接中词汇标签概念,在 Google 搜索引擎检索词汇标签,最终确定网页合法或者仿冒。CANTINA 工作步骤如下:

对于给定的网页, 计算其中每一个术语的 TF-IDF 值;选取前5个最大 TF-IDF值的术语, 形成词汇标签;在 Google搜索引擎检索这个词汇标签;如果该网页域名与前N个搜索结果的域名匹配,则该网页是合法网。否则,认为是钓鱼网页。

为降低 CANTINA 系统的虚警率, Zhang 等人采用了 SpoofGuard 软件和 FILFER 算法的特征。这些特征包括,域名存在周期、著名图片、可疑 URL、可疑链接、IP 地址、URL 点号数、表单等。当这些特征包含表 3 中内容时,用来确定钓鱼网站。

只采用基本 TF-IDF 算法检测钓鱼网站的误报率达 30%, 而这些特征时, CANTINA 的检测正确率达到 95%, 虚警率降到 10%。

2) Xiang 系列算法。卡耐基-梅隆大学的 Xiang 等人通过长期研究钓鱼网站,尤其是对网络钓鱼黑色产业链和网络钓鱼工具套件制作的钓鱼网站,提出了一系列的钓鱼网站检测算法。这些算法,一方面利用人工识别的黑名单降低虚警率,另一方借助启发式检测识别未知的钓鱼网页。在介绍这些算法之前,先给出判别钓鱼网站的标准:

Heuristic	Suspected phishing			
Age of Domain	<=12 months			
Known Images	Page contains any known logos and not on a domain owned by			
Kilowii Iiilages	logo owner			
Suspicious URL	URL contains @ or -			
Suspicious Links	Link on page contains @ or -			
IP Address	URL contains IP address			
Dots in URL	>=5 dots in URL			
Forms	Page contains a text entry field			

- (1) 钓鱼网站通过复制整个或者部分著名目标网站文件 和结构,以实现从外观上与目标网站视觉相似。
 - (2)钓鱼网站内含的域名与其仿冒目标网站的域名不一致。
 - (3) 钓鱼网站有需要用户提交机密信息的登录表单。

钓鱼网站为更好迷惑用户,通常采用复制目标网站文件,以实现外观上视觉相似。而且,很少钓鱼网站仅使用一次,许多钓鱼者会使用同样的文件,仿冒某一特定网站。钓鱼网站被打包成钓鱼工具包,以方便下一次进行新的攻击。对于这一类钓鱼网站,最初提出应对策略的是深层次 MD5 匹配策略。随后,Xiang 等人提出基于哈希近乎复制的钓鱼网页检测算法。此算法中采用 SHA1 哈希算法,计算钓鱼网站源文件哈希校验和。哈希校验和用来匹配其他钓鱼网站源文件的哈希校验和,检测钓鱼网站工具套件部署的钓鱼网站。

使用钓鱼网站工具套件部署的钓鱼网站还有一个特点是: URL 地址中文件路径一样。如下面 URL 地址中黑体部分。

 $http://ww5.chase.com.bank84.com/ccp/ \mbox{{\bf clientconfrm.jsp/}} \\ ?siteid = 17xwhhjde Dwcycu Okhb \\$

http://ww2.chase.com.id746.com/ccp/clientconfrm.jsp/?site=18bzkwdydDjImrdnOrdn

http://ww8.chase.com.cert83.com/ccp/clientconfrm.jsp/?taskid=14zrohDprhbhOkhb

http://ww6.chase.com.sid36.com/ccp/clientconfrm.jsp/ ?ref=22hknZZeoDdrdsatWrdnOhsa

即使钓鱼网站复制目标网站整个网站文件,但为实现钓鱼目的,钓鱼网站中存在与目标网站不一致地方,比如域名。因此,根据这一特点,分析网页内容,确定出其中异常,从而发现钓鱼网站。

另一种用于降低虚警率的方法是登录表单过滤(Login Form Filter)。钓鱼网页为获取用户机密信息,需要用户在钓鱼网站的登录表单上提交个人信息。因此,能根据网页中是否存在登录表单过滤钓鱼网页。如果,网页中没有登录表单,则网页为正常网页,否则,作为可疑的钓鱼网页候选对象进行深层次检测。3.2.3 基于网站身份的检测技术

