

doi:10.11835/j.issn.1000-582X.2021.03.014

基于集成学习的智能电网主机恶意软件检测方法

李旭阳¹, 牛鑫¹, 胡军星², 袁俊锋², 孟晗²

(1. 国网河南省电力公司 经济技术研究院, 河南 郑州 450000; 2. 大河南九域腾龙信息工程有限公司, 河南 郑州 450000)

摘要: 目前智能电网恶意软件检测系统主要基于特征库对已知恶意软件进行检测, 不适用检测恶意软件未知变种。而现有基于机器学习的恶意软件未知变种检测方法的准确性和鲁棒性有待进一步提升, 不足以满足智能电网实际需要。因此, 提出一种基于集成学习的恶意软件未知变种检测方法, 利用多源数据集和多种机器学习方法交叉构建单一检测模型, 并设计一种基于 Logistic 的集成学习方法, 构建恶意软件未知变种集成检测模型。实验对比分析表明, 构建的集成检测模型相较于传统单一检测模型在准确性和鲁棒性方面有着显著提升。

关键词: 智能电网; 恶意软件未知变种检测; 机器学习; 集成学习

中图分类号: TP309

文献标志码: A

文章编号: 1000-582X(2021)03-144-07

Ensemble learning based malware detection method for smart grid

LI Xuyang¹, NIU Xing¹, HU Junxing², YUAN Junfeng², MENG Han²

(1. State Grid Henan Economic Research Institute, Zhengzhou 450000, P. R. China;

2. Henan Jiuyu Tenglong Information Engineering Co., Ltd., Zhengzhou 450000, P. R. China)

Abstract: The traditional malware detection system of smart grid mainly detects known malware based on feature database, which is not applicable for detecting unknown malware variants. Although the machine learning based detection methods can detect unknown malware variants, but the accuracy and robustness of the existing methods need to be further improved, which is not enough to meet the actual needs of smart grid. Therefore, this paper proposes an ensemble learning based unknown malware variants detection method, which uses multi-source data and multiple machine learning methods to construct several single detection models respectively, and designs a hybrid detection model based on logistic. Compared with the traditional single detection models, the accuracy and robustness of the hybrid detection model are significantly improved.

Keywords: smart grid; malware variants detection; machine learning; ensemble learning

随着信息通信及互联网技术的发展, 智能电网逐步发展起来, 并给电力系统带来了巨大变革。智能电网是一种新型的电力系统, 它能够监测分析客户、电网设备及网络节点上电力流与信息流, 控制电力流与信息流双向流动, 实现电网自主优化运行。

收稿日期: 2020-10-12

基金项目: 国家自然科学基金项目(61572517)。

Supported by National Natural Science Foundation of China(61572517).

作者简介: 李旭阳(1979—), 男, 高级工程师, 主要从事电力工程技术经济方向研究, (Tel)15303832395, (E-mail) lxy0393@163.com。

智能电网是建立在集成的、高速双向通信网络的基础上,通过先进的传感和测量技术、先进的设备技术、先进的控制方法以及先进的决策支持系统技术,实现电网的可靠、安全、经济、高效、环境友好和使用安全的目标,其主要特征包括自愈、激励和包括用户、抵御攻击、提供满足 21 世纪用户需求的电能质量、容许各种不同发电形式的接入、启动电力市场以及资产的优化高效运行^[1-2]。

智能电网给人们生活带来巨大的便利的同时,也面临着一些安全威胁,特别是恶意软件威胁。在智能电网中,主机系统安全关系着信息系统的正确性,攻击者通过向主机植入恶意软件,以达到破坏电力监控系统或篡改、伪造信息流的目的,进而影响电力系统的稳定性。常见的基于恶意软件的攻击类型包括 C&C、APT 等。

为了能有效抵御恶意软件威胁,常通过部署恶意软件检测系统对智能电网主机操作系统进行扫描和检测。早期恶意软件检测方法常通过提取恶意软件可执行程序 Hash 码对其进行唯一标志,并提取可采集的所有已知恶意软件的 Hash 码构建恶意软件 Hash 码库,通过匹配待检测样本与库中的 Hash 码以判别待检测样本是否是恶意软件。然而,随着恶意软件未知变种的泛滥(恶意软件变种的 Hash 码已发生改变),传统恶意软件检测方法已不再适用于检测未知变种。近几年学术界一些学者提出基于机器学习的检测方法以应对恶意软件未知变种威胁,然而受限于单一训练数据集和单一机器学习方法的局限性,基于机器学习的检测方法的鲁棒性存在一定不足。

