

Escuela de Posgrado

MAESTRÍA EN CIENCIA DE DATOS

Tesis

**Desarrollo de un modelo de clasificación para la
detección de anomalías en redes lot utilizando el
Dataset Ciciot2023**

Raul Raico Gallardo

Para optar el Grado Académico de Maestra en Ciencia de Datos

Lima, 2025

Repositorio Institucional Continental
Tesis digital



Esta obra está bajo una Licencia "Creative Commons Atribución 4.0 Internacional" .

ANEXO 6

**INFORME DE CONFORMIDAD DE ORIGINALIDAD DEL
TRABAJO DE INVESTIGACIÓN**

A : Mg. Jaime Sobrados Tapia
: Director Académico de la Escuela de Posgrado

DE : **Kevin Rafael Palomino Pacheco**
: Asesor del Trabajo de Investigación

ASUNTO : Remito resultado de evaluación de originalidad de Trabajo de
Investigación

FECHA : 14 de febrero del 2025

Con sumo agrado me dirijo a vuestro despacho para saludarlo y en vista de haber sido designado Asesor del Trabajo de Investigación/Tesis/Artículo Científico titulado **“Desarrollo de un modelo de clasificación para la detección de anomalías en redes iot utilizando el dataset CICIOT2023”**, perteneciente a **Bach. RAICO GALLARDO RAUL**, de la **Maestría en Ciencia de Datos**; se procedió con el análisis del documento mediante la herramienta “Turnitin” y se realizó la verificación completa de las coincidencias resaltadas por el software, cuyo resultado es **7%** de similitud (informe adjunto) sin encontrarse hallazgos relacionados con plagio. Se utilizaron los siguientes filtros:

- Filtro de exclusión de bibliografía SÍ NO
- Filtro de exclusión de grupos de palabras menores (Máximo nº de palabras excluidas: < 40) SÍ NO
- Exclusión de fuente por trabajo anterior del mismo estudiante SÍ NO

En consecuencia, se determina que el trabajo de investigación constituye un documento original al presentar similitud de otros autores (citas) por debajo del porcentaje establecido por la Universidad.

Recae toda responsabilidad del contenido de la tesis sobre el autor y asesor, en concordancia a los principios de legalidad, presunción de veracidad y simplicidad, expresados en el Reglamento del Registro Nacional de Trabajos de Investigación para optar grados académicos y títulos profesionales – RENATI y en la Directiva 003-2016-R/UC.

Esperando la atención a la presente, me despido sin otro particular y sea propicia la ocasión para renovar las muestras de mi especial consideración.

Atentamente,



Kevin Rafael Palomino Pacheco
DNI: 1045711819

DECLARACIÓN JURADA DE AUTENTICIDAD

Yo, RAICO GALLARDO RAUL, identificado con Documento Nacional de Identidad N° 42155193 de la MAESTRÍA EN CIENCIA DE DATOS, de la Escuela de Posgrado de la Universidad Continental, declaro bajo juramento lo siguiente:

1. El Trabajo de Investigación titulado "DESARROLLO DE UN MODELO DE CLASIFICACIÓN PARA LA DETECCIÓN DE ANOMALÍAS EN REDES IOT UTILIZANDO EL DATASET CICIOT2023", es de mi autoría, el mismo que presento para optar el Grado Académico de MAESTRO EN CIENCIA DE DATOS.
2. El Trabajo de Investigación no ha sido plagiado ni total ni parcialmente, para lo cual se han respetado las normas internacionales de citas y referencias para las fuentes consultadas, por lo que no atenta contra derechos de terceros.
3. El Trabajo de Investigación es original e inédito, y no ha sido realizado, desarrollado o publicado, parcial ni totalmente, por terceras personas naturales o jurídicas. No incurre en autoplagio; es decir, no fue publicado ni presentado de manera previa para conseguir algún grado académico o título profesional.
4. Los datos presentados en los resultados son reales, pues no son falsos, ni duplicados, ni copiados, por consiguiente, constituyen un aporte significativo para la realidad estudiada.

De identificarse fraude, falsificación de datos, plagio, información sin cita de autores, uso ilegal de información ajena, asumo las consecuencias y sanciones que de mi acción se deriven, sometiéndome a las acciones legales pertinentes.



RAICO GALLARDO RAUL

DNI. N° 42155193

Lima, 18 de marzo de 2025



Huella

Arequipa

Av. Los Incas S/N,
José Luis Bustamante y Rivero
(054) 412 030

Calle Alfonso Ugarte 607, Yanahuara
(054) 412 030

Huancayo

Av. San Carlos 1980
(064) 481 430

Cusco

Urb. Manuel Prado - Lote B, N° 7 Av. Collasuyo
(084) 480 070

Sector Angostura KM. 10,
carretera San Jerónimo - Saylla
(084) 480 070

Lima

Av. Alfredo Mendiola 5210, Los Olivos
(01) 213 2760

Jr. Junín 355, Miraflores
(01) 213 2760

DESARROLLO DE UN MODELO DE CLASIFICACIÓN PARA LA DETECCIÓN DE ANOMALÍAS EN REDES IOT UTILIZANDO EL DATASET CICIOT2023

INFORME DE ORIGINALIDAD



FUENTES PRIMARIAS

1 Submitted to Universidad Continental 1%

Trabajo del estudiante

2 Submitted to Universidad Europea de Madrid 1%

Trabajo del estudiante

3 dl.dropboxusercontent.com 1%

Fuente de Internet

4 educa.fme.cl 1%

Fuente de Internet

5 www.mdpi.com <1%

Fuente de Internet

6 www.iso.org <1%

Fuente de Internet

7 netfritz-technology.online <1%

Fuente de Internet

8 repositorio.upci.edu.pe <1%

Fuente de Internet

9

Fuente de Internet

<1 %

10

Vita Santa Barletta, Danilo Caivano, Mirko De Vincentiis, Anibrata Pal, Michele Scalera.
"Hybrid quantum architecture for smart city security", Journal of Systems and Software, 2024

Publicación

<1 %

11

repository.icesi.edu.co

Fuente de Internet

<1 %

12

cdn.www.gob.pe

Fuente de Internet

<1 %

13

slidehtml5.com

Fuente de Internet

<1 %

14

Submitted to Universidad Peruana de Las Americas

Trabajo del estudiante

<1 %

15

revistas.ceeep.mil.pe

Fuente de Internet

<1 %

16

Submitted to Institución Tecnológica Metropolitana de Medellín

Trabajo del estudiante

<1 %

17

Submitted to Universidad Rey Juan Carlos

Trabajo del estudiante

<1 %

18

Submitted to The University of Wolverhampton

<1 %

19

Submitted to Universitat Politècnica de València

Trabajo del estudiante

<1 %

20

repositories.nust.edu.pk

Fuente de Internet

<1 %

Excluir citas

Apagado

Excluir coincidencias

< 40 words

Excluir bibliografía

Activo

Asesor

Dr. Kevin Rafael Palomino Pacheco

Agradecimiento

En primer lugar, quiero expresar mi más sincero y humilde agradecimiento a Dios y la Santísima Virgen María, cuya guía y apoyo han sido fundamentales en cada paso de mi vida, con su presencia y ayuda constante me han brindado la fortaleza necesaria en los momentos de mayor dificultad y me han guiado en mi crecimiento espiritual y profesional. Gracias Padre Dios por la buena y bonita vida que me has dado, ¡Gracias por tanto y perdón por tan poco!

A mis padres, María Eudocia y Amílcar (¡ojalá ya en el cielo!), les estoy eternamente agradecido por su inmenso amor y cuidado, la educación y los valores que me han inculcado, su ejemplo, dedicación y sacrificio han sido pilares en mi formación personal y profesional. Gracias por enseñarme la alegría de vivir y la importancia del esfuerzo y la perseverancia.

También quiero extender mi gratitud al profesor Kevin Rafael Palomino Pacheco, asesor del presente trabajo, cuya experiencia, conocimiento y orientación han sido esenciales en el desarrollo de este proyecto, su buena voluntad, paciencia, motivación, consejos y sugerencias han guiado y enriquecido enormemente este trabajo.

A todos mis compañeros de la maestría, que con las exposiciones de sus diferentes trabajos y demás participaciones en clase me ayudaron a comprender mejor los contenidos de los diferentes cursos.

Finalmente, agradezco todos los docentes de la Maestría en Ciencia de Datos por sus valiosas y oportunas enseñanzas y a todos los trabajadores de la Universidad Continental, que de una u otra manera han interactuado y me han apoyado a lo largo de toda la maestría.

Índice

Asesor	ii
Agradecimiento.....	iii
Índice	v
Índice de figuras	vii
Índice de cuadros	viii
Índice de tablas	ix
Resumen	x
Abstract	xi
Introducción.....	xii
Capítulo I: Planteamiento del estudio	14
1.1. Planteamiento y formulación del problema	14
1.1.1. <i>Planteamiento del problema.</i>	14
1.1.2. <i>Formulación del problema.</i>	16
1.2. Determinación de objetivos	16
1.2.1. <i>Objetivo general.</i>	16
1.2.2. <i>Objetivos específicos.</i>	16
1.3. Justificación e importancia del estudio.....	17
1.3.1. <i>Justificación teórica.</i>	17
1.3.2. <i>Justificación metodológica.</i>	17
1.3.3. <i>Justificación social.</i>	18
1.4. Limitaciones de la presente investigación.....	21
Capítulo II: Marco teórico	22
2.1. Antecedentes de la investigación	22
2.1.1. <i>Internacionales.</i>	22
2.1.2. <i>Nacionales.</i>	23
2.2. Bases teóricas	25
2.2.1. <i>Desarrollo histórico.</i>	25
2.2.2. <i>Fundamentación teórica.</i>	26
2.2.3. <i>Marco conceptual.</i>	27
A. <i>Definición conceptual de la variable 1: Características de la anomalía.</i>	27
B. <i>Definición conceptual de la variable 2: Tipo de Anomalías.</i>	29
2.3. Definición de términos básicos.....	30
Capítulo III: Hipótesis y variables	33
3.1. Hipótesis de investigación	33
3.1.1. <i>Hipótesis general.</i>	33
3.1.2. <i>Hipótesis específicas.</i>	33
3.2. Operacionalización de variables.....	33

3.2.1. <i>Variable Independiente</i>	33
3.2.2. <i>Variable dependiente</i>	35
3.3. Matriz de operacionalización de variables	36
Capítulo IV: Metodología del estudio	37
4.1. Enfoque, tipo y alcance de investigación	37
4.1.1. <i>Enfoque</i>	37
4.1.2. <i>Tipo y alcance</i>	37
4.2. Diseño de la investigación	38
4.3. Población y muestra	38
4.3.1. <i>Población</i>	38
4.3.2. <i>Muestra</i>	38
4.4. Técnicas e instrumentos de recolección de datos	39
4.4.1. <i>Técnicas e instrumentos</i>	39
4.4.2. <i>Validez y confiabilidad</i>	39
4.4.3. <i>Procedimiento de recolección de datos</i>	39
4.5. Técnicas de análisis de datos	39
Capítulo V: Resultados	56
5.1. Análisis de resultados	56
5.1.1. <i>Análisis exploratorio de datos</i>	56
5.1.2. <i>Modelos de clasificación para la detección de anomalías</i>	69
A. <i>Regresión Logística</i>	69
B. <i>Árbol de Decisión</i>	71
C. <i>Random Forest (Bosque Aleatorio)</i>	73
D. <i>AdaBoost (Adaptive Boosting: Potenciación Adaptativa)</i>	75
E. <i>Perceptrón</i>	77
5.2. Discusión de resultados	81
5.3. Conclusión general	91
Conclusiones	94
Recomendaciones	95
Referencias	96
Anexos	99

