

# 基于多因素认证的零信任网络构建

孙 瑞<sup>1</sup>, 张 正<sup>2\*</sup>

(1. 南京邮电大学计算机学院, 江苏 南京 210023; 2. 金陵科技学院网络与安全学院, 江苏 南京 211169)

**摘 要:**随着云计算技术的发展,传统基于边界的安全策略主要面临着两大问题:安全边界被突破后攻击者在系统内畅通无阻;攻击可能来自内部可信用户。零信任网络架构技术通过完全去除这种信任假设解决了边界为中心带来的安全隐患。其核心是最小权限访问控制与控制策略的动态更新,这就需要系统对用户可访问资源进行严格的控制并对用户行为保持持续关注。根据大数据与机器学习算法,设计了一种基于多因素身份认证的零信任网络构建模型,实现了基于用户动态行为的访问控制策略。实验结果表明,该模型能够通过实时分析用户行为模式进行身份认证,并剥夺异常用户的访问权限,实现零信任网络架构。

**关键词:**零信任;身份认证;访问控制;多因素分析

中图分类号:TP393

文献标识码:A

文章编号:1672-755X(2020)01-0021-06

## Building Zero Trust Network Based on Multi-factor Authentication

SUN Rui<sup>1</sup>, ZHANG Zheng<sup>2\*</sup>

(1. Nanjing University of Posts and Telecommunications, Nanjing 210023, China;

2. Jinling Institute of Technology, Nanjing 211169, China)

**Abstract:** With the development of cloud computing technology, the traditional boundary-based security policy is mainly faced with two problems: The attacker is unobstructed in the system after the security boundary is broken, and the attack may come from the internal trusted user. Zero trust network architecture technology solves the security risk of boundary-centric by completely removing this trust hypothesis. The core of zero trust network technology is the minimum privilege access control and dynamic update of control strategy, which requires the system to strictly control the user's accessible resources and maintain continuous attention to user behavior. Based on the big data and machine learning algorithm, this paper designs a zero trust network model based on multi-factor authentication, and realizes the access control strategy based on the user's dynamic behavior. Experimental results show that the model can authenticate the user behavior pattern through real-time analysis, deprive the abnormal user of access rights, and realize the zero-trust network architecture.

**Key words:** zero trust; authentication; access control; multi-factor analysis

随着云计算的兴起,越来越多的企业将数据与应用部署在云端,与此同时,以内外网为划分的安全边界变得模糊。将企业应用数据置于虚拟专用网中,并通过防火墙进行保护的模式现已难以为继。采用边界为中心的安全策略所依赖的是内部网络上一切都是可以信任的,然而这种假设已经不再是安全的。在

收稿日期:2020-01-04

基金项目:国家重点研发计划网络空间安全重点专项(2017YFB0802800,2017YFB0802802)

作者简介:孙瑞(1997—),男,江苏沐阳人,硕士研究生,主要从事信息安全研究。

通信作者:张正(1973—),男,江苏宝应人,研究员,主要从事信息安全研究。

内外网技术中,多采用口令认证,然而根据调查显示<sup>[1]</sup>,超过 30%的用户在不同系统中采用同一口令,一旦攻击者盗取用户 VPN 账户密码<sup>[2]</sup>,便可长驱直入获取内网全部信息。这一安全隐患对访问权限较高的账户影响尤为明显。

针对上述问题,Google 公司的工程师 John Kindervag 于 2010 年提出了零信任网络<sup>[3]</sup>。零信任网络是一种新型安全概念,其核心原则为“永远不要相信”和“始终验证”<sup>[4-5]</sup>。本文针对这两项原则,提出一种基于多因素身份认证的零信任网络模型,对用户日常使用行为进行抽象建模,实时分析用户行为,检测可疑账户并剥夺其访问权限。为了实现灵活的认证授权机制,本文还利用软件定义网络(Software Defined Network,SDN)<sup>[6]</sup>实现访问控制策略动态更新,降低了系统部署的难度。

