



Improving the speed of the intrusion detection system and maintaining its performance by reducing the data volume using kernel-based DBSCAN.

S.A. Shahzadeh Fazeli*, A. Ghoveh Nodoushan, J. Zarepour Ahmadabadi

*Associate Professor, Yazd University, Yazd, Iran

Received: 2024/10/05, Revised: 2024/12/12, Accepted: 2024/01/12, Published: 2024/02/01

DOR: <https://dor.isc.ac/dor/20.1001.1.23224347.1403.12.4.8.2>

ABSTRACT

The Internet of Things (IoT) is a rapidly evolving technology that connects physical devices through networked systems. However, as IoT continues to expand, it poses various security challenges that require appropriate solutions to protect sensitive information and user privacy. This paper focuses on improving the speed of intrusion detection systems (IDS) as a critical solution for IoT security. In IDS, the large volume of data can slow down the learning process. In this paper, the DBSCAN clustering algorithm is modified by adding a minimum neighborhood parameter to reduce data samples in a targeted manner, aiming to enhance the speed of IDS and reduce learning time and costs. The parameters of the modified DBSCAN are tuned using a genetic algorithm. Experimental results on the Kaggle and NSL_KDD datasets demonstrate that the proposed model can maintain classification accuracy above 96% for the Kaggle dataset and above 92.51% for the NSL_KDD dataset, even with up to an 80% reduction in data volume. Additionally, computation time for the Kaggle dataset decreased from 458.09 ms to 47.21 ms, and for the NSL_KDD dataset from 995.2 ms to 223.60 ms. Thus, despite improvements in speed and reductions in time and cost, the model's optimal performance is maintained.

Keywords: : Internet of Things, Intrusion detection system, Clustering, Classification and regression tree algorithm, DBSCAN, Data reduction, RN_DBSCAN, Genetic algorithm, Kaggle dataset, NSL_KDD dataset.

* Corresponding Author Email: fazeli@yazd.ac.ir

بهبود سرعت سیستم تشخیص نفوذ و حفظ عملکرد آن از طریق کاهش حجم داده‌ها با استفاده از DBSCAN مبتنی بر هسته

سید ابوالفضل شاهزاده فاضلی^{۱*}، اعظم قوه ندوشن^۲، جمال زارع پور احمدآبادی^۳

۱-دانشیار، ۲-دانشجوی کارشناسی، ۳-استادیار، دانشگاه یزد، یزد، ایران

(دریافت: ۱۴۰۳/۰۷/۱۴، بازنگری: ۱۴۰۳/۰۹/۲۲، پذیرش: ۱۴۰۳/۱۰/۲۳، انتشار: ۱۴۰۳/۱۱/۱۳)

DOR: <https://dor.isc.ac/dor/20.1001.1.23224347.1403.12.4.8.2>



* این مقاله یک مقاله با دسترسی آزاد است که تحت شرایط و ضوابط مجوز Creative Commons Attribution (CC BY) توزیع شده است.

© نویسندهان

ناشر: دانشگاه جامع امام حسین (ع)

چکیده

اینترنت اشیاء یک فناوری به سرعت در حال تکامل است که دستگاه‌های فیزیکی را از طریق سیستم‌های شبکه‌ای به هم متصل می‌کند. با این حال، همان‌طور که اینترنت اشیاء به گسترش خود ادامه می‌دهد، چالش‌های امنیتی مختلفی را ایجاد می‌کند که نیازمند راحلهای مناسب برای محافظت از اطلاعات حساس و حریم خصوصی کاربران است. این مقاله بر روی بهبود سرعت سیستم تشخیص نفوذ به عنوان یک راحلهایی برای امنیت اینترنت اشیاء تمرکز دارد. در سیستم‌های تشخیص نفوذ، وجود حجم زیاد داده موجب کاهش سرعت یادگیری می‌شود. در این مقاله، الگوریتم خوشبندی DBSCAN با افزودن پارامتر حداقل همسایگی جهت کاهش هدفمند نمونه‌ها اصلاح شده است که سعی در افزایش سرعت سیستم تشخیص نفوذ و کاهش زمان و هزینه یادگیری دارد. تنظیم پارامترهای DBSCAN اصلاح شده با الگوریتم ژنتیک انجام می‌شود. نتایج آزمایش‌ها بر روی مجموعه‌داده NSL_KDD و Kaggle نشان می‌دهد که مدل پیشنهادی قادر است با کاهش تا ۸۰٪ از حجم داده‌ها، دقت طبقه‌بندی را برای مجموعه‌داده Kaggle بالای ۹۶٪ و برای مجموعه‌داده NSL_KDD بالای ۵۱٪ حفظ نماید. همچنین، زمان محاسبات برای مجموعه‌داده NSL_KDD از ۴۷/۲۱ ms به ۴۵۸/۰۹ ms و برای مجموعه‌داده RN_DBSCAN از ۹۹۵/۲ ms به ۶۰ ms کاهش یافته است. به این ترتیب، با وجود بهبود در سرعت و کاهش زمان و هزینه، عملکرد مطلوب مدل حفظ شده است.

کلیدواژه‌ها: اینترنت اشیاء، سیستم تشخیص نفوذ، الگوریتم درخت طبقه‌بندی و رگرسیون، DBSCAN، کاهش داده،

.NSL_KDD، Kaggle، مجموعه‌داده RN_DBSCAN

حریم خصوصی این اطلاعات برای جلوگیری از سوءاستفاده و نفوذ به اطلاعات شخصی افراد بسیار حائز اهمیت است. همچنین سیستم‌های IoT به طور مداوم در معرض حملات امنیتی مانند حمله توزیع شده انکار از سرویس^۲، نفوذ به شبکه و سرقت اطلاعات هستند. بدون توجه به امنیت مناسب، دستگاه‌های IoT می‌توانند به دروازه‌های ورودی برای حملات به سایر سیستم‌ها تبدیل شوند. از طرفی در بعضی موارد مانند صنایع نفت و گاز، حمل و نقل هوایی و صنایع برق، سیستم‌های IoT نقش حیاتی در عملکرد و امنیت این صنایع را ایفا می‌کنند. نقص امنیتی در این دستگاه‌ها می‌تواند منجر به تخریب سیستم‌های حیاتی، ایجاد خسارات‌های جدی و حتی تهدید جان افراد شود [۳، ۲].

به طور کلی، امنیت در اینترنت اشیاء باید در همه جوانب توسعه و استفاده از دستگاه‌های متصل به اینترنت مدنظر قرار

۱. مقدمه

اینترنت اشیاء^۱ به شبکه‌ای از اشیاء فیزیکی اطلاق می‌شود که از طریق اینترنت به یکدیگر متصل می‌شوند [۱]. امنیت در اینترنت اشیاء بسیار مهم است؛ زیرا اینترنت اشیاء شامل دستگاه‌هایی است که به شبکه اینترنت متصل می‌شوند و اطلاعات حساس را به اشتراک می‌گذارند. اگر این دستگاه‌ها مورد حملات نفوذ قرار گیرند، می‌تواند منجر به دسترسی غیرمجاز به اطلاعات شخصی، کنترل غیرمجاز بر روی دستگاه‌ها و حتی خطرات فیزیکی مانند حمله به خودروهای هوشمند یا سیستم‌های خانه هوشمند شود. با توجه به تعداد روزافزون دستگاه‌های متصل به اینترنت، تهدیدات امنیتی نیز روبرو باشند. برای مثال، دستگاه‌های متصل به IoT اطلاعات حساسی را درباره کاربران، رفتارها، و محیط‌ها جمع‌آوری می‌کنند. حفظ

² distributed denial-of-service (DDOS)

¹ Internet of things (IOT)

* رایانه‌نامه نویسنده مسئول: fazel@yazd.ac.ir

۴. حفظ حریم خصوصی: با کاهش حجم داده‌ها، میزان اطلاعات حساسی که باید در دسترس سیستم‌های تشخیص نفوذ قرار گیرد کمتر می‌شود. این امر به حفظ حریم خصوصی کاربران و موجودیت‌های سازمان کمک می‌کند.

به طور خلاصه، کاهش حجم داده‌ها در سیستم‌های تشخیص نفوذ بهبود کارایی، کاربردی تر شدن و افزایش دقت و قابلیت تشخیص سیستم را به همراه دارد [۸]. در ادامه، کارهای مرتبط انجام شده توسط محققان دیگر در بخش ۲ شرح داده می‌شود و در بخش ۳ پیش‌نیازهای تحقیق بیان می‌شود. در بخش ۴ روش کاهش پیشنهادی به همراه شبه کد آن توضیح داده می‌شود و تنظیم پارامترهای آن توسط الگوریتم زنگین، انجام می‌شود و همچنین مدل استفاده شده به همراه روش‌های هرس و بهبود آن، ارائه می‌شود. در بخش ۵ نتایج عملی روش پیشنهادی به همراه نمودارهای بررسی عملکرد و همچنین نمودار مقایسه روش پیشنهادی با کارهای مشابه قبل، قرار داده شده است و در نهایت، در بخش ۶، نتایج کلی بیان می‌شود.

۲. مروری بر کارهای مرتبط

فرایند کاهش داده‌ها را می‌توان با استفاده از تکنیک‌های خوشه‌بندی انجام داد که هدف آن تقسیم داده‌ها به خوشه‌های کوچک برای یادگیری است. ویهارت و همکاران [۹]، خوشه‌بندی را انجام داده‌اند و سپس از هر خوشه برای الگوریتم یادگیری لونبرگ - مارکوارت و الگوریتم شبیه نیوتون استفاده نموده‌اند تا برای نتیجه‌گیری با الگوریتم بیز ساده، ترکیب شود. استفاده از خوشه‌بندی برای کاهش داده‌ها نیز توسط جان و همکاران انجام شد [۱۰] که ترکیبی از کاهش داده و کاهش ویژگی است. در تحقیقات فوق، با استفاده از خوشه‌بندی قصد دارند داده‌ها را به تعدادی خوشه تقسیم کنند، به طوری که مقدار داده برای هر خوشه کمتر از کل داده‌ها باشد. این رویکرد، خوشه‌بندی کل داده‌ها را کاهش نمی‌دهد، بلکه داده‌ها را به تعداد ثابتی از خوشه‌ها تقسیم می‌کند و اگر همه داده‌های هر خوشه جمع شوند، مانند داده‌های اصلی باقی می‌مانند و در نتیجه پیچیدگی محاسبات را افزایش می‌دهد. استفاده از خوشه‌بندی، علاوه بر تقسیم داده‌ها، می‌تواند مبنای برای کاهش داده‌ها باشد تا از حجم داده‌ها کاسته شود.

