

Detection of Malicious Application Based on Improved Naive Bayesian Algorithm Android

Ruzhen Shi

School of Communication and Information Engineering, Nanjing University of Posts and Telecommunications,
Nanjing Jiangsu
Email: shiruzhen0@126.com

Received: Aug. 19th, 2016; accepted: Sep. 9th, 2016; published: Sep. 12th, 2016

Copyright © 2016 by author and Hans Publishers Inc.
This work is licensed under the Creative Commons Attribution International License (CC BY).
<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

With the popularity of smart phones, Android mobile applications are becoming more widely and the attendant problem about the Android security issues becomes more critical. Malicious applications are endless. It is a serious threat to the user's information security. This paper describes the current status of mobile security and the current Android malicious application detection techniques are analyzed. As the situation of Android applications becomes more and more serious, an improved Naive Bayes algorithm is proposed. The malicious application detection system is designed by using this algorithm. The system uses a reverse analysis method to produce the malicious application classifier. The experimental results are analyzed.

Keywords

Android, Naive Bayes, Reverse Analysis, Classifier

基于改进朴素贝叶斯算法Android恶意应用的检测研究

石汝振

南京邮电大学通信与信息工程学院, 江苏 南京

Email: shiruzhen0@126.com

收稿日期: 2016年8月19日; 录用日期: 2016年9月9日; 发布日期: 2016年9月12日

摘要

随着智能手机的普及, Android手机上的应用也越来越广泛, 随之而来的问题是Android安全问题越来越严峻。各种恶意应用也层出不穷, 严重威胁到用户的信息安全。本文首先介绍了当前的手机安全状况, 并对当前的Android恶意应用的检测技术进行了分析。由于Android应用越来越严峻的形势, 本文提出了一种基于改进的朴素贝叶斯算法, 并应用此算法设计了恶意应用检测系统。本文中的系统应用了逆向分析的方式来产生恶意应用分类器。然后, 本文对实验结果进行了分析。

关键词

Android, 朴素贝叶斯, 逆向分析, 分类器

1. 引言

根据 360 发布的《2015 年中国手机安全状况报告》 [1]显示:

2015 全年, 360 互联网安全中心累计截获 Android 平台新增恶意程序样本 1874.0 万个。分别是 2013 年, 2014 年的 27.9 倍、5.7 倍。平均每天截获新增恶意程序样本也高达 51,342 个。

2015 全年, 360 互联网安全中心累计监测到 Android 用户感染恶意程序 3.7 亿人次, 分别是 2013 年、2014 年的 3.8 倍和 1.1 倍; 平均每天感染量达到了 100.6 万人次。

在所有手机恶意程序中, 资费消耗类恶意程序的感染量仍然保持最多, 占比高达 73.6%; 其次为恶意扣费(21.5%)和隐私窃取(4.1%)。

恶意应用的开发者们不断的开发着新的恶意代码,新的木马病毒样本不断被发现。从恶意扣费、窃取隐私、隐式安装恶意软件、窃取网银密码等到远程控制操作用户的手机。Android 智能手机不仅仅给用户带来了方便也导致了 Android 安全问题越来越严重。Android 手机上恶意应用的检测越来越被人们所重视。

机器学习算法可应用于 Android 平台上恶意代码检测, 例如朴素贝叶斯算法(Native Bayes NB 算法)。但朴素贝叶斯算法存在着较高误判率的缺陷[2], 本文提出了应用改进的朴素贝叶斯算法进行机器学习, 对 Android 软件进行字节码级别的静态分析后, 提取出软件的特征值信息作为样本特征。本文基于该原理实现了对 Android 应用进行恶意应用检测的框架, 对 Android 软件进行特征码扫描后, 通过基于改进的朴素贝叶斯算法分类器的判定区分是否是恶意应用, 从而有效提高了恶意代码检测的精确度, 降低了误判率。

2. Android 恶意应用检测技术

Android 恶意代码的检测一般分为静态检测和动态检测, 在此基础上一般采用特征值检测技术和基于启发式的检测技术这两种方式[3]。

2.1. 特征值匹配技术

特征值检测是目前大多数反病毒产品普遍使用的检测技术。该技术的主要工作原理是根据特定的规则对目标文件进行扫描, 并提取出特征值, 与病毒库中已有的特征按照某种匹配算法进行完全匹配或者