针对上述问题,提出一种基于集成学习的恶意软件未知变种检测方法,采用集成学习方法,融合由多个数据集和多种机器学习方法训练得到的模型,以提升恶意软件未知变种检测的准确性和鲁棒性。研究的贡献包括:1)提出一种基于集成学习的恶意软件检测方法,集成多源数据和多种机器学习方法实现恶意软件未知变种检测;2)提出一种基于 Logistic 的集成学习方法,集成由支持向量机和卷积神经网络训练得到的单一检测模型;3)实验结果表明,提出的基于集成学习的恶意软件未知变种检测方法的精度约为 93%,召回率约为 97%。

1 国内外研究现状

随着机器学习等人工智能技术的发展,为检测恶意软件未知变种提供新的技术方法支撑。近年来,一些学者纷纷提出基于机器学习的恶意软件未知变种检测方法,通过提取恶意软件变种共性语义特征,并采用机器学习算法,如支持向量机^[3]、卷积神经网络^[4]、深度信念网络^[5]等,训练恶意软件变种检测模型以判别未知样本是否为恶意软件。

McLaughlin 等人^[6]对恶意软件操作码 n-gram 模型进行分布式表征,并采用卷积神经网络自适应修正分布式表征的连接权值,以训练恶意软件检测模型。Zhang 等人^[7]通过 n-gram 模型和统计表征表示恶意软件操作码序列和 API 调用,分别采用卷积神经网络和多层感知器提取恶意软件操作码特征和 API 调用特征,通过合并操作码特征和 API 调用特征并采用分类器训练恶意软件检测模型。Zhang 等人^[8]通过计算操作码 bi-gram 边缘概率矩阵,并以此构建操作码图像,最后采用 K 近邻方法通过比较操作码图像的相似性来检测恶意软件。Yan 等人^[9]提取若干段等长的恶意软件操作码序列以构建等宽二维操作码图像,并以此为特征,通过采用卷积神经网络训练恶意软件检测模型。Cesare 等人^[10]基于操作码控制流图表示恶意软件,并采用 EditString 算法比较待检测的未知样本与已知恶意样本间的相似性来判别未知样本是否为恶意样本。Joshua 等人^[11]采用基于 bi-gram 模型提取恶意软件字节码特征,并采用深度神经网络算法训练恶意软件检测模型。陈志锋等人^[12]采用动态分析技术,提出了一种基于数据特征的内核恶意软件检测方法。该方法首先通过分析内核数据生命周期内的访问行为以构建内核数据对象访问模型,并利用扩展页表监控内核对象访问操作,提取恶意软件内核数据特征,最后基于该特征进行检测。Huang 等人^[13]通过分析用户接口与程序顶层函数之间的交互,基于相似度方法检测隐藏的恶意行为。Bai 等人^[14]通过分析恶意软件运行时执行路径中触发的恶意行为,提取行为特征并采用机器学习分类算法训练恶意软件检测模型。

尽管基于机器学习的恶意软件未知变种检测方法已取得一定进展,但仍存在一些问题,主要体现在:1)训练样本有限,导致模型覆盖的特征分布不全面;2)模型拟合性能有限,利用有限的样本数据进行训练,容易产生过拟合问题。因此,提出一种基于集成学习的恶意软件未知变种检测方法,采用集成学习方法进一步

提升恶意软件未知变种检测的准确性和鲁棒性。

集成学习方法主要分为 bagging 和 boosting 两类, bagging 方法通过对若干子模型的结果进行加权平均以综合多模型的决策结果; boosting 方法通过利用上一模型的错误分类样本重新训练子模型并对所有子模型的结果进行加权平均以综合多模型的决策结果, 如 Adaboost^[15]、随机森林^[16]等。

2 基于集成学习的恶意软件检测方法

2.1 问题描述

基于机器学习的恶意软件未知变种检测方法利用有限的已标记的恶意/良性样本训练检测模型, 以检测恶意软件未知变种。然而, 一方面受限于有限的训练样本及其特征分布, 另一方面受限于单一模型拟合性能的局限性, 恶意软件未知变种检测模型的准确性有待进一步提升, 特别在开放环境下现有恶意软件未知变种检测模型的鲁棒性存在明显不足, 在检测多来源样本数据时, 检测结果存在较严重的偏斜, 制约模型实际应用效果。

