

# Study on the unsupervised APT attack detection method based on adaptive fuzzy clustering

Jinhu Huang

Fujian CITIC network security information Technology Co., Ltd., Fuzhou, Fujian, 350000, China

## Abstract

In order to solve the problem of concealment and complexity of APT (advanced persistent threat) attacks, this paper proposes an adaptive fuzzy clustering algorithm, which can dynamically adjust the clustering parameters according to the attack characteristics to adapt to different types and strengths of APT attacks. The method is applied to the APT attack detection scenario, and the network traffic data through unsupervised learning. The results show that the proposed method can identify APT attacks with higher detection rate and lower false alarm rate. The unsupervised APT attack detection method based on adaptive fuzzy clustering can effectively improve the detection performance of APT attack, and provides new ideas and means for network security protection.

## Keywords

adaptive fuzzy clustering; unsupervised APT attack detection method; application

# 基于自适应模糊聚类的无监督 APT 攻击检测方法研究

黄金虎

福建中信网安信息科技有限公司, 中国·福建 福州 350000

## 摘要

为解决APT(高级持续性威胁)攻击的隐蔽性和复杂性问题,本文提出了一种自适应模糊聚类算法,该算法能够根据攻击特征动态调整聚类参数,以适应不同类型和强度的APT攻击。将该方法应用于APT攻击检测场景,通过无监督学习方式对网络流量数据进行聚类分析,识别出异常行为模式。结果表明,与传统的无监督攻击检测方法相比,本文提出的方法能够更有效地识别出APT攻击,具有较高的检测率和较低的误报率。本文提出的基于自适应模糊聚类的无监督APT攻击检测方法能够有效提高APT攻击的检测性能,为网络安全防护提供了新的思路 and 手段。

## 关键词

自适应模糊聚类; 无监督APT攻击检测方法; 应用

## 1 引言

随着互联网技术的飞速发展,网络安全问题日益突出。APT(Advanced Persistent Threat, 高级持续性威胁)攻击作为一种隐蔽性强、持续时间长、攻击目标明确的网络攻击手段,给我国网络安全带来了严重威胁。APT攻击往往针对特定目标,通过长期潜伏、窃取敏感信息等手段,对国家安全、经济和社会稳定造成极大危害。本文提出了一种基于自适应模糊聚类的无监督APT攻击检测方法。该方法通过自适应模糊聚类算法对网络流量进行聚类,提取出潜在的安全威胁特征,从而实现对APT攻击的检测。

## 2 基于自适应模糊聚类的无监督 APT 攻击检测方法

### 2.1 自适应模糊聚类算法原理

#### 2.1.1 模糊聚类的概念

模糊聚类是一种将数据集中的对象按照一定的相似度关系划分成若干个类别的无监督学习方法。与传统的硬聚类方法不同,模糊聚类允许一个对象可以同时属于多个类别,其成员度表示对象对某个类别的隶属程度<sup>[1]</sup>。这种灵活性使得模糊聚类在处理复杂、不确定的领域问题时具有独特的优势。

#### 2.1.2 自适应机制的实现

自适应模糊聚类算法的核心思想是在聚类过程中根据聚类结果动态调整聚类参数,从而提高聚类效果。常见的自适应机制实现方法主要包括以下四类:(1)动态调整聚类数目:根据聚类过程中的聚类中心变化,实时调整聚类数目,使聚类结果更加合理。(2)动态调整聚类半径:根据聚类

【作者简介】黄金虎(1977-),男,中国福建南安人,硕士,工程师,从事数据安全,网络安全研究。

中心的距离变化,动态调整聚类半径,使聚类结果更加紧凑。(3)动态调整隶属度权重:根据聚类过程中对象的分布情况,动态调整隶属度权重,使聚类结果更加均衡。(4)自适应调整聚类算法参数:根据聚类结果和目标数据集的特性,自适应调整聚类算法的参数,如迭代次数、学习率等。

## 2.2 在 APT 攻击检测中的应用

### 2.2.1 数据特征提取

在 APT 攻击检测中,首先需要收集到的网络流量、系统日志、用户行为等数据进行特征提取。数据特征提取是 APT 攻击检测的基础,其目的是从大量原始数据中提取出具有代表性的特征,以便后续的聚类分析和攻击检测。主要特征提取方法包括统计特征、频率特征、时序特征、异常检测特征等。统计特征包括平均连接时间、数据包大小、数据包到达速率<sup>[2]</sup>。频率特征包括 IP 地址、端口、协议类型、域名等频繁出现的特征。时序特征包括时间序列分析、滑动窗口等。异常检测特征包括基于机器学习的异常检测方法,如孤立森林、支持向量机等。

