

IoT Dataset Validation Using Machine Learning Techniques for Traffic Anomaly Detection

Laura Vigoya * , Diego Fernandez , Victor Carneiro and Francisco J. N  voa

Centre for Information and Communications Technology Research (CITIC), Campus de Elvi  a s/n, 15071 A Coru  a, Spain; dfernandez@udc.es (D.F.); victor.carneiro@udc.es (V.C.); fjnovoa@udc.es (F.J.N.) * Correspondence: l.v.vigoya@udc.es; Tel.: +34-881-011-213

Abstract: With advancements in engineering and science, the application of smart systems is increasing, generating a faster growth of the IoT network traffic. The limitations due to IoT restricted power and computing devices also raise concerns about security vulnerabilities. Machine learning-based techniques have recently gained credibility in a successful application for the detection of network anomalies, including IoT networks. However, machine learning techniques cannot work without representative data. Given the scarcity of IoT datasets, the DAD emerged as an instrument for knowing the behavior of dedicated IoT-MQTT networks. This paper aims to validate the DAD dataset by applying Logistic Regression, Naive Bayes, Random Forest, AdaBoost, and Support Vector Machine to detect traffic anomalies in IoT. To obtain the best results, techniques for handling unbalanced data, feature selection, and grid search for hyperparameter optimization have been used. The experimental results show that the proposed dataset can achieve a high detection rate in all the experiments, providing the best mean accuracy of 0.99 for the tree-based models, with a low false-positive rate, ensuring effective anomaly detection

Keywords: IoT; sensors; dataset validation; machine learning; intrusion detection systems; analysis; metric; algorithm design

اعتبارسنجی مجموعه‌داده اینترنت اشیا با استفاده از روش‌های یادگیری ماشین برای تشخیص ناهنجاری ترافیک

لائورا ویگویا*، دیگو فرناندز، ویکتور کارنیرو و فرانسیسکو خ. نووا

استناد: ویگویا، ل؛ فرناندز، د؛ کارنیرو، و؛ نووا، ف.خ. اعتبارسنجی مجموعه‌داده اینترنت اشیا با استفاده از روش‌های یادگیری ماشین برای تشخیص ناهنجاری ترافیک. الکترونیکس ۲۰۲۱، ۱۰، ۲۸۵۷.

<https://doi.org/10.3390/electronics10222457>

سردبیر علمی: دیمیترای. کاکلامانی

دریافت: ۲۹ زوئن ۲۰۲۱

پذیرش: ۱۶ نوامبر ۲۰۲۱

انتشار: ۱۹ نوامبر ۲۰۲۱

یادداشت ناشر: انتشارات MDPI در مورد ادعاهای حوزه‌ای موجود در نقشه‌های منتشرشده و وابستگی‌های مؤسسه‌ای بی‌طرف باقی می‌ماند.

حق چاپ: © ۲۰۲۱ توسط نویسنندگان.

ناشر تحت پروانه MDPI، بازل، سوئیس.

این مقاله یک مقاله دسترسی‌آزاد است که تحت شرایط و ضوابط مجوز Creative Commons Attribution (CC BY) منتشر شده است. (<https://creativecommons.org/licenses/by/4.0/>)

مرکز تحقیقات فناوری اطلاعات و ارتباطات (CITIC)، کمپوس ال وینیا، شماره‌ای ۱۵۰۷۱ آ کورونیا، اسپانیا؛ آسیب‌پذیری‌های امنیتی ایجاد کرده است. روش‌های مبتنی بر یادگیری ماشین (ML) اخیراً در کاربردهای موفقیت‌آمیز برای تشخیص ناهنجاری‌های شبکه، از جمله شبکه‌های IoT، معتبر شده‌اند. با این حال، روش‌های یادگیری ماشین بدون داده‌های نماینده نمی‌توانند عمل کنند. با توجه به کمبود مجموعه‌داده‌های اختصاصی IoT، مجموعه‌داده DAD به عنوان ابزاری برای شناخت رفتار شبکه‌های اختصاصی IoT-MQTT مطرح شده است. این مقاله با هدف اعتبارسنجی مجموعه‌داده DAD با به کارگیری رگرسیون لجستیک (Logistic Regression)، ناوی بیز (Naive Bayes)، جنگل تصادفی (Random Forest)، آدابوست (AdaBoost) و ماشین بردار پشتیبان (Support Vector Machine) برای تشخیص ناهنجاری‌های ترافیکی در IoT ارائه شده است. برای به دست آوردن بهترین نتایج، از تکنیک‌های مدیریت داده‌ای نامتوازن، انتخاب ویژگی و جستجوی شبکه‌ای (grid search) برای بهینه‌سازی هایپرپارامترها استفاده شده است. نتایج آزمایش‌ها نشان می‌دهد که مجموعه‌داده پیشنهادی می‌تواند نرخ تشخیص بالایی را در تمام آزمایش‌ها به دست آورد و بهترین دقت میانگین ۰.۹۹ را برای مدل‌های مبتنی بر درخت، همراه با نرخ کم خطای مشبت کاذب-positive rate، فراهم آورد که این امر یک تشخیص مؤثر از ناهنجاری‌ها را تضمین می‌کند.

وازگان کلیدی: اینترنت اشیا(IoT)؛ حسگرها؛ اعتبارسنجی مجموعه‌داده؛ یادگیری ماشین؛ سیستم‌های تشخیص نفوذ؛ تحلیل؛ معیار؛ طراحی الگوریتم

۱. مقدمه

با توجه به رشد و پیاده‌سازی گستردگی شبکه‌های اینترنت اشیا (IoT) در حوزه‌های مختلف، این فناوری نقش چشمگیری در فعالیت‌های روزمره‌ی ما ایفا می‌کند و به عنوان نیروی محركه‌ی خانه‌های هوشمند، شهرهای هوشمند، سیستم‌های سلامت نوین و تولید پیشرفته تحول یافته است [۱]. با این حال، با پیشرفت‌های حاصل در مهندسی و علوم، ابعاد کاربردی سیستم‌های هوشمند در حال افزایش است که رشد سریع‌تری در ترافیک شبکه‌ی IoT ایجاد می‌کند و همین امر دغدغه‌هایی در مورد آسیب‌پذیری‌ها و محدودیت‌های امنیتی نیز به همراه دارد. اینترنت و کاربران آن از پیش تحت مداومت حملاتی هستند که می‌توانند به صورت ناهنجاری‌های ترافیکی خود را نشان دهند و تهدیدی برای IoT به عنوان مجموعه‌ای از دستگاه‌های محدود محسوب شوند. این واقعیت می‌تواند دیدگاه‌های متفاوتی را ایجاد کند که به تولید الگوهای مخرب جدید و خلاقانه منجر شود. چالش اصلی جلوگیری از گسترش این الگوها یا حداقل کاهش و محدود کردن تأثیر آن‌هاست [۲].

Nahنجاری‌های شبکه همیشه قابل دسته‌بندی به عنوان یک حمله نیستند و لزوماً عناصر مضر محسوب نمی‌شوند؛ با این وجود، بینش‌های مهمی را درباره‌ی رفتار ترافیک فراهم می‌کنند و می‌توانند در شناسایی اطلاعات مهم و حیاتی در کاربردهای مختلف کمک کننده باشند [۳]. یکی از راه‌های تشخیص تغییرات در رفتار شبکه، استفاده از سیستم‌های تشخیص نفوذ (IDSs) است که به کشف، تعیین و شناسایی استفاده‌های غیرمجاز، تکثیر، دستکاری و تخرب سیستم‌های اطلاعاتی کمک می‌کنند. روش‌های سنتی IDS به دلیل ویژگی‌های خاص IoT—مانند انرژی محدود، همه‌گیری (ubiquity)، ناهمگونی، پهنای باند محدود و اتصال‌پذیری جهانی—برای امنیت سیستم‌های IoT کمتر مؤثر یا ناکافی هستند [۴]. در واقع، IoT نیازمند استانداردها و پروتکل‌های ارتباطی تخصصی برای مقابله با چالش‌های ناشی از این ویژگی‌هاست. در لایه انتقال، پروتکل‌های UDP و TCP برای اکثر کاربردها غالب هستند. با این حال، این توابع باید به شیوه‌های استاندارد MQTT اغلب برای همکاری پیاده‌سازی شوند و بسته به نیازهای کاربردی IoT، توابع توزیع پیام‌های مختلفی مورد نیاز است. پروتکل MQTT برای انتشار و/یا اشتراک‌گذاری پیام‌های دریافتی از دستگاه‌هایی با پهنای باند پایین در شبکه‌های غیرقابل اعتماد به کار می‌رود. از طریق TCP برای ارسال و دریافت داده‌های لایه‌ی واسطه (middleware) به مراکز IoT عمل می‌کند.

یک جنبه‌ی حیاتی این است که یادگیری ماشین بدون داده‌های نماینده نمی‌تواند عمل کند و صحت داده‌های جمع‌آوری‌شده از دستگاه‌ها برای ساخت سیستم‌های تصمیم‌گیری هوشمند و همچنین مدیریت مؤثر محیط‌های IoT ضروری است. روش‌های نظارت‌شده (که همچنین به عنوان روش‌های طبقه‌بندی شناخته می‌شوند) به یک مجموعه آموزشی برچسب‌دار نیاز دارند که شامل نمونه‌های عادی و ناهنجار باشد تا مدل پیش‌بینی‌کننده ساخته شود. از نظر نظری، روش‌های نظارت‌شده نرخ تشخیص بهتری را نسبت به روش‌های نیمه‌نظارت‌شده و بدون نظارت ارائه می‌دهند، زیرا به اطلاعات بیشتری دسترسی دارند. رایج‌ترین الگوریتم‌های نظارت‌شده، شبکه‌های k-Nearest Neighbors (k-NN)، مارکوفیانه‌ای بردار پشتیبان (Support Vector Machines)، شبکه‌های نوری (Neural Networks)، شبکه‌های بیزی (Bayesian Networks) و درخت‌های تصمیم (Decision Trees) هستند [۵].

۱.۱. مانگیزه

به دلیل تنظیمات درهم‌تنیده و وابسته‌ی اینترنت اشیا (IoT)، این محیط مستعد مسائل امنیتی و حریم خصوصی مختلفی است که ممکن است منجر به انجام وظایف غیرمجاز توسط کاربران مخرب از راه دور شود [۶] و در نتیجه، رفتار غیرعادی در شبکه ایجاد کند. در موارد دیگر، ترافیک ناهنجار در یک شبکه‌ی IoT ممکن است ناشی از پیکربندی اشتباه، نصب غیرمعمول یا خرابی سخت‌افزاری حسگرها باشد.

در نتیجه، نیاز به تحقیق درباره‌ی روش‌های تشخیص و پیشگیری خاص و طراحی و توسعه‌ی راه حل‌های امنیتی هوشمند برای محافظت از دستگاه‌های IoT آسیب‌پذیر و در برابر ناهنجاری‌های ناشی از دستگاه‌های IoT فشرده شده وجود دارد. ایجاد سیستم‌های تشخیص ناهنجاری (IDSS) نیازمند دانش نماینده از محیط است.

با توجه به کمبود مجموعه‌داده‌های IoT، مجموعه‌داده‌ی DAD [۷] به عنوان یک مکانیزم برای شناخت رفتار دستگاه‌های سبک وزن و اختصاصی در شبکه‌های IoT-MQTT مطرح شده است. مجموعه‌داده‌ی DAD [۷] یک مجموعه‌داده‌ی کامل و برچسب‌گذاری شده برای تشخیص ناهنجاری ترافیک در دنیای واقعی است که در شرایط مناسب، با حجم کافی ردپا (trace)، سناریوهای متنوع ناهنجاری و تحلیل پیشین ارائه شده و قرار است در الگوریتم‌های یادگیری ماشین (ML) به کار گرفته شود. با این حال، این مجموعه‌داده باید تحت کاربرد مدل‌های یادگیری ماشین قرار گیرد تا قابلیت استفاده‌ی آن برای هدفی که برای آن طراحی شده است، تأیید شود.

