

# 基于机器学习的校园网恶意网页检测方法

杨洪娇

(天津外国语大学, 天津 300204)

**摘要:**随着校园网网页逐渐增多, 恶意网页已成为校园网用户的主要安全威胁之一。笔者提出了一种基于 One-class SVM 机器学习算法的恶意网页检测方法, 首先通过规则匹配方法检测恶意网页疑似页面, 然后通过 One-class SVM 算法进行恶意网页检测。One-class SVM 仅通过正常网页数据进行训练, 克服了恶意网页数据集收集困难的问题。仿真结果表明该方法的准确率较高, 能够应用于校园网恶意网页检测工作中。

**关键词:**恶意网页; 机器学习; 校园网

**中图分类号:** TP181 **文献标识码:** A **文章编号:** 1003-9767 (2016) 11-175-02

## 1 概述

随着互联网与高校信息化的不断发展, 越来越多的服务和网站部署在校园网中, 很多页面可以被校园网用户及校外互联网用户访问到。导致恶意网页对网络用户的攻击不断增多<sup>[1]</sup>, 校园网具有用户群相对固定、用户防护能力相对有限、发生恶意网页安全事件后传播较快等特点。所以, 对校园网部署的网站进行恶意网页检测并定期监控, 已成为校园网亟需解决的安全问题之一。

专业人员提出了很多技术进行恶意网页检测, 其检测技术可以分为基于黑名单的方法、基于动态分析的方法、基于静态分析的方法三类<sup>[2]</sup>。基于黑名单的方法通过列出已知的恶意网页的 URL、IP 地址等信息来屏蔽用户访问, 从而简单地检测出恶意网页, 该方法能够快速、准确地识别出已被确认的恶意网页, 但对未对外网开放访问的校园网环境具有局限性。基于动态分析的方法即根据网页的动态特征进行分析, 常见的动态特征包括: 注册表、文件夹及文件的变化情况、网页间的跳转关系和浏览器行为等。通过捕获网页异常行为的动态特征来判断恶意网页, 该方法多借助蜜罐技术或虚拟机技术进行分析。该方法分析的难点在于获取网页的动态特征, 这些动态特征需要人们对恶意网页长时间的 analysis 才能够获得, 其分析过程也相对耗时和耗费资源。基于静态检测技术的方法对页面中的 HTML 信息、JavaScript 代码特征等静态信息进行分析, 该方法能够及时发现恶意网页并采取相应措施, 但对于隐蔽性强的恶意网页容易出现漏判现象。Suyeon Yoo、李洋等都提出了利用机器学习技术检测恶意网页的方法<sup>[2-3]</sup>, Suyeon Yoo 的方法为采用两阶段的算法综合检测恶意网页。侯冰楠等提出了一种 JSFEA 检测方法, 采用预过滤的方法结合动态检测方法对恶意网页进行检测<sup>[4]</sup>, 其预过滤阶段采用四个静态条件判断页面是否为可疑页面, 然后利用动态检测方法模拟脚本运行来判定恶意网页。

本文提出了结合规则匹配与机器学习算法的恶意网页检测方法, 通过规则匹配方法预先筛选可疑页面, 且采用合理

的规则减少了漏报的概率, 然后利用 One-class SVM 算法判断可疑页面是否为恶意页面, 由于不采用动态模拟运行脚本, 使得网页特征较易获取, 同时由于当前恶意网页的数据集较难获取, 采用 One-class SVM 仅需正常样本训练即可, 提高了数据采集与训练的效率。实验表明, 在校园网环境下采用该方法能够有效地完成恶意页面的检测分析工作。

## 2 One-class SVM 算法

支持向量机 (SVM) 是一种基于统计学习理论发展的机器学习算法, 支持向量机算法把原始数据映射到高维空间, 通过寻找高维空间中的一个超平面, 使得该超平面既能保证分类的精度, 又能使超平面两侧的空白区域最大化, 从而使支持向量机可以用于模式识别、预测、判别等领域。同时, 支持向量机在理论上可以实现对现行可分数据的最优分类<sup>[5]</sup>。

支持向量机技术具有良好的泛化能力, 在多个领域得到广泛使用, 在传统 SVM 算法中, 对样本数量较少或样本类内差异较小时设置不同的惩罚因子, 以防止训练的支持向量机出现分类错误。

One-class SVM 是由 Bernhard Scholkopf 等人提出的用于解决一类分类问题的方法<sup>[6]</sup>。作为一类分类器, 其核心思想是将样本数据映射到高维空间, 寻找高维空间中样本数据密度较大的区域, 如果某一数据处于较大密度区域之内, 则认为其是正常数据, 否则就认为其是异常数据。由此可见, One-class SVM 的训练只需要一类正常样本即可进行, 值得一提的是, One-class SVM 算法处理有噪声的样本时也有较好的表现。

## 3 网页静态特征

利用机器学习算法需要有效的特征来区分恶意网页和非恶意网页, 恶意网页检测技术中的静态特征较多, Hou 等人将网页特征分为三类, 即 HTML 文档特征、Javascript 特征和 ActiveX 特征, 其方法共使用了 171 个特征进行检测<sup>[7]</sup>。李洋等的方法结合有效特征和数学统计的特征, 采用了 13