2.2 总体方案

笔者提出一种基于集成学习的恶意软件未知变种检测方法, 采用 bi-gram 模型提取恶意/良性软件特征, 并采用 2 种机器学习方法分别利用 3 个来源不同的样本数据集进行训练, 构建 6 个恶意软件未知变种单一检测模型, 并基于 bagging 策略, 设计一种 Logistic 集成学习方法通过集成 6 个单一模型最终形成统一的检测模型, 进一步提升恶意软件未知变种检测的准确性和鲁棒性, 并利用该集成模型对未知恶意软件变种进行检测, 如图 1 所示。

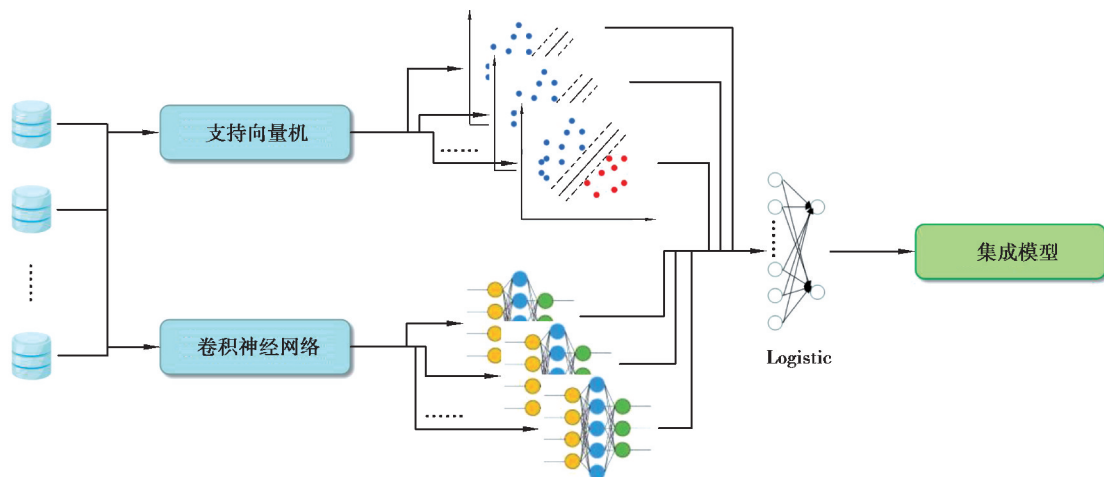


图 1 总体方案

Fig. 1 The architecture of approach

2.3 特征提取

笔者采用 IDA Pro 提取恶意/良性软件的操作码序列, 并采用 bi-gram 模型通过计算相邻两两操作之间的局部语义关系。令操作码序列 $OP = \{op_1, op_2, \dots, op_n\}$, 其中 op_i 表示序列中第 i 操作码, 采用 bi-gram 模型提取操作码序列中相邻两操作码之间的概率以表示恶意/良性代码局部语义 $bi\text{-}gram = \{\langle op_1, op_2 \rangle, \langle op_2, op_3 \rangle, \dots, \langle op_{n-1}, op_n \rangle\}$, 其中 $\langle op_{i-1}, op_i \rangle$ 表示相邻两操作码对, $p(\langle op_{i-1}, op_i \rangle)$ 表示相邻两操作码对 $\langle op_{i-1}, op_i \rangle$ 在序列中的共生概率, 该共生概率作为恶意/良性代码序列特征, 用于训练检测模型。

2.4 模型训练

基于上述提取的恶意/良性操作码序列局部语义特征 $p(\text{bi-gram})$, 分别采用支持向量机 (SVM, support vector machine) 和卷积神经网络 (CNN, convolution neural network) 2 种机器学习算法, 利用来源于 VxHeaven、Microsoft 和企业云安全中心的样本数据集, 交叉训练 6 种恶意软件检测模型: SVM_{vxheaven} 、 $SVM_{\text{Microsoft}}$ 、 $SVM_{\text{security_center}}$ 、 CNN_{vxheaven} 、 $CNN_{\text{Microsoft}}$ 、 $CNN_{\text{security_center}}$ 。其中支持向量机算法核函数采用线性核, 训练一组权重