۱.۲. سهم‌های این پژوهش

این پژوهش با هدف نشان دادن این موضوع انجام شده است که مجموعه‌داده‌ی DAD را می‌توان برای تشخیص ناهنجاری‌ها در شبکه‌های ترافیکی MQTT-IoT به کار برد و برای این منظور از پنج روش رایج یادگیری سطحی (shallow learning) برای تشخیص ناهنجاری ترافیک در محیط‌های IoT استفاده شده است. سهم‌های اصلی این پژوهش عبارتند از:

- تحلیل و استخراج ویژگی‌ها برای شبکه‌های حسکر بی‌سیم IoT؛
- پیاده‌سازی تکنیک‌های آماده‌سازی داده، شامل گروه‌بندی (binning)، مدیریت داده‌های نامتوازن و استخراج ویژگی در محیط‌های واقعی؛
- به کار گیری الگوریتم‌های مختلف یادگیری ماشین برای اعتبارسنجی مجموعه‌داده‌ی انتخاب شده؛
- استفاده از معیارها و مقایسه‌ی عملکرد روش‌های مختلف یادگیری ماشین در تشخیص ناهنجاری‌های ترافیکی.

۱.۳. ساختار مقاله

این پژوهش به صورت زیر سازمان‌دهی شده است: مروری بر سیستم‌های تشخیص نفوذ (IDSS) مبتنی بر روش‌های یادگیری ماشین برای شبکه‌های اینترنت اشیا (IoT) در بخش ۲ ارائه می‌شود. بخش ۳ تعریف مختصه‌ی از الگوریتم‌های انتخاب شده را همراه با پارامترهای در نظر گرفته شده برای تنظیم های پر پارامترها ارائه می‌دهد. بخش ۴ شرح مختصه‌ی از سناریوی مورد استفاده را ارائه می‌کند و داده‌های مربوط به مجموعه‌داده را توصیف، بررسی و کیفیت آن‌ها را تأیید می‌نماید. سپس، ارتقای کیفیت داده‌ها به سطح مورد نیاز توسط روش‌های یادگیری ماشین انتخاب شده ضروری است. بخش ۵ تکنیک‌ها و فرآیندهای لازم برای این‌که داده‌ها به صورت کامل برای الگوریتم‌های یادگیری ماشین قابل فهم باشند را معرفی می‌کند. بخش ۶ نتایج بدست‌آمده برای هر یک از مدل‌ها در مرحله آموزش را تحلیل می‌کند. در ادامه، بخش ۷ نتایج بدست‌آمده در مرحله آزمون را بررسی کرده و همچنین مقایسه‌ای میان مدل‌ها انجام می‌دهد. در نهایت، بخش ۸ نتیجه‌گیری‌های اصلی این پژوهش را ارائه می‌دهد.

۲. کارهای مرتبط

پیاده‌سازی نسبتاً ساده و نتایج حاصل از سیستم‌های تشخیص نفوذ (IDSS) مبتنی بر الگوریتم‌های یادگیری ماشین، باعث شده است که بسیاری از پژوهشگران این روش‌ها را به عنوان ابزاری برای توسعه پروژه‌های خود انتخاب کنند. با این حال، کاربرد IDSS در محیط IoT به تازگی به صورت گستردگای مورد توجه قرار گرفته است. این بخش مروری مختصه بر مهم‌ترین روش‌های یادگیری ماشین به کار رفته برای اعتبارسنجی مجموعه‌داده‌های IoT در سال‌های اخیر ارائه می‌دهد.

با توجه به کمبود اولیه مجموعه‌داده‌های عمومی اختصاصی IoT، کاربرد IDS در این محیط با تطبیق مجموعه‌داده‌های عمومی که توسط جامعه پذیرفته شده بودند، بر روی رده‌پاهای شبکه IoT آغاز شد. متداول‌ترین مجموعه‌داده به کاررفته در IoT، مجموعه‌داده UNSW-NB15 [۸] بوده است که پژوهشگران متعددی کاربردهای یادگیری ماشین خود را بر این مجموعه‌داده آزمایش کرده‌اند. کورونیوتیس و همکاران [۹] از چهار روش طبقه‌بندی، شامل درخت تصمیم (DT)، داده‌کاوی مبتنی بر قواعد انجمنی (ARM)، شبکه عصبی مصنوعی (ANN) و ناوی بیز (NB) استفاده کردند تا توانایی الگوریتم‌ها در شناسایی بردارهای حمله را اعتبارسنجی نمایند. AD-IoT [۱۰] نیز بخشی از مجموعه‌داده UNSW-NB15 را انتخاب کرده و سیستمی برای تشخیص ناهنجاری جهت شناسایی سایبر‌حمله‌ها در گره‌های فوگ یک شهر هوشمند با استفاده از روش‌های گروهی (ensemble) مبتنی بر جنگل تصادفی (RF) و درخت اضافی (ET) پیشنهاد داده است. علاوه بر این، مجموعه‌داده‌های NIMS botnet UNSW-NB15 و حاوی داده‌های حسگرهای IoT شبیه‌سازی شده توسط مصطفی و همکاران [۱۱] برای پیاده‌سازی یک سیستم تشخیص نفوذ شبکه (NIDS) مبتنی بر الگوریتم یادگیری گروهی AdaBoost با به کارگیری سه روش یادگیری ماشین—DT، NB و ANN—مورد استفاده قرار گرفت. از سوی دیگر، پاجووه و همکاران [۱۱] از NSL-KDD استفاده کرده و مدلی را پیشنهاد دادند که از یک مازول طبقه‌بندی دو سطحی مشکل از ناوی بیز (NB) و نسخه ضریب اطمینان (Certainty Factor) الگوریتم k نزدیک‌ترین همسایه (K-Nearest Neighbor) برای شناسایی فعالیت‌های مخرب نظیر حملات کاربر به ریشه (U2R) و راه دور به محلی (R2L) بهره می‌برد. با گسترش کاربرد محیط‌های IoT، تقاضا برای سیستم‌های تشخیص نفوذ تخصصی افزایش یافت و در نتیجه، نیاز به انتشار مجموعه‌داده‌های اختصاصی که امکان اعتبارسنجی سیستم‌ها را فراهم کند، بیشتر شد. سیستم کیتسون (Kitsune) [۱۲] یک سیستم تشخیص نفوذ شبکه است که از یک گروه از شبکه‌های عصبی به نام اتوانکندر (autoencoder) برای تمایز الگوهای ترافیک عادی و غیرعادی استفاده می‌کند و با به کارگیری مجموعه‌داده N-BaIoT [۱۳] توسعه یافته است. مجموعه‌داده N-BaIoT ترافیک را از دو شبکه جمع‌آوری می‌کند: شبکه نظارت بروی‌دئوی با دوربین‌های IP که در آن هشت نوع مختلف حمله علیه در دسترس بودن و بکار چگی پیوندهای ویدئویی انجام شده است؛ و شبکه IoT شامل سه رایانه شخصی و نُه دستگاه IoT که با بدافزار باتنت میرای (Mirai) آلوده شده‌اند. علاوه بر این، عباسی [۱۴] از این مجموعه‌داده برای طبقه‌بندی با استفاده از رگرسیون لجستیک (LR) و شبکه عصبی مصنوعی (ANN) استفاده کرده است. دوشی و همکاران [۱۵] با تکیه بر این واقعیت که ترافیک IoT اغلب متفاوت از دستگاه‌های دیگر متصل به اینترنت است، یک خط‌لوله یادگیری ماشین را توسعه داده‌اند که ویژگی‌هایی را طراحی می‌کند تا از رفتارهای شبکه‌ای خاص IoT بهره ببرد. این خط‌لوله برای اجرا روی میان‌جعبه‌های شبکه (مثل مسیریاب‌ها، فایروال‌ها یا سوئیچ‌های شبکه) طراحی شده تا ترافیک ناهنجار و دستگاه‌های مرتبط که ممکن است بخشی از یک باتنت فعال باشند را شناسایی کند. در این پژوهش، نویسنده‌گان طیفی از رده‌بندان را برای تشخیص حمله مقایسه کرده‌اند که شامل الگوریتم k نزدیک‌ترین همسایه مبتنی بر درخت (KDTTree)، ماشین بردار پشتیبان با کرنل خطی (LSVM)، درخت تصمیم با معیار ناخالصی جینی (DT)، جنگل تصادفی با معیار ناخالصی جینی (RF) و شبکه عصبی (NN) می‌شود و اثربخشی جنگل تصادفی را نشان می‌دهند. حسن و همکاران [۱۷] با استناد به مجموعه‌داده خاص IoT ارائه شده توسط پاهل و همکاران [۱۶]، یک سیستم یادگیری ماشین هوشمند، امن و قابل اعتماد را پیشنهاد داده‌اند که از LR، DT، SVM، RF و ANN استفاده می‌کند. زیرساخت مبتنی بر IoT از چندین رده‌بند استفاده می‌کند که می‌تواند سیستم را هنگامی که در حالت غیرعادی قرار دارد تشخیص دهد و از آن محافظت نماید. بهترین عملکرد در دقت آموزش و آزمون با RF و ANN به دست آمده است. متأسفانه، این مجموعه‌داده تنها شامل یک روز از ترافیک است. کورونیوتیس و همکاران [۱۸] که پیش‌تر در زمینه کاربرد یادگیری ماشین در سیستم‌های تشخیص نفوذ فعالیت داشتند، مجموعه‌داده خاصی را برای باتنت در IoT با نام Bot-IoT طراحی کردند تا اثربخشی روش‌ها در این محیط را تأیید کنند. ساختار این مجموعه‌داده هم ترافیک عادی مرتبط با IoT و هم ترافیک شبکه‌های دیگر را دربرمی‌گیرد و همچنین انواع مختلف ترافیک حمله را در باتنت‌ها را شامل می‌شود. این مجموعه‌داده دارای دسته‌ها و زیردسته‌های حمله برای امکان طبقه‌بندی چندکلاسه است. آن‌ها از الگوریتم‌های ماشین بردار پشتیبان (SVM)، شبکه عصبی بازگشتی (RNN) و شبکه عصبی بازگشتی حافظه بلندمدت کوتاه‌مدت