## 1 相关工作

### 1.1 传统网络向零信任网络的转变

通常企业通过部署防火墙、入侵检测、漏洞扫描构建内网安全防护体系。谷歌公司在 2009 年经历高度复杂的 APT(极光行动,Operation Aurora)攻击后,开始尝试重新设计员工与设备访问内部应用的安全架构<sup>[7]</sup>,零信任架构 Beyond Corp 应运而生。与传统的边界安全模式不同,零信任网络将所有应用部署到公网上,用户通过认证与授权进行访问<sup>[2]</sup>。其主要设计理念包括:假定所有网络设备都是不可信的,因为安全威胁同时可能来自于外部网络或内部网络;需要基于受控设备和合法用户进行资源访问控制;任何对服务的访问都必须进行身份验证、授权和加密。由此,员工可以实现在任何地点的安全访问,无需传统的 VPN。图 1 为传统网络结构向零信任网络的转变。

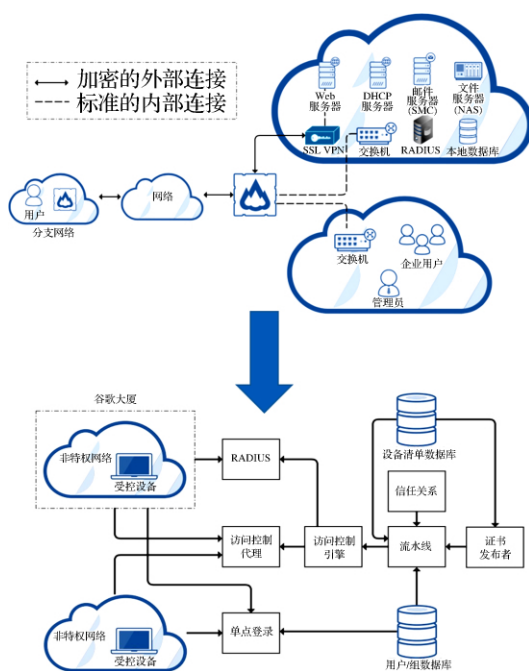


图 1 传统网络结构向零信任网络的转变

### 1.2 基于用户行为的认证技术

身份认证是指在计算机以及计算机网络系统中确认操作者身份的过程,从而确定该用户是否具有对某种资源的访问和使用权限,进而使计算机和网络系统的访问策略能够可靠、有效的执行。传统的认证方式有静态密码、智能卡、短信密码、动态口令、生物识别等<sup>[8-9]</sup>,以及上述方法组成的双因素甚至多因素认证方法。由于传统的身份认证方式为一次性认证,无法持续保障系统安全,此时大数据和机器学习技术通过对用户行为的动态分析,为身份认证提供了新的思路<sup>[4,6-7]</sup>。利用大数据和机器学习技术可以具体解决的问题有:分类、聚类、估计、预测、相关性分组或关联规则等。实验证明通过用户行为建模,机器学习算法可高效、准确地判断当前用户行为是否可疑,为动态访问控制提供强有力的依据。

## 2 身份认证方法

### 2.1 认证流程

根据用户行为建模进行身份认证的流程如图 2 所示。

我们提出基于多因素的用户行为认证算法,对用户在日常系统中的日常行为进行抽象建模。由于恶意用户需要完成篡改、窃取等攻击

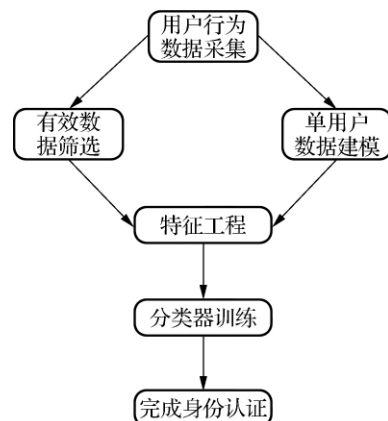


图 2 认证流程图

行为,他们在行为习惯上与合法用户有着较大区别,因此在数据特征上也与合法用户有着一定区别。利用人工智能技术从冗杂的数据中抽取有效的特征,并利用这些特征对用户画像进行建模。随后根据分类技术对恶意用户以及被盗用户进行动态身份识别,从而能够有效保护用户数据的安全。