برخی از مدل‌های کاهش داده با استفاده از الگوریتم‌های خوشه‌بندی، در ادامه توصیف می‌شوند. لی خاک و همکاران [۱۱]، یک مدل کاهش داده را با ترکیب الگوریتم شباهت نزدیک‌ترین همسایه مشترک^۲ و الگوریتم خوشه‌بندی بر اساس

گیرد. استفاده از سیستم‌های تشخیص نفوذ^۱ مبتنی بر یادگیری ماشین می‌تواند به عنوان یک راه حل قوی و کارآمد برای شناسایی و پیشگیری از حملات نفوذ در اینترنت اشیاء مورد استفاده قرار گیرد.

سیستم تشخیص نفوذ، یک سیستم امنیتی هست که طراحی شده است تا به طور خودکار فعالیت‌های ناهنجار و هجمه‌های امنیتی را در یک شبکه یا سیستم کامپیوتری تشخیص دهد و آن‌ها را گزارش کند. هدف اصلی سیستم تشخیص نفوذ، شناسایی هجمه‌های امنیتی و محافظت از سیستم‌ها و شبکه‌ها در برابر تهدیدهای امنیتی است [۴، ۵].

یادگیری ماشین و سیستم تشخیص نفوذ دو حوزه مرتبط هستند که در کنار یکدیگر استفاده می‌شوند. یادگیری ماشین، به عنوان یک فرایند که به وسیله آن ماشین‌ها و سیستم‌ها قادر به یادگیری الگوها و رفتارهای مشکوک و استنتاج از داده‌ها هستند، در سیستم‌های تشخیص نفوذ به کار می‌رود [۶، ۷].

در سیستم‌های تشخیص نفوذ، کاهش حجم داده‌ها یکی از موارد مهم برای بهبود عملکرد و کارایی سیستم است. دلیل اصلی این امر، مربوط به حجم بسیار زیاد داده‌های شبکه و سیستم است که باید تحلیل شوند تا فعالیت‌های مشکوک و نفوذ را تشخیص دهند. زمانی که یک IDS فعال است، در هر لحظه داده‌های مختلفی از سیستم‌ها و شبکه‌ها را بررسی می‌کند. این داده‌ها شامل پروتکل‌های شبکه، لگ‌ها، رویدادها و اطلاعات مختلف دیگر هستند. با توجه به حجم بالای این داده‌ها، تحلیل و پردازش آن‌ها ممکن است زمان برا و پرهزینه باشد؛ بنابراین، کاهش حجم داده‌ها به منظور تسهیل فرایند تحلیل و تشخیص نفوذ اهمیت دارد. به طور کلی، با کاهش حجم داده‌ها در سیستم‌های تشخیص نفوذ، مزایای زیر به دست می‌آید:

۱. کارایی بالا: با کاهش حجم داده‌ها، فرایند تحلیل و تشخیص نفوذ سریع‌تر انجام می‌شود. زمان واکنش به فعالیت‌های مشکوک کاهش می‌یابد و امکان تشخیص بهموقع تهدیدات امنیتی افزایش می‌یابد.

۲. صرفه‌جویی در منابع: با کاهش حجم داده‌ها، نیاز به منابع سخت‌افزاری و نرم‌افزاری کمتر می‌شود. این به معنی کاهش هزینه‌ها و استفاده بهینه از ظرفیت سیستم است.

۳. کاهش اشتباہات: با انتقال و تحلیل حجم کمتری از داده‌ها، احتمال بروز خطاهای و اشتباہات انسانی کاهش می‌یابد. تمرکز بر داده‌های مهم‌تر و مشکوک‌تر، دقت و قابلیت تشخیص بیشتری را فراهم می‌کند.

² Sharing Nearest Neighbor (SNN)

¹ Intrusion Detection System (IDS)

نمونه‌گیری^۴ و برش داده^۵ استفاده نمودند. نمونه‌گیری یک روش ناپارامتری کاهش تعداد است که با نمایش مجموعه داده‌های بزرگ به عنوان زیر مجموعه داده‌های تصادفی کوچک‌تر، داده‌ها را کاهش می‌دهد. با این حال تکنیک نمونه‌گیری هنوز مشکلات متعددی در نمایش داده‌ها دارد، یعنی نمونه‌گیری می‌تواند برخی از داده‌هایی را که ممکن است با داده‌های گرفته شده همگن نباشد را حذف کند. این بر میزان دقت در نتایج طبقه‌بندی تأثیر می‌گذارد.

در سیستم تشخیص نفوذ که حاوی داده‌های پرت و نویز زیادی است، انتخاب الگوریتم برای خوشه‌بندی بسیار مهم است. از بین روش‌های خوشه‌بندی، خوشه‌بندی DBSCAN به دلیل توانایی در تشخیص و حذف داده‌های پرت و کشف خوشه‌های پیچیده، برای سیستم تشخیص نفوذ که بسیار حجمی و نویزی می‌باشد، مناسب‌تر است. در [۱۹] کاهش داده با استفاده از الگوریتم DBSCAN_m انجام شده است این کاهش با اصلاح الگوریتم DBSCAN بوده است که با درنظرگرفتن نقاط همگن برای هر نقطه بازدید نشده و حذف این نقاط انجام می‌شود. این روش نسبت به روش‌های کاهش Slice Data و Sampling بهتر است؛ اما به دلیل اینکه نقاط همگن با هر نقطه را کاهش می‌دهد، باعث ایجاد نقاط نویزی می‌شود.

۳. پیش نیاز تحقیق

داده‌های مورداستفاده: این تحقیق بر اساس داده‌های ارزیابی تشخیص نفوذ که نتیجه شبیه‌سازی محیط شبکه نظامی LAN نیروی هوایی ایالات متحده است، انجام شده است. شرکت دارپا که متعلق به وزارت دفاع ایالات متحده می‌باشد، اولین داده‌های استاندارد را برای بررسی و ارزیابی سیستم‌های تشخیص نفوذ، جمع‌آوری نمودند. در این پژوهش از مجموعه داده‌های Kaggle و NSL_KDD استفاده شده است. مجموعه داده شامل ۲۵,۱۹۲ نمونه برای داده‌های آموزش و ۲۲,۵۴۴ نمونه برای داده‌های آزمون است. همچنین مجموعه داده NSL_KDD نسخه جدید از مجموعه داده KDD-99 است که برای حل برخی از مشکلات ذاتی مجموعه داده KDD-99 پیشنهاد شده است و از توزیع مناسبی برخوردار است.

الگوریتم DBSCAN: یکی از الگوریتم‌های بسیار قدرتمند در زمینه خوشه‌بندی، الگوریتم «خوشه‌بندی» بر اساس تراکم

تراکم فاصله‌ها برای برنامه‌های کاربردی نویزدار^۱ ایجاد کردند. نقطه ضعف این مدل این است که در صورت وجود حجم زیادی از داده‌ها، وجود دو فرایند یعنی SNN و DBSCAN، سرعت محاسبات را به شدت تحت تأثیر قرار می‌دهند. توسعه انجام شده توسط وانگ و همکاران [۱۲]، کاهش در مقدار داده‌ها را با کاهش نمونه تصادفی ترکیب می‌کند و از روش تحلیل مؤلفه اصلی برای خوشه‌بندی بیشتر با استفاده از الگوریتم خوشه‌بندی C_میانگین استفاده می‌کند. با اشاره به الگوریتم خوشه‌بندی مورداستفاده، نتیجه حاصل از تحقیق [۱۱] بهتر از [۱۲] بوده است؛ به دلیل اینکه عملکرد الگوریتم خوشه‌بندی DBSCAN نسبت به C_میانگین، بهویژه برای داده‌های نویزدار بهتر است [۱۳، ۱۴].

الگوریتم کاهش مبتنی بر خوشه‌بندی دیگری وجود دارد که از طریق خوشه‌های همگن^۲ کاهش می‌یابد. RHC، بر اساس مفهوم همگنی است، اما از الگوریتم خوشه‌بندی k_میانگین استفاده می‌کند [۱۵-۱۷]. الگوریتم برای داده‌های غیرهمگن خوشه‌بندی را انجام می‌دهد به طوری که همه داده‌ها همگن می‌شوند. نقطه ضعف این الگوریتم این است که استفاده از k_میانگین با مجموعه داده‌های پرت و نویز به خوبی کار نمی‌کند. DBSCAN در مقایسه با k_میانگین از قابلیت حساسیت خوبی برخوردار است [۱۴]. بر اساس تحقیقات انجام شده توسط اوجباروگلو و همکاران [۱۷]، کاهش داده‌ها با استفاده از الگوریتم کاهش خوشه‌های همگن برای تقریباً ۸۰٪ کاهش داده‌ها انجام شده است و از الگوریتم‌های طبقه‌بندی شبکه عصبی، مانشین بردار پشتیبان و k_نزدیک‌ترین همسایگان استفاده شده است که دقت حاصل زیر ۹۰٪ می‌باشد.

الگوریتم RHC، روی داده‌های نویزی ضعیف عمل می‌کند. با توجه به این ضعف، الگوریتم کاهش خوشه‌های همگن با عنوان «ویرایش و کاهش از طریق خوشه‌های همگن»^۳، توسط اوجباروگلو و همکاران، توسعه یافته [۱۸]. در الگوریتم eRHC، وقتی یک نقطه داده تشکیل یک خوشه می‌دهد، به عنوان نویز در نظر گرفته شده و دور انداخته می‌شود. در مطالعه [۱۸]، نشان می‌دهد که کاهش داده با الگوریتم eRHC می‌تواند عملکردی بدقت کمتر از ۹۰٪، زمانی که با الگوریتم طبقه‌بندی مانشین بردار پشتیبان ترکیب شود، به دست آورد.