(LSTM-RNN) برای مقایسه اعتبار و ارزیابی دقت مجموعه‌داده استفاده کرده‌اند. به‌طور مشابه، نویسنده‌گان دیگری مانند سوسیلو و همکاران [۱۹] و السامیری [۲۰] آزمایش‌های خود را بر اساس مجموعه‌داده Bot-IoT انجام داده‌اند. سوسیلو سیستم تشخیص ناهنجاری را با به‌کارگیری MLP، SVM، RF و CNN ساخته است، در حالی که السامیری از MLP، RF، NB، KNN استفاده کرده است. آنتی و همکاران [۲۱] تکرارشونده ۳ (ID3)، AdaBoost، تحلیل تشخیصی درجه دوم (QDA) و KNN استفاده کرده است. آنتی و همکاران [۲۲] یک سیستم تشخیص نفوذ سه‌لایه (IDS) ارائه داده‌اند که مجموعه‌داده‌ای نماینده تولید می‌کند و از روش نظارت شده برای تشخیص طیفی از سایر حمله‌های رایج مبتنی بر شبکه در شبکه‌های IoT استفاده می‌کند. برای اعتبارسنجی سیستم‌های تشخیص نفوذ بالقوه، آن‌ها NB شبکه بیزی، OneR، Zero R، J48، RF و MLP، SVM را در نظر گرفته‌اند. با این حال، مجموعه‌داده به‌کاررفته از پروتکل‌های خاص IoT استفاده نکرده و تنها شامل دو روز از ترافیک ناهنجار است. MedBiOT [۲۲] یک مجموعه‌داده برچسب‌گذاری شده رفتاپی از IoT ارائه می‌دهد که حاوی ترافیک عادی و مخرب باتنت است. سناریوی پیشنهادی یک زیرساخت IoT با اندازه متوسط (۸۳ دستگاه IoT) است که از ترکیب دستگاه‌های واقعی و شبیه‌سازی شده رایج IoT و سه بدافزار باتنت برجسته (RF، DT، SVM و KNN) پیاده‌سازی Torii و BashLite، Mirai) تشکیل شده است. مدل‌های طبقه‌بندی یادگیری ماشین (RF، DT و SVM) و شده‌اند تا قابلیت استفاده از مجموعه‌داده پیشنهادی تأیید شود. تامارایسلوی و همکاران [۲۳] از الگوریتم‌هایی مانند RF، NB و DT برای تشخیص ناهنجاری در شبکه‌های IoT استفاده کرده‌اند. برای انجام آزمایش‌های خود، از مجموعه‌داده IoT-23 [۲۴] بهره برده‌اند که مجموعه‌داده بزرگی از بیست و سه ردپای مختلف ترافیک شبکه IoT ارائه می‌دهد. این سناریوها به بیست ردپای شبکه از دستگاه‌های IoT آلووه و سه ردپای شبکه از ترافیک واقعی دستگاه‌های IoT تقسیم شده‌اند. ضبطهای بدافزار برای بازه‌های طولانی اجرا شده‌اند و سناریوهای حمله متنوعی از جمله Mirai، Torii و Gagfyt را اجرا کرده‌اند. در [۲۵] کاربردی برای IIoT یافت شده است. الگوریتم‌های تشخیص ناهنجاری مبتنی بر یادگیری ماشین برای یافتن ترافیک مخرب در یک مجموعه‌داده تولیدشده مصنوعی از ارتباط Modbus/TCP در یک سناریوی صنعتی ساختگی به کار گرفته شده‌اند. این مجموعه‌داده شامل حملات نفوذ رایج مبتنی بر خانه و دفتر است که رفتار زمان‌بندی و نرخ بسته‌ها در واحد زمان را شبیه‌سازی می‌کند که این امر عامل تمایز خوبی برای حملات محسوب می‌شود. چهار الگوریتم یادگیری ماشین مختلف برای تشخیص ناهنجاری به کار گرفته شده‌اند: KNN، RF، SVM و خوشبندی k-means. بهترین نتایج با ماشین بردار پشتیبان (SVM) بدست آمده است. لیو و همکاران [۲۶] یک مجموعه‌داده عمومی با استفاده از حسگرهای هوشمند در یک شبکه IoT ایجاد کرده‌اند. داده‌ها در محیط آزمایشگاه هوشمند و خانه هوشمند با استفاده از بوردهای حسگر Rainbow HAT نصب شده روی رزبری‌پایه جمع‌آوری شده‌اند. این مجموعه‌داده حاوی پنج نوع سایر حمله است: مسموم‌سازی ARP (ARP Poisoning)، حمله سرکوب سرویس (ARP Flood)، سیل (ARP DoS)، حمله فورس‌بروت (Flood) و پروتکل SlowLoris و Asterisk. آن‌ها روشی ترکیبی پیشنهاد داده‌اند که یک مدل جاسازی شده برای انتخاب ویژگی را با RCNN یک شبکه عصبی کانولوشنی (CNN) برای طبقه‌بندی حمله ترکیب می‌کند. روش تشخیص نفوذ پیشنهادی دو مدل دارد: (الف) که در آن جنگل تصادفی با CNN ترکیب شده است؛ و (ب) که در آن XCNN با XGBoost ترکیب شده است. یک تکنیک ترکیبی که تلفیقی از رایانش فوگ و رایانش ابری است، به کار گرفته شده است. همچنین، آن‌ها طبقه‌بندی‌ها را روی CCD-INID-V1 و دو مجموعه‌داده دیگر IoT، یعنی N_BaIoT [۱۳] و CIRA-CIC-DoHBrw-2020 (DoH20) [۲۷] بررسی کرده‌اند تا اثربخشی این مدل‌های امنیتی مبتنی بر یادگیری را بررسی کرده و کارایی روش پیشنهادی خود را با الگوریتم‌های سنتی یادگیری ماشین مانند KNN، NB، LR، SVM و SlowITe مقاریسه نموده و نتایج مقایسه‌ای طبقه‌بندی ناهنجاری و چندکلاسه را ارائه دهنده. واکاری و همکاران [۲۸] را ارائه داده‌اند که بر پروتکل MQTTset متمرکز است و از دستگاه‌های IoT با ماهیت‌های مختلف (مانند حسگرهای دما، رطوبت، حرکت و غیره) تشکیل شده است تا محیط خانه/دفتر/ساختمان هوشمند را شبیه‌سازی کند. MQTTset قانونی را با انواع مختلف ترافیک مخرب/حمله‌ای نظیر حملات سرکوب سرویس با سیل، سیل انتشار MQTT، داده‌های NAMENAS و احراز هویت با فورس‌بروت ترکیب می‌کند تا زمینه‌های مختلفی از جمله اتوماسیون خانگی، نظارت بر زیرساخت‌های حیاتی یا

محیط‌های صنعتی را شبیه‌سازی نماید. بر اساس این سناریو، آن‌ها یک سیستم تشخیص نفوذ پیاده‌سازی کرده‌اند که از MLP، RF، DT و NN و ناوی بیز گاوی برای شناسایی رفتارهای مخرب و در نتیجه محافظت از سیستم در برابر حملات بهره می‌برد. بیشتر کارهای انجام‌شده در سیستم‌های تشخیص نفوذ شبکه (NIDS) به دنبال استخراج ویژگی‌ها و مجموعه‌ای از ترکیب الگوریتم‌ها هستند که نتایج را بهینه کنند. سرهان و همکاران [۲۹] در پژوهش خود تلاش کردند تا روش‌ها را استانداردسازی کنند تا بتوان آن‌ها را بر هر مجموعه‌داده‌ای اعمال نمود. برای آزمایش‌های خود، آن‌ها از شش مدل یادگیری ماشین استفاده کردند: شبکه پیشخور عمیق (DFF)، شبکه عصبی کانولوشنی (CNN)، شبکه عصبی بازگشتی (RNN)، درخت تصمیم (DT)، رگرسیون لجستیک (LR) و ناوی بیز (NB). همچنین آن‌ها سه الگوریتم استخراج ویژگی—تجزیه و تحلیل مؤلفه‌های اصلی (PCA)، تحلیل تشخیصی خطی (LDA) و اتوانکر (AE)—را بر روی سه مجموعه‌داده معیار [۳۰] و CSE-CIC-IDS2018 [۳۱] بررسی نمودند. نتایج بهدست آمده نشان می‌دهند که روش استخراج ویژگی یا مدل یادگیری ماشین خاصی وجود ندارد که بتواند بهترین نتایج را برای تمام مجموعه‌داده‌ها کسب کند و انتخاب مجموعه‌داده به طور قابل توجهی عملکرد روش‌های به کاررفته را تحت تأثیر قرار می‌دهد. در ابتدا، روش‌های یادگیری ماشین بر مجموعه‌داده‌هایی اعمال شدند که متعلق به شبکه‌های اینترنت اشیا (IoT) نبودند و تنها با تعدیل ردپاها (traces) به این محیط، نشان داده شد که این روش‌ها می‌توانند برای تشخیص حملات در چنین سناریوهایی به کار گرفته شوند. ظهور مجموعه‌داده‌های اختصاصی IoT امکان مدل‌سازی محیط‌های واقعی را فراهم کرده و کاربرد بهینه الگوریتم‌های یادگیری ماشین را برای اعتبارسنجی مجموعه‌داده‌های موجود تضمین می‌کند. پس از مرور روش‌های اصلی به کاررفته برای تشخیص ناهنجاری ترافیک با استفاده از یادگیری ماشین در محیط‌های IoT، تصمیم گرفتیم تا پنج مدل مختلف یادگیری ماشین با ماهیت‌های متفاوت را برای اعتبارسنجی مجموعه‌داده DAD انتخاب کنیم: رگرسیون لجستیک (Logistic Regression)، ناوی بیز (Naive Bayes)، جنگل تصادفی (Naïve Bayes)، آدابوست (AdaBoost) و ماشین بردار پشتیبان (Support Vector Machine) (Random Forest).

۳. روش‌های یادگیری ماشین

یادگیری ماشین به استخراج دانش از داده‌ها می‌پردازد. موفق‌ترین انواع الگوریتم‌های یادگیری ماشین، الگوریتم‌هایی هستند که با تعمیم از مثال‌های شناخته‌شده، فرآیندهای تصمیم‌گیری را خودکار می‌کنند. این روش به عنوان یادگیری نظارت‌شده شناخته می‌شود [۳۲]. از میان روش‌های متعدد یادگیری ماشین، برخی کمتر انعطاف‌پذیر یا محدودتر هستند؛ مدل‌های سطحی (shallow models) در این دسته قرار می‌گیرند، زیرا تنها می‌توانند طیف نسبتاً محدودی از اشکال را برای تخمین تولید کنند. به طور کلی، برای ساخت یک مدل طبقه‌بندی با دقت بالا، انتخاب یک الگوریتم یادگیری ماشین قدرتمند و همچنین تنظیم مناسب پارامترهای آن بسیار مهم است. بیشتر الگوریتم‌های یادگیری ماشین در صورتی که پارامترهایشان به درستی تنظیم شوند، نتایج بهینه‌ای حاصل می‌کنند. تنظیمات نامناسب پارامترها منجر به نتایج ضعیف طبقه‌بندی می‌شود. بهینه‌سازی پارامترها در صورت انجام دستی—به ویژه زمانی که الگوریتم یادگیری دارای پارامترهای زیادی باشد—می‌تواند بسیار زمان بر بارد. از این رو، جستجوی شبکه‌ای (grid search) در اصل یک جستجوی جامع بر پایه زیرمجموعه‌ای از فضای هایپرپارامترهاست. جستجوی شبکه‌ای با استفاده از تکنیک اعتبارسنجی متقطع (CV) به عنوان معیار بدقت پیش‌بینی کند. تکنیک اعتبارسنجی متقطع می‌تواند از مشکل بیش‌برازش (overfitting) جلوگیری کند [۳۳]. همان‌طور که پیش‌تر اشاره شد، برای انجام آزمایش‌ها، پنج الگوریتم طبقه‌بندی نظارت‌شده و سطحی با ماهیت‌های مختلف انتخاب شدند: رگرسیون لجستیک، ناوی بیز، جنگل تصادفی، آدابوست و ماشین بردار پشتیبان. در ادامه توضیح مختصی از هر الگوریتم ارائه می‌شود، همراه با شرح کوتاهی از هایپرپارامترهای در نظر گرفته شده برای تنظیم آن‌ها. اولین مدل یادگیری ماشین در نظر گرفته شده، رگرسیون لجستیک است که یک روش آماری برای تحلیل و پیش‌بینی یک خروجی دودویی محسوب می‌شود. این مدل یادگیری ماشین حالت خاصی از مجموعه‌ای از مدل‌های خطی تعمیم‌یافته است. هدف رگرسیون لجستیک، ایجاد بهترین مدل برآشی برای برقراری رابطه وابستگی بین