### 2.2 用户画像建模

用户画像是对基于大量带标签的用户行为数据进行特征提取和分类器训练后得到的,目的是提取用户的历史行为特征,来判断用户当前行为是否异常。但目前用户行为记录不仅数量多,且行为特征多样,难以处理。因此,我们将用户行为特征共分为 6 个方面,如表 1 所示:

表 1 用户行为特征

项目	具体行为
登录方式	登录方式、登录时间戳
登录时长	使用总时长、活跃时段
登录设备	设备名称、设备编号、设备个数
登录 IP	IP 切换与操作内容、时间的对应关系
操作对象	访问资源名
键入信息	键入文字

根据数据集的不同,具体行为还可进行扩充或删减。值得注意的是,虽然进行了分类精简,六种用户行为特征仍包含不同的类型和数据结构。在进行特征工程之前需对数据进行预处理。这里将六种特征分为三种数据结构:1)类别型。该类特征包含登录方式、登录设备、登录 IP 和操作对象。这类数据以 string 方式进行存储,每条记录对应某种类别,本文采用 one-hot 编码,将不同特征映射至矩阵空间中,每一种特征对应唯一向量。2)数值型。这类特征包含登录时长、登录时间戳。这类数据的值具有具体的意义,且为连续分布。若以具体数值来进行建模,则特征工程会面临维度爆炸的问题。本文中登录时长按小时进行划分;登录时间戳以星期+小时进行划分,以达到降维的目的。3)文本型。这类特征主要为键入文字。本文对文本信息采用词袋模型进行处理。

### 2.3 分类器训练

在机器学习中有众多分类方法,如决策树分类、朴素贝叶斯分类、支持向量机、神经网络等。本文选择使用 XGBoost(eXtreme Gradient Boosting)<sup>[10]</sup>算法来实现非法用户的分类功能。XGBoost 是一种 tree boosting 可扩展机器学习系统,这个系统可以作为开源的软件包使用。XGBoost 的优异表现<sup>[11-12]</sup>使其在大量的机器学习和数据挖掘挑战中被认可,并被广泛地应用到各领域。

XGBoost 的基本思想是把成百上千个分类准确率较低的树模型组合起来,成为一个准确率很高的模型。图 3 为特征树构建方式。这个模型不断迭代,每次迭代都会生成一棵新的树。XGBoost 算法源起于 Boosting 集成学习方法,在演化过程中又融入了 Bagging 集成学习方法的优势,通过 Gradient Boosting 框架自定义损失函数提高了算法解决通用问题的能力,同时引入更多可控参数可针对问题场景进行优化。总体算法如算法 1 所示。不同于其他 boosting 方法,XGBoost 的特点是计算速度快,模型表现好。

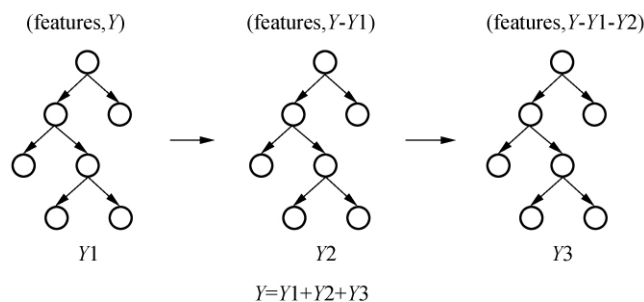


图 3 特征树构建

### 算法 1 树划分贪心算法

输入: I(当前节点)