و بهارتو و همکاران [۱۹]، جهت کاهش داده، از الگوریتم‌های

¹ Density Based Spatial Clustering algorithm for Applications with Noise (DBSCAN)

² Reduced Homogeneous Clusters (RHC)

³ Editing and Reduction through Homogeneous Clusters (eRHC)

حجم داده‌ها به منظور بهبود زمان و فضا و همچنین کاهش هزینه‌های موجود، استفاده می‌گردد.

روند کلی کار انجام شده، در شکل (۱) نشان داده شده است. همان‌طور که در قسمت (الف) نمودار گردشی شکل (۱) مشخص است، بعد از جمع‌آوری داده‌ها و انجام عملیات پیش‌پردازش، داده‌ها براساس کلاس تفکیک شده و سپس هر گروه داده به صورت مجزا کاهش می‌یابد. بعد از کاهش بر روی هر گروه، داده‌های کاهش یافته، ادغام می‌شوند و به دو بخش برای یادگیری و آزمون روش پیشنهادی، تقسیم می‌شوند. سپس داده‌های نرمال شده با الگوریتم CART طبقه‌بندی می‌شوند. در نهایت، عملکرد سیستم تشخیص نفوذ بهوسیله معیارهای ارزیابی مختلف، مورد بررسی و آزمایش قرار می‌گیرد و نتیجه نهایی به دست می‌آید.

مرحله پیش‌پردازش به سه بخش؛ پاک‌سازی داده‌ها، عددی‌سازی داده‌ها و جداسازی داده‌ها تقسیم می‌شود. عددی‌سازی داده‌ها با استفاده از روش «وان هات»^۶ انجام شده است. در روش کدبندی وان‌هات، بهازای هر مقدار از ویژگی غیر عددی، ستونی به مجموعه‌داده اضافه می‌گردد، در صورتی که نمونه، آن ویژگی را داشته باشد، مقدار یک و در غیر این صورت، مقدار صفر قرار می‌گیرد. همچنین از روش نرمال‌سازی «نمره استاندارد»^۷ استفاده شده است. این روش سعی دارد، برای یک مجموعه‌داده، مقدارهایی را به دست آورد که دارای میانگین صفر و واریانس یا انحراف معیار یک باشد؛ بنابراین اگر میانگین داده‌های اصلی برابر با μ و انحراف معیار آن‌ها نیز σ باشد، مقدار Z را براساس رابطه (۱) می‌توان به دست آورد.

$$z = \frac{x - \mu}{\sigma} \quad (1)$$

به‌این‌ترتیب مشخص است که داده‌های تبدیل یافته Z دارای میانگین صفر و واریانس یک هستند.

باتوجه به شکل (۱) قسمت (ب)، مرحله کاهش داده‌ها از دو بخش تشکیل می‌شود، بخشی از کاهش ابعادی از نظر مقدار داده می‌باشد که با استفاده از الگوریتم اصلاح شده DBSCAN انجام CART می‌شود و بخشی دیگر از کاهش ابعاد، مربوط به الگوریتم است. این الگوریتم علاوه بر آموزش و طبقه‌بندی داده‌ها، کاهش ویژگی نیز انجام می‌دهد.

⁶ One hot encoding
⁷ Z-score

فاصله‌ها برای برنامه‌های کاربردی نویزدار^۱ است. این الگوریتم بر اساس چگالی داده‌ها، خوشبندی را انجام می‌دهد [۲۰]. الگوریتم DBSCAN نیاز به دو پارامتر نقطه حداقل^۲ و اپسیلون^۳ دارد. هر نقطه از گردشی دارد، اگر این فاصله بین دو نقطه مفروض، کمتر از اپسیلون باشد، به عنوان همسایه آن نقطه در نظر گرفته می‌شود. هر نقطه مفروض که به تعداد نقطه حداقل یا بیشتر، همسایه داشته باشد، یک نقطه مرکزی محاسب می‌شود. در واقع با استفاده از این دو پارامتر، حداقل چگالی یک خوشبندی می‌شود [۲۱].

الگوریتم ژنتیک: الگوریتم ژنتیک یک الگوریتم بهینه‌سازی است که از انتخاب طبیعی و از مفهوم تکامل الهام‌گرفته شده است. این الگوریتم یک جستجوی مبتنی بر جمعیت هست و از نظریه داروینی در مورد بقای بهترین‌ها در طبیعت تقلید می‌کند [۲۲].

الگوریتم CART: یکی از محبوب‌ترین و در عین حال ساده‌ترین الگوریتم‌های درخت‌های تصمیم، درخت طبقه‌بندی و رگرسیون^۴ است که کاربردهای فراوانی در طبقه‌بندی و رگرسیون دارد. این الگوریتم بر اساس درخت‌های دودویی (باینری) بنا نهاده شده است. الگوریتم درخت طبقه‌بندی و رگرسیون، برای ساخت درخت تصمیم، داده‌ها را به قسمت‌های دوتایی تقسیم کرده و بر اساس آن‌ها درخت دودویی را می‌سازد. این درخت (و البته درخت‌های دیگر) می‌تواند پایه‌ای برای الگوریتم‌های پیچیده‌تر مانند جنگل تصادفی^۵ باشد [۲۳].

۴. روش پیشنهادی

مسئله موردنظر در این پژوهش، بهینه‌سازی سیستم‌های تشخیص نفوذ است. در واقع هدف تعیین راهبردی برای حداقل‌سازی عملکرد سیستم تشخیص نفوذ می‌باشد؛ به طوری که بتوان زمان و هزینه را تاحدامکان کاهش داد.

برای بهبود الگوریتم یادگیری ماشینی که قادر به کاوش اطلاعات و در پی آن تعیین الگویی برای تشخیص نفوذ باشد، نیازمند به انتخاب برخی نمونه‌هایی که بتوانند الگوریتم طبقه‌بندی را به خوبی آموزش داده و از حجم مناسبی برخوردار باشند، حس می‌شود؛ بنابراین در این مقاله از بهبود خوشبندی DBSCAN جهت حذف سروصدایی‌های موجود و همچنین کاهش

¹ Density Based Spatial Clustering algorithm for Applications with Noise (DBSCAN)

² Min Point (MinPts)

³ Epsilon (Eps)

⁴ Classification And Regression Tree (CART)

⁵ Random forest

$$d_{ij} = \sqrt{\sum_{v=1}^n (x_{vi} - x_{vj})^2} \quad (2)$$

شبیه کد مربوط به فرایند کاهش داده‌ها با روش RN در الگوریتم‌های ۱ و ۲ نشان داده شده است. در این الگوریتم‌ها، مواردی که پررنگ‌تر نوشته شده‌اند، تفاوت الگوریتم پیشنهادی این مقاله و الگوریتم DBSCAN است.

RN-DBSCAN : ۱

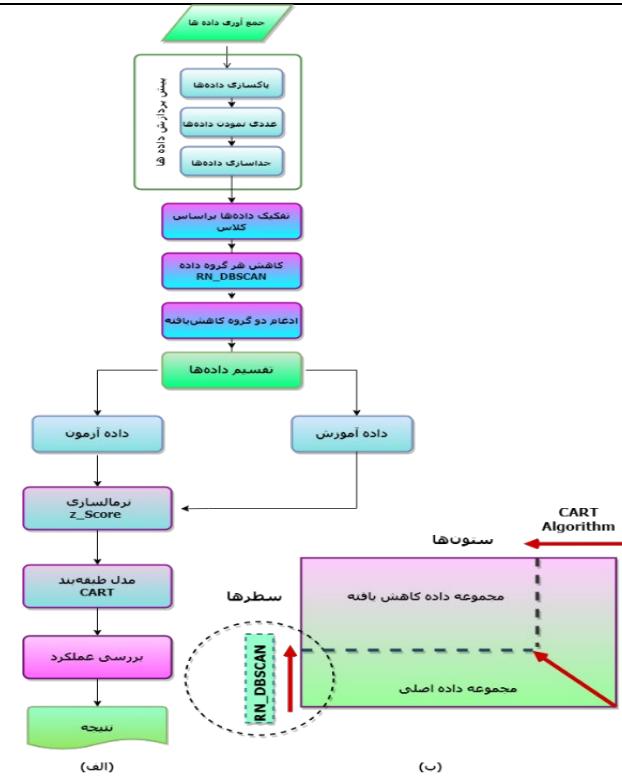
ورودی: پارامترهای $D, Eps, MinPts, MinNr$ را وارد کنید.
خروجی: تشکیل خوشه و یافتن نقاط مرکزی.

- ۱ $C = 1$ قرار دهد.
- ۲ برای هر نقطه p ملاقات نشده در مجموعه داده D مراحل زیر را انجام دهد.
- ۳ p را به عنوان بازدید شده علامت‌گذاری کنید.
- ۴ همسایه‌های p به شعاع Eps را در N قرار دهد.
- ۵ اگر تعداد N کوچک‌تر از $MinPts$ آنگاه $Noise$ را به عنوان p علامت‌گذاری کنید.
- ۶ در غیراینصورت p را در خوشه C جدید قرار دهد.
- ۷ همسایه‌های p به شعاع $MinNr$ را در X قرار دهید و حذف نماید.
- ۸ $ExpandCluster(P, N, C, Eps, MinPts, MinNr, X)$
- ۹ خوشه را با p گسترش دهد.
- ۱۰ پایان آنگاه C است.
- ۱۱ پایان حلقه.
- ۱۲ پایان حلقه.

الگوریتم ۲: ExpandCluster

- ورودی: $P, N, C, Eps, MinPts, MinNr, X$ را وارد کنید.
خروجی: گسترش خوشه‌ها و یافتن نقاط همگن با نقاط مرکزی جهت کاهش.
- ۱ p را به عنوان خوشه C اضافه کنید.
 - ۲ برای هر نقطه p' در N قرار دهد.
 - ۳ اگر p' ملاقات نشده آنگاه p' را به عنوان ملاقات شده علامت‌گذاری کنید.
 - ۴ همسایه‌های p' به شعاع Eps را در N' قرار دهد.
 - ۵ اگر تعداد N' بزرگ‌تر ساواحی با $MinPts$ آنگاه N' را با N ترتیب کرده در N قرار دهد.
 - ۶ همسایه‌های p' به شعاع $MinNr$ را در X قرار دهید و حذف نماید.
 - ۷ پایان آنگاه C است.
 - ۸ پایان حلقه.
 - ۹ اگر p' در هیچ خوشه‌ای نیست آنگاه p' را به خوشه C اضافه کنید.
 - ۱۰ پایان آنگاه C است.
 - ۱۱ پایان حلقه.
 - ۱۲ پایان آنگاه C است.
 - ۱۳ پایان آنگاه C است.
 - ۱۴ پایان حلقه.