基于网站身份的钓鱼网站检测是根据钓鱼网站所宣称的 目标网站身份与真实身份之间的差异判断。网络钓鱼是一种 身份窃取攻击,仿冒目标网站的身份,欺骗用户泄漏个人机密 信息。虽然钓鱼网站仿冒合法网站的身份,但为了实现非法目的,不可能与目标网站一模一样。这种差异体现在结构特征、HTTP 处理和搜索引擎结果。

1) 机器学习分类技术。最初提出基机器学习检测的是Pan 等人提出的 Anomaly 算法。Anomaly 算法的检测器由身份抽取器和页面分类器组成。对于输入的网页,身份抽取器构建与网站身份相关的 DOM 对象集合,记为 D。这些网页 DOM 对象包括:Title、Description、Copyright、ALT/title、Address和 Body。处理 D,得到一个候选网站身份的单词集合 W,表示网站身份的单词将从这个集合中获取。根据观察可知,表示网站身份的单词分布偏离普通单词分布,采用 x^2 分布表示单词分布,具有最大 x^2 分布值的单词作为网站身份。网站身份表示为 $I=\{s_0,s_1,...,s_k\}$, s_i 是第 i 个网站身份单词。

页面分类器采用支持向量机,输入具有 10 个网站结构特征的向量 输出标志" 1 "表示是钓鱼网页", -1 "表示合法网页。 Anomaly 算法的特征向量 $V_P = \langle F_1, F_2, F_3, F_3, F_3, F_4, F_5, F_5, F_6, F_7 \rangle$,其中特征 F_1 :URL 地址,特征 F_2 :DNS 记录,特征 F_{31} :空锚点(Nil Anchor), F_{32} :ID 锚点, F_{33} :域名锚点, F_4 :服务器表单处理, F_{51} :ID 请求 URL, F_{52} :域名请求 URL, F_6 :cookie 中域名, F_7 :SSL 认证。Anomaly 算法的取得了较低的漏报率和虚警率。

与 Anomaly 算法相类似的是 He 等人提出的算法 $^{[40]}$ 。对可疑网页进行预处理后,得到 12 个特征向量,利用支持向量机训练正常和钓鱼网页,最后进行分类决策。网站身份由术语身份集合和 URL 身份,12 个特征组成。特征向量 $V_P=<F_1$, F_2 , F_3 , F_4 , F_5 , F_6 , F_7 , F_8 , F_9 , F_{10} , F_{11} , $F_{12}>$,包括 URL 地址、锚点 URLs、请求 URLs、cookies、SSL 认证、URL 点号数、搜索引擎结果等。其中,9 个特征采用了 Anomaly 方法,2 个特征是 PILFER 方法,1 个 CANTINA 方法。

2) 启发检测技术。Xiang 等人提出一种身份发现与关键词检索的混合钓鱼网站检测系统(简称 Hybrid 系统),系统结构如图 1 所示。

Hybrid 系统流程为:对于可疑网页,先通过URL白名单过滤,不在URL白名单的网页,解析网页DOM,检测登录表单,如FORM标记、INPUT标记等,通过基于身份的钓鱼网页检测部件和关键词检索部件检测,标记网页是否属于钓鱼网页。URL白名单过滤机制和登录表单检测机制是为了减少系统的虚警率。基于身份检测部件直接利用信息抽取技术发现钓鱼网站身份与模仿的身份不一致确定钓鱼网站。关键词检索部件利用信息检索算法,借助搜索引擎结果确定钓鱼网站。下面详细介绍基于身份检测部件。

网站身份是指网站商标名,一般出现在网页 Title 和 Copyright 位置。基于网站身份检测部件最初定位网页 DOM 文

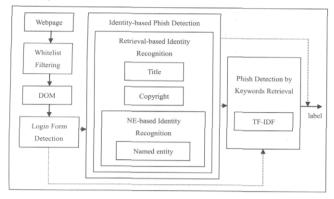


图1 系统结构图

本节点,查找标识网站商标名节点和属性。然后,借助搜索引擎查找此网站身份相应的网站域名,与网页中域名比较。根据预先设定的策略,如果域名匹配,则网页标识为"正常网页",否则网页标识为"钓鱼网页"。网站商标名出现在 Title和 Copyright 元素时,采用基于检索的身份识别算法;而对于网站商标名不在这些元素时,采用基于命名实体算法识别网站商标名。

