

Application research of network traffic analysis data based on frequent pattern mining

Daming Meng

Wuzhou Branch, China Mobile Communications Group Guangxi Co., Ltd., Wuzhou, Guangxi, 543000, China

Abstract

With the continuous growth of network traffic data, how to extract effective information from it to improve network management and security has become a hot topic of current research. As an important method in data mining field, frequent pattern mining can be used to reveal hidden patterns and association rules in network traffic. Based on the basic principle of frequent pattern mining, this paper discusses its application and advantages in network traffic analysis. A set of network traffic analysis model based on frequent pattern mining is designed, and the performance and practical application effect of the model are evaluated. The results show that this method has significant advantages in abnormal traffic detection, traffic classification and optimization strategy formulation. In this paper, the potential challenges and future development direction are further analyzed, in order to provide theoretical support and practical guidance for network traffic analysis.

Keywords

frequent pattern mining; Network traffic; Data analysis; Anomaly detection; Association rule

基于频繁模式挖掘的网络流量分析数据应用研究

蒙达明

中国移动通信集团广西有限公司梧州分公司, 中国·广西 梧州 543000

摘要

随着网络流量数据规模的不断增长, 如何从中提取有效信息以提升网络管理和安全性成为当前研究的热点。频繁模式挖掘作为数据挖掘领域的重要方法, 可用于揭示网络流量中的隐藏模式与关联规则。本文从频繁模式挖掘的基本原理出发, 探讨其在网络流量分析中的应用与优势, 设计了一套基于频繁模式挖掘的网络流量分析模型, 并对模型的性能与实际应用效果进行评估。研究表明, 该方法在异常流量检测、流量分类与优化策略制定等方面具有显著优势。本文对其潜在挑战及未来发展方向进行了进一步分析, 以期能为网络流量分析提供理论支持与实践指导。

关键词

频繁模式挖掘; 网络流量; 数据分析; 异常检测; 关联规则

1 引言

随着互联网技术的飞速发展, 网络流量数据量呈指数级增长。这些数据中蕴含了大量有价值的信息, 包括用户行为模式、网络性能状态以及潜在的安全威胁。然而, 由于数据规模庞大、结构复杂和动态性强, 传统的网络流量分析方法在处理效率和精度上存在一定局限性。如何高效挖掘网络流量数据中的潜在模式和规律, 以提升网络管理效率、优化资源配置并提高安全防护水平, 成为当前亟须解决的重要问题。

频繁模式挖掘是数据挖掘中的核心技术之一, 主要用于发现数据集中出现频率较高的模式。其在商业推荐系统、市场分析等领域已得到广泛应用, 但在网络流量分析领域的

研究尚处于初级阶段。通过将频繁模式挖掘引入网络流量分析, 可以有效揭示网络行为的规律性, 为流量优化、异常检测和网络安全管理提供新的视角。

本文旨在基于频繁模式挖掘技术, 探索其在网络流量分析中的应用潜力。具体研究目标包括: 设计基于频繁模式挖掘的网络流量分析模型, 验证其在实际数据集上的性能, 并分析其应用中的潜在挑战与解决思路。研究的最终目的是为网络管理和安全提供一种高效、可扩展的数据分析工具。

2 频繁模式挖掘的基本理论

2.1 频繁模式挖掘的定义与基本方法

频繁模式挖掘是指从数据集中发现频繁出现的项目集、序列或模式的过程。其主要目的是揭示数据中隐藏的关联关系或规律性。例如, 在网络流量数据中, 可以通过挖掘分析高频访问的 IP 地址、常见的协议组合或频繁的用户行为路径, 为网络优化和安全管理提供支持。Apriori 算法是频繁

【作者简介】蒙达明(1977-), 本科, 高级工程师, 从事计算机网络通信技术研究。

模式挖掘的经典方法之一，其通过迭代生成候选项集并筛选满足支持度阈值的频繁项集，逐步构建完整的频繁模式集合。此外，FP-Growth 算法通过构建频繁模式树（FP-tree）有效减少了计算开销，为大规模数据的频繁模式挖掘提供了高效解决方案。

2.2 频繁模式挖掘的特点

频繁模式挖掘具有以下特点：首先，针对数据集中出现频率较高的模式，能够有效挖掘显著特征，尤其适用于高维度数据分析。其次，其过程具有良好的可扩展性，可应用于各种数据结构，包括事务型数据、时间序列数据等。此外，频繁模式挖掘方法能够结合领域知识，通过设置适当的支持度和置信度阈值过滤噪声，提高结果的可靠性和适用性。