سازوکار این الگوریتم مشابه الگوریتم خوشبندی DBSCAN است، با این تفاوت که نقطه P علاوه بر شعاع Eps ، شعاع $MinNr$ را جهت یافتن نقاط همگن با نقطه مرکزی، یعنی نقاطی که بسیار نزدیک و مشابه نقاط مرکزی هستند و به اصطلاح همپوشانی دارند، محاسبه می‌نماید، تا در نهایت از این نقاط برای کاهش استفاده کند. در واقع کاری که انجام شده است، استفاده از الگوریتم DBSCAN، جهت خوشبندی و یافتن نقاط مرکزی می‌باشد، تا بتوان از این نقاط برای کاهش بهره‌مند شد، به این صورت که نقاطی که در شعاع $MinNr$ به مرکز نقاط مرکزی



شکل (۱): نمودار گردشی روش پیشنهادی تحقیق

۱-۴. کاهش ابعاد با الگوریتم RN_DBSCAN

الگوریتم DBSCAN برای انجام فرایند کاهش، مورد اصلاح قرار می‌گیرد و با عنوان «کاهش بهوسیله همسایگان در الگوریتم خوشبندی بر اساس تراکم فاصله‌ها برای برنامه‌های کاربردی نویزدار^۱» مطرح می‌شود.

اصلاح الگوریتم DBSCAN با افزودن پارامتری به نام حداقل همسایگی^۲ انجام شده است که این پارامتر برای تعیین فاصله چگالی تا نقطه مرکزی مورداستفاده قرار می‌گیرد و سپس برای حذف علامت‌گذاری می‌شود. در واقع سطح چگالی اطراف نقطه مرکزی به صراحت به عنوان مقدار حداقل همسایگی تعریف می‌شود. کاهش داده‌ها در این روش، با حذف برخی از داده‌ها که تراکم بیشتری نسبت به داده‌های دیگر دارند، انجام می‌شود. این تراکم با توجه به توزیع داده‌ها تعیین می‌شود.

جهت تعیین پارامترهای Eps , $MinNr$ و $MinPts$ محاسبه ماتریس فاصله طبق رابطه (۲) برای هر داده انجام می‌گیرد، که در آن d_{ij} فاصله بین نقطه i و نقطه j و x_{vi} مختصات نقطه i برای بعد v و x_{vj} مختصات نقطه j برای بعد v می‌باشد [۲۴].

¹ Reduction by Neighboring in DBSCAN (RN_DBSCAN)

² MinNeighborhood (MinNr)

تابع برازش $fitness_1$ و $fitness_2$ طبق رابطه (۴) تعريف شده اند:

$$fitness_1 = |N_{main} - N_{remain}|, \quad (4)$$

$$fitness_2 = |1 - accuracy|$$

که در آن، $fitness_1$ تفاضل N_{main} تعداد مطلوب دادهای کاهش یافته و N_{remain} تعداد دادهای کاهش یافته با RN_DBSCAN است و $fitness_2$ از تفاضل دقت به دست آمده با دقت کل، که یک درنظر گرفته شده است، به دست می‌آید.

۲. با درنظر گرفتن برچسب داده: در این حالت تابع برازش

طبق رابطه (۵) محاسبه می‌شود:

$$TotalFitness = w_1 fitness_1 + w_2 fitness_2 + w_3 fitness_3 \quad (5)$$

که در آن w_1 و w_2 و w_3 به ترتیب وزن‌های $fitness_1$ و $fitness_2$ و $fitness_3$ هستند.

مقادیر $fitness_1$ و $fitness_2$ و $fitness_3$ از روابط ۶ تا ۸ بدست می‌آیند.

$$fitness_1 = |N_{main_n} - N_{remain_n}|, \quad (6)$$

$$fitness_2 = |N_{main_an} - N_{remain_an}|, \quad (7)$$

$$fitness_3 = |1 - accuracy|. \quad (8)$$

در رابطه (۶)، $fitness_1$ از تفاضل مقدار دادهای کاهش یافته با روش RN_DBSCAN از کلاس نرمال (N_{remain_n})، با تعداد کاهش مطلوب از این گروه (N_{main_n}) به دست می‌آید. همچنین با توجه به رابطه (۷)، $fitness_2$ از تفاضل مقدار دادهای کاهش یافته با روش RN_DBSCAN از کلاس ناهنجار (N_{remain_an}) با تعداد کاهش مطلوب از این گروه (N_{main_an}) حاصل می‌شود. $fitness_3$ در رابطه (۸)، نیز بهترین دقت حاصل از درصد کاهش مدنظر را بررسی می‌کند.

نتایج تنظیم پارامترهای RN_DBSCAN با استفاده از الگوریتم ژنتیک، جهت کاهش داده‌ها برای مجموعه داده Kaggle، بدون درنظر گرفتن برچسب داده در جدول (۱) و با درنظر گرفتن برچسب داده در جدول (۲)، قرار گرفته است. همچنین تنظیم پارامترهای RN_DBSCAN با استفاده از الگوریتم ژنتیک، جهت کاهش داده‌ها برای مجموعه داده NSL-KDD، با درنظر گرفتن برچسب داده انجام گرفت که نتایج آن در جدول (۳) نشان داده شده است.

هستند، کاهش یابد.

دلیل استفاده از نقاط مرکزی جهت کاهش، تراکمی هست که در اطراف این نقاط وجود دارد، بنابراین حذف نقاط همگن به شعاع $MinNr$ از نقاط مرکزی، اطلاعاتی را از دست نداده و می‌توان مدل یادگیر را به وسیله نقاط باقیمانده به طور مناسبی آموزش داد. همچنین به دلیل داشتن تراکمی که در اطراف نقاط مرکزی وجود دارد، می‌توان مطمئن بود که با حذف نقاط اطراف نقاط مرکزی، آن نقطه تبدیل به نیز نخواهد شد.

۴-۲. تنظیم پارامترهای RN_DBSCAN با استفاده

از الگوریتم ژنتیک

همواره یکی از چالش‌های موجود در الگوریتم $DBSCAN$ تنظیم پارامترهای آن است. حال با افزودن پارامتر جدید به این الگوریتم جهت کاهش، این چالش بیش از پیش مشهود است؛ بنابراین بهمنظور استفاده از روش کاهش RN_DBSCAN و به دست آوردن نتیجه مطلوب، تنظیم پارامترهای آن از اهمیت ویژه‌ای برخوردار است.

در این پژوهش برای تنظیم پارامترهای الگوریتم RN_DBSCAN از الگوریتم ژنتیک بهره گرفته شده است. الگوریتم ژنتیک علاوه بر تنظیم خودکار پارامترها قادر است پارامترهای بهینه را به دست آورد. در واقع با استفاده از این الگوریتم، برای هر درصد کاهش، نمونه‌هایی که با اطلاعاتی کمتری دارند حذف می‌شوند؛ بنابراین، می‌توان گفت کارایی و دقت طبقه‌بندی نیز علاوه بر تعداد داده‌های باقیمانده از کاهش، مدنظر قرار گرفته می‌شود.

برای تعیین اثر کاهش داده‌ها، فرایند کاهش داده‌ها، با ایجاد دو رویکرد، یعنی کاهش داده‌ها بدون درنظر گرفتن برچسب داده‌ها و کاهش داده‌ها با درنظر گرفتن برچسب‌های داده، انجام می‌شود. کاهش با درنظر گرفتن برچسب‌های داده با گروه‌بندی داده‌ها بر اساس برچسب‌ها انجام می‌گیرد؛ به این منظور که درصد کاهش مدنظر، بر روی هر گروه از کلاس‌ها به صورت مجزا انجام می‌شود. با توجه به این دو حالت، برای الگوریتم ژنتیک، تابع برازش متفاوتی تعریف می‌شود که در ادامه توضیح داده می‌شود.

۱. بدون درنظر گرفتن برچسب داده: تابع برازش در این حالت، طبق رابطه (۳) محاسبه می‌شود. این تابع برازش کلی، از مجموع دو تابع برازش وزن دار به دست می‌آید. در رابطه (۳)، w_1 و w_2 وزن‌هایی هستند که به هر یک از تابع برازش $fitness_1$ و $fitness_2$ نسبت داده می‌شوند. w_1 وزن تعداد نمونه‌های کاهش یافته و w_2 وزن دقت می‌باشد.

$$TotalFitness = w_1 fitness_1 + w_2 fitness_2 \quad (3)$$

جدول (۱). نتایج کاهش داده‌ها بدون درنظر گرفتن کلاس (Kaggle)

پارامترهای RN_DBSCAN				
کاهش داده (%)	Eps	MinPts	MinNr	تعداد داده‌ها
۰	-	-	-	۲۵۱۹۲
۵	۲/۷	۲۱	۲/۴۵	۲۳۹۳۲
۲۰	۵/۵۹	۵۹	۳/۶	۲۰۱۵۴
۳۰	۱۱/۳۵	۹۷	۱۰/۲۷	۱۷۶۳۴
۴۹	۲۰/۰۲	۱۲	۱۱/۱۷	۱۲۸۴۸
۶۰	۲۰	۲	۱۴/۹۰۵	۱۰۰۷۷
۸۰	۷۶/۰۴	۲	۴۲/۶۲	۵۰۳۸

جدول (۲). نتایج کاهش داده‌ها با درنظر گرفتن کلاس (Kaggle)