متغیر کلاس و ویژگی هاست [۳۴]. این مدل در واقع هایپرپارامترهای حیاتی زیادی برای تنظیم ندارد. با این حال، پارامترهای L1 و L2 و C را می‌توان به عنوان تکنیک‌های نظمدهی (regularization) برای مقابله با بیش‌برازش تنظیم کرد. تکنیک نظمدهی L1 مربوط به مدل رگرسیون لاسو (Lasso Regression) است، در حالی که L2 بر مدل رگرسیون ریج (Ridge Regression) است. در تخمین میانگین داده‌ها دارد. به طور خلاصه، نظمدهی L1 سعی می‌کند میانه داده‌ها را تخمین بزند، در حالی که نظمدهی L2 سعی در تخمین میانگین داده‌ها دارد. پارامتر C معکوس پارامتر نظمدهی است؛ این یک متغیر کنترلی است که با قرارگیری معکوس نسبت به ضریب لاندا (Lambda)، میزان قدرت نظمدهی را تنظیم می‌کند. مقادیر بالاتر C متناظر با نظمدهی کمتری هستند که از بیش‌برازش مدل جلوگیری می‌کند. مدل ناوی بیز یک مدل احتمالی بیزی بسیار ساده شده است که بر پیش‌فرض استقلال قوی بین ویژگی‌ها عمل می‌کند. این بدین معناست که احتمال یک ویژگی بر احتمال دیگری تأثیر نمی‌گذارد. با فرض داشتن مجموعه‌ای از n ویژگی، رده‌بند ناوی بیز $n/2$ فرضیه استقلال ایجاد می‌کند. با این وجود، نتایج رده‌بند ناوی بیز اغلب صحیح است [۳۵]. رده‌بندهای ناوی بیز خانواده‌ای از رده‌بندها هستند که از نظر ساختاری بسیار شبیه به رگرسیون لجستیک و رده‌بند بردار پشتیبان خطی (Linear SVC) می‌باشند. با این حال، آن‌ها عموماً در فرآیند آموزش حتی سریع‌تر عمل می‌کنند. قیمت پرداختی برای این کارایی این است که مدل‌های ناوی بیز اغلب عملکرد تعیین‌پذیری را ارائه می‌دهند که کمی ضعیفتر از رده‌بندهای خطی است. دلیل کارایی بالای مدل‌های ناوی بیز این است که آن‌ها پارامترها را با بررسی هر ویژگی به صورت جداگانه یاد می‌گیرند و آمار ساده‌ای را برای هر کلاس از هر ویژگی جمع‌آوری می‌کنند. برای طبقه‌بندی داده‌های دودویی، از ناوی بیز برنولی (Bernoulli Naive Bayes) استفاده می‌شود. این رده‌بند فرض می‌کند که داده‌ها دودویی هستند و تعداد دفعاتی را می‌شمارد که هر ویژگی در هر کلاس مقدار غیرصرف دارد. ناوی بیز برنولی تنها یک پارامتر، یعنی α ، دارد که پیچیدگی مدل را کنترل می‌کند. نحوه عملکرد α بدین صورت است که الگوریتم α نقطه داده مجازی به داده‌ها اضافه می‌کند که تمام ویژگی‌های آن‌ها مقادیر مثبت دارند. این امر منجر به «هموارسازی» (smoothing) آمار می‌شود. مقدار α بزرگ‌تر به معنای هموارسازی بیشتر و در نتیجه مدل‌هایی با پیچیدگی کمتر است. عملکرد الگوریتم نسبت به تنظیم α نسبتاً قوی (robust) است، به این معنا که تنظیم دقیق α برای دستیابی به عملکرد خوب ضروری نیست [۳۶]. از سوی دیگر، مدل‌های مبتنی بر درخت در کاربردهای خود در حوزه سیستم‌های تشخیص نفوذ مبتنی بر رفتار، عملکرد چشمگیری از خود نشان داده‌اند. جنگل تصادفی (RF) یا Random Forest چندین درخت تصمیم بدون هرس می‌سازد که با روش نمونه‌برداری با جایگذاری (bootstrapping) از داده‌های آموزشی و انتخاب تصادفی زیرمجموعه‌ای از ویژگی‌ها القا می‌شوند. پیش‌بینی با تجمعی (aggregating) نتایج مجموعه انجام می‌شود (در طبقه‌بندی، رأی اکثریت و در رگرسیون، میانگین‌گیری). جنگل تصادفی عموماً بهبود قابل توجهی در عملکرد نسبت به یک درخت واحد از خود نشان می‌دهد و نرخ خطای پایینی همراه با مقاومت بر جسته در برابر نویز فراهم می‌کند [۳۲]. پارامتر $n_{\text{estimators}}$ تعداد درخت‌هایی است که در مدل RF ساخته می‌شوند. تعداد درخت‌های بیشتر، مدل‌هایی پایدارتر تولید می‌کند اما نیازمند حافظه بیشتر و زمان اجراطولانی‌تر است. پارامتر \max_{features} تعداد ویژگی‌هایی است که هنگام جستجوی بهترین تقسیم (split) در نظر گرفته می‌شوند. در اینجا مقادیر \log_2 (یعنی $\lceil \log_2(\text{تعداد ویژگی‌ها}) \rceil$) و 2 (یعنی $\log_2(\text{تعداد ویژگی‌ها})$) را به عنوان مقادیر ممکن برای این پارامتر در نظر گرفته‌ایم. آدبوست (Adaboost) یک الگوریتم نماینده از خانواده روش‌های تقویت‌کننده (boosting) است که ایده اصلی آن انتخاب و ترکیب گروهی از رده‌بندهای ضعیف برای ساخت یک رده‌بند قوی است. الگوریتم تقویت‌کننده به صورت یک رویکرد تکرارشونده برای تولید یک رده‌بند قوی به کار گرفته می‌شود که با ترکیب رده‌بندهای ضعیفی که تنها از حدس‌های تصادفی ساخته شده‌اند، خطای آموزشی بسیار کوچکی را به دست می‌آورد. در این روش، از تکنیک گروهی (ensemble) برای استخراج سوابق آموزشی با به روزرسانی مکرر توزیع نمونه‌های داده‌های قبل‌آموزش دیده استفاده می‌شود [۳۷]. در این پیاده‌سازی، رده‌بند درخت تصمیم (Decision Tree) به عنوان تخمین‌گر پایه (base estimator) مورد استفاده قرار گرفته است. برای تنظیم هایپرپارامتر، پارامتر $n_{\text{estimators}}$ حداقل تعداد تخمین‌گرهایی را تعیین می‌کند که در آن فرآیند تقویت متوقف می‌شود. در صورت برآش کامل (perfect fit)، فرآیند یادگیری زودتر متوقف می‌شود. پارامتر learning_rate برای کاهش سهم هر رده‌بند (یعنی وزنی) که در هر تکرار تقویت به هر رده‌بند

اختصاص داده می‌شود) در نظر گرفته شده است. نرخ یادگیری (learning rate) بالاتر، سهم هر رده‌بند را افزایش می‌دهد. بین پارامترهای $n_{estimators}$ و learning_rate یک رابطه مبادله (trade-off) وجود دارد.

جدیدترین روش نظارت شده در این حوزه، ماشین بردار پشتیبان (SVM) یا Support Vector Machine است. SVM به فضاهای با بعد بالاتر تبدیل می‌کند و ابرصفحه‌ای (hyperplane) را پیدا می‌کند که بهترین شکل داده‌ها را از هم جدا کند. یک ماشین بردار پشتیبان با یافتن بهترین ابرصفحه‌ای که تمام نقاط داده یک کلاس را از نقاط داده کلاس دیگر جدا می‌کند، داده‌ها را طبقه‌بندی می‌نماید. SVM‌ها حول مفهوم «حاشیه» (margin) به دو سوی ابرصفحه‌ای که دو کلاس داده را از هم جدا می‌کند، می‌چرخدن. حاشیه، عرض حداقلی نواری موازی با ابرصفحه است که هیچ نقطه داده‌ای در داخل آن وجود ندارد. حداقل کردن این حاشیه و ایجاد بیشترین فاصله ممکن بین ابرصفحه جداً کننده و نمونه‌های هر دو سوی آن، ثابت شده است که کران بالای خطای تعمیم‌پذیری مورد انتظار را کاهش می‌دهد [۳۸]. کرنل (kernel)، نحوه تصویرسازی متغیرهای ورودی را کنترل می‌کند. گزینه‌های متعددی وجود دارند، اما کرنل‌های خطی، چندجمله‌ای و RBF (تابع پایه شعاعی) رایج‌ترین آن‌ها هستند. هنگامی که یک رده‌بند بردار پشتیبان (SVC) برای طبقه‌بندی از کرنل خطی استفاده می‌کند، به عنوان یک SVC خطی (Linear SVC) تعریف می‌شود. مدل‌های خطی در فضاهای کم‌بعد می‌توانند بسیار محدود کننده باشند، زیرا خطوط و ابرصفحه‌ها انعطاف‌پذیری محدودی دارند. یکی از راههای افزایش انعطاف‌پذیری مدل خطی، افزودن ویژگی‌های بیشتر است. یک پارامتر حیاتی، پارامتر نظم‌دهی جریمه‌ای C است که می‌تواند طیف وسیعی از مقادیر را پیزدیر و تأثیر چشمگیری بر شکل نواحی نهایی هر کلاس دارد. استفاده از مقیاس لگاریتمی می‌تواند نقطه شروع مناسبی باشد. شدت نظم‌دهی regularization) به صورت معکوس با C متناسب است؛ بنابراین C باید کاملاً مثبت باشد. همان‌گونه که در مدل‌های خطی دیده می‌شود، مقدار کوچک C به معنای مدلی بسیار محدود است که در آن هر نقطه داده تنها می‌تواند تأثیر بسیار محدودی داشته باشد [۳۶].

۴. مروری بر مجموعه‌داده

برای به کارگیری روش‌های یادگیری ماشین، وجود مجموعه‌داده‌ای با تعداد نمونه‌های قابل توجه، زمینه‌سازی شده و برچسب‌گذاری شده بضروری است. این بخش توضیح مختصری از سناریویی ارائه می‌دهد که در آن مجموعه‌داده انتخاب شده توسعه یافته است. سپس، تحلیل خلاصه‌ای از داده‌های مرتبط با رده‌بند ارائه می‌شود. تحلیل جامع تر این مجموعه‌داده در منبع [۳۹] قابل یافتن است.

۴.۱. توصیف سناریو

برای نزدیک کردن مجموعه‌داده به یک محیط واقعی، داده‌ها از حسگرهای دمای مرکز داده CITIC [۴۰] جمع‌آوری شده‌اند. این مجموعه‌داده که DAD نامیده می‌شود، هفت روز از وضعیت روزانه‌ای را نشان می‌دهد که در یک محیط واقعی رخ می‌دهد و دارای حجم کافی ردپا (trace) و استخراج ویژگی‌های مشخص است. حسگرها دمای مرکز داده را روی سه مؤلفه نظارت می‌کنند: رک‌ها (racks)، نوارهای برق (PDU) و دستگاه‌های سرمایشی (InRow). ما تنها از حسگرهای InRow استفاده کردیم، زیرا مقادیر حاصل از سایر حسگرها ثابت و بی‌معنی هستند. مهم‌ترین دستگاه‌های سیستم، واحدهای InRow (شماره‌های ۱۳، ۱۵، ۲۳ و ۲۵) هستند که مسئول سرمایش هوا در مرکز داده از طریق یک سیستم خنک‌کننده مایعی می‌باشند. این واحدهای InRow دارای چهار حسگر مرتبط هستند:

- (i) دمای هوای ورودی واحد (TAS)؛
- (ii) دمای هوای بازگشته واحد (TAR)؛
- (iii) دمای سیال ورودی واحد (TFEU)؛
- (iv) دمای سیال خروجی واحد (TFSU).

در مورد محیط مجازی‌سازی شده، پنج ماشین مجازی با سیستم‌عامل Ubuntu Server 18.04 ایجاد شد که هر کدام به شبکه داخلی اینترنت اشیا (IoT) متصل بودند، به طوری که ترافیک بین آن‌ها کاملاً ایزوله بود.

در هر یک از چهار واحد سرمایشی، چهار حسگر وجود دارد. هر گره مشتری (client node) فرآیندی را اجرا می‌کند که هر حسگر را با یک شناسه مرتبط با شناسه واحد سرمایشی (InRow) شبیه‌سازی می‌کند. به عنوان مثال، واحد سرمایشی ۱۳ برای TAS شناسه ۱۳۱، برای TAR شناسه ۱۳۲، برای TFEU شناسه ۱۳۳ و برای TFSU شناسه ۱۳۴ دارد. هنگامی که حسگر برای انتشار موضوع (topic) مربوطه به بروکر متصل می‌شود، پیامی با پروتکل MQTT و با شناسه گره خود به عنوان ClientId ارسال می‌کند. همه حسگرهای دارای یک آدرس IP یکسان هستند، اما در هر روز، پورت‌های مورد استفاده برای انتقال پیام برای هر حسگر متفاوت است. به همین دلیل، در مجموع از ۱۱۲ پورت TCP استفاده می‌شود. هر یک از این حسگرهای هر پنج دقیقه یکبار داده را به بروکر ارسال می‌کند، یعنی ۲۸۸ نمونه در ساعت و در مجموع ۴۰۳۲ نمونه در روز برای هر حسگر. این حسگرهای ارتباطی مستقیم با یکدیگر ندارند. مجموعه‌داده شامل سناریوهای متنوعی از ناهنجاری است که در آن ترافیک غیرعادی از نظر آماری با ترافیک عادی متفاوت است و اکثر نمونه‌های ترافیک شبکه، عادی هستند. این مجموعه‌داده در سطح بسته (packet) برچسب‌گذاری شده است و نوع توکن (عادی یا ناهنجار) را مشخص می‌کند.