Índice de figuras

Figura 1: Dispositivos IoT conectados a nivel mundial (en miles de millones)	14
Figura 2: Regiones consideradas y Tecnologías interconectadas	19
Figura 3: Archivos de la muestra	43
Figura 4: Resultado de ejecuciones para el archivo “combinado.pkl”	43
Figura 5: Gráfico de barras de las etiquetas	46
Figura 6: Matriz de correlación de dataframe con sus 40 características	49
Figura 7: Matriz de correlación 1 de las diez características seleccionadas.	51
Figura 8: Matriz de correlación 2 de las diez características seleccionadas	52
Figura 9: Gráfico de barras de Label en DF Preparado	55
Figura 10: Gráfico de barras de la característica Number	57
Figura 11: Gráfico de barras de la característica Time_To_Live	59
Figura 12: Gráfico de barras de la característica HTTPS	60
Figura 13: Gráfico de barras de la característica Std	61
Figura 14: Gráfico de barras de la característica Max	62
Figura 15: Gráfico de barras de la característica ack_flag_number	63
Figura 16: Gráfico de barras de la característica Tot size	64
Figura 17: Gráfico de barras de la característica Variance	65
Figura 18: Gráfico de barras de la característica Header_Length	66
Figura 19: Gráfico de barras de la característica DNS	68
Figura 20: Rendimiento de los modelos de clasificación	80

Índice de cuadros

Cuadro 1: Características del tráfico de red del dataset CICIoT2023	28
Cuadro 2: Matriz de operacionalización de variables	36
Cuadro 3: Cantidad de ataques.....	45
Cuadro 4: Sustitución de valores por categoría de ataque.....	47
Cuadro 5: Sustitución de valores por tráfico.....	48
Cuadro 6: Selección de características con SelectKBest, Correlación de Pearson y Chi2.....	50
Cuadro 7: Selección de 10 características con Random Forest.....	51
Cuadro 8: Diez características seleccionadas con la Correlación de Pearson.....	52
Cuadro 9: Características con valores originales (vista parcial).....	53
Cuadro 10: Estadísticas de las características seleccionadas con valores originales	53
Cuadro 11: Características escaladas (Vista parcial).....	54
Cuadro 12: Estadísticas de las características seleccionadas escaladas	54
Cuadro 13: DS Original y DF Preparado	55
Cuadro 14: Matriz de confusión de la Regresión Logística	71
Cuadro 15: Matriz de confusión del Árbol de decisión	73
Cuadro 16: Matriz de confusión del Random Forest	75
Cuadro 17: Matriz de confusión de AdaBoost	77
Cuadro 18: Matriz de confusión del Perceptrón	78
Cuadro 19: Rendimiento de los modelos de clasificación	80

Índice de tablas

Tabla 1: Estadísticas descriptivas de la característica Number.....	57
Tabla 2: Estadísticas descriptivas de la característica Time_To_Live.....	58
Tabla 3: Estadísticas descriptivas de la característica HTTPS.....	60
Tabla 4: Estadísticas descriptivas de la característica Std	61
Tabla 5: Estadísticas descriptivas de la característica Max	62
Tabla 6: Estadísticas descriptivas de la característica ack_flag_number	63
Tabla 7: Estadísticas descriptivas de la característica Tot size.....	64
Tabla 8: Estadísticas descriptivas de la característica Variance	65
Tabla 9: Estadísticas descriptivas de la característica Header_Length.....	66
Tabla 10: Estadísticas descriptivas de la característica DNS.....	67
Tabla 11: Reporte de clasificación de la Regresión Logística	71
Tabla 12: Reporte de clasificación del Árbol de decisión	73
Tabla 13: Reporte de clasificación de Random Forest.....	74
Tabla 14: Reporte de clasificación de AdaBoost.....	77
Tabla 15: Reporte de clasificación del Perceptrón	78

Resumen

En la presente investigación se evaluó cinco modelos de clasificación, Regresión Logística, Árbol de Decisión, Random Forest, Ada Boost y Perceptrón, para la detección de anomalías en redes IoT, utilizando una muestra del dataset CICIoT2023, las anomalías incluidas están agrupadas en siete categorías de ataque (DDoS, DoS, Reconocimiento, Basados en la web, Fuerza bruta, Suplantación y Mirai). Los resultados revelaron que todos los modelos lograron una exactitud superior al 97%, confirmando que modelos simples y que no requieren de mucha capacidad de cómputo son efectivos en la clasificación de tráfico malicioso y benigno. El modelo Random Forest destacó con una exactitud del 98.63%, seguido por Ada Boost con 98.55% y Árbol de Decisión con 98.28%, estos resultados fueron posibles mediante la preparación de los datos que incluyó la limpieza de datos faltantes, duplicados u otros irrelevantes y la selección de características mediante la Correlación de Pearson, que mejoraron la calidad del dataset. Con el análisis exploratorio de datos (EDA) se identificó patrones relevantes del tráfico malicioso y benigno, que dan una idea del modo de trabajo de las siete categorías de ataque incluidas en el dataset. Finalmente, los resultados obtenidos se compararon con los alcanzados por los antecedentes de la presente investigación.

Palabras clave: Detección de anomalías, redes IoT, CICIoT2023, modelos de clasificación, exactitud.

Abstract

In this research, five classification models, Logistic Regression, Decision Tree, Random Forest, Ada Boost, and Perceptron, were evaluated for the detection of anomalies in IoT networks, using a sample of the CICIoT2023 dataset. The anomalies included are grouped into seven attack categories (DDoS, DoS, Recon, Web-based, Brute Force, Spoofing, and Mirai). The results revealed that all the models achieved an accuracy of over 97%, confirming that simple models that do not require much computing capacity are effective in classifying malicious and benign traffic. The Random Forest model stood out with an accuracy of 98.63%, followed by Ada Boost with 98.55% and Decision Tree with 98.28%. These results were possible by preparing the data, which included cleaning missing, duplicate, or other irrelevant data, and feature selection using Pearson's Correlation, which improved the quality of the dataset. With the exploratory data analysis (EDA), relevant patterns of malicious and benign traffic were identified, which give an idea of the modus operandi of the seven attack categories included in the dataset. Finally, the results obtained were compared with those achieved by the background of this research.

Keywords: Anomaly detection, IoT networks, CICIoT2023, classification models, accuracy.

Introducción

Los entornos IoT son redes de dispositivos IoT interconectados a través de comunicaciones digitales, en nuestra vida diaria utilizamos dispositivos IoT de diversas formas, como en las casas inteligentes que controlan la iluminación y temperatura automáticamente, vehículos conectados que nos ayudan a navegar con facilidad, aplicaciones móviles para el monitoreo de nuestra salud, etc. Estos dispositivos nos simplifican la vida al mejorar aspectos importantes como la salud, mediante monitoreo constante, el transporte, con sistemas de navegación avanzados, el ahorro de tiempo y energía al automatizar tareas domésticas, la seguridad a través de sensores de presencia, cámaras de videovigilancia, alarmas inteligentes y demás, en resumen, nos brindan confort en general.

Dicha tecnología seguirá avanzando rápidamente en el futuro, esto significa que varios aspectos importantes de nuestra vida dependerán cada vez más del IoT, por ejemplo, podríamos ver hogares completamente automatizados donde todo funciona con solo dar órdenes por voz o con gestos y quizá en algún momento ya nada de eso, los vehículos podrían volverse aún más autónomos para reducir accidentes viales. Sin embargo, es necesario tomar conciencia de que en la tecnología IoT no todos son beneficios y que están libres de riesgos, sino que también existen peligros asociados con ataques cibernéticos a todo dispositivo conectado a internet. Los ciberdelincuentes pueden aprovecharse de las vulnerabilidades existentes en estos sistemas para robar información personal o causar daños significativos directamente a los usuarios, a infraestructuras críticas como sistemas de agua potable, centros de salud, plantas de generación de energía eléctrica, etc. Esto requiere el urgente y continuo desarrollo de sistemas seguros capaces de protegernos en el ciberespacio contra amenazas emergentes. Para abordar este desafío se plantea el siguiente problema, ¿Cómo desarrollar un modelo de clasificación para la detección de anomalías en redes IoT utilizando el dataset CICIoT2023?, este es un paso pertinente para la mejora de la seguridad en estos entornos interconectados.

Consecuentemente se tiene como objetivo, desarrollar un modelo de clasificación para la detección de anomalías en redes IoT utilizando el dataset CICIoT2023.

Con la ejecución del presente trabajo se comprobará, si el modelo de clasificación desarrollado es capaz de detectar anomalías en redes IoT con una exactitud superior al 90% utilizando el dataset CICIoT2023. Adicional a ello se presenta este estudio con la intención de comprender y aplicar modelos de aprendizaje automático para detectar anomalías en entornos IoT, así como también contribuir en el estado del arte en la seguridad IoT.

El presente trabajo contiene cinco capítulos, el primer capítulo, presenta la realidad problemática, según su influencia y alcance en nuestra vida cotidiana, además también se presenta los objetivos, la justificación, importancia y limitaciones de la investigación.

El segundo capítulo, comprende los antecedentes internacionales y nacionales de la investigación, las bases teóricas, el marco conceptual, definición de variables y términos básicos.

El tercer capítulo, muestra las hipótesis de investigación, la definición y operacionalización de variables, así como sus dimensiones, indicadores e instrumentos.

El cuarto capítulo, expone el enfoque, tipo, alcance, diseño, población y muestra de la investigación, técnicas e instrumentos de recolección de datos, técnicas de análisis de datos, limpieza de datos y la selección de las características más relevantes del dataset y su estructura.

El quinto capítulo, contiene el análisis y discusión de resultados, descripción y análisis de las características seleccionadas, sus estadísticas descriptivas y sus respectivos gráficos, hallazgos de patrones observados e interpretaciones. También se presentan los modelos de clasificación, Regresión Logística, Árboles de decisión, Random Forest, ADA Boost y Perceptrón, con sus respectivos componentes y modo de funcionamiento. Se incluyen también los resultados de rendimiento de cada modelo a través de reportes de clasificación y matrices de confusión, se discuten los resultados obtenidos, se presentan los resultados relevantes, la validación de objetivos e hipótesis, el contraste de los resultados de la presente investigación con los de los antecedentes, las conclusiones y recomendaciones de la investigación.