کاهش داده (%)	طبیعی				ناهنجر				تعداد کلی داده‌ها
	Eps	MinPts	MinNr	تعداد داده‌ها	Eps	MinPts	MinNr	تعداد داده‌ها	
۰	-	-	-	۱۲۱۹۲	-	-	-	۱۳۰۰۰	۲۵۱۹۲
۵	۱۲/۵۵	۲۳	۵/۷۵	۱۱۵۸۲	۲	۱۲	۱/۳۲	۱۲۳۵۰	۲۳۹۳۲
۲۰	۱۵/۳	۴	۱۳/۵۵	۹۷۵۴	۵	۵۸	۲/۰۲	۱۰۴۰۰	۲۰۱۵۴
۳۰	۲۲/۴	۴	۲۰/۷۹	۸۵۳۴	۸/۰۱	۲۷	۱/۹۴	۹۱۰۰	۱۷۶۳۴
۴۹	۵۶/۶۲	۱	۲۹/۱۹	۶۲۱۸	۹/۰۳	۲۱	۳/۰۱	۶۶۳۰	۱۲۸۴۸
۶۰	۵۶/۰۹	۲	۳۸/۹	۴۸۷۷	۲۳/۷	۴۰	۴/۱۱	۵۲۰۰	۱۰۰۷۷
۸۰	۸۵/۰۷	۲	۷۴/۸۷	۲۴۳۸	۴۵/۰۵	۳۴	۱۰/۲۳	۲۶۰۰	۵۰۳۸

جدول (۳). نتایج کاهش داده‌ها با درنظر گرفتن کلاس (NSL_KDD)

کاهش داده (%)	طبیعی				ناهنجر				تعداد کلی داده‌ها
	Eps	Min Pts	MinNr	تعداد داده‌ها	Eps	MinPts	MinNr	تعداد داده‌ها	
۰	-	-	-	۶۷۳۴۳	-	-	-	۵۸۶۳۰	۱۲۵۹۷۳
۵	۳۶/۸۳	۴۹۰	۳/۲۵	۶۳۹۷۶	۴/۸۴	۴۹۰	۳/۲۷	۵۵۶۹۹	۱۱۹۶۷۵
۲۰	۳۷/۴۸	۲۰۸	۳۶/۹۱	۴۸۷۴۵	۴/۸۵	۲۹۶	۲	۴۶۹۰۴	۱۰۰۷۷۸
۳۰	۷۰/۰	۱۴۳	۱۰/۸	۴۷۱۴۰	۴۸/۸۶	۲۹۶	۴	۴۱۰۴۱	۹۴۰۴۴
۴۹	۷۸/۷۵	۱۳۲	۲۴/۵۱	۴۴۳۴۵	۵/۵۵	۱۵۰	۲/۸۱	۲۹۹۰۱	۶۴۲۴۶
۶۰	۹۰/۹۳	۲۰۸	۷۱/۰۶	۴۶۹۳۷	۱۲/۱۱	۴۹۰	۳/۳۱	۲۳۴۵۲	۵۰۳۸۹
۸۰	۴۳۵/۰۸	۵۳۱	۶۸/۴۸	۱۳۴۶۹	۱۹/۰۱	۱۷۲	۵/۲۵	۱۱۷۲۶	۲۵۱۹۵

مقدار پارامتر $MinPts$ افزایش یابد، کاهش کمتری رخ می‌دهد. دلیل این امر به این خاطر است که با افزایش Eps ، شعاع بزرگ‌تر شده، بنابراین تعداد همسایه‌های اطراف نقطه مورد نظر بیشتر

باتوجه به پارامترهای به دست آمده برای هر مجموعه داده در جداول (۱) تا (۳)، می‌توان مشاهده نمود، هراندازه مقدار پارامترهای $MinNr$ و Eps افزایش یابد، کاهش بیشتر می‌شود و هر اندازه

سمت راست است.

$$N\emptyset(s, t) = 2P_L P_R Q(s|t), \quad (9)$$

$$\emptyset(s, t) = 2P_L P_R Q(s|t), \quad (10)$$

$$\emptyset(s|t) = \sum_{j=1}^{The number of catgorics} |p(j|t_{left}) - p(j|t_{right})|. \quad (11)$$

گره شاخه مرتب‌کننده‌ای که به مقادیر $\emptyset(s, t)$ بالاتر منجر می‌شود، مرتب‌سازی بهتری است و به عنوان تقسیم‌گر شاخه انتخاب می‌شود. گره t مانند تنها یک مشاهده در هر گره فرزند، تبدیل به یک گره پایانی می‌شود، اگر کاهش قابل توجهی در ناهمگنی/ناخالصی یا حداقل حد n وجود نداشته باشد. حداقل تعداد محدودیت‌ها در ترمینال انتهایی به طور کلی، پنج است. هنگامی که گره‌های ترمینال پیدا می‌شوند، تشکیل درخت تصمیم خاتمه می‌یابد.

۴-۳-۲ هرس درخت تصمیم

هدف اصلی از ساختن درخت طبقه‌بندی، ایجاد یک مدل طبقه‌بندی‌کننده به شکل درخت تصمیم است که کمترین مقدار خطای طبقه‌بندی یا نزدیک به صفر را داشته باشد. با این حال، درخت تصمیم در توصیف ساختار داده، بسیار پیچیده است. یکی از راه‌های تعیین درخت طبقه‌بندی بهینه، بدون کاهش دقیق مدل طبقه‌بندی‌کننده، هرس نمودن درخت طبقه‌بندی است. سطح اهمیت درخت با بهره آن، بر اساس حداقل پیچیدگی هزینه فرموله شده در معادله (۱۲) اندازه‌گیری می‌شود [۱۹].

$$D_\alpha(T_K) + \alpha |\tilde{T}_K| \quad (12)$$

که در آن:

$$(T_K) = D_\alpha(T_K)$$

(T_K) = مجموعه‌ای از گروه‌های ترمینال در

$$|\tilde{T}_K| = \text{تعداد گره‌های ترمینال به } \tilde{T}_K$$

$$\alpha = \text{پارامتر هزینه - پیچیدگی}$$

هستند. هرس حداقل هزینه پیچیدگی، یکی از انواع هرس درختان تصمیم است.

به طور کلی می‌توان گفت، هرس نمودن درخت طبقه‌بندی، باعث جلوگیری از خطای بیش برازش شده و عملکرد و اعتبار سیستم را افزایش می‌دهد و موجب کاهش هزینه‌های محاسباتی می‌شود. همچنین هرس درخت طبقه‌بندی، برای نمونه‌های دیده نشده، دقیق‌تری را رقم می‌زند.

۴-۳-۳ معیارهای ارزیابی عملکرد

خواهد شد و در نتیجه با فرض ثابت بودن مقدار $MinPts$ ، نقاط بیشتری به عنوان نقاط مرکزی، درنظر گرفته می‌شوند و با افزایش $MinNr$ شعاع اطراف نقطه مرکزی بزرگ‌تر شده بنابراین نقاط بیشتری حذف خواهند شد. در نتیجه با افزایش مقدار این دو پارامتر، میزان کاهش داده‌ها افزایش می‌یابد. اما با افزایش پارامتر $MinPts$ ، به دلیل اعمال سخت‌گیری بیشتر در انتخاب نقاط به عنوان نقطه مرکزی، تعداد کمتری نقطه مرکزی به دست می‌آید، بنابراین با افزایش پارامتر $MinPts$ با فرض ثابت بودن دو پارامتر دیگر، نقاط کمتری کاهش می‌یابند.

هدف اصلی، کاهش نقاطی است که بسیار به نقطه مرکزی نزدیک بوده و در واقع مشابه و همگن با آن است؛ بنابراین به مرکز نقطه مرکزی و به شعاع کوچک‌ترین همسایگی، نقاط کاهش می‌یابند. نقاط موجود در این شعاع همان نقاط همگن با نقطه مرکزی بوده و باید از مجموعه‌داده کاهش یابند. پر واضح است که مقدار اپسیلون می‌بایست از کوچک‌ترین همسایگی، بیشتر باشد. همچنین مقدار نقطه حداقل که تعداد همسایه‌ها است صحیح در نظر گرفته می‌شود.

CART ۴-۳ مدل

در این پژوهش از مدل CART جهت آموزش و آزمون عملکرد سیستم تشخیص نفوذ استفاده شده است. الگوریتم CART یک الگوریتم طبقه‌بندی است که از رویکرد طبقه‌بندی درخت تصمیم استفاده می‌کند. الگوریتم CART برای سیستم‌های تشخیص نفوذ، دارای عملکرد بسیار مناسبی است [۲۵]. در کل می‌توان گفت این الگوریتم دارای برتری‌های نسبت به الگوریتم‌های دیگر در این مسئله می‌باشد [۲۶-۲۸]. برای استفاده از مدل CART، داده‌ها با نسبت ۷۰٪ برای آموزش و ۳۰٪ برای اعتبارسنجی به اشتراک گذاشته می‌شوند و سپس با استفاده از الگوریتم CART طبقه‌بندی انجام می‌شود.

۴-۳-۴ تشکیل درخت طبقه‌بندی

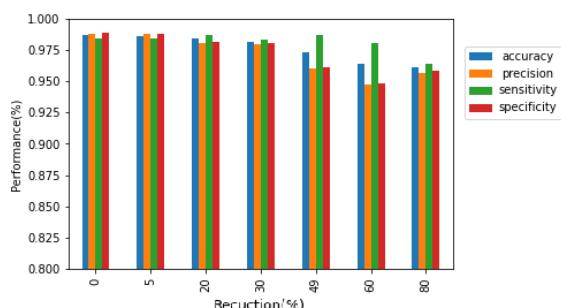
داده‌های نمونه برای یادگیری ناهمگن به عنوان مبنای برای تشکیل یک درخت طبقه‌بندی استفاده می‌شوند. در این مرحله، با رعایت قوانین مرتب‌سازی و معیارهای تقسیم‌بندی، مرتب‌سازی از داده‌هایی برای یادگیری انتخاب می‌شوند که منجر به بیشترین کاهش در سطوح ناهمگونی شوند.