رفتار گره ناهنجار به یکی از روش‌های زیر تغییر یافته است:

- اعتراض (Interception): حذف تصادفی برخی از بسته‌های ارسالی.
- تغییر (Modification): تغییر دمای ارسالی بدون رعایت الگوی تعیین شده.
- تکثیر (Duplication): ارسال توکن‌هایی بیشتر از تعداد برنامه‌ریزی شده اولیه.

۴.۲. توصیف محتوا

داده‌های عددی به منظور درک مدل مفهومی مقادیر مشاهده شده تحلیل شده‌اند. ترافیک ناهنجار از واحد ۱۳ IP با آدرس ۱۰.۶.۵۶.۴۱ ایجاد شده است و بسته‌های ناهنجار در روزهای خاصی از هفته (چهارشنبه، پنجشنبه، جمعه و شنبه) ارسال شده‌اند. مجموعه‌داده DAD در مجموع دارای ۱۰۱,۵۸۳ بسته است که ترافیک UDP و TCP غالب هستند. از این بسته‌های ۶۳.۳٪ متعلق به ترافیک MQTT هستند که ۱۶٪ از آن‌ها به عنوان ناهنجار علامت‌گذاری شده‌اند. تعداد بایت‌های منبع و مقصد، بسته‌های منبع و مقصد، بسته‌های TCP، UDP و MQTT در طول روزهای هفت‌های یک‌نواخت است و به طور میانگین روزانه شامل ۱۴۰۲۰ بسته است که از این میان ۹,۲۶۷ بسته متعلق به MQTT هستند. ناهنجاری‌های ناشی از «اعتراض» در مجموعه‌داده برچسب‌گذاری نشده‌اند، زیرا این بسته‌ها حذف شده‌اند. بسته‌های مربوط به اعتراض وجود ندارند و به همین دلیل در آمار کلی به عنوان بسته‌های غیرعادی منعکس نمی‌شوند. به عبارت دیگر، در صد ترافیک ناهنجار، ناهنجاری‌های ناشی از اعتراض را در نظر نمی‌گیرد. هر حسگر هر ۵ دقیقه یکبار پیام به بروکر ارسال می‌کند. یک زمان انتظار بیکاری (idle timeout) برابر با ۳۰ ثانیه پیکربندی شده است، بنابراین اگر گره‌ای برای این مدت بسته‌ای ارسال نکند، بسته بعدی به عنوان بخشی از یک جریان (flow) جدید در نظر گرفته می‌شود. جریان‌ها به صورت یک‌سویه (unidirectional) تعریف شده‌اند. با توجه به نوع اتصال معمول در شبکه‌های MQTT، تغییر در پیام یکی از حسگرهای بر کل جریان MQTT تأثیر می‌گذارد. در نتیجه، اگر بسته‌ای متعلق به جریانی باشد که حداقل یکی از بسته‌های آن به عنوان ناهنجار برچسب‌گذاری شده باشد، تمام بسته‌های آن جریان غیرعادی در نظر گرفته می‌شوند. مجموعه‌داده در مجموع دارای ۲۲۴ جریان اتصال TCP است که از این میان ۲۰۸ جریان عادی و ۱۶ جریان ناهنجار هستند. علاوه بر این، این مجموعه‌داده شامل ۶۷,۸۴۸ جریان MQTT است که از

آن‌ها ۵۴۴ جریان مربوط به ناهنجاری هستند. تعداد جریان‌ها در طول روزهای هفته یکنواخت است، اما در روز یکشنبه تعداد کمتری دیده می‌شود که به دلیل قطع اتصال‌هایی است که روز بعد انجام می‌شود. تمام ناهنجاری‌ها بر روی جریان‌های TCP اعمال شده‌اند.

۵. آماده‌سازی داده‌ها

به عنوان مرحله‌ای پیش از به کارگیری الگوریتم‌های یادگیری ماشین، لازم است داده‌ها، نمونه‌ها و ویژگی‌های مورد استفاده در تحلیل مشخص شوند. این وظیفه شامل ارتقای کیفیت داده‌ها به سطح مورد نیاز توسط روش‌های تحلیلی انتخاب شده است و در برگیرنده تبدیل‌های داده و ساخت ویژگی‌های جدید از یک یا چند ویژگی موجود برای اهداف پاکسازی و همچنین بررسی تأثیر احتمالی آن‌ها بر نتایج تحلیل می‌باشد.

شناخت پارامترهایی که الگوریتم‌های یادگیری ماشین برای اجرای بهینه نیاز دارند، امری حیاتی است. با توجه به ماهیت و نحوه عملکرد رده‌بند، ۱۴ ویژگی به عنوان مرتبط شناسایی و برای پردازش انتخاب شدند:

«ip.src»، «tcp.srcport»، «frame.len»، «frame.time»،
 «ip.dst»، «tcp.dstport»، «tcp.flags»، «tcp.flags.ack»،
 «tcp.flags.fin»، «tcp.flags.res»، «tcp.flags.reset»،
 .««label»، «protocol»، «tcp.flags.syn»»

ویژگی‌های جدیدی برای بهبود تخمین مدل ایجاد شدند:

- (روز هفته): روز هفته استخراج شده از ویژگی Weekday
- (نوع پروتکل): مشخص می‌کند که آیا پروتکل مربوطه MQTT-IoT Protocol_Type است یا خیر.
- (شناسه مشتری): ویژگی که بر اساس ip.src و tcp.srcport مشخص می‌کند که هر سنسور به کدام واحد ClientID تعلق دارد.

برخی داده‌ها به طور ذاتی چرخه‌ای هستند، این موضوع در مورد ساعات، دقایق و ثانیه صدق می‌کند. اگر ساعت را به عنوان یک متغیر خطی در نظر بگیریم، رده‌بند متوجه نخواهد شد که ساعت ۲۳ قبل از ساعت ۰ قرار دارد. برای حل این مشکل، باید مقادیر چرخه‌ای به الگوریتم یادگیری ماشین ارائه شوند؛ این کار با جایگزینی ویژگی موجود با دو ویژگی جدید با استفاده از نمایش سینوسی انجام می‌شود. هر دو متغیر ضروری هستند، در غیر این صورت حرکت صحیح در طول زمان از بین می‌رود. این امر به دلیل تغییر مشتق توابع سینوس یا کسینوس در طول زمان است که در آن موقعیت (y, x) به صورت پیوسته در اطراف دایره واحد تغییر می‌کند.

در مورد ساعت‌شمار ۲۴ ساعته، این کار به ترتیب با تبدیل‌های سینوسی و کسینوسی مطابق معادله (۱) انجام می‌شود:

$$x = \sin(2 * \pi * hour / 24)$$

$$y = \cos(2 * \pi * hour / 24).$$
(۱)

فاصله بین دو نقطه در این فضای جدید، معادل اختلاف زمانی بین آن‌هاست که دقیقاً همان چیزی است که از یک چرخه ۲۴ ساعته انتظار داریم.

همان‌طور که پیش‌تر اشاره شد، در مجموعه‌داده از ۱۱۲ پورت TCP استفاده شده است. این امر در فرآیند نگاشت منجر به ساخت ۱۱۲ ویژگی می‌شود که اطلاعات تکراری ضمنی در ClientId را دارند. با توجه به این موضوع، ویژگی‌های tcp.dstport و tcp.srcport حذف شدند.

بهترین روش نمایش داده‌ها نه تنها به معنای داده‌ها، بلکه به نوع مدل مورد استفاده نیز بستگی دارد. مدل‌های خطی و مدل‌های مبتنی بر درخت (مانند درخت‌های تصمیمی، درخت‌های تقویت‌شده گرادیانی و جنگل تصادفی)، دو خانواده بزرگ و بسیار رایج، خواص بسیار متفاوتی در برخورد با نمایش‌های مختلف ویژگی دارند. مدل‌های خطی تنها می‌توانند روابط خطی را مدل کنند که در مورد یک ویژگی منفرد به صورت خطوط ساده ظاهر می‌شوند. یکی از راههای افزایش توان مدل‌های خطی در مواجهه با داده‌های پیوسته، استفاده از روش گروه‌بندی (binning) یا گسسته‌سازی (discretization) است که در آن یک ویژگی پیوسته به چندین ویژگی گسسته تبدیل می‌شود. در اینجا، به عنوان فرآیندی مشترک برای تمام مدل‌ها، یک ویژگی ورودی پیوسته از مجموعه‌داده به یک ویژگی گسسته تبدیل می‌شود که نشان می‌دهد هر نمونه داده به کدام دسته (bin) تعلق دارد [۳۶]. ویژگی‌های گسسته در مجموعه‌داده باید به نمایش عددی تبدیل شوند. این فرآیند با استفاده از کدگذاری دودویی معمول انجام می‌شود، به طوری که هر متغیر گسسته با m مقدار ممکن، با $1m - 1$ متغیر مجازی (dummy variable) جایگزین می‌شود. یک متغیر مجازی برای یک دسته خاص مقدار ۱ و برای سایر دسته‌ها مقدار ۰ دارد [۴۱]. در مورد ویژگی برچسب (label)، کلاس عادی مقدار عددی ۰ و کلاس ناهمجارت مقدار عددی ۱ دریافت می‌کند. ویژگی frame.time به یک متغیر تاریخ تبدیل می‌شود تا امکان استخراج اطلاعات مرتبط دیگر از آن فراهم شود. ویژگی d_flow نشان‌دهنده مدت زمان جریان (flow) است و frame.len_mean معادل میانگین حسابی طول فریم تمام بسته‌های متعلق به یک جریان است. تمام این ویژگی‌ها به فرم گسسته تبدیل شده‌اند.

ویژگی‌های نهایی به صورت زیر هستند:

```
[«frame.len», «hora_cos», «hora_sin», «Npackperflow», «d_flow»,
«tcp.flags.ack», «tcp.flags.fin», «tcp.flags.reset», «tcp.flags.syn»,
«frame.len_mean», «tcp.flags.res», «label_fl»,
«ip.src_10.6.56.1», «ip.src_10.6.56.34», «ip.src_10.6.56.36»,
«ip.src_10.6.56.41», «ip.src_10.6.56.50», «ip.dst_10.6.56.1»,
«ip.dst_10.6.56.34», «ip.dst_10.6.56.36», «ip.dst_10.6.56.41»,
«ip.dst_10.6.56.50», «ClientId_131», «ClientId_132», «ClientId_133»,
«ClientId_134», «ClientId_151», «ClientId_152», «ClientId_153»,
«ClientId_154»,
«ClientId_231», «ClientId_232», «ClientId_233», «ClientId_234»,
«ClientId_251», «ClientId_252», «ClientId_253», «ClientId_254»,
«ClientId_Broker», «protocol_TIP_MQTT», «protocol_TIP_TCP»,
«weekday_Friday», «weekday_Monday», «weekday_Saturday»,
«weekday_Sunday», «weekday_Thursday», «weekday_Tuesday»,
[«weekday_Wednesday»]
```