### 2.2.2 聚类分析过程

聚类分析是 APT 攻击检测的关键步骤,通过将具有相似性的数据点归为一类,可以揭示 APT 攻击的潜在模式。基于自适应模糊聚类的无监督 APT 攻击检测方法,采用以下步骤进行聚类分析:(1)初始化聚类中心:选择一定数量的数据点作为初始聚类中心。(2)计算相似度:根据数据特征,计算每个数据点与聚类中心的相似度。(3)分配数据点:将每个数据点分配到与其最相似的聚类中心所在的类别。(4)更新聚类中心:根据分配到每个类别的数据点,计算新的聚类中心。(5)迭代优化:重复步骤(2)至(4),直到满足终止条件(如聚类中心变化较小或达到最大迭代次数)。

### 2.2.3 攻击检测的判定依据

在完成聚类分析后,若某个类别中的数据点数量远小于其他类别,或具有明显的异常特征,则存在 APT 攻击。通过分析聚类结果,找出具有相似特征的攻击模式,如恶意软件传播、数据窃取等。根据聚类结果,分析攻击者的行为模式,如攻击时间、攻击目标、攻击手段等。挖掘聚类结果中的关联规则,发现潜在的攻击路径和攻击手段。

## 3 应用场景

### 3.1 企业网络安全防护

#### 3.1.1 内部网络监控

企业内部网络监控是 APT 攻击检测的重要手段。通过实时监测网络流量、系统日志等信息,及时发现异常行为,为企业提供安全预警。本文提出的方法可以应用于内部网络监控,提高检测 APT 攻击的准确性和效率。

#### 3.1.2 防止数据泄露

数据泄露是 APT 攻击的常见目的之一。企业通过实施有效的安全策略,加强对数据访问和传输的监控,可以降低

数据泄露风险。本文提出的自适应模糊聚类方法可以帮助企业识别潜在的数据泄露风险,提高数据安全防护水平。

## 3.2 关键基础设施保护

### 3.2.1 能源领域应用

通过自适应模糊聚类方法,对电力系统中的海量数据进行实时分析,识别出异常行为,从而实现对电力系统的实时监控和预警。例如,识别电网中异常的电流、电压等参数,以及设备故障等,保障电力系统的安全稳定运行。自适应模糊聚类方法可以帮助电力系统进行负荷预测、设备状态评估和能源调度优化。通过对历史数据的分析,识别出负荷变化的规律,为电力调度提供依据,提高能源利用效率。在新能源并网过程中,自适应模糊聚类方法可用于检测新能源发电设备的异常状态,确保新能源发电设备的稳定运行,同时保障电网的稳定。

### 3.2.2 交通领域应用

自适应模糊聚类方法可应用于交通监控系统中,对交通流量、车辆行驶状态等数据进行实时分析,识别出异常情况,如拥堵、交通事故等,为交通管理部门提供预警信息。通过分析公共交通系统的运行数据,自适应模糊聚类方法可帮助优化公交线路、调度方案,提高公共交通的运行效率和服务水平。自适应模糊聚类方法在智能交通系统中,可用于车辆识别、道路拥堵分析、交通流量预测等方面,为智能交通系统的建设提供有力支持。

## 3.3 移动终端安全

### 3.3.1 智能手机、平板电脑的安全检测

针对智能手机和平板电脑等移动终端,对移动终端中应用的行为、权限、流量等信息进行聚类分析,可以发现具有潜在威胁的恶意应用。对移动终端的访问记录、网络流量、电池使用情况等进行聚类分析,可以发现异常行为,如数据泄露、非法访问等。对移动终端的安全漏洞、系统配置等进行聚类分析,可以评估安全风险,为安全加固提供依据。

### 3.3.2 防范移动 APT 攻击

针对移动 APT 攻击,通过对移动终端的异常行为、恶意代码、网络流量等进行聚类分析,可以发现攻击特征,提高检测精度。通过对历史攻击数据进行分析,可以预测未来的攻击趋势,为安全防护提供预警。根据聚类分析结果,可以为移动终端制定更加合理的安全策略,提高安全防护效果。

## 4 存在的问题

### 4.1 数据质量和复杂性

#### 4.1.1 数据噪声和缺失值的影响

在基于自适应模糊聚类的无监督 APT 攻击检测方法中,数据质量和完整性是影响检测效果的关键因素。数据噪声和缺失值的存在会对聚类结果产生不利影响,噪声数据会干扰聚类过程,导致聚类结果不精确,进而影响 APT 攻击的检测效果。噪声数据来源于网络传输、传感器采集等多种途径,

使得聚类算法难以准确识别攻击模式<sup>[3]</sup>。缺失值的存在会导致聚类过程中某些样本特征信息不完整,影响聚类结果的可靠性。在APT攻击检测中,缺失值导致攻击模式识别不准确,降低检测效果。

#### 4.1.2 处理大规模数据的挑战

随着网络攻击的日益复杂,APT攻击数据量呈现爆炸式增长,大规模数据处理需要更多的计算资源,包括CPU、内存等。对于资源有限的系统,大规模数据处理成为瓶颈,影响检测效果。随着数据量的增加,聚类算法的性能会逐渐下降,导致检测效果降低。此外,算法的复杂度也会随着数据量的增加而增加,使得算法在实际应用中难以有效运行。在处理大规模数据时,如何从海量的特征中选择对APT攻击检测最有价值的特征成为一大难题。特征选择不当导致聚类效果不佳,进而影响APT攻击检测的准确性。