تعیین مناسب‌بودن تقسیم $(s, t) \emptyset$ شاخه کاندید s در نقطه تصمیم t به صورت معادلات (۹) تا (۱۱) محاسبه می‌شود [۱۹]. که در آن t_{left} کاندید شاخه چپ از نقطه تصمیم t ، t_{right} کاندید شاخه سمت راست از نقطه تصمیم t ، $p(j|t_{left})$ احتمال j در گره شاخه سمت چپ، $(p(j|t_{right}))$ احتمال j در شاخه

داده در نظر گرفته نشود، عمل کاهش بر روی تمامی داده‌ها بدون توجه به برچسب آن‌ها، انجام می‌شود و دراین‌بین، امکان دارد تنها داده‌هایی از یک کلاس حذف گردد و یا بهصورت غیرمساوی از هر کلاس کاهش اعمال شود که این امر موجب عدم تعادل در مجموعه‌داده خواهد شد. اما با درنظرگرفتن برچسب داده‌ها طی فرایند کاهش، داده‌ها بر اساس برچسب به دو گروه (باتوجه به اینکه دو کلاس داده وجود دارد) تقسیم شده و برای هر گروه، درصد کاهش بهصورت جدا انجام می‌شود. جهت بررسی این که آیا تغییرات در تعادل داده‌ها بر عملکرد نتایج طبقه‌بندی تأثیر بهصورت Kaggle می‌گذارد یا خیر، این مورد بر روی داده‌های مجزا انجام‌گرفته است که در ادامه به بررسی نتایج پرداخته می‌شود.

۱-۵ نتایج بر روی مجموعه Kaggle بدون تفکیک کلاس

شکل (۲) نشان می‌دهد که درصدهای مختلف کاهش، بر عملکرد سیستم تشخیص نفوذ اثر می‌گذارد. با توجه به شکل، دقت بهطور یکنواخت با افزایش درصد کاهش، کم شده است. معیارهای خاصیت و صحت با افزایش درصد کاهش داده‌ها، سیر تغییر مشخص و یکنواخت نداشته‌اند. همچنان معيار حساسیت نسبت به بقیه معیارها از عملکرد بهتری برخوردار بوده است. با توجه به تمامی پارامترها، هر چه درصد کاهش داده‌ها افزایش می‌یابد، عملکرد سیستم تشخیص نفوذ، کاهش می‌یابد. اما این کاهش عملکرد نسبت به درصد کاهشی که اعمال می‌شود، کم است و همان‌طور که مشخص است، برای ۸۰٪ کاهش، دقت ۹۶/۱۲٪ حفظ شده است.



شکل (۲): تأثیر کاهش داده‌ها بر عملکرد سیستم تشخیص نفوذ، بدون درنظر گرفتن برچسب‌ها، در مجموعه داده Kaggle. نمودار بعدی که مورد ارزیابی قرار می‌گیرد، بررسی و مقایسه زمان با عملکرد سیستم تشخیص نفوذ است.

محاسبه زمان برای اندازه‌گیری اثر کاهش درصد داده‌ها بر

عملکرد مدل سیستم IDS توسعه‌یافته با اصلاح DBSCAN با استفاده از ماتریس ابهام، اندازه‌گیری شد. پارامترهای عملکرد مورداستفاده شامل دقت^۱، صحت^۲، حساسیت^۳ و خاصیت^۴ و امتیاز F^۵ بود. پارامترهای عملکرد را می‌توان با استفاده از معادلات (۱۳) تا (۱۷) محاسبه کرد.

ماتریس ابهام		منفی	مثبت
منفی	TN	FP	
مثبت	FN	TP	

$$Accuracy = \frac{TP + TN}{TP + FN + FP + TN}, \quad (13)$$

$$Sensitivity = TP_{rate} = \frac{TP}{TP + FN}, \quad (14)$$

$$Specificity = TN_{rate} = \frac{TN}{TN + FP}, \quad (15)$$

$$Precision = \frac{TP}{TP + Fp}, \quad (16)$$

$$F_Score = \frac{2 * TP}{2 * TP + Fp + FN}. \quad (17)$$

۵. نتایج عملی

نتایج حاصل از ترکیب کاهش RN-DBSCAN و مدل CART بر روی دو مجموعه‌داده Kagg le و NSL-KDD با استفاده از نمودارهای عملکرد و زمان برای درصد کاهش‌های ۵٪، ۲۰٪، ۳۰٪، ۴۹٪، ۶۰٪ و در نهایت ۸۰٪، مورد بررسی و تحلیل قرار گرفته است. برای بررسی این نتایج به دلیل نیاز به محیط یکسان جهت بررسی زمان حاصل شده، از محیط colab، بهره گرفته شده است.

در این مقاله، دو مجموعه‌داده بهصورت مجزا مورد بررسی جهت بررسی تعیین اثر فرایند کاهش گرفتند. داده‌های کاهش با دو سناریو، بدون درنظر گرفتن برچسب داده و با درنظر گرفتن برچسب داده و گروه‌بندی داده‌ها بر اساس برچسب، برای انجام NSL-KDD بررسی شدند. همچنان از مجموعه‌داده فرایند کاهش از رویکرد با درنظر گرفتن برچسب داده و گروه‌بندی داده‌ها بر اساس برچسب، استفاده نموده است. زمانی که برچسب

¹ Accuracy

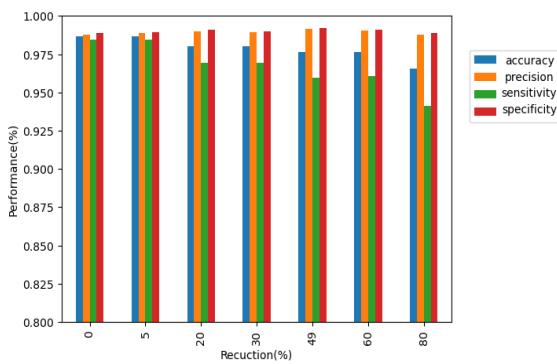
² Precision

³ Sensitivity

⁴ Specificity

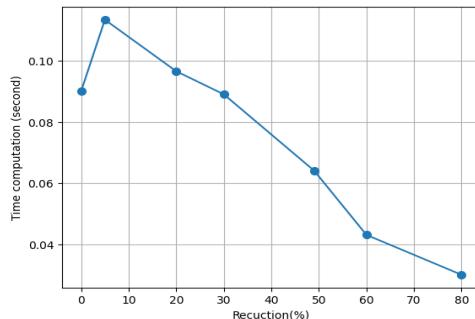
⁵ F_Score

عملکرد سیستم تشخیص نفوذ، بهویژه برای پارامتر خاصیت نسبتاً کم بود. پارامتر حساسیت، بیشترین کاهش را دارد، اما به طور کلی، دقت عملکرد سیستم تشخیص نفوذ با نرخ کاهش داده تا ۸۰٪ هنوز نسبتاً خوب است. پارامتر صحت، که توانایی سیستم تشخیص نفوذ در انتقال حمله را می‌سنجد [۲۶]، درصد بالای داشت. به طور کلی با افزایش میزان درصد کاهش داده‌ها، الگوی عملکرد به صورت خطی کاهش یافته است.



شکل (۴): تأثیر کاهش داده‌ها بر عملکرد سیستم تشخیص نفوذ با تفکیک کلاس، در مجموعه داده Kaggle

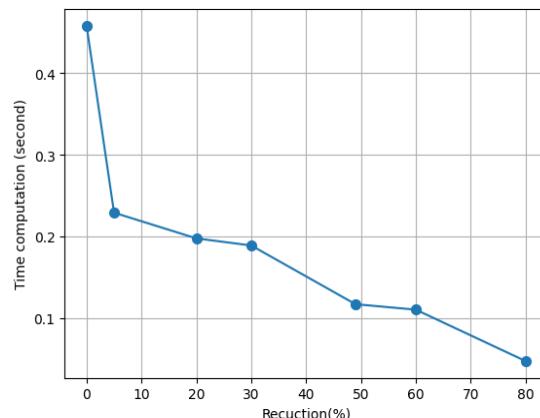
آزمون بعدی شامل محاسبه زمان، با درنظر گرفتن برچسب‌های داده‌ها می‌باشد. نتایج آزمون در شکل (۵) ارائه شده است.



شکل (۵): زمان محاسبات سیستم تشخیص نفوذ، با درنظر گرفتن برچسب‌ها، در مجموعه داده Kaggle

استفاده از کاهش داده‌ها با RN_DBSCAN می‌تواند زمان محاسبات را به طور قابل توجهی کاهش دهد و در عین حال عملکرد سیستم تشخیص نفوذ را حفظ کند. با توجه به شکل (۵) مشاهده می‌شود که با ۸۰٪ کاهش داده، زمان از ۹۰ ms به ۳۰ ms کاهش می‌یابد، در صورتی که دقت از ۹۸/۶۶ به ۹۶/۵۶٪ کاهش داشته است که نسبت به کاهش زمان، کم بوده و همچنان بالای ۹۵ درصد حفظ شده است. مقدار دقیق نتایج به دست آمده در جدول (۵) قرار داده شده است.

زمان انجام می‌شود. نتایج ارزیابی پارامترهای عملکرد محاسبات زمانی، برای طبقه‌بندی در سیستم تشخیص نفوذ، در شکل (۳) نشان داده شده است. همان‌طور که از شکل مشخص است؛ هر چه درصد کاهش داده‌ها بیشتر باشد، زمان محاسبه کمتر است. از منظر عملکرد، کاهش نسبتاً کم است، اما کاهش زمان محاسبات بسیار قابل توجه است. برای کاهش ۸۰ درصدی داده، زمان محاسبه از ۴۷/۲۱ ms به ۴۵۸/۰۹ ms کاهش یافت.



شکل (۳): زمان محاسبات سیستم تشخیص نفوذ، بدون درنظر گرفتن برچسب‌ها، در مجموعه داده Kaggle

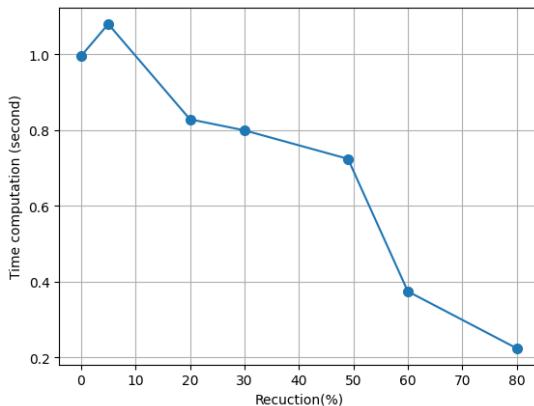
مقدار دقیق نتایج به دست آمده در جدول (۴) قرار داده شده است.