$W = \{w_1, w_2, \dots, w_n\}$, 使得 $W = \{w_1, w_2, \dots, w_n\}$ 满足公式(1), $p(\text{bi-gram})_{\text{malw}}$ 表示恶意软件特征, $p(\text{bi-gram})_{\text{benign}}$ 表示良性代码特征。卷积神经网络采用损失函数为最小平方误差, 如公式(2)所示, 通过训练一组权重 $W = \{w_1, w_2, \dots, w_n\}$, 使得损失函数最小。

$$W \cdot p(\text{bi-gram})_{\text{malw}} + b \geq 1 \text{ 或 } W \cdot p(\text{bi-gram})_{\text{benign}} + b \leq -1, \quad (1)$$

$$\text{Loss} = (h(W \cdot p(\text{bi-gram})) - y)^2. \quad (2)$$

2.5 模型集成

基于上述6种恶意软件检测模型, 基于 bagging 策略, 设计一种 Logistic 集成学习方法集成6个模型。令6个模型的判别结果(输出的置信度向量)为 $\mathbf{X} = \{x_1, x_2, \dots, x_6\}$, 并为每个模型分别分配一个随机初始化的决策权重 $W' = \{w_1', w_2', \dots, w_6'\}$, 根据样本标签 y , 基于 Logistic 函数进行有监督训练, 如公式(3)所示, 自适应优化每个模型的决策权重。每个模型输出的置信度作为 Logistic 函数的输入, Logistic 函数的权重采用随机初始化, 利用样本标签进行有监督训练, 经过若干次迭代后决策权重收敛, 此时 Logistic 函数中的权重即为每个模型的决策权重。当决策权重收敛时, 6个模型通过该 logistic 函数集成为统一的判别模型。其中损失函数采用交叉熵, 如公式(4)所示, 采用梯度下降法进行训练, 权重 W' 根据公式(5)进行迭代更新。

$$h(W' \cdot \mathbf{X}) = \text{logistic}(W' \cdot \mathbf{X}) = \frac{1}{1 + e^{-\sum w_i' \cdot x_i}}, \quad (3)$$

$$\text{Loss} = -(y \cdot h(W' \cdot \mathbf{X}) + (1 - y) \cdot (1 - h(W' \cdot \mathbf{X}))), \quad (4)$$

$$\Delta W' = \alpha \cdot \frac{\partial \text{Loss}}{\partial h(W' \cdot \mathbf{X})} \cdot \frac{\partial h(W' \cdot \mathbf{X})}{\partial W'} = \alpha \cdot (y - h(W' \cdot \mathbf{X})) \cdot \mathbf{X}. \quad (5)$$

2.6 模型检测

基于上述统一的判别模型对未知样本进行检测, 通过提取未知样本的操作码 bi-gram 特征输入至模型中, 模型对未知样本的 bi-gram 特征进行计算得到恶意/良性分类的置信度, 如果恶意类别的置信度大于良性类别的置信度, 则说明未知样本属于恶意样本, 否则属于良性样本。

3 实验分析

3.1 实验环境

所有实验均在相同计算机软硬件环境和配置中进行, 其中 CPU 为 Intel i5-3470 @ 3.20 GHz, 内存容量为 16.0 GB, 操作系统为 Ubuntu 16.04, 编译器采用 Eclipse 3.5/JRE 1.8。

3.2 数据集

实验所采用的恶意软件数据集来源于 VxHeaven、Microsoft 和企业云安全中心, 其中 VxHeaven 包括 17 000 余个恶意样本, Microsoft 包括 5 000 余个恶意样本, 企业云安全中心收集 4 000 余个恶意样本, 包含 6 类样本家族, 其数量分布如表 1 所示。训练阶段和测试阶段均采用了恶意样本和良性样本, 其中 3 个恶意样本数据集中的 50% 样本用于分别训练单一检测模型, 合并 3 个恶意样本数据集中剩下 50% 样本形成新的数据集用于模型测试。良性样本从若干用户计算机采集, 约 11 000 余个, 其中随机选取 50% 的良性样本用于训练单一检测模型, 剩余 50% 的良性样本用于模型测试。

表 1 恶意软件数据集

Table 1

恶意软件家族	数量
Virus	1 049
Worm	2 922
Trojan	17 911
Backdoor	4 835
DoS	353
Flooder	200
共计	27 270