Capítulo I: Planteamiento del estudio

1.1. Planteamiento y formulación del problema

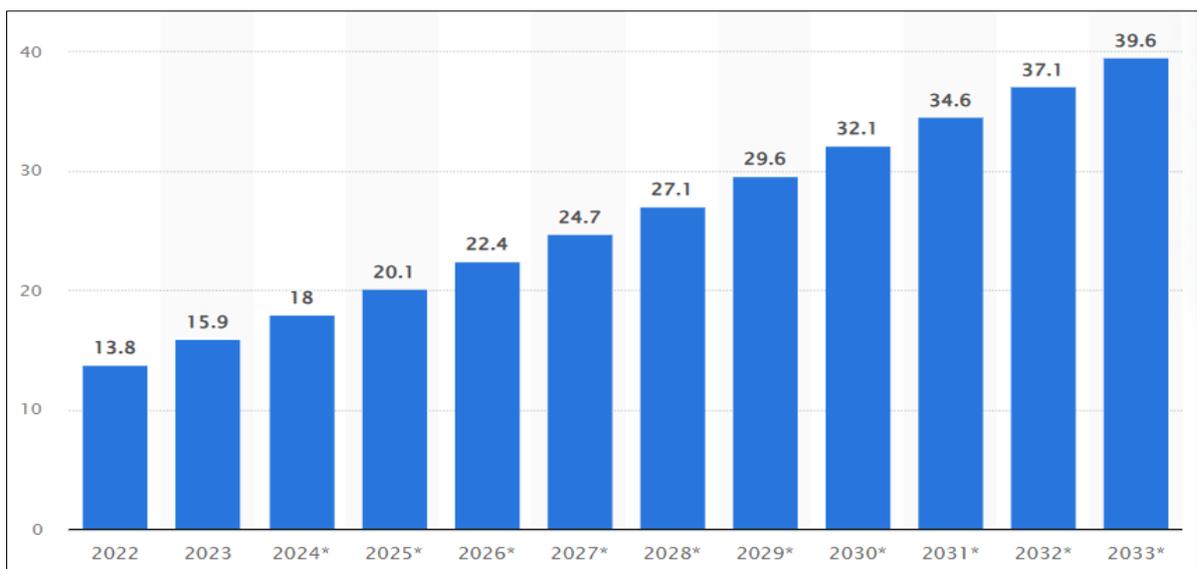
1.1.1. Planteamiento del problema.

A lo largo de la historia la tecnología ha influenciado en la vida de las personas, generando cambios positivos como negativos, dichos cambios principalmente se dan según el modo en que se utilice la tecnología en general. Una de las tecnologías actuales más utilizadas en la vida cotidiana es internet.

Del mismo modo que la Internet transformó radicalmente la sociedad, la Internet de las cosas (Internet of Things o IoT por sus siglas en inglés) repercutirá en todos los ámbitos de la vida humana: desde nuestros hogares, vehículos, lugares de trabajo y fábricas, hasta nuestras ciudades y pueblos, la agricultura y los sistemas sanitarios. También afectará a todos los niveles de la sociedad (individuos, empresas y estado), desde lo urbano hasta lo rural, pasando por el mundo natural y más allá. Por ello, es fundamental tener un conocimiento adecuado de la IoT y de los retos que conlleva. (Bacelli, 2021, p. 3).

La presencia y el alcance del IoT es global. Sujay (2024) pronosticó que el número de dispositivos IoT, en todo el mundo puede llegar a casi duplicarse, de 15.9 mil millones en 2023, a más de 32.1 mil millones de dispositivo IoT en 2030.

Figura 1: Dispositivos IoT conectados a nivel mundial (en miles de millones)



Fuente: <https://www.statista.com/statistics/1183457/iot-connected-devices-worldwide/>

En la figura 1 se puede apreciar el incremento considerable del uso de la tecnología IoT, lo cual implica el surgimiento de nuevas oportunidades, así como también desafíos ya deducidos y aún no previstos.

Latinoamérica también está experimentando el crecimiento notable del IoT, lo cual genera oportunidades sin precedentes en términos de ventas y crecimiento empresarial, se estima que en la región el mercado de IoT alcance los 62 mil millones de dólares para el año 2025, pero este crecimiento prometedor se enfrenta a desafíos significativos, principalmente la falta de regulaciones claras en materia de seguridad y privacidad de datos, así como la falta de infraestructura de conectividad en algunas áreas rurales. (Moabits, 2024).

La Asociación Peruana de Telecomunicaciones (APTC) indicó que el alcance de internet ha llegado al 74% de la población y para el 2032 se espera que el alcance llegue al 85% y que, con iniciativas del gobierno, así como de entidades privadas, se extienda la cobertura de internet a áreas rurales y urbanas marginales. Además, el Perú está en el proceso de implementación de redes 5G, que permitirán el impulso de nuevas aplicaciones y servicios, desde el IoT hasta soluciones avanzadas de telemedicina y educación a distancia. (APTC, 2024).

El acceso a internet y la adopción de tecnologías digitales han generado un impacto positivo y bienestar social en el país, pero como afirmó Villón (2024) la digitalización de la sociedad peruana ha abierto nuevas puertas a diversas amenazas cibernéticas que pueden afectar directamente a la seguridad nacional y la vida cotidiana de sus ciudadanos. Estas amenazas cibernéticas pueden tener como objetivo robar información, causar daños o interrumpir servicios.

Se reconoce que los avances tecnológicos benefician a la sociedad que los utiliza, pero es muy importante considerar que dicho uso también presenta riesgos que deben ser abordados oportunamente, para la seguridad y bienestar de la población en general. El Centro Nacional de Planeamiento Estratégico (CEPLAN) indicó que la combinación de tecnologías, entre ellas la inteligencia artificial (IA) y el IoT, ofrece una oportunidad para mejorar la seguridad digital. Con la IA se puede identificar

patrones y anomalías en los datos, lo cual permite la detección temprana de amenazas cibernéticas y a la vez una respuesta automatizada a ellas. (CPLAN, 2024).

1.1.2. Formulación del problema.

Problema general

¿Cómo desarrollar un modelo de clasificación para la detección de anomalías en redes IoT utilizando el dataset CICIoT2023?

Problemas específicos

- ¿Cómo analizar y procesar los datos del dataset CICIoT2023?
- ¿Cómo diseñar un modelo de clasificación para la detección de anomalías en el dataset CICIoT2023?
- ¿Cómo evaluar el modelo de clasificación para la detección de anomalías en el dataset CICIoT2023?

1.2. Determinación de objetivos

1.2.1. Objetivo general.

Desarrollar un modelo de clasificación para la detección de anomalías en redes IoT utilizando el dataset CICIoT2023.

1.2.2. Objetivos específicos.

- Realizar el análisis exploratorio y el procesamiento de datos del dataset CICIoT2023.
- Diseñar un modelo de clasificación para la detección de anomalías en el dataset CICIoT2023.
- Evaluar el modelo de clasificación para la detección de anomalías en el dataset CICIoT2023.

1.3. Justificación e importancia del estudio

1.3.1. Justificación teórica.

Con el creciente número de dispositivos interconectados, Baccelli (2021) indicó que con la combinación de tecnologías se puede transformar el modo en que operan las empresas, así como la producción de bienes a través de la automatización de procesos, pero que a la vez surge el riesgo de ciberataques, por tal motivo, el IoT asistido por la IA puede ayudar a detectar y prevenir ciberataques al analizar el tráfico de la red, identificar anomalías y responder a las amenazas en tiempo real.

El presente trabajo de investigación contribuirá en el campo de la ciberseguridad y el aprendizaje automático aplicado a redes IoT. La investigación explora la interacción de la IA con la seguridad de IoT a través del estudio del dataset CICIoT2023, en el dataset (conjunto de datos) mencionado se ejecutaron 33 ataques en condiciones reales a 105 dispositivos IoT, lo cual le brinda relevancia y pertinencia para abordar retos actuales. El uso del dataset CICIoT2023 ofrece una valiosa oportunidad de analizar y modelar algoritmos de aprendizaje automático para clasificar y evaluar patrones de tráfico en redes IoT actuales. Este dataset, que incluye como ya se mencionó una amplia gama de dispositivos IoT comerciales, permite una comprensión más profunda de las características y vulnerabilidades específicas de los dispositivos IoT, lo cual puede generalizarse a otros dispositivos de comportamiento y características similares. (Pinto et al., 2023).

1.3.2. Justificación metodológica.

Para abordar el presente estudio se iniciará con el análisis exploratorio de datos (EDA), que según IBM (2023), los EDA se utilizan para analizar conjuntos de datos y resumir sus características principales, a menudo empleando métodos de visualización de datos. En ese sentido el EDA será útil para comprender las características, el modo de organización, posibles errores y patrones, tipos de datos, entre otros aspectos relevantes del dataset CICIoT2023.

También se debe atender la necesidad de comprender y aplicar técnicas de aprendizaje automático, Alpaydin (2020) indicó que a partir de programas o algoritmos de aprendizaje se pueden procesar grandes cantidades de datos, de diferentes áreas del conocimiento, para crear modelos descriptivos y predictivos, además se deben optimizar dichos modelos mediante técnicas de entrenamiento para hacerlos más eficientes en el reconocimiento de patrones, así como predicciones futuras. Por lo mencionado se entiende que los modelos que aprendizaje automático nos serán útiles para clasificar y en consecuencia detectar anomalías en el dataset CICIoT2023.

