

文章编号: 1672-2892(2011)02-0234-04

数据流异常检测及其在僵尸网络检测中的应用

邓 军

(西南交通大学 电气工程学院, 四川 成都 610031)

摘 要: 目前大多关于 P2P 僵尸网络检测的研究都采用传统的逆向工程方法, 这些方法检测都比较准确, 但其工程实施难度太大, 效率较低, 且对于变种病毒, 该类检测方法无能为力。本文尝试通过数据流异常检测技术的应用, 找到一种适合数据流应用场景的异常检测方法, 并尝试将其应用于 P2P 僵尸病毒的检测当中, 通过监控网络数据流, 能够有效地发现 P2P 僵尸病毒在传播过程当中的特殊行为, 并通过捕获这些行为来实现发现僵尸主机的目的。

关键词: 僵尸网络; 数据流异常检测; 聚类建模

中图分类号: TP393.08

文献标识码: A

Data flow anomaly detection technique and its application in Botnet detection

DENG Jun

(School of Electrical Engineering, Southwest Jiaotong University, Chengdu Sichuan 610031, China)

Abstract: Most of the current detection of P2P(Peer to Peer) Botnet adopts traditional reverse engineering method, which is very accurate, but difficult to be implemented and shows low efficiency. It becomes ineffective for variants. This paper attempts to find a data stream anomaly detection method suitable to the data stream application cases, and tries to apply it to P2P Zombie Virus detection. By monitoring network data stream, the special behaviors of P2P Zombie Virus in their spreading can be found. The locating of the zombie host can be realized by capturing those behaviors.

Key words: Botnet; data stream anomaly detection; clustering model

计算机网络安全作为数据流异常检测技术的一个重要应用领域, 正随着互联网和个人计算机的发展越来越受到人们的重视。僵尸网络(Botnet)的出现对网络和数据的安全构成极大的威胁, 已引起了各方面的高度重视。

1 僵尸网络背景介绍

随着 P2P 技术的发展, 采用 P2P 协议特点构建的僵尸网络成为近几年出现的一种新型僵尸网络。与传统的基于 IRC(Internet Relay Chat)协议的僵尸网络相比, P2P 僵尸网络克服了其必须具有固定集中控制节点(IRC 服务器)的缺陷, 每个被 Bot 感染的机器都可作为 P2P 僵尸网络的服务器(控制节点)或是客户端(僵尸机), 其典型结构如图 1 所示。

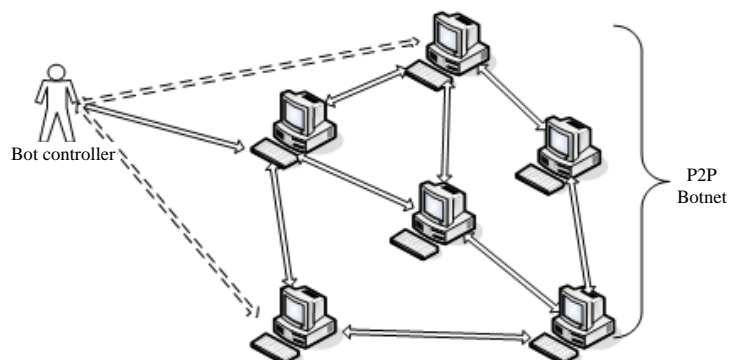


Fig.1 Typical structure of P2P Botnet
图 1 典型僵尸网络结构图

2 P2P 僵尸网络检测技术

相比 P2P 僵尸网络的快速蔓延, P2P 僵尸网络检测技术只能算刚刚起步。现有大多成熟的僵尸网络检测技术都是针对 IRC 僵尸网络提出的, 而对于检测隐蔽性和灵活性都更好的 P2P 僵尸网络, 则很少有研究者提出更有效的检测方法。