### 4.2 聚类算法的局限性

#### 4.2.1 对初始参数的敏感性

聚类算法在执行过程中需要设定一系列参数,如聚类数量、隶属度阈值等。这些参数的选择对聚类结果有着显著影响。然而,在实际应用中,由于缺乏对攻击数据的深入理解,很难准确确定这些参数的取值。过多的聚类会导致数据过度分割,影响攻击检测的准确性;过少的聚类则导致攻击样本被错误地归入同一类,降低检测效果。隶属度阈值是衡量数据点与聚类中心相似度的标准。阈值设置过高会导致攻击样本难以被正确识别;阈值设置过低则导致正常流量被错误地归为攻击流量,增加误报率。

#### 4.2.2 聚类错误

聚类算法在处理复杂、高维数据时,容易受到噪声、异常值等因素的影响,导致聚类结果出现错误。具体表现为攻击数据中存在噪声,导致聚类结果不准确。攻击数据中存在异常值,这些异常值会扭曲聚类结果,使得攻击样本被错误地归入其他类别。由于攻击数据分布不均匀,聚类边界存在模糊现象,导致攻击样本被错误地归入其他类别。

## 5 策略与建议

### 5.1 数据预处理和优化

#### 5.1.1 数据清洗和增强技术

在数据预处理阶段,对采集到的原始数据进行清洗,去除重复、异常、噪声等不良数据,提高数据质量。为了提高模型的泛化能力,可以对原始数据进行线性变换、非线性变换等,增加数据的变化范围<sup>[4]</sup>。通过复制、旋转、缩放等方法,生成新的数据样本,丰富数据集。

#### 5.1.2 特征选择和提取方法的改进

在特征提取阶段,对原始特征进行筛选,保留对攻击检测具有较高贡献的特征,降低模型复杂度。对筛选后的特征进行提取,提高特征的表达能力。具体方法包括基于距离的方法,如K-means、层次聚类等,通过聚类分析提取特征。

还可以采用基于变换的方法,如主成分分析(PCA)、线性判别分析(LDA)等,通过降维提取特征。

### 5.2 算法改进和优化

#### 5.2.1 结合其他机器学习算法

将自适应模糊聚类与其他机器学习算法(如决策树、支持向量机、神经网络等)结合,构建集成学习模型。通过集成学习模型的优势互补,提高APT攻击检测的准确率和鲁棒性。在原始数据的基础上,通过特征工程提取更具代表性的特征,为自适应模糊聚类提供更好的输入。例如,可以结合时间序列分析、主成分分析等方法进行特征提取。将自适应模糊聚类与其他异常检测算法(如孤立森林、One-Class SVM等)结合,提高对APT攻击的检测能力。通过对比不同算法的检测结果,可以降低误报率。

#### 5.2.2 自适应参数调整策略

根据攻击数据的特征,动态调整聚类中心,使聚类结果更贴近实际情况。例如,可以采用自适应调整聚类中心的策略,如基于距离的聚类中心调整、基于相似度的聚类中心调整等。根据攻击数据的分布,自适应调整隶属度阈值,提高检测的准确性<sup>[5]</sup>。例如,可以采用基于聚类密度、聚类轮廓系数等指标来确定隶属度阈值。根据攻击数据的分布,自适应调整聚类数量,使聚类结果更符合实际情况。例如,可以采用基于信息增益、轮廓系数等指标来确定聚类数量。根据攻击数据的分布,自适应调整模糊系数,提高聚类结果的精确度。例如,可以采用基于聚类轮廓系数、聚类熵等指标来确定模糊系数。

## 6 结论

自适应模糊聚类算法在APT攻击检测中具有较好的应用前景,能够有效提取网络流量中的潜在威胁特征。通过自适应模糊聚类算法对网络流量进行聚类,可以降低特征提取的复杂度,提高检测效率。针对APT攻击检测过程中存在的问题,本文优化聚类算法参数,提高聚类效果。结合多种特征提取方法,提高检测精度。引入异常检测技术,增强对APT攻击的检测能力。本文提出的基于自适应模糊聚类的无监督APT攻击检测方法在理论研究和实际应用中具有较好的前景,有助于提高我国网络安全防护水平。

### 参考文献

- [1] 刘笑笑.基于溯源图长期特征提取的APT攻击检测方法研究[D].北京交通大学,2023.
- [2] 李雪鸥.基于样本特征强化的APT攻击多阶段检测方法研究[D].中国民航大学,2023.
- [3] 谢丽霞,李雪鸥,杨宏宇,等.基于样本特征强化的APT攻击多阶段检测方法[J].通信学报,2022,43(12):66-76.
- [4] 王培森.基于图神经网络表示的APT攻击溯源识别和横向移动检测[D].华中科技大学,2023.
- [5] 林昌建.基于深度学习的APT攻击检测与溯源技术研究与实现[D].国防科技大学,2021.