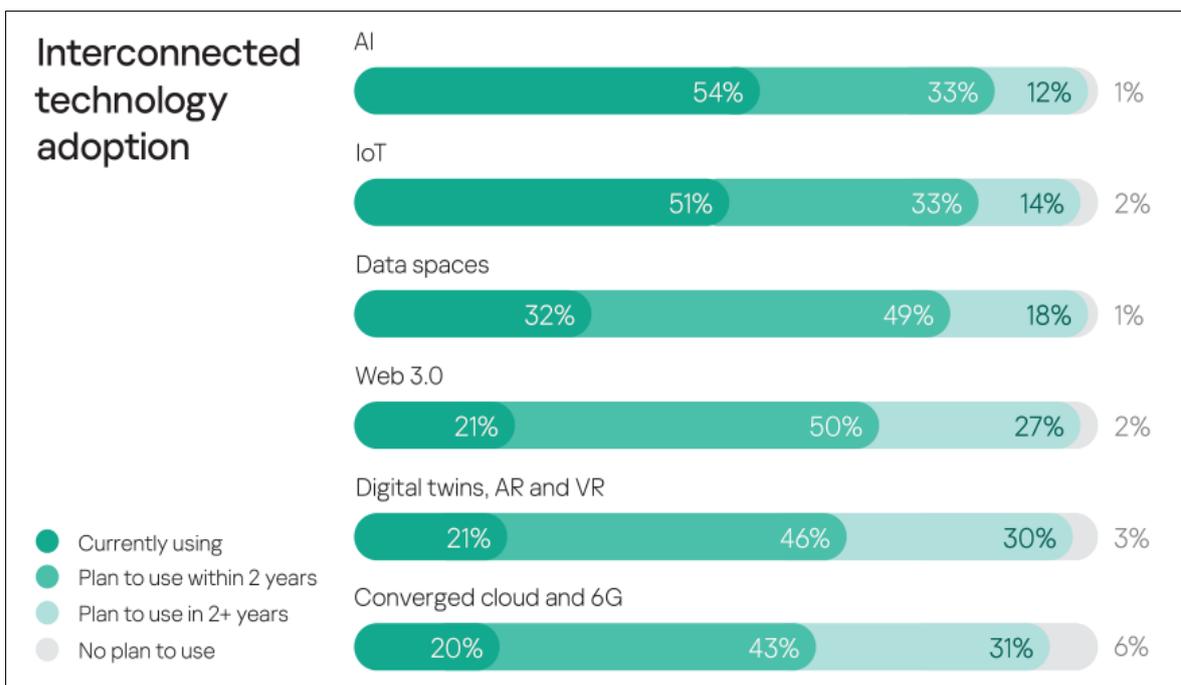
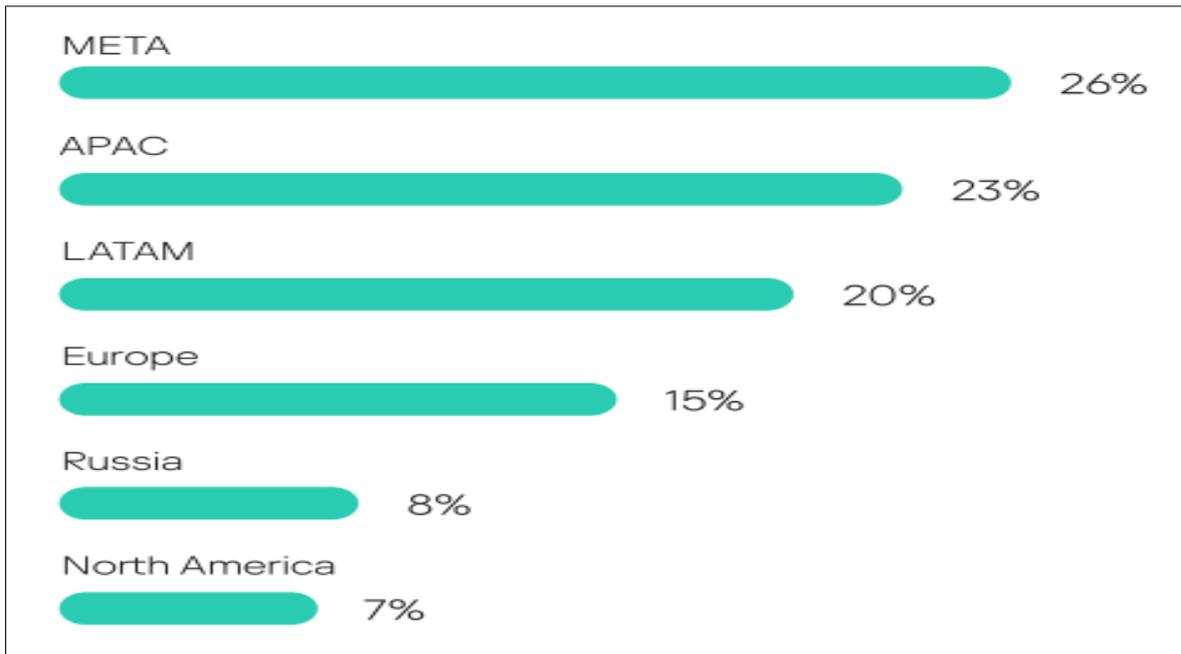
1.3.3. Justificación social.

Junto al desarrollo y avance del IoT aumentará también nuestra dependencia vital, según Baccelli (2021), está surgiendo una tendencia de implantes inteligentes, como es el caso de los sistemas de páncreas artificiales (APS: Artificial Pancreas System) que, mediante el ajuste del bombeo de insulina basal, se mantiene la glucosa de la sangre en un rango seguro durante la noche y entre comidas. Si el APS se avería, la vida del paciente corre peligro inmediatamente. Además, desde un enfoque más general se han producido incidentes de seguridad y fiabilidad de la IoT, que afectan a una gran variedad máquinas, automóviles conectados e incluso aviones, debido a fallas de la IoT, las vidas de los usuarios corrieron peligro o fueron víctimas de accidentes fatales. En consecuencia, la aplicación del IoT se está extendiendo a diferentes áreas y aspectos de la vida cotidiana, que genera beneficio a la población, pero aún se debe seguir investigando para que el IoT sea más confiable y seguro.

Actualmente la IA y el IoT tienen presencia y participación en los procesos empresariales, Kaspersky (2024) informó, a partir de una encuesta realizada a 560 líderes senior de seguridad informática de América del Norte, América Latina, Europa, Oriente Medio y África, Rusia y Asia-Pacífico, que más del 50% de empresas ya han implementado la IA y el IoT dentro de sus infraestructuras y que el 33% de empresas ya están planeando adoptar estas tecnologías interconectadas dentro de dos años (Ver figura 2). Por lo cual los líderes empresariales deben asegurarse de hacer inversiones tecnológicas adecuadas contando para ello con

personas y sistemas necesarios que garanticen el éxito. La adopción de tecnologías interconectadas, brinda inmensas oportunidades comerciales, impulsan cambios en todas las industrias y también traen consigo nuevos riesgos y desafíos que deben ser atendidos oportunamente para asegurar a las empresas y para proteger a sus clientes brindándoles seguridad y confianza.

Figura 2: Regiones consideradas y Tecnologías interconectadas



Fuente: Kaspersky (2024).

Ya es notorio en nuestras vidas y a la vez se reconoce que el IoT, está revolucionando sectores diversos, como los hogares inteligentes, el transporte y las ciudades, según Ouaisa et al. (2025), en el ámbito del hogar, dispositivos como termostatos inteligentes, luces, cámaras y sistemas de seguridad permiten a los propietarios gestionar sus hogares de manera eficiente y a distancia. En el transporte, la tecnología IoT mejora la seguridad y la eficiencia a través de infraestructuras conectadas y vehículos autónomos. Asimismo, la creación de ciudades inteligentes permite optimizar recursos y reducir la contaminación, promoviendo un desarrollo urbano sostenible.

También afirma que, en cuanto a la salud y la atención médica, el IoT actualmente ya desempeña un papel fundamental en dicho sector, al permitir el monitoreo remoto de pacientes a través de dispositivos de monitoreo como sensores vestibles (llevados sobre el cuerpo) que permiten el seguimiento de signos vitales y niveles de actividad, así como plataformas de telemedicina donde se realizan consultas virtuales. La gestión de los medicamentos se optimiza a través de dispensadores inteligentes y aplicaciones que recuerdan a los pacientes, la ingesta de sus respectivos medicamentos, la frecuencia y las dosis recetadas. Además, la monitorización de salud y bienestar se realiza mediante dispositivos que analizan hábitos de actividades físicas y del sueño, ofreciendo recomendaciones personalizadas según lo requiera el usuario. Importante también es la respuesta a emergencias que se ven favorecidas por dispositivos IoT que brindan información en tiempo real acerca del estado de los pacientes y las condiciones del entorno.

Finalmente se concluye que el impacto del IoT en diversos sectores es innegable, destacando su potencial en la mejora de la calidad de vida y en hacer más eficientes los sistemas existentes, pero es vital abordar cuestiones de privacidad, seguridad y accesibilidad para maximizar sus beneficios, puesto que a medida que avanzamos hacia un futuro cada vez más conectado y digital, la integración del IoT en la vida cotidiana tiende a revolucionar la forma en que interactuamos entre nosotros y con nuestro entorno, ofreciéndonos soluciones innovadoras y personalizadas que mejoren nuestras condiciones de vida, la atención médica e incluso la gestión de recursos y el cuidado medio ambiente.

1.4. Limitaciones de la presente investigación

Se han identificado y considerado las siguientes limitaciones:

- Calidad y representatividad del dataset, pues, aunque CICIoT2023 es extenso y diverso, su calidad puede verse afectada por la representación limitada de ciertos tipos de ataques o comportamientos anómalos.
- Complejidad de los modelos de aprendizaje automático, porque la complejidad excesiva puede generar sobreajuste durante el entrenamiento, dicha limitación puede comprometer la capacidad de los modelos para detectar anomalías en entornos dinámicos.
- Dinamismo de las amenazas cibernéticas que evolucionan rápidamente, siendo necesario actualizar frecuentemente los modelos de aprendizaje.
- Recursos computacionales, las limitaciones de este aspecto restringen el alcance y profundidad del entrenamiento y evaluación de los modelos.

Aspectos legales y regulatorios, el estudio no se centra en normas legales de ninguna jurisdicción.

Capítulo II: Marco teórico

2.1. Antecedentes de la investigación

2.1.1. Internacionales.

Las investigaciones internacionales sobre la detección de anomalías en redes de IoT han presentado diversos enfoques y metodologías innovadoras. La investigación realizada por Tseng et al. (2024) en Taiwán propone un método de detección de intrusiones en redes IoT utilizando el dataset CICIoT2023. Este estudio aplicó siete modelos de aprendizaje profundo, destacando el uso del modelo Transformer y logró alcanzar una exactitud del 99.40% en clasificación multiclase. Por otro lado, en Canadá, Pinto et al. (2023), desarrollaron un extenso dataset sobre ataques en entornos IoT, conocido como CICIoT2023, fue desarrollado en condiciones reales, que incluye 33 tipos de ataques clasificados en siete categorías y un tráfico benigno, este dataset es crucial para fomentar el desarrollo de aplicaciones de análisis de seguridad y para evaluar los métodos de aprendizaje automático en diversos escenarios, los desarrolladores del dataset también le ejecutaron cinco modelos clasificación, siendo Random Forest el que alcanzó la mayor exactitud de 99.68%. En otro contexto, Ragab et al. (2023), en Arabia Saudita, se centraron en la detección de ataques DDoS mediante un algoritmo de optimización y un clasificador de aprendizaje profundo (PHHO-ODLC), logrando una exactitud del 99.20%, lo que subraya la eficacia de su enfoque en la mejora de la seguridad IoT. En España, Vigoya (2023) estudió la detección de anomalías en el tráfico de IoT, utilizando varios algoritmos de aprendizaje automático, su trabajo generó dos datasets en escenarios virtuales y se basa en la implementación de modelos de clasificación, logrando una exactitud de 95.24% con SVM y 99.99% con Random Forest. Finalmente, Plata (2022) en Colombia, desarrolló un sistema de detección de anomalías físicas en redes IoT empleando el algoritmo K-Nearest Neighbors (KNN), obteniendo una exactitud del 99.53% en la identificación de anomalías. Estos cinco estudios destacan la variedad de enfoques implementados en la detección de anomalías y su relevancia en el contexto actual del IoT.

Las investigaciones internacionales mencionadas son altamente relevantes para el presente estudio. En primer lugar, el estudio de Tseng et al. (2024) se relaciona directamente al abordar la detección de intrusiones en redes IoT, lo que puede proporcionar valiosas referencias metodológicas para la construcción del modelo de aprendizaje automático a desarrollar en esta investigación. Además, el dataset CICIoT2023, propuesto por Pinto et al. (2023), es fundamental para la validación del modelo, al ser una base extensa que incluye diversos tipos de ataques que se pueden utilizar para evaluar eficazmente el rendimiento del modelo. Por otro lado, el enfoque de Ragab et al. (2023) en la detección DDoS podría ser considerado al desarrollar un modelo que no solo identifique intrusiones, sino también otros comportamientos anómalos en el tráfico IoT. Esto permitirá ampliar el rango de detección más allá de ataques específicos. La investigación de Vigoya (2023) ofrece un enfoque alternativo, proporcionando ejemplos de cómo técnicas de aprendizaje automático pueden ser aplicadas para detectar anomalías de tráfico, sugiriendo que diferentes algoritmos podrían ser probados en el modelo que se desarrollará. Finalmente, el estudio de Plata (2022) complementa este marco al resaltar la importancia del análisis del tráfico de red y cómo los sistemas de detección deben ser capaces de identificar y reaccionar ante diversas anomalías, lo cual es crucial para la seguridad en un entorno IoT. En conjunto, estas investigaciones ofrecen un sólido trasfondo que justificará y guiará el desarrollo del modelo propuesto en la tesis.