پس از نگاشت ویژگی‌های نمادین به مقادیر عددی، در صورت وجود واریانس قابل توجه، نیاز به مقیاس‌بندی ویژگی‌ها احساس شد. مقیاس‌بندی ویژگی‌ها از طریق نرمال‌سازی میانگین (mean normalization) انجام گرفت. از سوی دیگر، کاهش تعداد ویژگی‌ها می‌تواند اجرای الگوریتم‌های یادگیری ماشین را کارآمدتر (با پیچیدگی فضایی یا زمانی کمتر) و مؤثرتر کند. برخی الگوریتم‌های یادگیری ماشین ممکن است توسط ویژگی‌های ورودی نامرتبط گمراه شوند که این امر منجر به عملکرد پیش‌بینی ضعیفتر می‌شود. با توجه به همگن‌بودن پروتکل‌های MQTT و نحوه عملکرد الگوریتم‌های مختلف یادگیری ماشین، انتخاب ویژگی‌ها امکان حذف ویژگی‌های هم‌خطی Recursive Feature (correlated) و بهبود کارایی محاسباتی الگوریتم‌ها را فراهم می‌آورد.تابع حذف بازگشتی ویژگی‌ها (RFE) یا Elimination روش انتخاب معکوس از پیش‌بین‌کننده‌ها بر اساس رتبه‌بندی اهمیت آن‌ها پیاده‌سازی می‌کند. در این روش، پیش‌بین‌کننده‌ها رتبه‌بندی شده و کم‌اهمیت‌ترین آن‌ها قبل از مدل‌سازی به صورت متوالی حذف می‌شوند. هدف یافتن زیرمجموعه‌ای از پیش‌بین‌کننده‌های است که بتوان از آن برای تولید یک مدل دقیق استفاده کرد [۴۲]. مدل انتخاب‌شده به عنوان تخمین‌گر RFE، یک رده‌بند درخت تصمیم با اعتبارسنجی متقطع k-fold استراتیفاید با پنج تقسیم بود. پس از انتخاب مهم‌ترین پیش‌بین‌کننده‌ها توسط الگوریتم، دقت بدست‌آمده در تعدادهای مختلفی از پیش‌بین‌کننده‌های مورد استفاده برآورد شد. همان‌طور که در شکل ۱ مشاهده می‌شود، بهترین دقت با استفاده از ۱۲ ویژگی حاصل شد. ویژگی‌های انتخاب‌شده عبارت بودند از: d_flow, hour_sin, hour_cos, weekday_Friday, ClientId_134, ClientId_133, ClientId_132, ClientId_131, ip.src_10.6.56.41, weekday_Wednesday, weekday_Thursday, weekday_Saturday

شکل ۱. دقت در برابر اهمیت متغیر.

در محیط تشخیص نفوذ شبکه، معمولاً با این وضعیت مواجه هستیم که یکی از کلاس‌ها بخش بسیار کوچکی از داده‌ها را تشکیل می‌دهد، در حالی که همین کلاس اقلیت، موارد مهمی برای تشخیص محسوب می‌شوند. هنگام یادگیری از داده‌های بسیار نامتوازن، احتمال این که نمونه انتخاب‌شده حاوی تعداد بسیار کم یا حتی فاقد نمونه‌های کلاس اقلیت باشد، به‌طور قابل توجهی بالاست؛ این امر منجر به ایجاد الگوریتمی با عملکرد ضعیف در پیش‌بینی کلاس اقلیت می‌شود [۴۳]. این مشکل باعث می‌شود که الگوریتم‌های یادگیری ماشین به راحتی دسته‌بندی‌های نادرستی انجام دهند. اعتبارسنجی متقطع (CV) یک روش آماری برای ارزیابی عملکرد تعمیم‌پذیری است که در مقایسه با تقسیم ساده به مجموعه آموزش و آزمون، پایدارتر و جامع‌تر است. در اعتبارسنجی متقطع، داده‌ها به صورت مکرر تقسیم شده و مدل‌های متعددی آموزش داده می‌شوند. اعتبارسنجی متقطع k-fold cross-validation (stratified k-fold cross-validation) نوعی تغییریافته از روش k-fold است که چینش‌های استراتیفاید (stratified folds) تولید می‌کند. در اعتبارسنجی متقطع استراتیفاید، چینش‌ها به گونه‌ای ساخته می‌شوند که درصد نمونه‌های هر کلاس در تمام چینش‌ها حفظ شود و رابطه تعادل‌یافته بین کلاس‌ها در داده‌ها حفظ گردد. معمولاً استفاده از اعتبارسنجی متقطع k-fold استراتیفاید به جای اعتبارسنجی متقطع k-fold ساده برای ارزیابی یک رده‌بند، ایده خوبی است، زیرا تخمین‌های قابل اطمینان‌تری از عملکرد تعمیم‌پذیری ارائه می‌دهد [۳۶]. رویکرد دیگری که امکان کاهش مشکل نامتوازن‌بودن داده را فراهم می‌کند، استفاده از تکنیک‌های نمونه‌برداری است. SMOTE تابعی برای مدیریت مسائل طبقه‌بندی نامتوازن است. این روش با کاهش تصادفی نمونه‌های کلاس اکثریت (down-sampling) همراه است و به جای نمونه‌برداری با جایگذاری از کلاس اقلیت، نمونه‌های مصنوعی از کلاس اقلیت ایجاد می‌کند تا تعداد آن را افزایش دهد [۴۳, ۴۴]. در مجموعه داده DAD، ۱۰۸۸ بسته ناهمجارت و ۹۶,۸۹۴ بسته عادی وجود دارد. این امر نسبت ۱:۸۹ را نشان می‌دهد که وضعیت آشکار نامتوازن‌بودن داده را بیان می‌کند. برای تضمین این که نمونه‌هایی از هر دو کلاس در نظر گرفته شوند، آزمایش‌ها با به‌کارگیری SMOTE در درون یک حلقه تو در توی اعتبارسنجی متقطع k-fold استراتیفاید با k=5 انجام شد. به صورت پیش‌فرض، SMOTE از پنج همسایه (-

(neighbors) استفاده می‌کند. در ابتدا، کلاس اقلیت تا نسبتی حدود ۰.۴۰ افزایش نمونه‌برداری (over-sampled) شد. سپس، کلاس اکثریت تا نسبت ۰.۵ کاهش نمونه‌برداری (under-sampled) گردید، بهطوری که هر چینش (fold) در تقسیم‌بندی اعتبارسنجی متقطع، توزیع کلاسی یکسانی مطابق با پیکربندی SMOTE داشته باشد.

۶. آموزش مدل‌ها

در آزمایش‌ها، تمام مدل‌ها در چارچوب یک اعتبارسنجی متقطع k-fold استراتیفاید تو در تو ارزیابی شدند؛ حلقه بیرونی با $k=5$ و حلقه درونی نیز با $k=5$ در حالی که SMOTE در حلقه درونی اعمال شد. در این بخش، روش‌های برآش یادگیری ماشین روی داده‌های باقی‌مانده از مجموعه آموزش k-fold cv تشریح می‌شود و مقدار «ناحیه زیر منحنی عملکرد دستگاه دریافت‌کننده» (ROC AUC) از نمره‌های پیش‌بینی محاسبه می‌گردد. نتایج ارائه‌شده متناظر با میانگین مقادیر بهدست‌آمده از حلقه درونی هستند. برای رگرسیون لجستیک، از حل‌کننده liblinear استفاده شد. مجازات‌های L1 و L2 در نظر گرفته شدند و پارامتر C در مقادیر ۱، ۰.۰۰۱، ۰.۰۰۰۱ و ۰.۰۰۰۰۱ تنظیم گردید. به دلیل ماهیت تصادفی الگوریتم، نتایج بهدست‌آمده بسته به تکرارهای انجام‌شده متفاوت است. بهترین نتیجه بهدست‌آمده، مقدار ۰.۹۷۹۰ با مجازات L2 و مقدار C برابر با ۱۰۰۰ است. کمترین مقدار C با مجازات L1 و مقدار C برابر با ۰.۹۶۷۶ است. میانگین کلی در تمام تکرارها برابر با ۰.۹۷۵۶ و انحراف معیار آن $10 \times 2.06 \times 10^{-3}$ است. معمولاً مقادیر بالای C آزادی بیشتری به مدل می‌دهد، بنابراین انتظار عملکرد بهتری در مقادیر بالای C وجود دارد. با این حال، همان‌طور که در شکل ۲ مشاهده می‌شود، تغییر پارامتر C به تنها‌ی بدون انتخاب مناسب نوع مجازات، به صورت قابل توجهی عملکرد مدل را بهبود نمی‌بخشد.

شکل ۲. میانگین نمره رگرسیون لجستیک بر اساس تکرار.

جدول ۱ بهترین پارامترهای انتخاب‌شده برای بهترین نمره میانگین در هر یک از حلقه‌های درونی CV و همچنین انحراف معیار بهدست‌آمده را نشان می‌دهد. پارامترهای تنظیم‌شده که در اکثر موارد بهترین نتایج را ارائه داده‌اند، C برابر با ۱ و مجازات L1 هستند. پارامترهای ارائه‌شده در جدول، همان پارامترهایی هستند که برای ارزیابی مدل در مجموعه آزمون انتخاب شده‌اند.

جدول ۱. بهترین مدل رگرسیون لجستیک بر اساس تکرار.

در میان مدل‌های ناوی بیز که پیش‌تر ذکر شدند، ناوی بیز برنولی (Bernoulli Naive Bayes) انتخاب شد. برای یافتن ترکیب بهینه هایپرپارامترها کهتابع از دستداد (loss function) از پیش‌تعریف‌شده را کمینه کند، نتایج با تغییر پارامتر α در مقادیر ۰.۱، ۰.۰۱ و ۰.۰۰۵ و ۰.۰۰۰۵ محاسبه شدند. میانگین مقدار ROC AUC در تمام تکرارها برابر با ۰.۹۶۵۴ و انحراف معیار آن $10 \times 3.23 \times 10^{-3}$ است. کمترین مقدار ROC AUC برابر با ۰.۹۵۸۷ با $\alpha = 0.1$ و بهترین مقدار ROC AUC برابر با ۰.۹۷۲۱ با $\alpha = 0.01$ بهدست آمد. شکل ۳ میانگین نتایج بهدست‌آمده از تنظیم هایپرپارامترها را بر اساس تکرار نشان می‌دهد. همان‌طور که در مورد رگرسیون لجستیک مشاهده شد، ناوی بیز برنولی نیز عملکرد متناسب با مقادیر α را در هر تکرار بهبود نمی‌بخشد.

شکل ۳. میانگین نمره ناوی بیز بر اساس تکرار.

جدول ۲ پارامترهایی را ارائه می‌دهد که در هر تکرار درونی بهترین نمره را ارائه داده‌اند. این پارامترها در ارزیابی مدل مورد استفاده قرار خواهند گرفت. بهترین مقدار بهدست‌آمده، ۰.۹۷۲۱ با $\alpha = 0.1$ و کمترین انحراف معیار برابر با $10 \times 3.22 \times 10^{-3}$ است.

جدول ۲. بهترین مدل ناوی بیز بر اساس تکرار.

برای تنظیم هایپرپارامترهای جنگل تصادفی (Random Forest)، مهم‌ترین تنظیمات عبارتند از تعداد درخت‌ها در جنگل (n_estimators) و تعداد ویژگی‌هایی که برای تقسیم در هر گره برگ (leaf node) در نظر گرفته می‌شوند (max_features). در طول اجرا، ارزیابی با مقادیر ۱۰۰، ۲۰۰ و ۳۰۰ درخت و مقادیر sqrt و 2 برای max_features آزمایش شد. میانگین کلی بهدست آمده برابر با ۰.۹۹۹۹ و انحراف معیار آن 5.57×10^{-6} بود که در برخی تکرارها به ROC AUC برابر با ۱ دست یافت. شکل ۴ میانگین نتایج ROC AUC را بر اساس تکرار نشان می‌دهد. نقاط بالا در نمودار نشان‌دهنده مقادیر ۱ در آن تکرارهای است. جنگل تصادفی بهترین عملکرد را در مرحله آموزش از خود نشان داده است.

شکل ۴. میانگین نمره جنگل تصادفی بر اساس تکرار.

جدول ۳ بهترین پارامترها را در هر تکرار درونی نشان می‌دهد. چندین پیکربندی مختلف وجود دارند که در تمام تکرارها بهترین پاسخ را ارائه می‌دهند. با این حال، تعداد درخت‌ها و ویژگی‌ها می‌توانند زمان و عملکرد مدل را در حین اجرا تحت تأثیر قرار دهند.