大多关于 P2P 僵尸网络检测的研究，都是采用传统的逆向工程方法：对某个特定的 P2P 僵尸病毒进行分析，然后通过得到的静态特征对该种僵尸病毒进行检测；或通过静态分析的结果，修改僵尸病毒源码并伪装僵尸机加入僵尸网络，将跟踪获得的僵尸网络信息发回到指定的终端进行分析，从而实现僵尸网络的识别^[1-4]。这 2 种逆向工程方法的检测都比较准确，但其工程实施难度太大，效率较低，并且只对所分析的僵尸病毒有效，如果僵尸的制造者修改了僵尸病毒(变种)，或出现了其他的僵尸病毒，则该类检测方法无能为力。

3 基于聚类建模的数据流异常检测方法

相较于有监督异常检测方法，无监督或者半监督的异常检测方法^[5]的实用性更高。无监督的异常检测方法不需要带有准确标签的训练数据集；半监督的异常检测方法只需要正常数据集作为训练集，这就为某些可能的应用放宽了限制，如在网络恶意行为的检测上，对未知入侵行为(或病毒的传播行为等)进行检测时，是不可能事先得到异常数据行为作为训练集的。本文选择无监督的聚类算法作为异常检测的基础算法，主要是因为对包含未知异常对象的数据流进行检测时，训练数据集(包括正常数据对象和异常数据对象)是很难完整定义和收集的。但由于数据流具有动态数据集的特点，原始的聚类算法并不能很好地适应基于数据流的异常检测，其解决办法可以分为 2 种：

- a) 采用动态的聚类算法。当输入新的数据对象时，将该数据对象进行聚类，同时更新之前保存的聚类结果；
- b) 利用聚类结果作为训练集建立分类器，将新输入的数据对象通过分类器直接进行分类得到结果。

基于网格聚的异常检测就是一种快速的动态聚类方法，其本身的效率很高。但基于网格的聚类方法只能识别出球形的聚类结果，对于复杂的、具有不规则类边界数据集，则效果不理想。动态局部异常因子(Local Outlier Factor, LOF)算法是将静态 LOF 算法进行改进而得到的数据流异常检测算法，在保证准确率的情况下，其效率远远高于重复执行相应的静态 LOF 算法，但因每得到一个新输入的数据对象都要动态地更新原有部分数据对象的特征信息，其计算复杂度仍然很高，且该方法继承了静态 LOF 的缺点，即无法发现凝聚成团的异常数据^[6-8]。

本文基于后一种方式，提出了一种通过聚类结果进行建模的数据流异常检测方法，该方法主要包括聚类、分类器训练 2 个部分。聚类算法通过收集数据流中的数据对象进行聚类，给出聚类结果，然后对聚类结果数据集中带有标签的数据进行分类器训练。在建成分类器之后，所有新输入的数据对象就会通过分类器直接进行判断，给出最终分类结果。采用先聚类后训练分类器的方式进行数据流的异常检测方法，解决了反复进行聚类效率低下的问题。聚类建模具有一定的时间复杂度，但一旦分类器训练完成，就能够实时进行数据流的异常检测。该方法不需要事先知道待检测数据流中异常数据和正常数据信息，整体上是一种无监督的异常检测方法，这也放宽了本文尝试将该方法应用到僵尸网络检测上的限制。