2.3 频繁模式挖掘的挑战与局限

尽管频繁模式挖掘方法在理论上具有较强的普适性，但在实际应用中仍面临一些挑战。首先，数据规模和维度的快速增长可能导致算法的计算复杂度急剧上升，对存储和计算资源提出更高要求。其次，如何合理设置支持度阈值以平衡结果的覆盖率和精度是一项关键问题。此外，对于动态数据或实时流数据，传统频繁模式挖掘方法的效率和效果可能受到限制，需引入在线挖掘和增量挖掘技术进行优化。

3 网络流量分析中的频繁模式挖掘

3.1 网络流量的特性与分析需求

网络流量数据是互联网运行过程中生成的重要信息，其主要特性包括海量性、时序性和多样性。首先，随着网络用户数量和应用场景的不断增加，网络流量数据呈现爆炸式增长，其规模动辄以 TB 甚至 PB 计算。这种海量性要求数据分析技术具备高效处理大规模数据的能力。其次，网络流量具有明显的时序特性，即数据生成和传输都以时间为关键变量，流量的时变特征对分析用户行为和网络性能具有重要意义。例如，不同时间段的流量模式可能反映用户的访问习惯或网络负载的变化。最后，网络流量数据的多样性体现在数据类型和内容的复杂性上，包括传输协议、访问 IP 地址、端口号、数据包大小以及传输时延等多种信息。

3.2 基于频繁模式挖掘的网络流量分析模型

本文提出了一种基于频繁模式挖掘的网络流量分析模型，该模型通过挖掘网络流量数据中的频繁模式，为流量优化和异常检测提供科学依据。模型的设计分为三个主要步骤：数据预处理、频繁模式挖掘和结果解释。

在数据预处理阶段，首先需要网络流量数据进行清洗和格式化。由于原始数据可能包含噪声和冗余信息，如丢失的数据包或重复的流量记录，数据清洗的目的是去除这些无效信息，以确保分析结果的准确性。接下来，数据格式化包括对关键字段的提取和统一，例如协议类型、数据包大小和时间戳等关键信息。这些特征为后续的模式挖掘提供了必要的输入。此外，数据预处理阶段还包括数据分组和时间窗

口划分，以便识别不同时间段内的流量特征和模式变化。

4 模型应用效果评估

4.1 实验数据与设置

为了验证基于频繁模式挖掘的网络流量分析模型的有效性，实验选取了某企业的实际网络流量数据集。数据集包含为期一个月的网络访问记录，总数据量达到数百万条，记录的内容包括时间戳、源 IP 地址、目标 IP 地址、协议类型、端口号及数据包大小等多种信息。这些数据涵盖了企业网络的日常运行情况，具有高度代表性和真实性。

实验采用 FP-Growth 算法进行频繁模式挖掘，同时设置支持度阈值为 0.01，这意味着仅挖掘在整个数据集中出现频率达到 1% 以上的模式。FP-Growth 算法通过构建频繁模式树（FP-tree）减少了传统频繁项集生成过程中对候选项集的依赖，大幅提升了计算效率。为对比验证 FP-Growth 算法的性能，实验还采用了 Apriori 算法作为基准。Apriori 算法虽然是经典的频繁模式挖掘方法，但其在数据规模较大时的性能劣势较为显著，因此成为对比分析的合理选择。

实验平台的硬件配置为 Intel i7 处理器、32GB 内存和 1TB SSD 存储，操作系统为 Linux。软件环境采用 Python 作为实现语言，并利用 Scikit-learn 库完成数据处理和算法实现。通过统一的实验环境和数据集，确保了实验结果的可重复性和对比分析的公平性。

4.2 挖掘结果与分析

实验结果表明，基于 FP-Growth 算法的网络流量分析模型能够快速挖掘出多组高频模式，这些模式不仅反映了网络流量的整体特征，还揭示了某些特定时间段内的行为规律。例如，某些高频访问的协议组合（如 HTTP 和 HTTPS）以及频繁出现的端口号组合显示了用户在办公时间段的主要网络行为。这些模式能够帮助网络管理员识别流量的时序特征，从而为优化资源分配和制定网络安全策略提供参考依据。

此外，实验还发现了一些异常模式，例如某些 IP 地址在短时间内频繁访问多个端口，这种行为可能与潜在的网络攻击有关。通过进一步分析这些异常模式，管理员可以及时采取措施，如阻止异常 IP 地址的访问或加强该时间段的网络监控。相比之下，使用 Apriori 算法的挖掘结果虽然同样能揭示这些模式，但其计算时间明显更长，同时对复杂模式的覆盖率相对较低，这表明 FP-Growth 算法在处理大规模数据时具有明显优势。