3.3 验证方法

所有实验均采用 K -折交验证法,其中 $K=10$ 。评估指标包括精度 PRE、召回率 REC 和 F_1 值。精度表示模型判别结果中正确结果的比例,如公式(6)所示,召回率表示恶意样本中被检测出来的比例,如公式(7)所示, F_1 值表示精度和召回率重要程度相同时的准确性,如公式(8)所示。其中, TP (true positive)表示验证集中被模型正确识别的高风险网络节点数量, FP (false positive)表示验证集中被模型错误识别的高风险网络节点数量, TN (true negative)表示验证集中被模型正确识别的低风险网络节点数量, FN (false negative)表示验证集中被模型错误识别的低风险网络节点数量。根据 TP 、 FP 、 TN 、 FN ,计算准确度精度 (PRE)、召回率(REC)和 F_1 值。

$$PRE = \frac{TP}{TP + FN}, \quad (6)$$

$$REC = \frac{TP}{TP + FP}, \quad (7)$$

$$F_1 = \frac{2 \cdot PRE \cdot REC}{(PRE + REC)}. \quad (8)$$

3.4 性能评估

通过测试集样本对比分析多个单一检测模型以及集成模型的精度、召回率和 F_1 值,从而验证集成模型相较于单一模型在准确性方面有显著优势,实验结果如图 2 所示,结果表明:

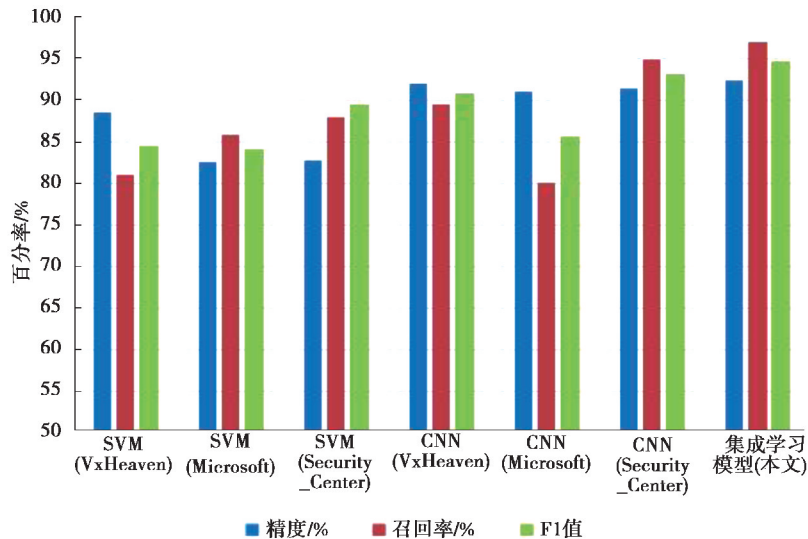


图 2 准确性对比实验结果

Fig. 2 The accuracy comparison results

1)对比分析采用相同机器学习算法、不同训练集的单一检测模型,单一检测模型的准确性受有限训练样本的影响较大,因此在准确性方面,利用不同数据集训练的模型准确性存在一定差异。

2)对比分析采用相同训练集、不同机器学习算法的单一检测模型,卷积神经网络的拟合性能相较于支持向量机的拟合性能更好。

3)对比分析集成模型和单一检测模型,集成模型的精度有一定提升,召回率呈现显著提升,说明集成模型的鲁棒性明显优于单一检测模型。

对比分析集成模型和单一检测模型的检测时间开销,如图 3 所示,结果表明:集成模型的检测时间略高于单一卷积神经网络模型,虽然单一支持向量机模型的检测时间开销相较单一卷积神经网络模型和集成模型,但在精度和召回率等指标方面存在明显不足。

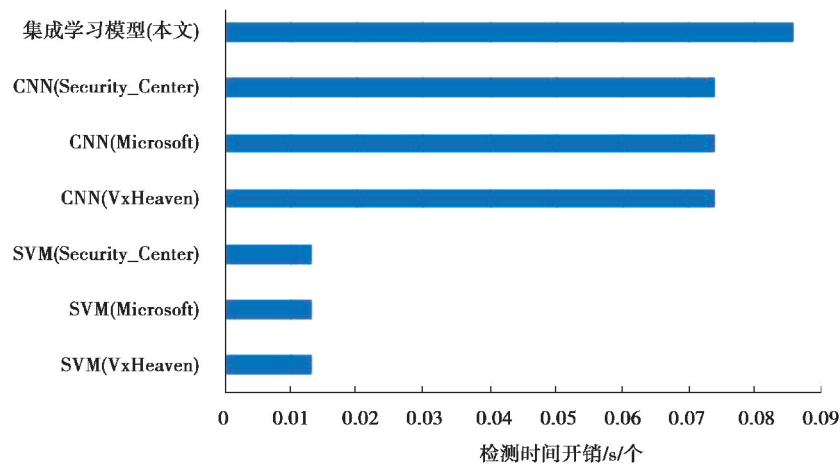


图 3 检测时间对比实验结果

Fig. 3 The time cost comparison results

4 总 结

提出一种基于集成学习的恶意软件未知变种检测方法。该方法实现多种恶意软件未知变种检测模型的集成,以进一步提升恶意软件未知变种的检测性能。基于集成学习的恶意软件未知变种检测方法可以综合多个检测模型的性能,覆盖更多的样本数据特征,以改善单一检测模型样本特征和收敛性能的局限性。最后对提出的模型集成方法进行了实验分析,结果表明方法相比于单一检测模型,集成模型可以达到更高的精度、召回率和 F_1 值。