数据流作为动态数据集还有一个显著的特性，就是其包含的数据对象可能存在概念漂移。为了解决这个问题，

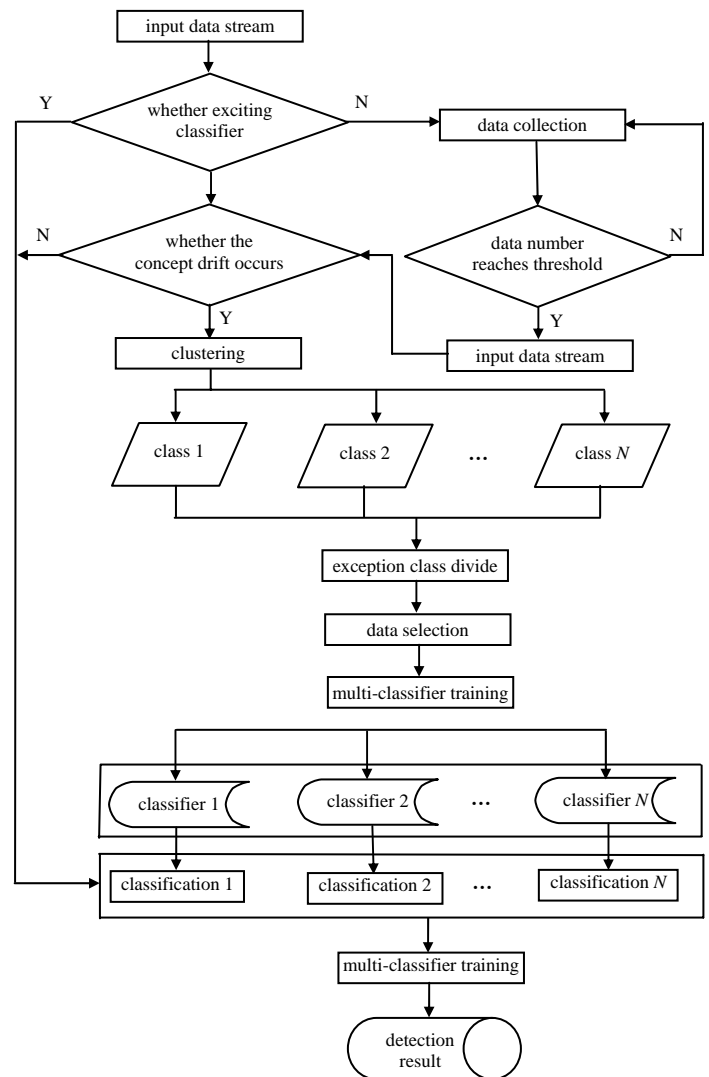


Fig.2 Procedure of improved data stream anomaly detection method
图 2 改进后的数据流异常检测方法流程

需要对基于聚类建模的数据流异常检测方法进行一些改进,改进之后的流程图如图2所示。本文采用了一种概念漂移检测结合概念漂移容忍的响应式方法,以捕获数据流概念漂移的发生为触发点,对异常检测模型进行更新,在某2个概念漂移触发点之间选择建立能够容忍概念漂移的分类器,以提高检测方法对数据漂移的健壮性。本文在分类器训练之前加入训练数据集平衡策略和多分类器集成来容忍一定程度的概念漂移,并且在数据流异常检测整体流程中加入概念漂移检测模块,以及时捕获数据流的概念漂移并触发分类器更新。改进之后的数据流异常检测方法,既能根据数据流的变化对异常检测模型进行更新,又能快速有效地进行数据流的异常检测。

本文提出的数据流异常检测方法的有效性基于如下2个前提条件:

- 待分析数据流中的正常数据对象和异常数据对象其本质上确实存在差异;
- 待分析的数据流中,正常数据对象是主要组成部分,其所占的比例远远超过异常数据对象所占比例。

4 基于数据流异常检测的僵尸网络发现

4.1 P2P 僵尸病毒传播状态模型

本文通过对大量僵尸病毒的分析,抽象出一种僵尸病毒感染过程的状态模型。该模型描述了一个潜在的目标计算机,从被感染到成为僵尸机并进一步作为僵尸网络的一部分,感染其他潜在受害者的变化过程。状态模型包含5个状态,代表了僵尸病毒一次感染扩散的全过程。原始的状态模型并没有针对P2P僵尸进行特殊的描述。本文针对P2P僵尸网络检测的最终目的,通过对典型僵尸病毒Phatbot的分析,改进了状态模型,最终得到描述P2P僵尸病毒的传播状态模型,如图3所示。假设状态模型为 $S = \{E_1, E_2, E_3, E_4, E_5\}$,其中各元素代表的含义为: E_1 —僵尸网络对潜在受害者发起的漏洞扫描; E_2 —利用漏洞和后门进行入侵感染; E_3 —被攻陷的主机通过内建信息(预置的PeerList)加入P2P僵尸网络; E_4 —与僵尸网络的命令控制节点或其他僵尸机进行通信和交互; E_5 —作为僵尸网络新的组成部分继续扫描下一个潜在受害者。