جدول (۴): مقادیر دقیق عملکرد سیستم تشخیص نفوذ، مجموعه داده Kaggle بدون درنظر گرفتن برچسب‌ها

کاهش داده (%)	دقت	خاصیت	حساسیت	صحبت	زمان (ms)
۰	۹۸/۶۶	۹۸/۸۶	۹۸/۴۴	۹۸/۸۰	۴۵۸/۰۹
۵	۹۸/۶۳	۹۸/۸۳	۹۸/۴۲	۹۸/۷۵	۲۲۹/۰۹
۱۰	۹۸/۳۷	۹۸/۰۹	۹۸/۶۶	۹۷/۹۹	۱۹۷/۶۱
۳۰	۹۸/۱۸	۹۸/۰۴	۹۸/۳۳	۹۷/۹۴	۱۸۸/۹۱
۴۹	۹۷/۳۴	۹۶/۰۹	۹۸/۶۶	۹۵/۹۸	۱۱۶/۹۲
۶۰	۹۶/۳۸	۹۴/۰۰	۹۸/۰۶	۹۴/۷۰	۱۱۰/۲۳
۸۰	۹۶/۱۲	۹۵/۸۸	۹۶/۳۷	۹۵/۶۸	۴۷/۲۱

۵-۲ نتایج بر روی مجموعه Kaggle با تفکیک کلاس

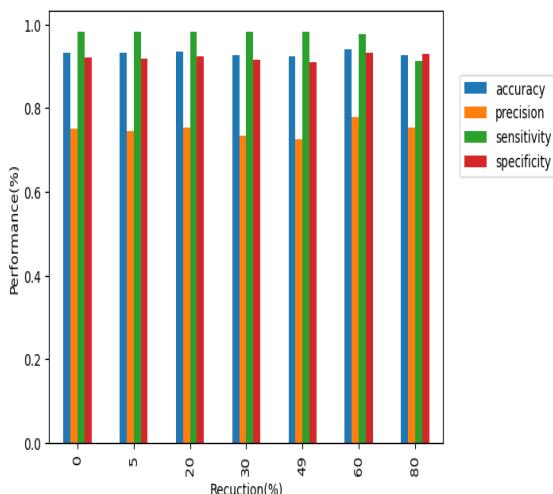
در شکل (۴)، مشخص است که افزایش درصد کاهش داده‌ها بر عملکرد سیستم تشخیص نفوذ تأثیر می‌گذارد. کاهش در



شکل (۷): زمان محاسبات سیستم تشخیص نفوذ، با درنظر گرفتن برچسب‌ها، در مجموعه داده NSL_KDD

همان‌طور که در شکل (۷) مشخص است، هر چه درصد کاهش داده‌ها بیشتر شود، پارامترهای عملکرد کاهش یافته و زمان محاسبات نیز کاهش می‌یابد. با توجه به جدول (۶) که مقادیر دقیق بدست آمده از ترکیب کاهش RN و RN_DBSCAN و اعمال مدل CART برای مجموعه داده NSL_KDD می‌باشد، مشاهده می‌شود که زمان محاسبات در ۰٪ کاهش، ۹۹۵/۲ ms بوده است که این مقدار برای ۸۰٪ کاهش به ۲۲۳/۶۰ ms رسیده است. همان‌طور که مشخص است کاهش در زمان محاسبات بسیار قابل توجه می‌باشد، درصورتی که کاهش در نتیجه عملکرد نسبتاً کم است و دقت بالای ۹۲/۳۶٪ حفظ شده است

شکل (۶): تأثیر کاهش داده‌ها بر عملکرد سیستم تشخیص نفوذ، در مجموعه داده مجموعه داده NSL-KDD



جدول (۵): مقادیر دقیق عملکرد سیستم تشخیص نفوذ، مجموعه داده Kaggle با تفکیک کلاس

کاهش داده (%)	دقت	خاصیت	حساسیت	صحت	زمان (ms)
۰	۹۸/۶۶	۹۸/۸۶	۹۸/۴۴	۹۸/۷۹	۹۰
۵	۹۸/۶۹	۹۸/۹۴	۹۸/۴۳	۹۸/۸۷	۱۱۲/۴۲
۲۰	۹۸/۰۴	۹۹/۰۸	۹۶/۹۵	۹۹	۹۶/۵۶
۳۰	۹۸	۹۹/۰۲	۹۶/۹۴	۹۸/۹۴	۸۹/۰۳
۴۹	۹۷/۶۴	۹۹/۲۳	۹۵/۹۶	۹۹/۱۶	۶۴/۰۶
۶۰	۹۷/۶۳	۹۹/۱۰	۹۶/۰۸	۹۹/۲	۴۳/۰۶
۸۰	۹۶/۵۶	۹۸/۹۰	۹۴/۱۱	۹۸/۷۷	۳۰

کاهش داده‌ها، با درنظر گرفتن تفکیک کلاس یا بدون درنظر گرفتن آن، تأثیری بر عملکرد سیستم ندارد، اما برای پیش‌بینی وقوع داده‌های نامتعادل تولید شده در فرایند کاهش، تفکیک کلاس بهتر است انجام شود. همچنین کاهش با درنظر گرفتن برچسب داده‌ها، مشکل عدم چگالی یکسان در سطح داده [29] که یکی از موانع تنظیم صحیح پارامترهای الگوریتم DBSCAN می‌باشد را تا حدی حل می‌نماید.

۵-۳ نتایج بر روی مجموعه NSL_KDD با تفکیک کلاس

باتوجه به نمودار عملکرد سیستم تشخیص نفوذ در شکل (۶)، پارامتر دقت که مهم‌ترین معیار ارزیابی عملکرد سیستم می‌باشد، برای درصد کاهش‌های مختلف تا ۸۰٪ کاهش، بالای ۹۰٪ حفظ شده است. همچنین معیار حساسیت، از صفر تا چهل و نه درصد کاهش، بالای ۹۸٪ بوده و عملکرد بسیار مناسبی دارد. معیار صحت نسبت به پارامترهای دیگر ارزیابی، از عملکرد ضعیفتری برخوردار بوده است، اما به طور کلی می‌توان گفت، در این مجموعه داده نیز با کاهش ۸۰ درصدی داده‌ها، عملکرد سیستم به خوبی حفظ شده است.

معیار مورد ارزیابی بعدی برای این مجموعه داده، بررسی زمان، برای درصدهای مختلف کاهش می‌باشد. محاسبه زمان برای این مجموعه داده در شکل (۷) نشان داده شده است.

۴-۵- کارهای مشابه در سال‌های اخیر

در سال‌های اخیر بسیاری از تحقیقات و مطالعات در حوزه سیستم‌های تشخیص نفوذ و کاهش داده‌ها انجام شده است. هدف این تحقیقات بررسی روش‌های مختلف برای کاهش داده‌ها در سیستم تشخیص نفوذ بوده است.

بر اساس نتایج جدول (۷)، روش پیشنهادی مورد بررسی در این مقاله نشان داده است که در مقایسه با روش‌های قبلی، عملکرد بهتری در کاهش داده‌ها در سیستم تشخیص نفوذ داشته است.

با توجه به جدول ارائه شده، این روش علاوه بر افزایش سرعت، توانسته است دقت و عملکرد را بهبود بخشد و نتایج بهتری را به دست آورد. با استناد به این مقاله و تحقیقات گذشته، می‌توان این نتیجه را گرفت که روش پیشنهادی در این مقاله، به عنوان یک جایگزین قوی برای روش‌های قبلی در حوزه کاهش داده‌ها در سیستم تشخیص نفوذ می‌تواند مورد استفاده قرار گیرد.

جدول (۶): مقادیر دقیق عملکرد سیستم تشخیص نفوذ

مجموعه داده NSL_KDD با تفکیک کلاس

زمان (ms)	صحت	حساسیت	خاصیت	دقت	کاهش داده (%)
۹۹۵/۲	۷۹/۹۵	۹۸/۲۲	۹۲/۰۶	۹۳/۲۶	۰
۱۰۸۰/۹	۷۴/۵۲	۹۸/۲۱	۹۱/۸۸	۹۳/۱۱	۵
۸۲۸/۳۱	۷۵/۳۲	۹۸/۱۵	۹۲/۲۲	۹۳/۳۷	۲۰
۷۹۹/۱۸	۷۳/۴۷	۹۸/۱۶	۹۱/۴۲	۹۲/۷۴	۳۰
۷۲۳/۹۶	۷۲/۴۰	۹۸/۱۹	۹۰/۹۵	۹۲/۳۶	۴۹
۳۷۳/۹۳	۷۷/۷۴	۹۷/۶۴	۹۳/۲۴	۹۴/۱	۶۰
۲۲۳/۶۰	۷۵/۴۲	۹۱/۲۱	۹۲/۸۰	۹۲/۵۱	۸۰

جدول (۷): مقایسه با تحقیقات قبلی

مطالعه	روش	رویکرد کاهش داده‌ها	مجموعه داده	امتیاز اف	دقت
باتاجار و همکاران [30]	SVM	PCA	Kaggle	-	95/2
ویهارت و همکاران [19]	NB	PCA	Kaggle	-	75/3
شهرآ و همکاران [32]	NB	PCA+Firefly	Kaggle	-	84/2
سارکر و همکاران [31]	NB	-	Kaggle	90/0	90/0
عبدالله و همکاران [33]	CART	m_DBSCAN	Kaggle	89/43	92/3
صبری و همکاران [34]	J48	CFS	NSL_KDD	-	86/1
کاسونگو [35]	Voting classifier	IG-Filters	NSL_KDD	82/3	86/67
پیشنهاد شده	DT	PIO	NSL_KDD	88/2	88/3
	RNN+ LSTM+ GRU	XGBoost	NSL_KDD	99/58	88/13
پیشنهاد شده	CART	RN_DBSCAN	Kaggle	96/56	96/39
			NSL_KDD	82/60	92/51

با این روش، داده‌های مهم برای آموزش مدل حفظ می‌شوند و داده‌های یکنواخت و همگن حذف می‌شوند.