**作者简介:** 杨洪娇 (1988-), 男, 河北石家庄人, 硕士, 助理工程师。研究方向: 图像处理与模式识别。

个特征进行检测分析。

本文根据 Hou、李洋等的方法,结合本文采用规则匹配的方法对网页进行可疑性检测,选取其中 11 个特征进行检测,见表 1。

表 1 恶意网页检测特征表

序号	所用特征
1	是否包含隐藏元素
2	Eval 字符串出现次数
3	Exec 字符串出现次数
4	Script 标签是否对称
5	Unicode 字符串出现的次数
6	特殊字符个数百分比
7	%u 个数百分比
8	长字符串个数
9	字符串平均长度
10	Iframe 元素的大小
11	空格字符百分比

#### 4 基于机器学习的恶意网页检测方法

采用机器学习方法对网页的静态特征进行分类,能够有效区分正常网页和恶意网页,本文通过采用 One-class SVM 方法对被规则匹配方法判断为疑似恶意的网页进行分类,对恶意网页进行检测。

##### 4.1 规则匹配方法

采用机器学习算法对网页静态特征进行检测往往需要提取较多特征并进行检测,校园网环境下的网站页面中大多数是不包含恶意代码的正常页面,如果在网站检测过程中对每个页面都进行特征提取将降低网页检测效率。JSFEA 方法中就采用了静态特征预过滤的手段判断可疑页面,对可疑的网页才进行下一步检测。但当前恶意页面多使用混淆技术等方法绕过安全检测<sup>[8]</sup>,面对这种问题,本文对 JSFEA 的预过滤条件进行完善,以减少漏报可能,提高混淆检测能力,在 JSFEA 原有的四个预过滤方法上增加两个规则匹配方法对可疑页面进行判断。

(1) URL 参数长度大于 10 个字符。当前恶意网页的逃逸技术中采用 URL 混淆技术来逃避检测,将对 URL 参数长度的判断作为判断其是否可疑的标准之一。

(2) 是否含有长字符串。恶意网页在对代码改变编码方式时,其编码的字符串通常较长,本文方法对长度超过 150 的字符串进行检测,以判断是否包含可疑编码信息。

通过运用规则匹配方法,恶意网页检测方法首先判断网页是否可疑,只有可疑的页面才进行基于机器学习方法的判断,凡是符合规则匹配阶段任何一个规则的都会被标记为可疑页面。该方法的应用可以提高系统检测效率,通过合理的规则设置又能降低漏报的概率。

##### 4.2 数据采集与训练

对 One-class SVM 进行训练需要有训练集,网页数据采集通常使用爬虫工具对网页进行抓取,在网页选择方面,正常网页通常选择 Alexa 网站排名靠前的网页进行抓取,这些网站通常访问量大、安全管控严格,被认为是正常数据。但当前恶意网页的数据采集较困难,原因是部分恶意网页的生命周期

持续缩短,同时恶意网页的逃逸技术不断升级<sup>[9]</sup>。本文采用的 One-class SVM 算法仅需正常数据样本即可进行训练,只需要采集正常网页数据,避免了恶意网页采集困难的情况。

本文针对校园网环境进行恶意网页检测,校园网部署的网站主要分为宣传类网站与业务系统,为采集正常网站,可选择校园部署的网站中源代码可控的网站进行采集,认为是正常网站数据。

One-class SVM 在恶意网页检测中分为两部分,即样本训练和检测应用。在样本训练部分,利用采集到的正常网页数据,对样本数据进行特征提取,对 One-class SVM 进行训练。在应用中通过对校园网网页与人工采集的恶意网页样本进行测试,最终应用到校园网恶意网页日常监测中。

#### 5 结 语

为了测试本文方法在恶意网页监测中的效果,对算法进行仿真实验分析,从仿真分析的结果看,小样本测试数据达到了较高的准确率。本算法从校园网恶意网页检测角度出发,首先采用规则匹配方法检测出疑似的恶意网页,然后对疑似的网页利用 One-class SVM 机器学习算法对校园网网页进行检测分析。仿真结果表明,利用该算法能较好地地进行校园网恶意网页的检测工作。但算法采用的训练样本较少,需进一步增加样本,提高检测的准确率。

#### 参考文献

- [1] 沙泓州,刘庆云,柳厅文,等.恶意网页识别研究综述[J].计算机学报,2015.
- [2] Yoo S, Kim S, Choudhary A, et al. Two-phase Malicious Web Page Detection Scheme Using Misuse and Anomaly Detection[J]. International Journal of Reliable Information and Assurance, 2014, 2(1).
- [3] 李洋,刘飏,封化民.基于机器学习的网页恶意代码检测方法[J].北京电子科技学院学报,2013,20(4):36-40.
- [4] 侯冰楠,俞研,吴家顺.基于预过滤的恶意 JavaScript 脚本检测与分析方法[J].计算机应用,2015,35(A01):60-62.
- [5] 贾彦茹.基于 SVM 的入侵检测方法[J].信息与电脑(理论版),2016(5).
- [6] B Schölkopf, J C Platt, J Shawetaylor, et al. Estimating the Support of a High-dimensional Distribution[J]. Neural Computation, 2001, 13(7):1443-1471.
- [7] YT Hou, Y Chang, T Chen, et al. Malicious Web Content Detection by Machine Learning[J]. Expert Systems with Applications, 2010, 37(1):55-60.
- [8] 马洪亮,王伟,韩臻.基于 JavaScript 的轻量级恶意网页异常检测方法[J].华中科技大学学报:自然科学版,2014,42(11):34-38.
- [9] 丁世飞,齐丙娟,谭红艳.支持向量机理论与算法研究综述[J].电子科技大学学报,2011,40(1):2-10.