在此模型的基础上,如果能够捕获僵尸病毒在传播过程中的这几个典型状态,则认为能够发现网络中的僵尸机,进而发现整个僵尸网络的存在。本文主要关注将数据流异常检测方法应用于 E_3 和 E_4 状态。

4.2 P2P 僵尸病毒通信的捕获

僵尸病毒为了更好地隐藏自己不被使用者发现,不同僵尸机之间的通信特征通常不会太明显,往往被淹没在大量正常的通信数据流中。这符合前文第3节中提到的假设条件,因此尝试采用本文提出的数据流异常检测方法,来检测通信数据流中混杂的小部分僵尸病毒。

首先需要收集P2P僵尸的通信数据来构造测试数据集,因此利用多台计算机和虚拟机软件Vmware搭建了一个可控的虚拟环境,其拓扑结构如图4所示。

图4的可控环境由5台虚拟机组成。首先在每个主机上运行网络嗅探程序,然后在每个虚拟机上运行P2P僵尸病毒样本PhatBot,通过网络嗅探器来捕获可控环境中各个主机之间的通信数据包。

每台主机上运行的嗅探器监控10个僵尸病毒(PhatBot)通信端口,并捕获这些端口上发送和接收的数据包。在捕获监控端口通信数据包的同时,分析数据包头信息,并对数据包的特征进行统计提取。设置时间间隔 $t=30\text{ s}$,每隔 t 时间将统计到的数据包特征保存到本地,作为1条僵尸病毒通信记录,然后将统计信息清零并重新统计下个 t 时间内的数据包特征信息。记录的每条数据包包含7维特征,每一维特征的值是10个端口的统计值之和,如某一维表示监控的10个僵尸通信端口发送的字节数,则最终记录的该维特征值为10个端口发送的字节数总和。

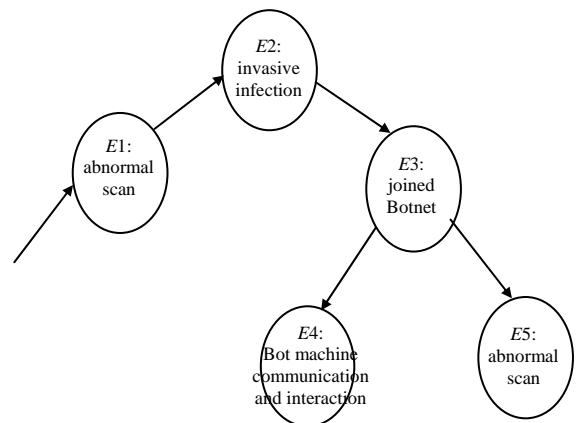


Fig.3 State model of once spread and infected P2P Zombie virus
图3 P2P 僵尸病毒一次传播感染状态模型