برای تنظیم پارامترهای الگوریتم RN_DBSCAN، از الگوریتم ژنتیک بهره برده شده است. همچنین برای طبقه‌بندی داده‌های کاهش یافته، از الگوریتم طبقه‌بندی CART استفاده شده است. نتایج حاصل از پژوهش نشان می‌دهد که سیستم تشخیص نفوذ با استفاده از الگوریتم RN_DBSCAN قادر است حجم داده‌ها را تا ۸۰٪ کاهش دهد، درحالی که عملکرد سیستم باقت بالای ۹۰٪ حفظ می‌شود. علاوه بر این، استفاده از این الگوریتم

پژوهش‌های گذشته نشان داده است که استفاده از روش‌های خوشه‌بندی می‌تواند به کاهش داده‌ها کمک کند و به عنوان پایه‌ای برای کاهش حجم داده‌ها عمل کند. در این پژوهش به‌منظور بهبود سیستم تشخیص نفوذ و در نتیجه آن بهبود امنیت شبکه‌های کامپیوتری، کاهش حجم داده‌ها با استفاده از روش RN_DBSCAN مطرح گردید. این روش بر اساس رویکرد چگالی کار می‌کند و با یافتن نقاط اطراف نقاط مرکزی که بسیار مشابه و نزدیک نقاط مرکزی هستند، این نقاط را کاهش می‌دهد.

۶. نتیجه‌گیری

تشخیص نفوذ می‌شود و در عین حال حجم داده‌ها را به طور آگاهانه کاهش می‌دهد تا داده‌های مهم برای آموزش مدل حفظ شوند. این پژوهش می‌تواند به توسعه و پیشرفت در حوزه امنیت شبکه‌های کامپیوتروی کمک کند.

[16] S. Ougiaroglou and G. Evangelidis, "RHC: a nonparametric cluster-based data reduction for efficient k-NN classification," *Pattern Analysis and Applications*, vol. 19, no. 1, pp. 93-109, 2016.

[17] S. Ougiaroglou, K. I. Diamantaras, and G. Evangelidis, "Exploring the effect of data reduction on neural network and support vector machine classification," *Neurocomputing*, vol. 208, pp. 101-110, 2018.

[18] S. Ougiaroglou and G. Evangelidis, "Efficient editing and data abstraction by finding homogeneous clusters," *Annals of Mathematics and Artificial Intelligence*, vol. 76, no. 3-4, pp. 327-349, 2016.

[19] A. K. Wicaksana and D. E. Cahyani, "Modification of a density-based spatial clustering algorithm for applications with noise for data reduction in intrusion detection systems," *International Journal of Fuzzy Logic and Intelligent Systems*, vol. 21, no. 2, pp. 189-203, 2021.

[20] F. O. Ozkok and M. Celik, "A new approach to determine Eps parameter of DBSCAN algorithm," *International Journal of Intelligent Systems and Applications Engineering*, vol. 5, no. 4, pp. 247-251, 2017.

[21] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 1996, pp. 226-231.

[22] M. Zbigniew, "Genetic algorithms + data structures = evolution programs," *Comput Stat*, 1996.

[23] L. Rutkowska, M. Jaworska, L. Pietruczuka, and P. Duda, "The CART Decision Tree for Mining Data Streams," *Information Sciences*, vol. 266, pp. 1-15, 2014.

[24] R. A. Johnson and D. W. Wichern, *Applied Multivariate Statistical Analysis*. Upper Saddle River, NJ: Pearson Prentice Hall, 2007.

[25] P. I. Radoglou-Grammatikis and P. G. Sarigiannidis, "An anomaly-based intrusion detection system for the smart grid based on CART decision tree," in *Proceedings of the Global Information Infrastructure and Networking Symposium (GIIS)*, 2018.

[26] H. H. Patel and P. Prajapati, "Study and analysis of decision tree based classification algorithms," *International Journal of Computer Sciences and Engineering*, vol. 6, no. 10, pp. 74-78, 2018.

[27] A. Priyam, G. R. Abhijeeta, A. Rathee, and S. Srivastava, "Comparative analysis of decision tree classification algorithms," *International Journal of Current Engineering and Technology*, vol. 3, no. 2, pp. 334-337, 2013.

[28] H. M. Sani, C. Lei, and D. Neagu, "Computational complexity analysis of decision tree algorithms," in *Proceedings of the Artificial Intelligence XXXV*, 2018, pp. 191-197.

[29] A. Zadehbalaei, A. Bagheri, and H. Afshar, "A study on DBSCAN Clustering algorithm issues and a survey on its improvements," *Soft Computing Journal*, vol. 6, 2021, pp. 2322-3707.

[30] S. Bhattacharya, P.K.R. Maddikunta, R. Kaluri, S. Singh, T.R. Gadekallu, M. Alazab, U. Tariq, "A novel PCA-firefly based XGBoost classification model for intrusion detection in networks using GPU," *Electronics*, vol. 9, no. 2, 2020.

[31] I.H. Sarker, Y.B. Abusark, F. Alsolami, A.I. Khan, "IntrudTree: a machine learning based cyber security intrusion detection model," *Symmetry*, vol. 12, no. 5, 2020.

[32] M.B. Shahbaz, X. Wang, A. Behnad, J. Samarabandu, "On efficiency enhancement of the correlation-based feature selection

منجر به کاهش زمان محاسباتی و کاهش پیچیدگی و هزینه می‌شود. همچنین، مقایسه الگوریتم‌های مختلف نشان داده است که RN_DBSCAN نسبت به الگوریتم‌های کاهش سطربی دیگر مانند Sampling و Slice و حتی الگوریتم m_DBSCAN کارایی بهتری دارد. این روش منجر به بهبود عملکرد و سرعت سیستم

۷. مراجع

- [1] E. A. Shammar and A. T. Zahary, "The Internet of Things (IoT): a survey of techniques, operating systems, and trends," *Library Hi Tech*, vol. 38, no. 1, pp. 5-66, 2020.
- [2] H. Tanha and M. Abbasi, "Traffic identification in Internet of Things infrastructure using neural networks and deep learning," *Scientific Journal of Electronic and Cyber Defense*, vol. 11, no. 2, pp. 1-13, Jul. 2023. (in Persian).DOR:20.1001.1.23224347.1402.11.2.1.4
- [3] Y. Zhang, Y. Zhang, T. Chen, and B. Xia, "Internet of Things (IoT) security: A survey," *Journal of Information Security and Applications*, vol. 50, p. 102419, 2020.
- [4] J. Mazloom and H. Bigdeli, "Optimized combination deep neural network integrated with feature selection for intrusion detection system in cyber-attacks," *Journal of Information Security and Applications*, vol. 10, no. 4, pp. 41-51, Feb. 2022. (in Persian).DOR:20.1001.1.23224347.1401.10.4.5.5
- [5] W. Zhang and S. Li, "A Deep Learning Approach for Intrusion Detection System," *IEEE Access*, vol. 9, pp. 35470-35479, 2021.
- [6] M. H. Nataj Salhaddar, "Investigating a New Hybrid Approach for Intrusion Detection System on Various Datasets," *Journal of Information Security and Applications*, vol. 10, no. 3, pp. 43-57, Dec. 2022. (in Persian).DOR:20.1001.1.23224347.1401.10.3.5.3
- [7] I. Ahmed, M. Mahfuzul Islam, and A. A. Adewole, "A survey of intrusion detection techniques in cloud computing," *Journal of Network and Computer Applications*, vol. 36, no. 1, pp. 42-57, 2013.
- [8] M. S. Farash and S. Samet, "Feature Selection for Intrusion Detection Systems: A Comprehensive Review," *Computer Networks*, vol. 74, pp. 443-460, 2014.
- [9] A. A. Wiharto and U. Permana, "Improvement of performance intrusion detection system (IDS) using artificial neural network ensemble," *Journal of Theoretical and Applied Information Technology*, vol. 80, no. 2, pp. 191-201, 2015.
- [10] D. Uhm, S. H. Jun, and S. J. Lee, "A classification method using data reduction," *International Journal of Fuzzy Logic and Intelligent Systems*, vol. 12, no. 1, pp. 1-5, 2012.
- [11] N. A. Le-Khac, M. Bue, M. Whelan, and M. T. Kechadi, "A clustering-based data reduction for very large spatiotemporal datasets," in *Advanced Data Mining and Applications*, 2010, pp. 43-54.
- [12] J. Wang, S. Yue, X. Yu, and Y. Wang, "An efficient data reduction method and its application to cluster analysis," *Neurocomputing*, vol. 238, pp. 234-244, 2017.
- [13] A. C. Benabdellah, A. Benghabrit, and I. Bouhaddou, "A survey of clustering algorithms for an industrial context," *Procedia Computer Science*, vol. 148, pp. 291-302, 2019.
- [14] J. M. Dudik, A. Kurosu, J. L. Coyle, and E. Sejdic, "A comparative analysis of DBSCAN, K-means, and quadratic variation algorithms for automatic identification of swallows from swallowing accelerometry signals," *Computers in Biology and Medicine*, vol. 59, pp. 10-18, 2015.
- [15] S. Ougiaroglou and G. Evangelidis, "Efficient dataset size reduction by finding homogeneous clusters," in *Proceedings of the 5th Balkan Conference in Informatics*, 2012, pp. 168-173.

for intrusion detection systems," in Proceedings of the 2016 IEEE 7th Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON), Vancouver, BC, Canada, 2016, pp. 1–7.

[33] M. Abdullah, A. Balamash, A. Alshannaq, S. Almabdy, "Enhanced Intrusion Detection System using Feature Selection Method and Ensemble Learning Algorithms," International Journal of Computer Science and Information Security (IJCSIS), vol. 16, no. 2, 2018, pp. 48–55.

[34] H. Alazzam, A. Sharieh, K.E. Sabri, "A feature selection algorithm for intrusion detection system based on Pigeon Inspired Optimizer," Expert Systems with Applications, vol. 148, 2020, 113249.

[35] S. M. Kasongo, "A deep learning technique for intrusion detection system using a recurrent neural networks based framework," Computer Communications, vol. 199, no. 1, pp. 113–125, 2023