参考文献:

- [1] 张东霞, 苗新, 刘丽平, 等. 智能电网大数据技术发展研究[J]. 中国电机工程学报, 2015, 35(01), 35: 2-12.
Zhang D X, Miao X, Liu L P, et al. Research on development strategy for smart grid big data[J]. Proceedings of the CSEE, 2015, 35(01), 35: 2-12. (in Chinese)
- [2] 李兴源, 魏巍, 王渝红, 等. 坚强智能电网发展技术的研究[J]. 电力系统保护与控制, 2009, 37(17), 37: 1-7.
Li X Y, Wei W, Wang Y H, et al. Study on the development and technology of strong smart grid[J]. Power System Protection and Control, 2009, 37(17), 37: 1-7. (in Chinese)
- [3] Cortes C, Vapnik V. Support-vector networks[J]. Machine Learning, 1995, 20(3): 273-297.
- [4] Krizhevsky A, Sutskever I, Hinton G E. ImageNet classification with deep convolutional neural networks [J]. Communications of the ACM, 2017, 60(6): 84-90.
- [5] Hinton G E, Salakhutdinov R R. Reducing the dimensionality of data with neural networks[J]. Science, 2006, 313(5786): 504-507.
- [6] McLaughlin N, Martinez del Rincon J, Kang B, et al. Deep android malware detection[C]// Proceedings of the Seventh ACM on Conference on Data and Application Security and Privacy. Scottsdale Arizona USA. New York, NY, USA: ACM, 2017: 301-308.
- [7] Zhang J X, Qin Z, Yin H, et al. A feature-hybrid malware variants detection using CNN based opcode embedding and BPNN based API embedding[J]. Computers & Security, 2019, 84: 376-392.
- [8] Zhang J X, Qin Z, Yin H, et al. Malware variant detection using opcode image recognition with small training sets[C]// 2016 25th International Conference on Computer Communication and Networks (ICCCN). August 1-4, 2016, Waikoloa, HI, USA. IEEE, 2016: 1-9.
- [9] Yan J P, Qi Y, Rao Q F. Detecting malware with an ensemble method based on deep neural network[J]. Security and

- Communication Networks, 2018, 2018: 1-16.
- [10] Cesare S, Xiang Y, Zhou W L. Control flow-based malware variant detection[J]. IEEE Transactions on Dependable and Secure Computing, 2014, 11(4): 307-317.
- [11] Saxe J, Berlin K. Deep neural network based malware detection using two dimensional binary program features[C]//2015 10th International Conference on Malicious and Unwanted Software (MALWARE). October 20-22, 2015, Fajardo, Puerto Rico. IEEE, 2015: 11-20.
- [12] 陈志锋, 李清宝, 张平, 等. 基于数据特征的内核恶意软件检测[J]. 软件学报, 2016, 27(12), 27: 3172-3191.
Chen Z F, Li Q B, Zhang P, et al. Data characteristics-based kernel malware detection[J]. Journal of Software, 2016, 27(12), 27: 3172-3191. (in Chinese)
- [13] Huang J J, Zhang X Y, Tan L, et al. AsDroid: detecting stealthy behaviors in android applications by user interface and program behavior contradiction[C]// Proceedings of the 36th International Conference on Software Engineering. [S.L.]: IEEE, 2014: 1036-1046.
- [14] Bai H, Hu C Z, Wang X Y, et al. Approach for malware identification using dynamic behaviour and outcome triggering [J]. IET Information Security, 2014, 8(2): 140-151.
- [15] Freund Y, Schapire R E. A decision-theoretic generalization of on-line learning and an application to boosting[J]. Journal of Computer and System Sciences, 1997, 55(1): 119-139.
- [16] Zhou Z H, Feng J. Deep forest: towards an alternative to deep neural networks[C]// Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence. August 19-26, 2017. Melbourne, Australia. California: International Joint Conferences on Artificial Intelligence Organization, 2017: 3553-3559.

(编辑 侯 湘)