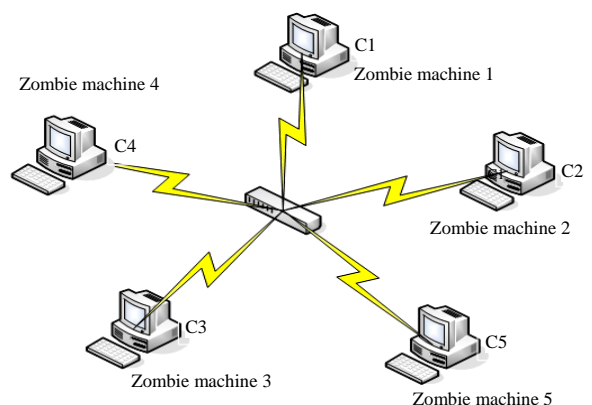


Fig.4 Virtual Botnet network environment
图4 虚拟僵尸网络环境

每条通信记录向量构成为： $\{S_{NCR}, Min_{ps}, Max_{ps}, S_{ub}, S_{db}, S_{ARP}, S_{ICR}\}$ 。其中， S_{NCR} 表示监控时间间隔内，主机上被监控的各端口建立连接总数； Min_{ps} 表示在监控时间间隔内，被监控各端口上数据包大小的下界平均值； Max_{ps} 表示在监控时间间隔内，被监控各端口上数据包大小的上界平均值； S_{ub} 表示在监控的时间间隔内，被监控的各端口所发送的字节总数； S_{db} 表示在监控的时间间隔内，被监控的各端口所接收的字节总数； S_{ARP} 表示在监控的时间间隔内，被监控的各端口发出 ARP 请求的次数； S_{ICR} 表示在监控的时间间隔内，被监控的各端口发出 ICMP 主机无法连接的次数。按照上述方式抓取僵尸通信数据包，提取特征并保存，最终从所有记录的僵尸病毒通信数据中选取特定条数，作为测试数据集中的僵尸病毒通信记录。

4.3 试验结果及分析

将4.2节中收集的测试数据集作为输入数据流，采用本文描述的聚类建模数据流异常检测方法进行异常检测。试验共进行2次，首次试验不对输入数据流进行概念漂移的检测，而是每分类3000个数据就对分类器进行更新。数据收集窗口大小为3000(同聚类数据块)，聚类参数设置为 $\alpha=0.96, \beta=6$ ，神经网络分类器训练个数为3个；在第2次试验中，加入概念漂移的检测，容忍度阈值设定为0.5，滑动窗口长度设定为1000，其余参数和第1次试验保持一致。检测结果如表1所示。

表1 可疑僵尸主机检测结果

	amount of suspected Bots	amount of actual Bots	amount omitted	amount misjudged
test 1	30	21	0	9
test 2	23	21	0	2

可以发现，在采用了概念漂移检测之后，基于数据流的异常检测方法对数据流中包含的异常数据(Phatbot的通信数据)检测准确率有明显提升。这是因为测试数据

集中正常数据对象本身的数据量较大，且相互之间较为相似，其特征分布并没有出现太大的偏差；而作为异常的僵尸通信信息存在着分布发生突变的情况。因此加入了概念漂移之后的数据流异常检测能够及时捕获到这个变化，重新聚类并且更新分类器，更新后的分类器能够更准确地表示当前异常数据的特征分布。

5 结论

本文提出了一种基于概念漂移检测的响应式数据流的异常检测方法，在建立异常检测模型时采用了容忍概念漂移的策略。通过概念漂移检测来捕获数据流中的分布变化，及时更新分类器，并配合能够容忍概念漂移的分类器训练策略，将数据流中的概念漂移对异常检测结果产生的影响大大降低。同时本文还尝试将该方法应用于网络安全领域中僵尸网络的发现，通过对P2P僵尸病毒传播过程中网络数据流的监控，实现僵尸机之间通信的检测。

参考文献：

- [1] 张晨. 数据流聚类分析与异常检测算法[D]. 上海:复旦大学, 2009.
- [2] 李人和. 数据流异常检测系统若干问题研究[D]. 上海:复旦大学, 2008.
- [3] 曹锋. 数据流聚类分析算法[D]. 上海:复旦大学, 2000.
- [4] 王明丽. 基于主机的P2P僵尸病毒检测技术研究[D]. 成都:电子科技大学, 2009.
- [5] 冯永亮. 结构化P2P僵尸网络检测技术的研究[D]. 武汉:华中科技大学, 2008.
- [6] 刘伟,刘嘉勇. 一种基于扩展角色的访问控制模型和方法[J]. 信息与电子工程, 2009,7(1):76-79.
- [7] Sadngi A Selhorst M Stuble C. TCG TPM Main Specification[Z]. Version 1.2 Revision 94,Part1,Design Principles, 2006.
- [8] Choi H, Lee H. Botnet Detection by Monitoring Group Activities in DNS Traffic[M]. Fukushima:IEEE Computer Society Press, 2007:715-720.

作者简介：



邓 军(1975-), 男, 成都市人, 助理工程师, 主要研究方向为电力系统及其自动化.email:dengj_mail@tom.com.