Hybrid 系统对 11449 份网页进行检测,检测的准确率为 90.06%, 虚警率为 1.95%。

3.3 基于视觉相似的检测技术

基于视觉相似的网络钓鱼检测是利用钓鱼网页与真实网页的视觉相似超过设定的阈值来检测网页是否为钓鱼网页。相比 URL 和启发式检测技术易被钓鱼者混淆,基于视觉相似的网络钓鱼检测技术以钓鱼网站的本质特征—外观上与被仿冒的目标网站视觉相似,作为判别标准,则更具鲁棒性。这类检测技术由3个步骤组成:

1)将待检测的可疑网页转换成图片文件;2)采用图像处理技术对转换的图片文件处理,得到表述图片的特征矢量;3)根据特征矢量,与目标网站的特征矢量匹配,以确定是否是钓鱼网站。

图像相似的研究已经取得较大成绩,但基于视觉相似的 钓鱼网站检测是一个难题,主要原因有:1)显示器分辨率会 影响网页在浏览器的解析效果和布局;2)字体和字体大小很 大程度上会影响网页外观;3)网页上经常更换的广告也会影 响网页外观;4)异步网络使得网页装载时版式发生变换;5)流行的网页模板设计增加检测的虚警率。

最初提出视觉相似判断钓鱼网站的是邓小铁等人,提出将网页分成关键区域、页面版式和整体样式衡量钓鱼网页与真实网页的视觉相似性;接着又提出基于 Earth Mover's Distance(EMD) 比较钓鱼网页与真实网页相似的方法;曹玖新等人先对 Web 图像进行分割,抽取子图特征并构建网页的ARG(Attributed Relational Graph),在计算不同 ARG 属性距离的基础上,采用嵌套 EMD(nested Earth Mover's Distance)

方法计算网页的相似度,实现对钓鱼网站的检测^[16,17]。Engin Kirda等人利用文本样式、嵌入的图片、整体视觉的相似性判断钓鱼网页;张卫丰等人首先提取渲染后网页的文本特征签名、图像特征签名以及网页整体特征签名,通过匈牙利算法计算二分图的最佳匹配来寻找不同网页签名之间匹配的特征对,客观地度量网页之间的相似性的方法检测钓鱼网页^[18]。The-chuang Chen等人借助格式塔理论,将网页当作不可分割的实体,并用超信号描述该实体,利用算法复杂理论计算网页间超信号的方法检测钓鱼网页。Kuanta Chen等人基于网页中区分关键点特征,提出一种新的钓鱼网站检测算法,通过不变的内容描述一具体内容直方图,计算可疑网页与目标网页的相似性^[19]。Matthew Dunlop等人首先获取网页图像,利用光学字符识别(OCR)将图像转换为文本,借助 Google 网页排名算法衡量当前网页是否为钓鱼网页。还有就是根据字符外观相似而 Unicode 码不同的域名检测算法 [^{20]}。

Masanori Hara 等人通过对 2,262 个钓鱼网页进行分析,确定出 224 个不同的钓鱼网页布局模拟相似点,并以此检测钓鱼网页。在没有被仿冒网页比较的情况下,该方法的检测钓鱼网页的检测率为 80%,而虚警率仅为 17.5%。

下面以 EMD 算法为例,介绍基于视觉相似的钓鱼网站检测技术思路。

3.3.1 EMD算法

最初提出用 EMD 衡量网页视觉相似是 Fu 等人,随后,Cao 等人改进 EMD 算法,提出 N-EMD 算法。EMD 将网页图像的特征相似性比较问题转化为"运输问题"。EMD 算法处理过程主要包括 3 个步骤:1)获取给定 URL 网页的图片;2)对图片进行正则化处理;3)将正则化的图片表示成视觉签名信息(visual signature),视觉签名信息由色块和坐标特征组成,用于计算网页对(钓鱼网页与目标网页)的视觉相似。

图片视觉签名信息考虑颜色的分布,将降质的色块按照颜色进行分类,求出降质颜色 dc 的质心 C_{dc} :

$$C_{dc} = \sum_{i=1}^{N_{dc}} \frac{C_{dc,i}}{N_{dc}}...$$
 (2)