جدول ۳. بهترین مدل جنگل تصادفی بر اساس تکرار.

ردهبند آدابوست (AdaBoost) به صورت پیش‌فرض از یک ردهبند درخت تصمیم به عنوان تخمین‌گر پایه استفاده می‌کند که با max_depth=1 مقداردهی اولیه شده است. برای برازش یک مدل خوب، ممکن است افزایش تعداد درخت‌های ردهبند مطلوب باشد، هرچند این امر ممکن است عملکرد سیستم را تحت تأثیر قرار دهد. در این مدل، تعداد تخمین‌گرها در مقادیر ۵۰۰، ۱۰۰۰ و ۲۰۰۰ و نرخ یادگیری (learning rate) در مقادیر ۰.۰۰۱، ۰.۰۰۱ و ۰.۱ تنظیم شد. میانگین کلی ROC AUC بهدست آمده برابر با ۰.۹۶۹ و انحراف معیار آن 1.69×10^{-3} بود. در شکل ۵ می‌توان مشاهده کرد که افزایش پارامتر نرخ یادگیری، عملکرد مدل را نسبت به تعداد درخت‌های انتخاب شده بهبود می‌بخشد؛ با این حال، نتایج عالی را می‌توان حتی با سطوح متوسط نرخ یادگیری و تعداد کم درخت‌ها در ردهبند بهدست آورد.

شکل ۵. میانگین نمره آدابوست بر اساس تکرار.

جدول ۴ پارامترهایی را نشان می‌دهد که در پیاده‌سازی اعتبارسنجی متقطع درونی، بهترین میانگین را بهدست آورده‌اند. با انحراف‌های معیار بسیار کم و دستیابی به نمرات ۱ در برخی تکرارها، بهترین مدل‌ها با مقادیر بالای نرخ یادگیری و تعداد زیاد درخت‌ها یافت شدند.

جدول ۴. بهترین مدل آدابوست بر اساس تکرار.

در نهایت، برای ماشین بردار پشتیبان (SVM)، از ردهبند بردار خطی استفاده شد، زیرا این مدل انعطاف‌پذیری بیشتری در انتخاب توابع جریمه وتابع زیان دارد و باید مقیاس‌پذیری بهتری نسبت به نمونه‌های زیاد داشته باشد. برای بهبود دقیق مدل، پارامتر C تنظیم شد. پارامتر جریمه C، تعادل بین مرز تصمیم و جمله سوء‌طبقه‌بندی را کنترل می‌کند و در مقادیر ۰.۱، ۰.۵ و ۱۰ ثابت شد. با در نظر گرفتن تمام آزمایش‌ها، میانگین نمره ROC AUC برابر با ۰.۹۷۵۹ و انحراف معیار آن 1.74×10^{-3} بهدست آمد. شکل ۶ میانگین نمره بهدست آمده در تکرارها را نشان می‌دهد و کاهش عملکرد مدل را در مقدار $C = 10$ بر جسته می‌کند. بهترین میانگین‌های نمره مدل در مقادیر C برابر با ۱ و ۵ مشاهده شد که در جدول ۵ قابل مشاهده است.

۷. نتایج و مقایسه مدل‌ها

برای انجام پیش‌بینی روی داده‌ها، لازم است از داده‌هایی استفاده شود که در فرآیند آموزش مدل مورد استفاده قرار نگرفته‌اند. تخمین عملکرد مدل یادگیری ماشین بر مبنای مجموعه آزمون حاصل از تقسیم‌بندی اعتبارسنجی متقطع انجام می‌شود. این بخش نتایج چهار معیار سنتی—دقت (accuracy)، صحت (precision)، فراخوانی (recall) و F_1 —را برای تمام مدل‌های یادگیری ماشین سطحی ارزیابی شده ارائه می‌دهد. علاوه بر این، مقایسه‌ای بین عملکرد الگوریتم‌ها انجام شده است. شکل ۲۷ میانگین مقادیر دقต را برای هر یک از مدل‌ها در تمام تکرارهای انجام‌شده نشان می‌دهد. در این شکل مشاهده می‌شود که کمترین عملکرد در این معیار مربوط به رده‌بند ناوی بیز است. مدل‌های مبتنی بر درخت تصمیم بهترین نتایج را کسب کرده‌اند و مقادیر بسیار نزدیک به نتیجه بهینه (1.0) را به دست آورده‌اند. در معیار صحت، باز هم رده‌بندهای مبتنی بر درخت تصمیم در سه مورد از پنج تکرار انجام‌شده، مقادیری برابر با 1 را کسب کرده‌اند که در شکل b7 نشان داده شده است. در فراخوانی (شکل C7)، بهترین نتایج متعلق به ناوی بیز است، در حالی که آدابوت و جنگل تصادفی نیز نتایج بسیار رضایت‌بخشی داشته‌اند. از سوی دیگر، در معیار F_1 (شکل d7)، باز هم آدابوت و جنگل تصادفی بهترین نتایج را در تمام تکرارها کسب کرده‌اند و عملکردی برجسته در مقایسه با سایر رده‌بندها از خود نشان داده‌اند. جدول ۶ میانگین‌ها برای تمام تکرارهای انجام‌شده برای تمام رده‌بندها همراه با انحراف معیار آن‌ها ارائه می‌دهد. رده‌بند ناوی بیز کمترین مقادیر دقت را در میان تمام مدل‌ها دارد، اما تنها مدلی است که فراخوانی کامل ($recall=1$) را به دست آورده است. همان‌طور که انتظار می‌رفت، به دلیل شباهت در ساختار مدل‌ها، رده‌بندهای رگرسیون لجستیک و ماشین بردار پشتیبان نتایج بسیار مشابهی ارائه داده‌اند. هر دو رده‌بند دققی حدود 0.96 ، صحبتی در حدود 0.1 و F_1 حدود 0.3 دارند. این نتایج، رده‌بندهای مبتنی بر درخت را به عنوان بهترین رده‌بندها برای این مجموعه‌داده تأیید می‌کنند؛ این مدل‌ها مقادیری تقریباً کامل (نزدیک به 1.0) را به دست آورده‌اند و همچنین کمترین انحراف معیار را نشان داده‌اند که این امر تضمین‌کننده پایداری رده‌بند است.

شکل ۷. معیارها بر اساس تکرار: (الف) دقت. (ب) صحت. (ب) فراخوانی. (ت) F_1

جدول ۶. بهترین مدل از نظر نمره بر اساس تکرار.

۸. نتیجه‌گیری و کارهای آینده

برای کاهش ناهنجاری‌های ترافیکی، توسعه‌دهندگان نیازمند روش‌های جدیدی برای تشخیص دستگاه‌های IoT آسیب‌دیده هستند. روش‌های یادگیری ماشین اخیراً به دلیل کاربردهای موفق در تشخیص ناهنجاری‌های شبکه، از جمله شبکه‌های IoT، اعتبار قابل ملاحظه‌ای کسب کرده‌اند. با توجه به کمبود مجموعه‌داده‌های IoT، مجموعه‌داده DAD به عنوان ابزاری برای شناخت رفتار دستگاه‌های سبک وزن و اختصاصی در شبکه‌های IoT-MQTT ایجاد شده است. DAD یک مجموعه‌داده کامل و برچسب‌گذاری شده با تحلیل‌های پیشین است که قرار است برای تشخیص ناهنجاری‌های ترافیکی با استفاده از الگوریتم‌های یادگیری ماشین به کار گرفته شود. بنابراین، این پژوهش با هدف نشان دادن این موضوع انجام شده است که مجموعه‌داده DAD را می‌توان برای تشخیص ناهنجاری در شبکه‌های ترافیکی-MQTT استفاده کرد. برای این منظور، داده‌ها باید به گونه‌ای پردازش شوند که برای رده‌بند کاملاً قابل فهم باشند. روش‌های پاک‌سازی و آماده‌سازی داده سعی می‌کنند نویز را هموار کرده و نقاط پرت را شناسایی نمایند. استفاده از ویژگی‌های چرخه‌ای، متغیرهایی یکتا برای رده‌بند ایجاد می‌کند، در حالی که تکنیک‌های گستره‌سازی رتبه‌بندی مرتبه بالایی از مقادیر را فراهم می‌آورند که می‌توانند روابط بین مشاهدات را هموار کنند. الگوریتم RFE پیاده‌سازی شد تا وابستگی‌ها و همخطی‌های بین ویژگی‌ها را حذف کرده و خطاهای محاسباتی و هزینه‌های محاسباتی را کاهش دهد. برای تضمین عملکرد بهینه مدل‌ها، تنظیم هایپرپارامترها اعمال شد که این امر مقادیر مناسبی را برای جلوگیری از بیش‌برازش (overfitting) تضمین می‌کند. علاوه بر این، استفاده از SMOTE در ترکیب با اعتبارسنجی متقطع k-

استراتیفاید به عنوان یک تکنیک مدیریت داده برای داده‌های نامتوازن، امکان حل ساده ولی مؤثر مشکل وجود کلاس اقلیت در مجموعه داده را فراهم آورد و به طور چشمگیری عملکرد ردهبند را بهبود بخشد. در مورد روش‌ها، ما پنج الگوریتم رایج یادگیری ماشین سطحی—رگرسیون لجستیک، ناوی بیز، جنگل تصادفی و ماشین بردار پشتیبان خطی—را انتخاب کردیم. سه مورد از این الگوریتم‌ها—رگرسیون لجستیک، ناوی بیز و SVM—بر توابع آماری خطی مبتنی هستند و همان‌طور که انتظار می‌رفت، نتایج بسیار مشابهی ارائه دادند. از سوی دیگر، ردهبندهای مبتنی بر درخت—جنگل تصادفی و آدابوست—رفتار مشابهی از خود نشان دادند. در مرحله آموزش، تمام مدل‌ها عملکرد استثنایی داشتند؛ کمترین نمره متعلق به ردهبند ناوی بیز با میانگین ROC AUC برابر با 0.9587 بود و بالاترین مقدار که برابر با 1 بود، توسط جنگل تصادفی و آدابوست کسب شد. در مرحله آزمون، نتایج به دست آمده از ردهبندهای مبتنی بر درخت، مقادیر عالی در معیارهای سنتی را نشان داد و عملکرد برجسته‌ای در کاربردهای سیستم‌های تشخیص نفوذ مبتنی بر رفتار ارائه کرد و بهبود چشمگیری در کارایی، نرخ خطای پایین و مقاومت در برابر نویز را به همراه آورد. مجموعه داده DAD به منظور تشخیص ناهنجاری در شبکه‌های ترافیکی MQTT-IoT ایجاد شده است. از این رو، این آزمایش‌ها تأیید می‌کنند که این مجموعه داده را می‌توان برای کاربرد تشخیص ناهنجاری ترافیک در IoT با به کارگیری تکنیک‌های بهینه‌سازی استفاده کرد. با توجه به جدید بودن این مجموعه داده، هنوز کار دیگری روی آن انجام نشده است. در نتیجه، ما اولین کاربرد الگوریتم‌های یادگیری ماشین را برای تشخیص ناهنجاری در شبکه‌های IoT بر اساس این مجموعه داده ارائه می‌دهیم. کارهای آینده می‌توانند شامل به کارگیری سایر انواع الگوریتم‌های یادگیری ماشین روی این مجموعه داده باشند تا با مقایسه نتایج به دست آمده، اثربخشی ردهبندهای ارائه شده تأیید گردد. مشارکت نویسنده‌گان: مفهوم پردازی، ل.و. و.ک؛ روش‌شناسی، و.ک. و ف.خ.ن؛ نرم‌افزار، ل.و. و.ک. و د.ف؛ اعتبارسنجی، و.ک. د.ف. و ف.خ.ن؛ تحلیل رسمی، ل.و. و.ک. و د.ف؛ پژوهش، ل.و. و.ک. و د.ف؛ منابع، ل.و. و.ک. د.ف. و ف.خ.ن؛ نگارش—پیش‌نویس اولیه، ل.و. و.ک. و د.ف؛ نگارش—بازبینی و ویرایش، ل.و. و.ک. د.ف. و ف.خ.ن؛ نظارت، و.ک؛ مدیریت پروژه، ل.و. تمام نویسنده‌گان نسخه منتشر شده مقاله را خوانده و با آن موافق هستند. تأمین مالی: این پروژه توسط برنامه «تأیید اعتبار، ساختاردهی و بهبود واحدهای پژوهشی تثبیت شده و مراکز منحصر به فرد» (ED431G/01) که توسط آموزش حرفه‌ای شورای زانتا د گالیسیا و با بودجه اتحادیه اروپا (صندوق FEDER) و وزارت علوم و نوآوری اسپانیا از طریق پروژه PID2019-111388GB-I00 تأمین شده است، حمایت مالی شده است.