计算出相似度，若得到完全匹配或者相似度超过规定的阈值，则表示检测为病毒文件[4]。所以特征值匹配需要分析大量的样本，并且不断的更新特征库，特征值匹配技术原理比较简单，而且在一定情况下具有容易实现，准确率高等众多优点。但是这种技术过分依赖于对特征值提取的精确程度，具有一定的不稳定性，而且对新出现的恶意代码的检测效果不理想。

2.2. 基于启发式的检测技术

启发式的检测技术包含动态和静态启发检测，静态启发是指对 Android apk 文件进行反编译生成 smile 代码和 java 代码，通过对文件外部静态信息分类，模拟跟踪代码执行流程来判断是否为恶意软件[5]。动态启发是指在系统中设置若干的特征点来监控软件行为，通过软件的恶意行为来判断是否是恶意软件。常用的启发式检测技术基于机器学习算法例如神经网络、朴素贝叶斯算法、决策树等。

3. 朴素贝叶斯算法分析与改进

3.1. 基本原理

朴素贝叶斯算法是目前公认的最简单而且有效的概率分类方法[6]。它是一种假定各个因素之间不存在任何联系，即完全独立得到的一种简化的贝叶斯算法。可以将朴素贝叶斯分类器应用于恶意代码的过滤，利用这种方法可以根据结果集自动学习，因此需要一定的恶意代码库来提取恶意代码的特征值来训练过滤器。

朴素贝叶斯方法是基于最小错误概率规则，尽量降低分类犯错误的概率。朴素贝叶斯方法通过在程序中的恶意代码的特征信息来计算程序中包含恶意代码的可能性。当判断某段代码是否是恶意代码的时候，利用提取出的特征集合 A 。定义 B 代表所有类别的一个随机取值， B 代表的是恶意样本或者普通代码，进行机器学习的目的就是计算 $P(B|A)$ 来判断 B 是恶意代码还是普通代码。由贝叶斯公式可得，

$$p(B|A) = \frac{P(A|B) \times P(B)}{P(A)} \quad \text{公式 2.1}$$

由于贝叶斯公式的条件是特征集合 A 中的各个特征值相互独立，如果软件样本包含多种属性 $A_1, A_2, A_3 \dots A_n$

$$P(B|A) = \frac{\prod_{i=1}^n P(A_i|B) \times P(B)}{\prod_{j=1}^n p(A_j)} \quad \text{公式 2.2}$$

从中选取最大概率作为其分类。

3.2. 朴素贝叶斯算法的改进

在恶意代码的实际分类中，由于只有两类即是否是恶意代码，所以相应的会出现两种分类错误：1) 将恶意代码判断成正常代码；2) 将正常代码判断成恶意代码；如果是将恶意代码判断成为正常代码，就会带来严重的后果，使用户信息丢失等等。在传统方法中一般当 $P(B_0|A) > P(B_1|A)$ 时就判定是否是恶意代码，但是这种方法并不精确存在很高的误判率和漏判率所以直接应用时分类偏差会比较大。

为了更准确的识别恶意代码，减少误判，设当 $\frac{P(B_1|A)}{P(B_0|A)} > \theta$ 时[7]，即当恶意代码的概率是

非恶意代码的 θ 倍时，将其判定为恶意代码。当 θ 越大，其为恶意代码的可能性就越大。由上述公式可以表示为，

$$\frac{P(B_1|A)}{1-P(B_1|A)} > \theta \tag{公式 2.3}$$

$$P(B_1|A) > \frac{\theta}{1+\theta} = h \tag{公式 2.4}$$

也就是当 $P(B_1|A) > h$ 时可以判定此应用为恶意 Android 应用。

4. Android 恶意代码静态检测系统的设计

4.1. 生成恶意代码检测分类器

通过对训练样本进行分析，得到样本的特征值集合，流程图如图 1 所示，首先收集 Android 的恶意软件的样本集，对样本集进行逆向分析，在分析出的信息中提取出分类器需要的信息特征，将这些特征作为分类器的样本输入，经过改进的朴素贝叶斯算法处理之后得到 Android 恶意软件的检测分类器[8]。