输入: d(特征维数)

gain  $\leftarrow$  0

$G \leftarrow \sum_{i \in I} g_i, H \leftarrow \sum_{i \in I} h_i$

FOR  $k = 0$  TO  $m$  DO

$G_L \leftarrow 0, H_L \leftarrow 0$

FOR  $j$  in *sorted*(I, by  $x_{jk}$ ) DO

$G_L \leftarrow G_L + g_j, H_L \leftarrow H_L + h_j$

$G_R \leftarrow G - G_L, H_R \leftarrow H - H_L$

$Gscore \leftarrow \max(score, \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{G^2}{H + \lambda})$

END

END

输出: 最大分值

## 2.4 实验结果与分析

2.4.1 数据集概述 实验所采用的数据集是京东提供的用户登录数据。其中选取了共 532 943 条有效数据, 每条数据包括登录 IP、登录时间戳、登录方式等 7 个字段。按照 2.2 节分为类别型、文本型及数值型。在输入 XGBoost 模型前已根据 2.2 节中的描述进行了预处理。本文中 70% 的数据作为训练数据, 10% 的作为验证数据, 20% 的作为测试数据。

2.4.2 模型构建 1) 特征构建。在 2.2 节中对特征进行的分类仅是针对预处理所做出的数据结构上的分类, 特征构建过程中仍需进行新的分类。在基于多因素的身份认证环境下, 用户画像的刻画不仅会受到单一特征的影响, 还有可能受到某几个特征组合的影响。因此, 本文运用了交叉特征方法, 对 7 个字段的特征(表 2)进行特征组合和特征选择, 共生成了 257 种特征组合。

2) 模型训练与预测。首先利用验证集数据进行模型训练, 将模型参数调整到最优。本文因篇幅限制, 直接给出该数据集上 XGBoost 模型的最优参数(表 3)。随后完成参数设置, 利用测试集对模型进行训练, 模型输出为 0~1 的数值, 根据数据集中的标签信息, 训练最佳阈值。本实验中得到的最佳阈值为 0.8, 高于 0.8 则认为该用户行为合法, 反之则行为异常。最后根据训练好的模型计算测试集上的结果, 并与标签进行比较。

表 2 单条用户记录对应字段

特征名	数据结构
登录 IP(I)	String
登录时长(L)	int
登录时间戳(S)	String
登录设备(E)	String
登录方式(W)	String
键入信息(T)	String
操作对象(O)	String

表 3 最优参数

参数	取值
max_depth	6
gamma	0.1
subsample	0.8
colsample_bytree	0.8
seed	32
nthread	10

2.4.3 实验结果分析 首先对特征的重要程度进行分析。本实验中一共包含 257 种特征组合, 逐一进行试验验证的计算开销较大, 因此利用 XGBoost 的随机森林分类器对所有特征进行权重分析, 表 4 列出了权重值在前 10 的特征组合。

表 4 中特征组合的特征名如表 2 所示,可以看到 IP 的影响权重最大,权重前 10 的特征组合中有 7 种组合包括 IP 信息,这说明在该数据集中,异常用户行为最有可能出现的特征是异常 IP 登录,这符合传统基于防火墙的防御理念;此外登录 IP、登录时间戳、设备的组合对异常用户预测的准确率最高,这符合零信任网络中“何人、何时、何地在哪种设备上登录”的要求。

随后对模型的预测准确率进行分析,测试集中共有 10 658 条记录,模型对异常用户的判定准确率为 82%,召回率为 79%。说明该模型基本能完成基于多因素的用户身份认证任务。

表 4 权重值前 10 的特征组合权重

特征组合	权重值
I+S+E	0.97
I+S	0.95
E+W	0.92
I+W	0.88
W+O+S	0.87
I+S+O	0.86
I+E+O	0.84
I+W+O	0.83
I+L+O	0.83
S+E+O	0.82