其中 C_{dc} , i 表示第 i 个像素在颜色 dc 的坐标 N_{dc} 表示具有颜色 dc 的像素和。

用 F_{dc} =<dc, C_{dc} > 作为降质颜色 dc 的特征 , N_{dc} 作为相应的权重 , 一个完整的图片视觉签名信息特征 S 表示为 :

有 n 个特征,特征对的 EMD 计算如下:

$$EMD(S_{s,a}, S_{s,b}, D) = \frac{\sum_{i=1}^{m} \sum_{j=1}^{n} f_{ij} \cdot d_{ij}}{\sum_{i=1}^{m} \sum_{j=i}^{n} f_{ij}}$$
(4)

技术研究

其中距离矩阵 $D=[d_{ij}](1-i-n,1-j-n)$,由色块距离 ND_{color} 和质心距离 ND_{color} 和质心理 ND_{color} 和质心理 ND_{color} 和质的 ND_{color} 和 ND_{col

色块距离 NDcolor 定义如下:

$$ND_{color}\left(d_{ci}, d_{cj}\right) = \frac{\sqrt{\left(dc_i - dc_j\right) \times \left(dc_i - dc_j\right)^T}}{MD_{color}} \dots (5)$$

MDcolor 表示最大的颜色距离。

质心距离 $D_{centroid}$ 定义如下:

$$ND_{centroid}\left(C_{dc_i}, C_{dc_j}\right) = \frac{\sqrt{\left(C_{dc_i} - C_{dc_j}\right) \times \left(C_{dc_i} - C_{dc_j}\right)^T}}{MD_{centroid}}.....(6)$$

MD_{centroid} 表示质心最大距离。

距离矩阵 $D=[d_n]$ 中,元素 d_n 计算如下:

当 $EMD(S_{s,a},S_{s,b},D)$ =0,图片是一致的 $EMD(S_{s,a},S_{s,b},D)$ =1,图片完全不一致。定义衡量图片视觉相似的基于 EMD 的指标 :

$$VS(S_{s,a}, S_{s,b}) = 1 - [EMD(S_{s,a}, S_{s,b}, D)]^{\alpha}$$
 (8)

设矢量 $T=<T_1,T_2,...,T_{Nprotected}>$ 表示可疑网页被分类成某一被保护网页的钓鱼网页的极值。 当检测可疑网页时,计算视觉相似矢量 $VS=<vs_1,vs_2,...,vs_{Nprotected}>$, $vs_i(1 i N_{protected})$ 是可疑网页与第 i 个保护网页的视觉相似性值。

钓鱼网页的分类标准定义如下:

$$isPhishing(VS) = \begin{cases} 1 & if \max(VS - T) \ge 0 \\ 0 & if \max(VS - T) < 0 \end{cases}$$
 (9)

当 isPhishing(VS)=1,可疑网页是钓鱼网页; isPhishing(VS)=0,可疑网页时正常网页。

3.3.2 存在的问题

基于视觉相似的网络钓鱼检测技术对钓鱼者的攻击更具 鲁棒,但也存在以下3个问题:1)钓鱼网页需要与真实网站的网页进行比较;2)需要额外存储设备存储被保护的目标网站网页的特征;3)表示图像特征较多,计算量大,检测算法的性能难以适应在线、实时响应用户请求。因此,需要新的基于视觉相似的钓鱼网页检测技术。

4 性能分析与未来发展方向

4.1 性能分析

网络钓鱼防御已取得不小的进展,但研究成果集中在用户端防御。为了验证提出的防御技术是否有效,很多研究人员对这些研究成果进行性能分析,指出现有的网络钓鱼防御技术也取得不小成绩,但也存在以下问题:

1) 用户无法分清钓鱼网站和合法网站,也经常忽略浏览器的安全警示;2) 让安全意识与技术处于弱势地位的用户最终判断当前网页是否为钓鱼网页;3) 检测的误报率和漏报率

较高,影响用户的使用意愿;4)需要额外安装插件,造成用户使用不便;5)很多网络钓鱼行为及其特征的分辨需要人工处理,此类处理耗时、费力。

4.2 未来发展方向

网络钓鱼防御已成为一个令全球关注并感到棘手的问题。 网络钓鱼防御不是一个国家、一个部门或者一种技术能解决 的问题,而是需要多方参与和多种防御技术联合起来共同防 御,才能取得成效。

随着对网络钓鱼认识,研究不断深入,通过分析网络钓鱼泛滥的原因,网络钓鱼防御的未来发展方向应从以下几个方面开展:

1) 改变目前网络中身份认证单向、静态的用户/密码的 认证体制,应该考虑动态口令、双重认证协议、手持设备认 证、生物信息认证等方式;2) 网站数字身份易于被窃取、复 制,应该从网站身份真实性加强对网站身份的认证、识别,主 动防御钓鱼网站的威胁;3) 国家立法部门、安全部门、金融 管理机构、互联网管理机构、互联网服务提供商等应该联合 起来打击网络钓鱼犯罪行为,取证钓鱼者犯罪行为,加大对 网络钓鱼犯罪行为的威慑力度;4) 加大网络钓鱼的危害宣传, 加强对用户识别网络钓鱼行为、防范网络钓鱼的危害进行培训,从而降低网络钓鱼的直接危害。

5 结束语

网络钓鱼除了带来经济上的损失,同时也使网民对电子商务产生不信任心理,减少甚至避免使用某些网络应用,从而阻碍我国互联网的深入发展。网络钓鱼防御不仅仅是技术问题,它还涉及提高网民的安全防范意识和技能;商家、各种支付系统和 Internet 服务提供商要对用户负责任,有很好的配套措施;法律、政府安全部门尽快制订相关管理细则,加大打击与惩罚力度,共同营造一个安全、可信的互联网环境。

本文分析了网络钓鱼定义、分类、攻击步骤及其发展趋势,探讨网络钓鱼攻击泛滥的原因,指出已有的防御技术。 重点介绍网络钓鱼检测技术的分类,介绍其中一些典型方法。 最后,分析已有防御技术的性能,指出下一步研究工作的方向。 (责编 程斌)

参考文献:

- [1] Anti-Phishing Working Group [EB/OL]. http://www.antiphishing.org, 2008-01/2012-02-05.
- [2] Aaron Emigh. Online Identity Theft: Phishing Technology, Chokepoints and Countermeasures [EB/OL].http://www.antiphishing.org/phishing-dhs-report.pdf, 2007-07/2012-02-05.
- [3] Markus Jakobsson. Modeling and Preventing Phishing Attacks [J]. Financial Cryptography and Data Security, 2005, (3570): 89-98.

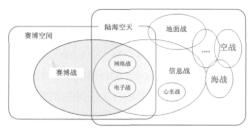


图5 赛博战与其他战术之间的关系

战超出其他任何常规武器作战的反应速度,没有平时和战时之分,随时发起随时结束。

- 4)作战目标广泛:包括军事、政治、经济、社会等领域,除了对敌方的用频设备、军用信息网络进行打击之外,也可以对一些事关国计民生的信息基础设施实施破坏,比如干扰正常的移动通信终端和电力供应。
- 5) 作战效果多样:可以是致命打击,比如传感器破坏、数据破坏、指挥中断等,也可以是非致命的攻击,比如雷达威力范围下降、通信误码率增大等。
- 6) 危害评估困难:赛博攻击具有很强的隐蔽性,导致被攻击方很难及时发现、定位、评估其危害,比如远程植入的木马和后门能够长期潜伏在敌方系统收集敏感信息。

4 结束语

本文以赛博空间一词的起源为出发点,在分析国内外相关研究的基础上,给出赛博空间的定义;然后讨论赛博空间与

传统的海、陆、空、天、赛的交互关系,给出赛博空间的组成结构;最后给出赛博战的概念及其与传统的网络战、电子战的关系,分析赛博战独有的特点。随着信息技术的不断发展及其在军事上的广泛应用,赛博空间作为一个全新的领域,已经成为各国发展和建设的热点,争夺赛博空间已提上日程,认识、研究、理解和构建赛博空间对于保护国家安全,获取未来作战优势具有重要意义。 (责编 张岩)

参考文献:

- [1] 维基百科[EB/OL]. http://en.wikipedia.org/wiki/Cyberspace, September 2011.
- [2] 维基百科 [EB/OL].http://en.wikipedia.org/wiki/Cybernetics, September 2011. [3] 维基百科 [EB/OL]. http://en.wikipedia.org/wiki/William_Gibson, September 2011.
- [4] Department of Defense. Department of Defense Dictionary of Military and Associated Terms[M]. Joint Publication 1 20, April 2001.
- [5] White House, United States. National Strategy of Secure Cyberspace[M]. February 2003.
- [6] Department of Defense. National Military Strategy for Cyberspace Operations[M]. December 2006.
- [7] White House, United States. National Security Presidential Directive 54/Homeland Security Presidential Directive 23 Cyberspace Policy[M]. January 2008.
- [8] Franklin D. Karmer, Stuart H. Starr, Larry Wentz. Cyberpower and National Security[M]. National Defense University Press, April 2009.
- [9] 石荣,李剑,黄鹏滔,李昊,贺岷珏.对信息战中赛博空间与赛博战的解析 [J]. 航天电子对抗,2010,26(04):44-46.