4.2. 应用分类器对 Android 应用进行检测

当测试是否是恶意应用时，需要先对测试的样本进行逆向分析，提取出特征信息，送入分类器进行检测。流程图如图 2 所示。

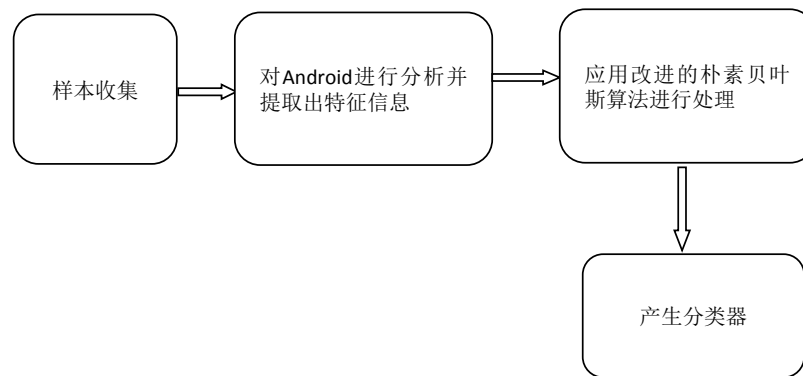


Figure 1. Malicious code detection classifier generation

图 1. 恶意代码检测分类器的生成

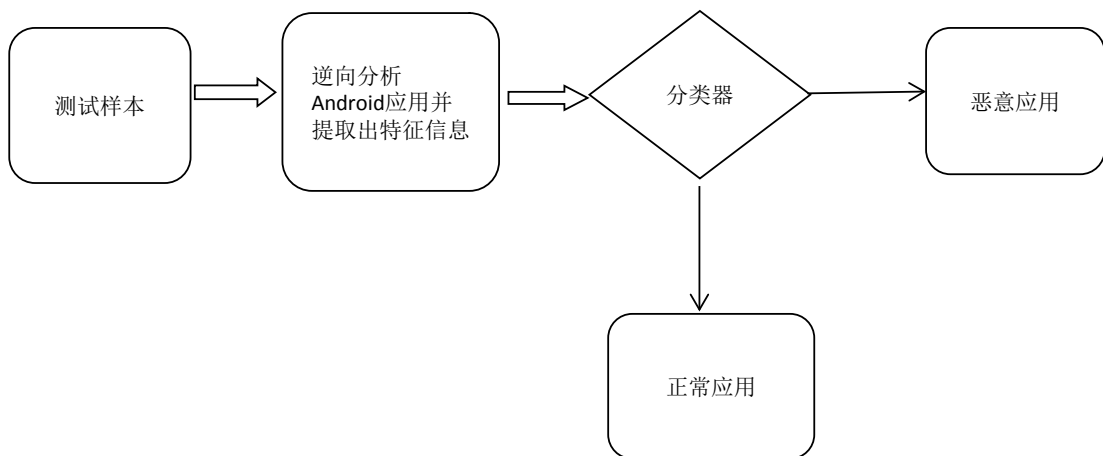


Figure 2. Classifier sample detector

图 2. 分类器样本的检测

4.3. Android 应用的静态逆向分析

APK 文件的反编译可以使用 dex2jar 把 classes.dex 文件反编译为 java 代码，反编译成功后的代码可以通过 jd-gui 进行阅读[9]。使转化后的 java 代码清晰易懂；或者把 classes.dex 文件反编译为 smali 代码，smali 代码是一种面向对象的汇编代码，其可读性较强，修改后的 smali 代码还可以重新打包和签名生成新的 APK 文件[6]。反编译成 smali 代码常用的工具是 apktool，此工具集合了解压缩 XML、文件解析、smali 正反编译、组装等各个功能。运用反编译工具可实现对 apk 文件的反编译，从而提取出所需要的特征信息，对其所申请的权限进行分析，运用字符串的比较算法比较签名等从而得到最终结果。