### 3 基于多因素身份认证的零信任网络

#### 3.1 网络的构建

目前有多种技术可以用来实现零信任网络架构,本文采用第 2 节中的身份认证模型实现认证与授权模块,并通过 SDN 实现中心化的控制策略和动态规则部署。SDN 作为一种新型的网络架构,实现了网络虚拟化,并通过 OpenFlow 技术将网络设备控制面与数据面分离开来,可以实现零信任网络中基于清单列表的访问控制。清单列表由多因素身份认证算法生成且动态更新,一旦发现已授权的用户有异常行为,立刻通过 SDN 控制器剥夺其所有访问权限,且网络中其他用户也不可访问该异常用户,以防止网络中存在一个以上的攻击者。此外,用户在没有访问权限时无法获得目标的任何信息,以防攻击者通过反馈信息窃取信息,该功能可通过 SDN 控制器完成包丢弃操作。

由于该模型基于多种因素完成访问策略的动态更新,因此为每个用户设计一个 32 位的 token,存放用户登录时的 IP、时间戳、设备和方式。该 token 被放在 TCP 报头中以完成后续的认证工作。SDN 控制器接收到该 token 后会传输到身份认证模块,记录用户访问时长、访问目标、键入信息等,并生成行为记录,通过多因素身份认证算法判断用户行为是否异常。一旦身份认证模块认为当前用户行为异常,则 SDN 控制器会拒绝该请求并进行丢包操作,用户将不会收到来自该网络的任何信息。该模式有着多种优势,首先它完成了不基于位置的零信任网络访问策略,可以避免来自内部的攻击;其次可运用于各种网络拓扑结构中,无论该网络使用了多少交换机连接设备,所有交换机都受到 SDN 的控制,因此可以通过软件编写信任列表,或设计动态访问策略以完成整个网络的访问控制;最后,该模型还具备 SDN 易于部署的优势。

#### 3.2 实验结果分析

3.2.1 实验环境与设置 本实验使用的平台为 ubuntu16.04;网络仿真平台为 Mininet;身份认证模块使用的环境为 Python3.6。

Mininet 作为 SDN 网络系统中的一种基于进程虚拟化的平台,它支持 OpenFlow、OpenvSwitch 等各种协议,可在一台计算机上模拟完整的网络<sup>[11]</sup>。我们利用 Mininet 模拟网络拓扑结构,验证在零信任网络架构下的认证与授权过程。实验利用 Mininet 构建的虚拟用户并赋予用户行为数据,判断系统能否在用户出现异常行为时剥夺用户访问权限。具体操作为将用户行为加入到 token 中输入网络,如果用户行为被标记为异常行为,检测该节点能否获取信息。

3.2.2 实验结果 实验中的网络拓扑结构如图 4 所示。利用前文提到的用户访问数据,逐一赋给图中主机,主机在接收到访问行为数据后,将静态信息加入 token 中输入给 SDN 交换机,SDN 交换机传输给控制器 c1,随后控

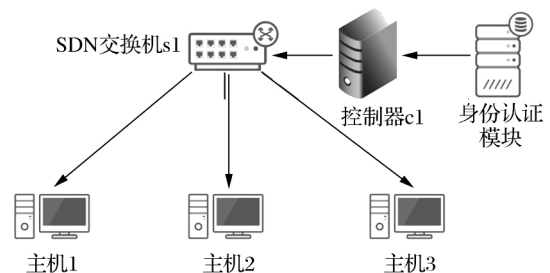


图 4 网络拓扑结构

制器将该信息传输给身份认证模块,身份认证模块将判断结果返回给控制器 c1,由控制器 c1 完成本次行为的授权或丢包处理。此处利用主机 2 输入异常用户行为判断控制器能否做出正确操作,输出结果如图 5 所示,主机 2 无法获取来自控制器的任何信息。