2.1.2. Nacionales.

En Perú también se han realizado investigaciones que enriquecen el contexto de la seguridad en las redes IoT. La investigación de Becerra et al. (2024) se centró en la evaluación del rendimiento de modelos de aprendizaje profundo, utilizando el dataset CICIoT2023 para clasificar ataques cibernéticos en redes IoT. Este estudio señala que el modelo basado en CNN obtuvo una exactitud del 99.10% en la clasificación multiclase y del 99.40% en la clasificación binaria, destacando la efectividad de esta arquitectura en la identificación de amenazas. Asimismo, Becerra et al. (2024) con un segundo equipo de investigadores también analizaron la mejora de la detección de ataques DDoS, pero en dicha investigación utilizaron otro dataset, el CICDDoS2019, donde el clasificador Random Forest superó otros

modelos con un rendimiento del 99.97%, subrayando la eficacia de los algoritmos de aprendizaje automático en este tipo de amenazas.

Además, la investigación de León et al. (2024) aportó un enfoque práctico al implementar un sistema que detecta anomalías en los logs de dispositivos IoT para prevenir accesos no autorizados. Este estudio demostró una efectividad que varió entre el 65% y el 100% dependiendo del volumen de logs analizados, lo cual destaca la importancia de la cantidad de datos en los sistemas de detección. En la misma línea, Bazan et al. (2024) se centraron en la clasificación de ataques ransomware, alcanzando el 99.4% de exactitud con Árboles de decisión y 99.6% con Random Forest, lo cual resalta la capacidad de los modelos de aprendizaje automático para identificar y clasificar amenazas específicas en entornos cibernéticos. Finalmente, Nolasco (2023) estudió la aplicación de aprendizaje automático en el pronóstico del desplazamiento de lluvias utilizando imágenes de radar. Aunque su enfoque se dirige a un área diferente, su investigación ilustra la versatilidad de las técnicas de aprendizaje automático y su impacto en la resolución de problemas complejos, lo que complementa el conocimiento necesario para el desarrollo del modelo de detección de anomalías en redes IoT.

Las investigaciones nacionales presentadas tienen una relación directa y significativa con la propuesta del presente trabajo de investigación. En particular, el trabajo de Becerra et al. (2024) es invaluable, ya que no solo evalúa el desempeño de diferentes modelos sobre el mismo dataset que se utilizará en la tesis, sino que establece un estándar que permitirá comparar los resultados del modelo propuesto. Además, su hallazgo de que el modelo CNN es altamente efectivo puede guiar la selección y el diseño del modelo a implementar en la investigación actual. Por otra parte, la investigación sobre la mejora de la detección de ataques DDoS en el dataset CICDDoS2019, de Becerra et al. (2024) y su segundo equipo de investigadores, enfatizó la importancia del preprocesamiento de datos y la selección de características, elementos que serán cruciales para la construcción de un modelo robusto que pueda detectar una variedad de anomalías en el tráfico IoT. La contribución de León et al. (2024) refuerza este enfoque al resaltar que un sistema autónomo de detección debe operar efectivamente sobre grandes volúmenes de logs, sugiriendo que la eficacia del modelo dependerá de la

calidad y cantidad de datos disponibles. Adicionalmente, el enfoque de Bazan et al. (2024) en ataques ransomware muestra que las técnicas de aprendizaje automático pueden ser efectivas en la identificación de diversas amenazas. Esto implica que el modelo a desarrollar no solo debe centrarse en intrusiones, sino también en clasificar otros comportamientos anómalos, reforzando la necesidad de un enfoque amplio que contemple todos los tipos de amenazas en un entorno interconectado. Por último, la investigación de Nolasco (2023), aunque distinta, confirma cómo las metodologías de aprendizaje automático pueden adaptarse y ser útiles en diferentes contextos, sugiriendo que un enfoque flexible y adaptable será clave en el desarrollo del modelo de detección propuesto.

2.2. Bases teóricas

2.2.1. Desarrollo histórico.

El Internet de las Cosas (IoT) ha evolucionado en varias décadas, transformando nuestra interacción con la tecnología y el entorno. Originalmente, se refería a la conexión de dispositivos físicos, pero ahora abarca un ecosistema más amplio que incluye sensores, software y conectividad (IBM, 2023). La idea de dispositivos interconectados comenzó en 1982, cuando en la Universidad de Carnegie Mellon, se modificó una máquina expendedora para que informara sobre su contenido y temperatura (González et al., 2020). Este avance sentó las bases para el desarrollo futuro del IoT. El término "Internet de las Cosas" fue acuñado por Kevin Ashton en 1999, quien propuso el uso de etiquetas RFID (Radio Frequency Identification - Identificación por Radiofrecuencia) para mejorar la eficiencia en las cadenas de suministro (Ashton, 2009). Su relevancia aumentó con la publicación de un informe de la Unión Internacional de Telecomunicaciones (UIT) en 2005, que abordaba las implicaciones sociales y económicas del IoT. Este informe, junto a las proyecciones de crecimiento en la conexión de dispositivos, mostró cómo el IoT estaba comenzando a cambiar industrias como la manufactura, el transporte y la agricultura (Ávila y Moreno, 2023). A medida que el IoT ha crecido, también lo han hecho los desafíos relacionados con la seguridad y la privacidad de los datos. Esta preocupación ha llevado al desarrollo de sistemas de detección de anomalías, que permiten identificar patrones inusuales y proteger la información (Muñoz, 2023).

Esto es crucial en la era de la Industria 4.0, donde la interconexión de sistemas es fundamental para optimizar procesos y crear modelos de negocio sostenibles (Ávila y Moreno, 2023). El análisis de datos generados por dispositivos IoT ha avanzado, utilizando técnicas de aprendizaje automático para mejorar la detección de anomalías. Estas técnicas se aplican en sectores como la gestión de residuos y la eficiencia energética, optimizando procesos y promoviendo la sostenibilidad (Osio et al., 2021). Los algoritmos de Machine Learning se convierten en herramientas clave para crear modelos predictivos que faciliten decisiones informadas en tiempo real (Muñoz, 2023). Mirando hacia el futuro, el IoT tiene un gran potencial, impulsado por el aumento de dispositivos conectados y su integración con tecnologías emergentes, como la inteligencia artificial (IA - Artificial Intelligence) y blockchain (cadena de bloques) (IBM, 2023). A medida que este ecosistema evoluciona, las oportunidades para mejorar la vida cotidiana y la eficiencia empresarial seguirán creciendo. La combinación de aprendizaje automático y sistemas de detección de anomalías permitirá que las empresas no solo protejan sus infraestructuras, sino que también aprovechen el amplio potencial de datos del IoT (González, L. et al., 2020).

2.2.2. Fundamentación teórica.

Internet de las Cosas (IoT - Internet of Things)

Es la creciente red de dispositivos conectados a Internet que pueden comunicarse y compartir información entre sí. Esto incluye desde electrodomésticos, como refrigeradores y bombillas inteligentes, hasta sistemas de control industrial. Gracias a la integración de sensores y tecnología avanzada, estos dispositivos recopilan datos en tiempo real y responden a las necesidades del usuario. El IoT ha permitido que miles de millones de objetos recopilen y transmitan datos, lo que crea un vasto conjunto de información que puede ser analizada para mejorar la eficiencia y la seguridad en diferentes entornos (Amazon Web Services, 2022). Este fenómeno destaca la importancia de contar con herramientas adecuadas para interpretar y actuar sobre los datos generados.

Inteligencia artificial (IA)

Es un campo de la ciencia que se enfoca en crear computadoras y máquinas que pueden razonar, aprender y actuar de manera similar a los humanos, especialmente en el manejo de grandes volúmenes de datos. Este campo abarca diversas disciplinas, como la informática, la estadística y la neurociencia. En el ámbito empresarial, la IA se basa principalmente en tecnologías de aprendizaje automático y profundo, utilizadas para análisis de datos, predicciones, procesamiento del lenguaje natural y más. Los sistemas de IA aprenden a través de datos y algoritmos, mejorando su rendimiento con el tiempo. La IA puede automatizar procesos, reducir errores humanos y trabajar sin descanso, lo que acelera la investigación y mejora la eficiencia en diversas aplicaciones, desde el reconocimiento de voz hasta la seguridad cibernética (Google Cloud, 2024).

Detección de anomalías

Es el proceso mediante el cual se identifican puntos de datos que se desvían significativamente del comportamiento esperado dentro de un conjunto de datos, este proceso es fundamental para garantizar la operación segura y eficiente de los sistemas, ya que las anomalías pueden indicar problemas serios, como fallas en equipos o intentos de intrusión maliciosa (IBM, 2024). En el contexto del IoT, la detección de anomalías permite a los sistemas discernir entre errores normales de datos y situaciones críticas que requieren atención inmediata, por lo tanto, desarrollar modelos de aprendizaje automático para esta tarea no solo mejora la calidad de los datos, sino que también optimiza la respuesta ante incidentes, creando un entorno más seguro y confiable en el que las decisiones se basan en análisis precisos y oportunos.

2.2.3. Marco conceptual.

A. Definición conceptual de la variable 1: Características de la anomalía.

Las características de la anomalía con las que se cuenta en el dataset CICIoT2023 son las propiedades o atributos específicos que describen el tráfico de red en un entorno IoT, Navaro (2024), les denomina columnas a estos atributos específicos del dataset. Estas características se utilizan para clasificar y distinguir entre

diversos tipos de tráfico en términos de benigno o malicioso, siendo fundamentales para la identificación de patrones y la detección de anomalías en entornos IoT.

Dimensiones de la variable 1.

- Definición conceptual de Cantidad de Características**

El dataset CICIoT2023 tiene 40 características de datos generados en condiciones reales, donde se ejecutaron 33 ataques, clasificados en 7 categorías por su similitud de comportamiento, dirigidos a 105 dispositivos IoT. Pinto et al. (2023).
- Definición conceptual de Tipos de Características**

Las características (ver Cuadro 1) incluyen métricas de tráfico, datos como, contadores de paquetes, duración del flujo, tipos de protocolos y flags que indican el estado de paquetes enviados o recibidos. Pinto et al. (2023).