5. 试验与结果分析

通过网络收集到的软件样本 600 个，其中有恶意应用 300 个，合法应用 300 个。得出基于朴素贝叶斯算法的实验结果如表 1 所示。

由表 1 得知，应用朴素贝叶斯算法时，300 个恶意应用中仍然有 35.5 个被判定为合法应用，300 个合法应用中有 37.7 个被判定为恶意应用。实验数据说明了朴素贝叶斯算法存在着一定的缺陷。

由公式 2.3 和 2.4 可得 h 的范围是 0~10，所以需要大量的实验来得到 h 值的大小。从中选取了部分 h 值实验结果如表 2 所示。可得当 $0 < h < 0.35$ 时分类器的判断正确概率接近 100%。

Table 1. Experiment based on naive Bayes algorithm

表 1. 基于朴素贝叶斯算法的实验

	系统判定为合法应用	系统判定为恶意应用	总数
实际为合法应用	262.3	37.7	300
实际为恶意应用	35.5	265.5	300
总数	297.8	302.2	600

Table 2. Experiment based on improved naive Bayes algorithm

表 2. 基于改进朴素贝叶斯算法的实验

	系统判定为合法应用	系统判定为恶意应用	数量	H 值
实际为合法应用	258.2	41.8	300	0.55
实际为恶意应用	40.1	259.9	300	0.55
实际为合法应用	265.5	35.5	300	0.5
实际为恶意应用	37.7	262.3	300	0.5
实际为合法应用	265.8	34.2	300	0.40
实际为恶意应用	10.4	289.6	300	0.40
实际为合法应用	272	28	300	0.35
实际为恶意应用	0	300	300	0.35
实际为合法应用	287	13	300	0.27
实际为恶意应用	0	300	300	0.27

6. 结束语

本文对现有的恶意应用的检测技术进行了分析，在一种朴素贝叶斯算法的基础上提出了一种改进型的朴素贝叶斯算法，并设计了一个检测系统对该改进算法的优势进行了检测。提高了恶意应用的检测概率，但是也存在着将合法应用判定为恶意应用的误差，这是算法需要进一步改进的地方。

参考文献 (References)

- [1] 2015 年中国手机安全报告[R]. <http://zt.360.cn/1101061855.php?dtid=1101061451&did=1101593997>
- [2] 王辉, 陈泓予, 刘淑芬. 基于改进朴素贝叶斯算法的入侵检测系统[J]. 计算机科学, 2014, 41(4): 111-119.
- [3] Kosmopoulou, A., Paliouras, G. and Androutsopoulos, I. (2008) Adaptive Spam Filtering Using Only Naïve Bayes Text Classifiers. *CEAS 2008 5th Conference on Email and Anti-Spam*, Mountain View, 21-22 August 2008, 3 p.
- [4] 蔡志平. 计算机病毒检测技术研究与应用[D]: [硕士学位论文]. 长沙: 国防科技大学, 2001: 50-56.
- [5] 彭国军, 李晶雯, 孙润康, 肖云倡. Android 恶意软件检测研究与进展[J]. 武汉大学学报(理学版), 2015, 61(1): 21-33.
- [6] 楼赞程, 施勇, 薛质. 基于逆向工程的 Android 恶意行为检测方法[J]. 信息安全与通信保密, 2015(4): 83-87.
- [7] 郑炜, 沈文, 张英鹏. 基于改进朴素贝叶斯算法的垃圾过滤器的研究[J]. 西北工业大学学报, 2010, 28(8): 622-628.
- [8] 郑吉飞. Android 恶意代码的静态检测研究[D]: [硕士学位论文]. 武汉: 华中师范大学, 2013.
- [9] 侯勤胜, 曹天杰. Android 恶意软件的分析与检测[J]. 河南科技大学学报: 自然科学版, 2015, 36(10): 54-59.

期刊投稿者将享受如下服务:

1. 投稿前咨询服务 (QQ、微信、邮箱皆可)
2. 为您匹配最合适的期刊
3. 24 小时以内解答您的所有疑问
4. 友好的在线投稿界面
5. 专业的同行评审
6. 知网检索
7. 全网络覆盖式推广您的研究

投稿请点击: <http://www.hanspub.org/Submission.aspx>