▶上接第35页 -

- [4] Chuang Yue, Haining Wang. BogusBiter:A Transparent Protection against Phishing Attacks [J]. ACM Transaction on Internet Technology, 2010, 10(02):1-31.
- [5] Akhilendra Pratap Singh, Vimal Kumar, Sandeep Singh Sengar, et al. Detection and Prevention of Phishing Attack Using Dynamic Watermarking [J]. Communication in Computer and Information Science, 2011, 147(01): 132 137.
- [6] Chad M.S. Steel, Chang-Tien Lu. Impersonator Identification through Dynamic Fingerprinting [J]. Digital Investigations, 2008, 5(02):60 70.
- [7] 郭敏哲, 袁津生, 王雅超. 网络钓鱼 Web 页面检测算法 [J]. 计算机工程, 2008, 34(20):161-164.
- [8] 陈涓,郭传雄.网络钓鱼攻击的在线检测及防治 [J].解放军理工大学学报(自然科学版),2007,8(02):134-138.
- [9] Steve Sheng, Bryant Magnien, Ponnurangam Kumaraguru, et al. Anti-Phishing Phil: The Design and Evaluation of a Game That Teaches People Not to Fall for Phish [C]In: Mads S, ed. Proc. of Symposium on Usable Privacy and Security. PA: ACM Press, 2007, 1-12.
- [10] Ponnurangam Kumaraguru, Steve Sheng, Alessandro Acquisti, et al. Teaching Johnny Not to Fall for Phish [J]. ACM Transactions on Internet Technology, 2009, 1(05):1-30.
- [11] Sujata Garera, Niels Provos, Monica Chew, et al. A Framework for Detection and Measurement of Phishing Attacks [A]. In: Christopher Kruegel, ed [C]. Proc. of the WORM '07. USA: ACM Press, 2007, 1-8.

- [12] Justin Ma, Lawrence K. Saul, Stefan Savage, et al. Beyond Blacklists: Learning to Detect Malicious Web Sites from Suspicious URLs [C]In: John Elder, ed. Proc. of the KDD ' 09. France: ACM Press, 2009, 1245-1254.
- [13] Neil Chou, Robert Ledesma, Yuka Teraguchi, et al. Client-side Defense against Web based Identity Theft[EB/OL].http://crypto.stanford.edu/SpoofGuard/webspoof.pdf, 2007 07/2012 02 05.
- [14] Yue Zhang, Jason Hong, Lorrie Cranor. CANTINA: A Content-Based Approach to Detecting Phishing Web Sites [C]In: Carey Williamson, Mary Ellen, eds. Proc. of the WWW 2007. Canada: ACM Press, 2007, 639-648.
- [15] He M, Horng S, Fan P, et al. An Efficient Phishing Webpage Detector [J]. Expert Systems with Application. 2011, 38(10): 12018 12027.
- [16] 曹玖新,毛波,罗军舟,等.基于嵌套 EMD 的钓鱼网页检测算法 [J]. 计算机学报,2009,32(05):922-929.
- [17] 李暄,刘莹.以视觉相似为基础的 phishing 检测方法 [J]. 清华大学学报(自然科学版), 2009, 49(01):146-148.
- [18] 张卫丰,周毓明,许蕾,等.基于匈牙利匹配算法的钓鱼网页检测方法 [J]. 计算机学报,2010,33(10):1963-1975.
- [19] Kuanta Chen, Chunrong Huang, Chusong Chen, et al. Fighting Phishing with Discriminative Keypoint Features [J]. IEEE Computer Society, 2009, 56-63.
- [20] 孙言,杜彦辉. 网络钓鱼防御中 UNICODE 字符相似度评估算法的研究 [J]. 计算机工程与应用, 2008, 44(26):86-87.