Cuadro 1: Características del tráfico de red del dataset CICIoT2023

Nº	CARACTERÍSTICA	DESCRIPCIÓN
1	Header Length	Longitud de cabecera
2	Protocol Type	Tipo de Protocolo: IP, UDP, TCP, IGMP, ICMP, Desconocido (Enteros)
3	Time_To_Live	Tiempo de vida (ttl)
4	Rate	Tasa de transmisión de paquetes en un flujo
5	fin flag number	Valor del flag Fin
6	syn flag number	Valor del flag Syn
7	rst flag number	Valor del flag Rst
8	psh flag numbe	Valor del flag Psh
9	ack flag number	Valor del flag Ack
10	ece flag numbe	Valor del flag Ece
11	cwr flag number	Valor del flag Cwr
12	ack count	Número de paquetes con el flag ack establecido en el mismo flujo
13	syn count	Número de paquetes con el flag syn establecido en el mismo flujo
14	fin count	Número de paquetes con el flag fin establecido en el mismo flujo
15	rst count	Número de paquetes con el flag rst establecido en el mismo flujo

16	HTTP	Indica si el protocolo de la capa de aplicación es HTTP
17	HTTPS	Indica si el protocolo de la capa de aplicación es HTTPS
18	DNS	Indica si el protocolo de la capa de aplicación es DNS
19	Telnet	Indica si el protocolo de la capa de aplicación es Telnet
20	SMTP	Indica si el protocolo de la capa de aplicación es SMTP
21	SSH	Indica si el protocolo de la capa de aplicación es SSH
22	IRC	Indica si el protocolo de la capa de aplicación es IRC
23	TCP	Indica si el protocolo de la capa de transporte es TCP
24	UDP	Indica si el protocolo de la capa de transporte es UDP
25	DHCP	Indica si el protocolo de la capa de aplicación es DHCP
26	ARP	Indica si el protocolo de la capa de enlace es ARP
27	ICMP	Indica si el protocolo de la capa de red es ICMP
28	IGMP	Indica si el protocolo de la capa de red es IGMP
29	IPv	Indica si el protocolo de la capa de red es IP
30	LLC	Indica si el protocolo de la capa de enlace es LLC
31	Tot sum	Suma total de las longitudes de los paquetes en flujo
32	Min	Longitud mínima del paquete en el flujo
33	Max	Longitud máxima del paquete en el flujo
34	AVG	Longitud promedio del paquete en el flujo
35	Std	Desviación estándar de la longitud del paquete en el flujo
36	Tot size	Tamaño total del paquete
37	IAT	La diferencia temporal con respecto al paquete anterior
38	Number	El número total de paquetes en el flujo
39	Variance	Varianza de las longitudes de los paquetes entrantes en el flujo / varianza de las longitudes de los paquetes salientes en el flujo
40	Label	Etiqueta del tráfico

Fuente: Pinto et al. (2023).

B. Definición conceptual de la variable 2: Tipo de Anomalías.

Es el proceso de identificar comportamientos o patrones inusuales en el tráfico de datos que pueden indicar problemas como el mal funcionamiento de los dispositivos, situaciones de riesgo, ataques cibernéticos o fallos en el sistema. Esto incluye la identificación de valores atípicos y cambios inesperados en las métricas de rendimiento. Además, complementa IBM (2024), que la detección de anomalías

requiere una atención inmediata para garantizar la operación segura y eficiente de los sistemas comprometidos.

Dimensiones de la variable 2.

- **Definición conceptual de Tráfico de Datos**
El tráfico de datos se determina, según sus métricas y comportamiento, como: malicioso o benigno, con algoritmos y técnicas de aprendizaje automático.
- **Definición conceptual de Métodos de Detección**
La detección de anomalías se ejecutará con los siguientes algoritmos de aprendizaje automático:
 - ✓ Regresión Logística
 - ✓ Árboles de decisión
 - ✓ Random Forest
 - ✓ ADA Boost
 - ✓ Perceptrón

2.3. Definición de términos básicos

- **Algoritmo:** Conjunto finito de instrucciones precisas y bien definidas que describe un proceso o procedimiento para resolver un problema o llevar a cabo una tarea específica. (Diccionario sobre inteligencia artificial, 2024).
- **Aprendizaje Automático (Machine Learning):** Es un componente fundamental de la inteligencia artificial que permite a las máquinas aprender patrones y mejorar su rendimiento sin intervención humana explícita. Utiliza algoritmos que permiten a las máquinas reconocer patrones en datos y tomar decisiones basadas en estos patrones, mejorando su desempeño con la experiencia. (Diccionario sobre inteligencia artificial, 2024).
Es la técnica utilizada para entrenar un modelo de aprendizaje automático. Hay diversos algoritmos (regresión lineal, máquinas de vectores de soporte, redes neuronales artificiales) y cada uno tiene sus propias formulaciones y complejidades. Sin embargo, todos ellos tienen como propósito reducir el

margen de error entre las predicciones de los modelos y el resultado deseado de los conjuntos de datos de entrenamiento. (ISO, 2024).

- **Ciberataque:** Intento deliberado de un ciberdelincuente de obtener acceso a un sistema informático sin autorización sirviéndose de diferentes técnicas y vulnerabilidades para la realización de actividades con fines maliciosos, como el robo de información, extorsión del propietario o simplemente daños al sistema. (Glosario de términos de ciberseguridad, 2020).
- **Clasificación binaria:** Es un tipo de clasificación que se basa en predecir una de dos clases, generalmente representando un estado normal y otro anormal. Algunos ejemplos comunes incluyen la detección de spam en correos electrónicos (spam o no), la predicción de abandono de clientes (churn o no) y la predicción de conversión (comprar o no). Para este tipo de tareas, se utilizan algoritmos como la regresión logística, k-vecinos más cercanos, árboles de decisión, máquinas de soporte vectorial y Naive Bayes. (Ouaissa, 2025).
- **Dato:** Unidad mínima de información (números, letras o símbolos) que representa un objeto, condición o situación y que requiere una interpretación para convertirse en información. (Glosario de términos TIC, 2021).
- **Detección de Anomalías (Outlier Detection):** Es el proceso de identificar observaciones inusuales o atípicas en un conjunto de datos. Estas observaciones pueden indicar errores en los datos, comportamientos anómalos o información relevante pero no típica. (Diccionario sobre inteligencia artificial, 2024).
- **Log:** Registros de eventos de la actividad de los usuarios y de los procesos asociados a dicha actividad, como pueden ser el inicio/salida de sesión, tiempo de actividad o conexiones, entre otros. Esta información ayuda a detectar fallos de rendimiento, mal funcionamiento, errores e intrusiones que permiten generar alertas en tiempo real gracias a los datos proporcionados a los sistemas de monitorización. (Glosario de términos de ciberseguridad, 2020).
- **Malware (Software malicioso):** Es un tipo de software que tiene como objetivo dañar o infiltrarse sin el consentimiento de su propietario en un sistema de información. Palabra que nace de la unión de los términos en

inglés de software malintencionado: malicious software. (Glosario de términos de ciberseguridad, 2020).

El malware representa riesgos sustanciales para la seguridad y funcionalidad de los sistemas informáticos. Puede resultar en el robo de información o datos privados, como información personal, financiera o confidencial. Además, el malware puede corromper archivos críticos, interrumpir operaciones y reducir el rendimiento general del sistema, causando daños extensos tanto a personas como a instituciones. (Ouaisa, 2025).

- **Modelo de aprendizaje automático:** Es una representación matemática que resulta de aplicar un algoritmo de aprendizaje automático a un conjunto de datos. Su principal objetivo es generalizar a partir de los datos de entrenamiento, lo que significa que debe ser capaz de hacer predicciones o tomar decisiones sobre nuevos datos que no se han visto antes, en lugar de simplemente memorizar los ejemplos con los que fue entrenado. Es la herramienta final que se utiliza para realizar predicciones o análisis en aplicaciones del mundo real. (ISO, 2024).
- **Redes Neuronales:** Modelos computacionales inspirados en el cerebro humano que están diseñados para realizar tareas de aprendizaje automático. Consisten en una red de unidades de procesamiento llamadas neuronas, que están organizadas en capas y conectadas entre sí mediante conexiones ponderadas. (Diccionario sobre inteligencia artificial, 2024).

Capítulo III: Hipótesis y variables

3.1. Hipótesis de investigación

3.1.1. Hipótesis general.

El modelo de clasificación desarrollado será capaz de detectar anomalías en redes IoT con una precisión superior al 90% utilizando el dataset CICIoT2023.

3.1.2. Hipótesis específicas.

- El análisis exploratorio y el procesamiento de datos mejorarán la calidad del dataset, lo que generará un rendimiento superior del modelo de clasificación en la detección de anomalías.
- El diseño de un modelo de clasificación permitirá identificar patrones anómalos en el tráfico de red con una tasa de detección significativamente alta.
- La evaluación del modelo mostrará que las métricas de rendimiento, Exactitud y Precisión alcanzan el rendimiento esperado para considerar al modelo efectivo en la detección de anomalías en redes IoT.

3.2. Operacionalización de variables

3.2.1. Variable Independiente.

Definición Operacional de las Características del Dataset.

Dimensiones:

- Cantidad de Características: 40
- Tipos de Características:
 - Métricas de tráfico: Número de paquetes (enviados y recibidos), duración de flujo (en milisegundos).
 - Protocolos: TCP, UDP, etc.
 - Flags de estado: Indicadores del estado del paquete (SYN, ACK, etc.).
 - Clasificación de ataques: Agrupados en DDoS, DoS, Reconocimiento, Basados en la web, Fuerza bruta, Suplantación y Mirai.

Indicadores:

- Número total de características: 40
- Métricas de las características: Número de paquetes, duración del flujo, etc.
- Clasificación de ataques: 7 categorías (DDoS, DoS, Reconocimiento, Basados en la web, Fuerza bruta, Suplantación y Mirai).

Escalas de valoración:

- Escala categórica: Clasificar el tráfico como malicioso o benigno.
- Escala numérica: Contar las ocurrencias eventos maliciosos o benignos.

Instrumentos:

- Herramientas de análisis de datos:
 - Pandas y NumPy para el manejo y manipulación de datos.
 - Scikit-learn para la implementación de técnicas de preprocesamiento.
- Análisis Exploratorio de Datos (EDA):
 - Utilizar técnicas de visualización (gráficas de barras, de dispersión, histogramas, boxplots) para entender la distribución y relaciones entre las características.
 - Métodos de correlación de Pearson para identificar relaciones entre las características.
- Técnicas de Preprocesamiento:
 - Selección de características: Aplicar herramientas de selección de características para identificar a las más relevantes para la variable objetivo.
 - Escalamiento: Ajustar las características a una misma escala para mejorar el rendimiento del modelo.
 - Eliminación de valores: Identificar y eliminar valores atípicos, duplicados, negativos, infinitos u otros que pueden afectar el rendimiento del modelo.
 - Eliminación de columnas duplicadas: Identificar las características que tienen totalmente iguales o redundantes.

3.2.2. Variable dependiente.

Definición Operacional de Detección de Anomalías.

Dimensiones

- Tipo de tráfico:
 - Malicioso.
 - Benigno.
- Modelos de Clasificación: Como herramientas para detectar las anomalías.
 - Regresión Logística
 - Árboles de decisión
 - Random Forest
 - ADA Boost
 - Perceptrón

Indicadores: Métricas de evaluación del rendimiento del modelo en términos de exactitud, precisión y capacidad de detección, según la presente investigación.