结合实验结果，频繁模式挖掘模型的应用价值主要体现在两个方面：一是能够快速获取网络流量的整体特征，为网络优化提供数据支撑；二是能够高效识别异常行为，为网络安全管理提供及时的预警信息。

4.3 性能评估

模型性能的评价主要从计算效率、结果覆盖率和实际

应用价值三个维度展开。实验数据显示,在支持度阈值为0.01的条件下,FP-Growth算法的运行时间仅为Apriori算法的40%。这是因为FP-Growth通过构建和利用频繁模式树,大幅减少了候选项集的生成次数,从而显著降低了计算复杂度。对于数百万条规模的网络流量数据,FP-Growth算法的运行时间大约为10分钟,而Apriori算法则需要超过25分钟,这种效率的提升对实时性要求较高的应用场景尤为重要。

在结果覆盖率方面,FP-Growth算法能够挖掘更多显著模式。实验结果显示,FP-Growth算法挖掘出的模式数量比Apriori算法高出约20%。这一结果表明,FP-Growth算法在确保效率的同时,能够提供更全面的数据洞察,为网络优化和异常检测提供更多可能性。

从实际应用价值来看,该模型能够结合实验数据的特性,将挖掘结果转化为实际应用。例如,通过高频模式优化网络带宽分配策略,减少拥塞情况;通过异常模式识别潜在威胁,提高网络安全防护能力。这些实际应用进一步验证了模型的实用性和可操作性。

4 频繁模式挖掘在网络流量分析中的应用前景

5.1 流量异常监测

在网络安全领域,频繁模式挖掘为流量异常检测提供了高效的工具。通过挖掘流量数据中的高频模式,可以快速发现异常行为模式,例如分布式拒绝服务(DDoS)攻击所表现出的高频请求特征。实验结果表明,通过将频繁模式挖掘应用于实时流量数据,能够在攻击发生的早期阶段检测到潜在威胁,并为管理员提供及时响应的依据。与传统的基于规则匹配的检测方法相比,频繁模式挖掘方法能够动态调整分析策略,适应不断变化的威胁场景,从而显著提高检测效率和准确性。

5.2 流量优化与资源分配

频繁模式挖掘在流量优化和资源分配方面的潜力同样

不容忽视。通过分析不同时间段的高频流量模式,可以识别网络使用的峰值时段和主要流量来源。这些信息为带宽分配和服务器资源调度提供了重要参考。例如,在访问高峰期间,可以优先保证关键服务的带宽分配,同时优化非关键流量的传输路径,最大限度地提升网络整体效率。此外,通过挖掘流量数据中常见的协议组合和端口使用情况,可以帮助网络管理员识别和清理冗余流量,提高资源利用率。

6 结语

本文基于频繁模式挖掘技术,研究了网络流量分析中的实际应用,并设计了一种高效的分析模型。实验结果验证了该模型在揭示网络行为规律、提高资源优化效率以及加强网络安全管理方面的显著作用。尤其是通过FP-Growth算法的应用,解决了传统方法在处理大规模数据时效率不足的问题,同时提高了模式挖掘的覆盖率和精度。

未来研究将聚焦于两个方向:一是进一步优化算法性能,包括提升算法的实时性和动态数据处理能力;二是探索模仿模式挖掘与其他数据分析技术的融合,例如机器学习和深度学习,以提升模型的智能化水平和适应性。通过不断的技术创新,频繁模式挖掘技术有望在网络管理和安全领域发挥更大的作用,为互联网的可持续发展提供强有力的技术支持。

参考文献

- [1] 陈威宇,王泷,何建锋,等.基于改进关联规则的报送信息大数据特征隐性加密算法[J].计算技术与自动化,2024,43(04):123-128.
- [2] 孙艳歌,邵罕,蒋明毅.基于闭合频繁模式的半随机森林数据流分类算法[J].信阳师范学院学报(自然科学版),2024,37(04):442-448.
- [3] 马亦晨,李改云,韩雪梅,等.基于数据挖掘的我国老年人慢性病共病关联规则分析[J/OL].海军军医大学学报,1-8[2025-01-15].
- [4] 徐实.医疗数据应用领域基于公共利益的同质豁免机制[J].网络法律评论,2024,26(00):237-262.