随后对多主机同时访问进行测试,添加 9 台虚拟主机至虚拟网络中,设置主机 3 与主机 4 注入异常用户行为数据,其余主机输入正常用户行为数据。实验结果表明,主机 3 与主机 4 无法完成访问而其余主机均可正常通信。

最后对全体数据进行测试,总计有 532 943 条用户行为数据。实验结果(表 5)表明,控制器可以基于身份认证模块的预测结果对访问进行控制,符合模型设计的初衷。

```
root@ubuntu:~# ping 10.0.0.2
PING 10.0.0.2 (10.0.0.2) 56(84) bytes of data:
From 10.0.0.1 icmp_seq=9 Destination Host Unreachable
From 10.0.0.1 icmp_seq=10 Destination Host Unreachable
From 10.0.0.1 icmp_seq=11 Destination Host Unreachable
```

图 5 输出结果

表 5 访问测试结果

输入条目	控制器判断	
	异常	正常
异常	53 295	0
正常	0	478649

## 4 结 语

在云服务越来越普及的今天,传统基于边界的网络防护模式无法满足对敏感数据的保护任务,本文在分析了当前零信任网络以及基于大数据的身份认证方法的研究现状后,提出了一种基于多因素身份认证的零信任网络构建方法。该方法对用户行为数据进行动态建模,通过特征过程及 XGBoost 算法完成对异常用户行为的判断,将判断结果作为零信任网络的动态访问策略,符合零信任网络中的“始终验证”原则,提高了系统的安全性。本文方法还有以下两个问题需要考虑:如何设置分类器阈值,使得识别准确率和用户体验之间达到平衡;实验中网络构建方式较为理想化,如何设计使方案更贴近实际应用需求。

### 参考文献:

- [1] Zhou C L, Lin Z Y. Study on fraud detection of telecom industry based on rough set[C]. As Vegas: Computing and Communication Workshop and Conference(CCWC), 2018
- [2] Abdullateef R, Mahmoud A A, Yaser J, et al. Authorship attribution of Arabic tweets[C]. Agadir: Computer Systems and Applications(AICCSA)2016 IEEE/ACS 13th International Conference, 2016: 1-6
- [3] Barckay O, Justin M, Betsy B, et al. BeyondCorp design to deployment at Google[EB/OL]. (2016-04-21)[2019-10-12]. <https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/44860.pdf>
- [4] 薛朝晖, 向敏. 零信任安全模型下的数据中心安全防护研究[J]. 通信技术, 2017, 50(6): 1290-1294
- [5] Casimer D C. Implementing zero trust cloud networks with transport access control and first packet authentication[C]. New York: The IEEE International Conference on Smart Cloud, 2016: 110-109
- [6] 左青云, 陈鸣, 赵广松, 等. 基于 OpenFlow 的 SDN 技术研究[J]. 软件学报, 2013, 24(5): 1078-1097
- [7] Rory W, Betsy B. BeyondCorp: A new approach to enterprise security[J]. Login, 2014, 39(6): 6-11
- [8] Murphy K P. Machine learning: a probabilistic perspective[J]. Chance, 2012, 27(2): 62-63
- [9] He H, Garcia E A. Learning from imbalanced data[J]. IEEE Transactions on Knowledge & Data Engineering, 2009, 21(9): 1263-1284
- [10] Fitriah N, Wijaya S K, Fanany M I, et al. EEG channels reduction using PCA to increase XGBoost's accuracy for stroke detection[C]. Depok: International Symposium on Current Progress in Mathematics and Sciences, 2016
- [11] Schapire R E, Singer Y. BoosTexter: A boosting-based system for text categorization[J]. Machine Learning, 2000, 39(2/3): 135-168
- [12] Chen T, Guestrin C. XGBoost: A scalable tree boosting system[C]. San Francisco: SIGKDD International Conference on Knowledge Discovery and Data Mining ACM, 2016: 785-794

(责任编辑: 湛 江)