- Accuracy (Exactitud): Porcentaje de detección correcta del tráfico malicioso y del benigno con respecto al tráfico total. Indica el rendimiento general del modelo.
- Precision (Precisión): Porcentaje de detección correcta de sólo del tráfico malicioso con respecto a la suma del tráfico malicioso más el tráfico benigno que fue detectado erróneamente como malicioso.
- Recall (Sensibilidad): Porcentaje de detección correcta de sólo del tráfico malicioso con respecto a la suma del tráfico malicioso más el tráfico malicioso que fue detectado erróneamente como benigno.
- F1-Score: Medida armónica de la precisión y el recall para evaluar modelos de clasificación con clases desequilibradas (desigualmente representadas).
- Support: Es el número total de instancias, elementos o datos por cada clase que se utilizan para entrenar, validar o evaluar un modelo de aprendizaje automático.

Escalas de valoración: Para Accuracy, Precisión, Recall y F1-score, valores entre 0 y 1, mientras más cercano al uno, mejor rendimiento del modelo clasificación.

Instrumentos:

- Modelos de aprendizaje automático: Implementar y ejecutar modelos de clasificación como Regresión Logística, Árboles de decisión, Random Forest, ADA Boost y Perceptrón, para la clasificación del tráfico.
- Herramientas de software: Para el análisis, la implementación, entrenamiento y evaluación de los modelos.
 - Anaconda Navigator 2.5.0
 - JupyterLab v3.6.3
 - Python

3.3. Matriz de operacionalización de variables

Cuadro 2: Matriz de operacionalización de variables

VARIABLE	DEFINICIÓN CONCEPTUAL	DEFINICIÓN OPERACIONAL	DIMENSIONES	INDICADORES	ESCALAS DE VALORACIÓN	INSTRUMENTOS
Independiente: Características del dataset	Propiedades o atributos específicos que describen el tráfico de red en entornos IoT, esenciales para clasificar tráfico benigno o malicioso.	Características generadas en condiciones reales durante 33 ataques a 105 dispositivos IoT, incluyendo métricas como contadores y protocolos.	Número de Características. Métricas de tráfico. Protocolos. Flags. Clasificación de ataques.	40 características. Métricas de las características. 7 categorías de ataques.	Categoría: Maliciosa o Benigna Numérica: Cantidad de ocurrencias.	Herramientas de Análisis: Pandas y NumPy. EDA: Visualizaciones y correlaciones. Técnicas de Preprocesamiento: Selección, escalamiento y eliminación de valores erróneos.
Dependiente: Detección de Anomalías	Proceso para identificar patrones inusuales en el tráfico de datos que indican problemas o ataques cibernéticos.	Uso de algoritmos y técnicas de aprendizaje automático para detectar valores atípicos en los datos.	Tráfico: Malicioso o Benigno. Métodos de detección: Regresión Logística, Árboles de decisión, KNN, Random Forest, ADA Boost y Perceptrón	Accuracy Precisión Recall F1-Score	Escalas de 0 a 1	Modelos: Regresión Logística, Árboles de decisión, KNN, Random Forest, ADA Boost y Perceptrón Herramientas: Python y Jupyter.

Fuente: Elaboración propia.

Capítulo IV: Metodología del estudio

4.1. Enfoque, tipo y alcance de investigación

4.1.1. Enfoque.

El enfoque de la presente investigación es cuantitativo, dado a que el objetivo principal es diseñar un modelo basado en el análisis de datos numéricos para clasificar las anomalías. La elección del enfoque cuantitativo se justifica porque se necesitan realizar tareas como, la recolección y análisis de datos numéricos, para ello se utilizarán métricas y estadísticas para evaluar el rendimiento del modelo de clasificación en la detección de anomalías, en ese sentido también se requiere aplicar algoritmos y técnicas cuantitativas para medir la efectividad del modelo.

4.1.2. Tipo y alcance.

Tipo de Investigación Aplicada: Se busca desarrollar un modelo práctico que pueda ser utilizado en entornos reales para detectar anomalías en redes IoT.

Alcance Explicativo: El alcance del presente trabajo de investigación es explicativo, por que busca desarrollar un modelo de clasificación para la detección de anomalías en el dataset CICIoT2023, para conseguirlo se debe entender e identificar las características de las anomalías. A través del análisis exploratorio de datos, se identifican características y patrones que contribuyen a la detección de comportamientos anómalos, proporcionando una base para explicar los resultados del modelo. Además, la evaluación del modelo no solo se enfoca en su precisión, sino también en comprender cómo y por qué funciona, analizando la relevancia de las características involucradas. Este enfoque, además de clasificar datos, también permite presentar y mostrar el variado y dinámico tráfico de entornos IoT.

4.2. Diseño de la investigación

Transversal: El diseño de la presente investigación es transversal porque se centra en la recopilación de datos, de un lapso de tiempo definido, lo que permite obtener una visión o perspectiva temporal (similar a una fotografía) de las características, entre ellas, la longitud de cabecera, el tipo de protocolo y demás en el tráfico IoT. Dado que los datos son específicos y se presentan como un conjunto de métricas en un intervalo tiempo determinado, el diseño transversal facilita la comparación y el análisis de estas características sin necesidad de observar cambios a lo largo del tiempo, lo cual sería el caso para un diseño longitudinal.

4.3. Población y muestra

4.3.1. Población.

La población está conformada por todos los registros del dataset CICIoT2023, específicamente de los 63 archivos MERGED en formato "csv", que en su totalidad contienen 45 019 243 de filas y 40 columnas (las 40 características previamente ya mencionadas, ver Cuadro 1. Este dataset incluye las respectivas etiquetas del tráfico de datos de los 33 tipos de ataque ejecutados y también el tráfico benigno o de normal funcionamiento de los 105 dispositivos IoT incluidos en el dataset de trabajo.

4.3.2. Muestra.

La muestra está conformada por seis archivos "csv", 10% de archivos de la población: Merged47, Merged49, Merged50, Merged51, Merged52 y Merged63, son los archivos más pequeños del dataset y como en el caso de la población, también incluyen el tráfico de datos de los 33 ataques y el tráfico benigno. Se han tomado los seis archivos de menor tamaño para agilizar la ejecución y comprobar la efectividad de los algoritmos que se han de implementar, pues la teoría ya presentada, nos afirma que el poder de los algoritmos de aprendizaje automático, es válido y se puede aplicar a conjuntos de datos más grandes. También se han considerado seis archivos para aplicar en ellos, uno de los estándares de división

de datasets, de 80% de datos para el entrenamiento y 20% de datos para la prueba. La muestra contiene 7 829 178 filas y 40 columnas.

4.4. Técnicas e instrumentos de recolección de datos

4.4.1. Técnicas e instrumentos.

Los datos se obtienen a través de la descarga directa del dataset CICIoT2023, publicado y compartido en el sitio web del Instituto Canadiense de Ciberseguridad.

4.4.2. Validez y confiabilidad.

Los datos han sido generados en el Laboratorio IoT del Instituto Canadiense de Ciberseguridad, equipado con diversas herramientas y software capaces de ejecutar ataques y capturar los datos del tráfico de ataque correspondiente, en dicho laboratorio, Pinto et al. (2023), implementaron una Topología de Red IoT, en la cual conectaron 105 dispositivos IoT, creando un entorno real de una casa inteligente, 7 dispositivos Raspberry Pi fueron los encargados de realizar el ataque a 67 dispositivos de la red. La topología mencionada, así como la comparación del dataset CICIoT2023 generado, con otros datasets similares, se adjuntan en la sección de anexos.

4.4.3. Procedimiento de recolección de datos.

Los datos están a libre disposición para todos los investigadores y personas interesadas en contribuir en el campo de la seguridad de entornos IoT. Los datos se pueden descargar del siguiente enlace:

http://205.174.165.80/IOTDataset/CIC_IOT_Dataset2023/Dataset/CSV/MERGED_CSV/

4.5. Técnicas de análisis de datos

La técnica aplicada es el Análisis Exploratorio de Datos (EDA), que se lleva a cabo a través de un análisis inicial del dataset CICIoT2023 para comprender su estructura, calidad y características. Esto incluirá visualizaciones gráficas y estadísticas descriptivas.

Antes de iniciar el EDA, se describe en que consiste el dataset CICIoT2023, es un conjunto de datos amplio e innovador, diseñado por el Laboratorio IoT del Instituto Canadiense de Ciberseguridad, para abordar una de las preocupaciones más apremiantes en el campo de la tecnología, la seguridad en el IoT. A medida que la conectividad de dispositivos IoT aumenta en diversos sectores, como el transporte, la salud, la educación, la industria y otros, surge la necesidad de actuar sobre los desafíos de seguridad existentes como la interoperabilidad, estándares de comunicación y mecanismos de defensa ante ataques cibernéticos.

Los creadores del dataset, realizaron un estudio exhaustivo al ejecutar 33 tipos de ataques diferentes en una infraestructura compuesta por 105 dispositivos IoT reales, los ataques ejecutados se clasifican en siete categorías principales:

1. DoS (Denial of Service: Denegación de Servicio)

Un ataque DoS se realiza enviando una gran cantidad de solicitudes a un servidor, saturando sus recursos, para hacer que el servicio sea inaccesible para los usuarios legítimos, lo cual provoca la interrupción del servicio, lo que puede resultar en pérdidas económicas y daño a la reputación de la empresa.

2. DDoS (Distributed Denial of Service: Denegación de Servicio Distribuido)

Un ataque DDoS utiliza múltiples dispositivos que envían solicitudes al mismo tiempo a un servidor, para saturar al servidor con tráfico proveniente de diversas fuentes, dificultando su defensa, esto puede causar interrupciones severas en el servicio y afectar a muchos usuarios simultáneamente.

3. Recon (Reconocimiento)

Ataque en el que se recopila información sobre un objetivo mediante escaneos y análisis para identificar vulnerabilidades que puedan ser explotadas en un ataque posterior. Un reconocimiento efectivo aumenta las posibilidades de éxito del ataque al proporcionar información crítica.

4. Web-Based (Basado en la Web)

Los ataques basados en la web se llevan a cabo a través de navegadores o aplicaciones web, aprovechando vulnerabilidades del software, para robar datos o inyectar malware en el sistema del usuario o en el servidor, esto puede generar

la pérdida y exposición de información sensible, causando daños a la reputación de las personas e instituciones atacadas.

5. Brute Force (Fuerza Bruta)

Un ataque de fuerza bruta consiste en probar todas las combinaciones posibles para adivinar contraseñas o claves criptográficas para acceder a cuentas protegidas, al haber descifrado mediante ensayo y error las contraseñas de los usuarios, lo cual puede llevar al acceso no autorizado a cuentas y sistemas críticos.

6. Spoofing (Suplantación)

El spoofing es la falsificación de la identidad del remitente en comunicaciones digitales, como correos electrónicos o direcciones IP, para engañar a los destinatarios con la intención de que crean que están interactuando con alguna persona o algún sistema confiable, con lo cual se pueden ejecutar fraudes y robo de información sensible.

7. Mirai

Mirai es un malware que infecta dispositivos IoT y los convierte en parte de una botnet (red maliciosa) controlada por atacantes, para lanzar ataques de suplantación, fuerza bruta y DDoS masivos utilizando los dispositivos infectados como fuentes de tráfico malicioso, esto puede causar interrupciones significativas en Internet al afectar muchos usuarios y servicios importantes

La creación de este dataset, ofrece un recurso valioso para guardar y analizar datos sobre ataques IoT, además también tiene aplicaciones en el desarrollo de soluciones de análisis de seguridad. Los datos recolectados se presentan en formatos accesibles, como “pcap” para la captura y “csv” para el análisis más detallado, permitiendo a los investigadores explorarlos de diversas maneras y formular nuevas herramientas de inteligencia en ciberseguridad.