1. Moustafa, N.; Turnbull, B.; Choo, K.R. An Ensemble Intrusion Detection Technique Based on Proposed Statistical Flow Features for Protecting Network Traffic of Internet of Things. *IEEE Internet Things J.* 2019, 6, 4815–4830.
2. Roman, R.; Najera, P.; Lopez, J. Securing the Internet of Things. *Computer* 2011, 44, 51–58.
3. Agrawal, S.; Agrawal, J. Survey on Anomaly Detection using Data Mining Techniques. *Procedia Comput. Sci.* 2015, 60, 708–713.
4. Asharf, J.; Moustafa, N.; Khurshid, H.; Debie, E.; Haider, W.; Wahab, A. A Review of Intrusion Detection Systems Using Machine and Deep Learning in Internet of Things: Challenges, Solutions and Future Directions. *Electronics* 2020, 9, 1177.
5. Omar, S.; Ngadi, M.; Jebur, H.; Benqdara, S. Machine Learning Techniques for Anomaly Detection: An Overview. *Int. J. Comput. Appl.* 2013, 79, 33–41.
6. Wazid, M.; Das, A.K.; K, V.; Vasilakos, A. LAM-CIoT: Lightweight authentication mechanism in cloud-based IoT environment. *J. Netw. Comput. Appl.* 2019, 150, 102496.
7. DAD: Dataset for Anomaly Detection. Available online: <https://github.com/dad-repository/dad> (accessed on 30 March 2021).
8. Moustafa, N.; Slay, J. UNSW-NB15: A comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set). In Proceedings of the 2015 Military

- Communications and Information Systems Conference (MilCIS), Canberra, Australia, 10–12 November 2015; pp. 1–6.
9. Koroniots, N.; Moustafa, N.; Sitnikova, E.; Slay, J. Towards Developing Network Forensic Mechanism for Botnet Activities in the IoT Based on Machine Learning Techniques. In Proceedings of the International Conference on Mobile Networks and Management (MONAMI), Melbourne, Australia, 13–15 December 2017.
 10. Alrashdi, I.; Alqazzaz, A.; Aloufi, E.; Alharthi, R.; Zohdy, M.; Ming, H. AD-IoT: Anomaly Detection of IoT Cyberattacks in Smart City Using Machine Learning. In Proceedings of the 2019 IEEE 9th Annual Computing and Communication Workshop and Conference (CCWC), Las Vegas, NV, USA, 7–9 January 2019; pp. 305–310.
 11. Pajouh, H.H.; Javidan, R.; Khayami, R.; Dehghantanha, A.; Choo, K.K.R. A Two-Layer Dimension Reduction and Two-Tier Classification Model for Anomaly-Based Intrusion Detection in IoT Backbone Networks. *IEEE Trans. Emerg. Top. Comput.* **2019**, *7*, 314–323.
 12. Mirsky, Y.; Doitshman, T.; Elovici, Y.; Shabtai, A. Kitsune: An Ensemble of Autoencoders for Online Network Intrusion Detection. *arXiv* **2018**, arXiv:1802.09089.
 13. Meidan, Y.; Bohadana, M.; Mathov, Y.; Mirsky, Y.; Shabtai, A.; Breitenbacher, D.; Elovici, Y. N-BaIoT—Network-Based Detection of IoT Botnet Attacks Using Deep Autoencoders. *IEEE Pervasive Comput.* **2018**, *17*, 12–22.
 14. Abbasi, F.; Naderan, M.; Alavi, S.E. Anomaly detection in Internet of Things using feature selection and classification based on Logistic Regression and Artificial Neural Network on N-BaIoT dataset. In Proceedings of the 2021 5th International Conference on Internet of Things and Applications (IoT), Isfahan, Iran, 19–21 May 2021; pp. 1–7.
 15. Doshi, R.; Aphorpe, N.; Feamster, N. Machine Learning DDoS Detection for Consumer Internet of Things Devices. In Proceedings of the 2018 IEEE Security and Privacy Workshops (SPW), Francisco, CA, USA, 24 May 2018; pp. 29–35.
 16. Pahl, M.; Aubet, F. All Eyes on You: Distributed Multi-Dimensional IoT Microservice Anomaly Detection. In Proceedings of the 2018 14th International Conference on Network and Service Management (CNSM), Rome, Italy, 5–9 November 2018; pp. 72–80.
 17. Hasan, M.; Islam, M.M.; Zarif, M.I.; Hashem, M.M. Attack and anomaly detection in IoT sensors in IoT sites using machine learning approaches. *Internet Things* **2019**, *7*, 100059.
 18. Koroniots, N.; Moustafa, N.; Sitnikova, E.; Turnbull, B.P. Towards the development of realistic botnet dataset in the Internet of Things for network forensic analytics: Bot-IoT dataset. *Future Gener. Comput. Syst.* **2019**, *100*, 779–796.
 19. Susilo, B.; Sari, R.F. Intrusion Detection in IoT Networks Using Deep Learning Algorithm. *Information* **2020**, *11*, 279.
 20. Alsamiri, J.; Alsubhi, K. Internet of Things Cyber Attacks Detection using Machine Learning. *Int. J. Adv. Comput. Sci. Appl.* **2019**, *10*, 627–663.
 21. Anthi, E.; Williams, L.; Słowińska, M.; Theodorakopoulos, G.; Burnap, P. A Supervised Intrusion Detection System for Smart Home IoT Devices. *IEEE Internet Things J.* **2019**, *6*, 9042–9053.
 22. Guerra-Manzanares, A.; Medina-Galindo, J.; Bahsi, H.; Nõmm, S. MedBIoT: Generation of an IoT Botnet Dataset in a Medium-sized IoT Network. In Proceedings of the 6th International Conference on Information Systems Security and Privacy—ICISSP, INSTICC, SciTePress, Valletta, Malta, 25–27 February 2020; pp. 207–218.
 23. Thamaraiselvi, D.; Mary, S. Attack and Anomaly Detection in IoT Networks using Machine Learning. *Int. J. Comput. Sci. Mob. Comput.* **2020**, *9*, 95–103.

24. Parmisano, A.; Garcia, S.; Erquiaga, M.J. Stratosphere Laboratory. A Labeled Dataset with Malicious and Benign IoT Network Traffic. Available online: <https://www.stratosphereips.org/datasets-iot23> (accessed on 3 September 2020).
25. Anton, S.D.; Kanoor, S.; Fraunholz, D.; Schotten, H.D. Evaluation of Machine Learning-based Anomaly Detection Algorithms on an Industrial Modbus/TCP Data Set. In Proceedings of the 13th International Conference on Availability, Reliability and Security, Hamburg, Germany, 27–30 August 2018.
26. Liu, Z.; Thapa, N.; Shaver, A.; Roy, K.; Siddula, M.; Yuan, X.; Yu, A. Using Embedded Feature Selection and CNN for Classification on CCD-INID-V1—A New IoT Dataset. *Sensors* **2021**, *21*, 4834.
27. MontazeriShatoori, M.; Davidson, L.; Kaur, G.; Lashkari, A.H. Detection of DoH Tunnels using Time-series Classification of Encrypted Traffic. In Proceedings of the 2020 IEEE IEEE International Conference on Dependable Autonomic, & Secure Computing International Conference on Pervasive Intelligence & Computing International Conference Cloud Big Data Computing International Conference Cyber Science and Technology Congress (DASC/Picom/Cbdcom/Cyberscitech), Calgary, AB, Canada, 17–22 August 2020; pp. 63–70.
28. Vaccari, I.; Chiola, G.; Aiello, M.; Mongelli, M.; Cambiaso, E. MQTTset, a New Dataset for Machine Learning Techniques on MQTT. *Sensors* **2020**, *20*, 6578.
29. Sarhan, M.; Layeghy, S.; Moustafa, N.; Gallagher, M.; Portmann, M. Feature Extraction for Machine Learning-based Intrusion Detection in IoT Networks. *arXiv* **2021**, arXiv:2108.12722.
30. Moustafa, N. New Generations of Internet of Things Datasets for Cybersecurity Applications based Machine Learning: TON_IoT Datasets. 2019. Available online: http://handle.unsw.edu.au/1959.4/resource/collection/resdatac_921/1 (accessed on 20 October 2021).
31. Sharafaldin, I.; Lashkari, A.H.; Ghorbani, A.A. Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization. In Proceedings of the 4th International Conference on Information Systems Security and Privacy (ICISSP 2018), Madeira, Portugal, 22–24 January 2018.
32. James, G.; Witten, D.; Hastie, T.; Tibshirani, R. *An Introduction to Statistical Learning: With Applications in R*; Springer Publishing Company: New York, NY, USA, 2014. Available online: <https://www.statlearning.com/> (accessed on 15 November 2021).
33. Syarif, I.; Prugel-Bennett, A.; Wills, G. SVM Parameter Optimization using Grid Search and Genetic Algorithm to Improve Classification Performance. *Telkomnika Telecommun. Comput. Electron. Control* **2016**, *14*, 1502.
34. Ghosh, P.; Mitra, R. Proposed GA-BFSS and logistic regression based intrusion detection system. In Proceedings of the 2015 Third International Conference on Computer, Communication, Control and Information Technology (C3IT), West Bengal, India, 7–8 February 2015; pp. 1–6.
35. Mukherjee, S.; Sharma, N. Intrusion Detection using Naive Bayes Classifier with Feature Reduction. *Procedia Technol.* **2012**, *4*, 119–128.
36. Muller, A.C.; Müller, A.C. *Introduction to Machine Learning with Python: A Guide for Data Scientists*, 1st ed.; O'Reilly Media: Sebastopol, CA, USA, 2016.
37. Mebawondu, O.J.; Alowolodu, O.D.; Adetunmbi, A.O.; Mebawondu, J.O. Optimizing the Classification of Network Intrusion Detection Using Ensembles of Decision Trees Algorithm. In Proceedings of the Third International Conference on Information and Communication

- Technology and Applications (ICTA 2020), Minna, Nigeria, 24–27 November 2020; pp. 286–300.
- 38. Hamid, Y.; Sugumaran, M.; Balasaraswathi, V. IDS Using Machine Learning- Current State of Art and Future Directions. *Br. J. Appl. Sci. Technol.* 2016, 15, 1–22.
 - 39. Vigoya, L.; Fernandez, D.; Carneiro, V.; Cacheda, F. Annotated Dataset for Anomaly Detection in a Data Center with IoT Sensors. *Sensors* 2020, 20, 3745.
 - 40. Centro de Investigación en Tecnologías da Información e as Comunicacións de Galicia. Available online: <https://www.citic-research.org/> (accessed on 30 January 2020).
 - 41. Hasan, M.A.M.; Nasser, M.; Pal, B.; Ahmad, S. Support Vector Machine and Random Forest Modeling for Intrusion Detection System (IDS). *J. Intell. Learn. Syst. Appl.* 2014, 6, 45–52.
 - 42. Recursive Feature Elimination with Cross-Validation. Available online: https://scikit-learn.org/stable/auto_examples/feature_selection/plot_rfe_with_cross_validation.html (accessed on 6 October 2021).
 - 43. Chen, C.; Breiman, L. Using Random Forest to Learn Imbalanced Data; University of California: Berkeley, CA, USA, 2004.
 - 44. Chawla, N.V.; Bowyer, K.W.; Hall, L.O.; Kegelmeyer, W.P. SMOTE: Synthetic Minority over-Sampling Technique. *J. Artif. Int. Res.* 2002, 16, 321–357.