Consideraciones importantes del dataset:

1. Relevancia: La creciente integración de dispositivos IoT en nuestra vida cotidiana y en procesos críticos de diversas industrias destaca la urgencia de abordar las vulnerabilidades en estos sistemas. Además, a medida que el

número de dispositivos IoT sigue creciendo, también lo hace el riesgo de ataques, lo que hace que la investigación en seguridad sea aún más esencial.

2. **Innovación Metodológica:** Las metodologías convencionales previas a menudo sopesan su análisis en entornos simulados o limitados en sus capacidades. El presente dataset, en contraste, examina un entorno de IoT real y presenta ataques provocados por dispositivos IoT maliciosos dentro de una red extensiva. Esta perspectiva no solo amplía el alcance de los estudios previos, sino que también proporciona un contexto más aplicable para futuros desarrollos en técnicas de defensa.
3. **Utilidad del dataset:** La disponibilidad del conjunto de datos en diferentes formatos es una gran ventaja; permite no solo el análisis inmediato sino también la reutilización y modificación de los datos para adaptarlos a diferentes modelos de aprendizaje automático. Esta flexibilidad es vital para los investigadores que buscan fortalecer su comprensión de los patrones de tráfico y ataques en entornos IoT.

El lanzamiento del dataset CICIoT2023 representa un progreso significativo en la investigación sobre la seguridad en IoT. Ayuda a cerrar lagunas en la comprensión de cómo los ataques pueden manifestarse en sistemas conectados y cómo realizar evaluaciones de seguridad efectivas. En el contexto de un mundo cada vez más interconectado, la capacidad de identificar y mitigar amenazas de manera efectiva es fundamental. El presente dataset de trabajo no solo proporciona recursos valiosos para académicos y profesionales de campo, sino que también enfatiza la necesidad de continuación en la investigación legal para abordar cuestiones emergentes en la seguridad de IoT, incluyendo el uso de protocolos futuros y características avanzadas de malware.

Dentro del dataset, cada fila tiene 40 características que ofrecen información sobre el tráfico, dentro de ellas se encuentra la variable objetivo, que es la etiqueta de cada tipo de ataque o también de tráfico benigno. Además, cada registro cuenta con una etiqueta que señala si el tráfico es benigno o corresponde a uno de los 33 tipos de ataque. Esta clasificación es crucial para entrenar modelos que puedan detectar y diferenciar entre tráfico normal y ataques.

A continuación, se inicia el Análisis Exploratorio de Datos en diferentes cuadernos de “JupyterLab”.

Carga de archivos

Para elaborar el presente trabajo de investigación se cargan los seis archivos mencionados, como se muestra en la figura 3.

Figura 3: Archivos de la muestra

Nombre	Tamaño	Tipo
 Merged47	132,658 KB	Archivo de valores separados por comas de Microsoft Excel
 Merged49	129,658 KB	Archivo de valores separados por comas de Microsoft Excel
 Merged50	121,259 KB	Archivo de valores separados por comas de Microsoft Excel
 Merged51	43,260 KB	Archivo de valores separados por comas de Microsoft Excel
 Merged52	13,260 KB	Archivo de valores separados por comas de Microsoft Excel
 Merged63	86,400 KB	Archivo de valores separados por comas de Microsoft Excel

Fuente: Elaboración propia.

Combinación de archivos

Luego de cargar los archivos, se procede a combinarlos y exportarlos en uno solo, en formato “pkl”, puesto que es un formato de archivo más pequeño que el tipo “csv”, aproximadamente la quinta parte del tamaño de un archivo “csv” y que contiene la misma información, pero de menor tamaño. El archivo “pkl” resultante es denominado el archivo “combinado” (ver figura 4), incluye 7 829 178 filas y 40 columnas. Con un tamaño de 788 MB.

Figura 4: Resultado de ejecuciones para el archivo “combinado.pkl”

```
Combinando el DataFrame en la ruta: D:\MERGED
...
¡Se han combinado 6 ARCHIVOS CSV en: 'combinado.pkl'!

# Filas y Columnas del DS COMBINADO, sin limpieza
print("Filas y Columnas del DS, sin limpieza:")
ds_data.shape

Filas y Columnas del DS, sin limpieza:
(7829178, 40)
```

Fuente: Elaboración propia.

Luego el archivo “combinado” se carga en otro cuaderno para liberar la RAM y empezar de nuevo con la mayor cantidad disponible de la misma.

Limpieza de datos

Proceso cuyo objetivo es mejorar la calidad de información del dataset, proceso que identifica, corrige errores en los datos o elimina los irrelevantes o duplicados, con ello se asegura que los datos sean más precisos, completos, consistentes y confiables.

En el archivo combinado se realizó la limpieza de la siguiente manera:

- 1° Se identificaron y eliminaron 96 filas con valores faltantes.
- 2° Se identificaron y eliminaron 6 153 347 filas duplicadas.
- 3° Se identificaron y eliminaron 14 filas con valores infinitos.
- 4° Se identificó y eliminó una fila con valor negativo.

A partir de lo ya indicado, el dataset (conjunto de datos en bruto, sin ningún tipo de limpieza) “DS Original” con 7 829 178 filas, se generó el dataframe (conjunto de datos con algún tratamiento de limpieza) “DF Limpio” resultante, con 1 675 720 filas, hasta este punto sólo se han procesado las filas, ambos conjuntos de datos aún conservan las 40 columnas, las cuales se procesarán más adelante.

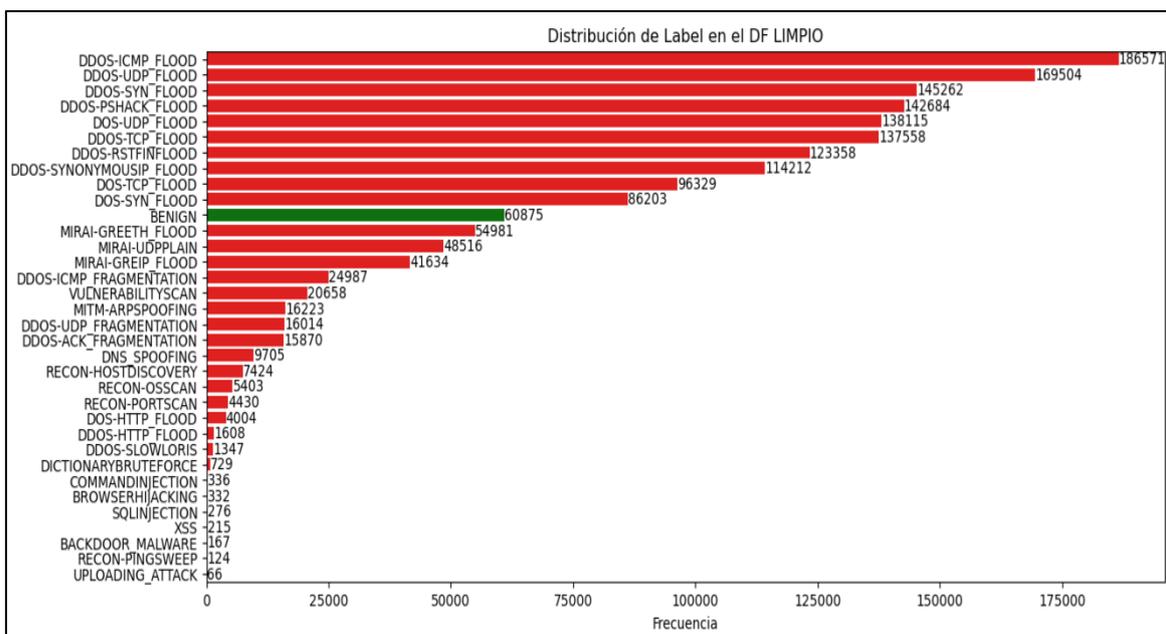
En el cuadro 3 y la figura 5 se pueden apreciar las diferentes etiquetas de “DF Limpio” con su respectivo número de ocurrencias.

Cuadro 3: Cantidad de ataques

Nº	LABEL	CANTIDAD
1	DDOS-ICMP_FLOOD	186571
2	DDOS-UDP_FLOOD	169504
3	DDOS-SYN_FLOOD	145262
4	DDOS-PSHACK_FLOOD	142684
5	DOS-UDP_FLOOD	138115
6	DDOS-TCP_FLOOD	137558
7	DDOS-RSTFINFLOOD	123358
8	DDOS-SYNONYMOUSIP_FLOOD	114212
9	DOS-TCP_FLOOD	96329
10	DOS-SYN_FLOOD	86203
11	BENIGN	60875
12	MIRAI-GREETH_FLOOD	54981
13	MIRAI-UDPPLAIN	48516
14	MIRAI-GREIP_FLOOD	41634
15	DDOS-ICMP_FRAGMENTATION	24987
16	VULNERABILITYSCAN	20658
17	MITM-ARPSPOOFING	16223
18	DDOS-UDP_FRAGMENTATION	16014
19	DDOS-ACK_FRAGMENTATION	15870
20	DNS_SPOOFING	9705
21	RECON-HOSTDISCOVERY	7424
22	RECON-OSSCAN	5403
23	RECON-PORTSCAN	4430
24	DOS-HTTP_FLOOD	4004
25	DDOS-HTTP_FLOOD	1608
26	DDOS-SLOWLORIS	1347
27	DICTIONARYBRUTEFORCE	729
28	COMMANDINJECTION	336
29	BROWSERHIJACKING	332
30	SQLINJECTION	276
31	XSS	215
32	BACKDOOR_MALWARE	167
33	RECON-PINGSWEEP	124
34	UPLOADING_ATTACK	66
TOTAL		1675720

Fuente: Elaboración propia.

Figura 5: Gráfico de barras de las etiquetas



Fuente: Elaboración propia.

Luego de que ya hemos obtenido el “DF Limpio”, se lo exporta nuevamente en formato “pkl”, ya con un menor tamaño de 519 MB (269 MB menos que el “DS Original”).

Sustitución de valores

Antes de empezar el entrenamiento de nuestro modelo, se carga “DF Limpio” ya creado en otro cuaderno, se procede a reemplazar el valor de las etiquetas por una de sus respectivas categorías, para agruparlos en categorías similares (DDoS, DoS, Reconocimiento, Basados en la web, Fuerza bruta, Suplantación, Mirai), proceso que se realiza para reducir la magnitud y complejidad del dataset y mejorar el posterior rendimiento de los modelos de clasificación, ver detalle en los cuadros 4 y 5.

Cuadro 4: Sustitución de valores por categoría de ataque

Nº	ATAQUE ORIGINAL	CATEGORÍA DE ATAQUE
1	DDOS-PSSHACK_FLOOD	DDOS
2	DDOS-ICMP_FLOOD	
3	DDOS-SYN_FLOOD	
4	DDOS-SYNONYMOUSIP_FLOOD	
5	DDOS-UDP_FLOOD	
6	DDOS-RSTFINFLOOD	
7	DDOS-TCP_FLOOD	
8	DDOS-UDP_FRAGMENTATION	
9	DDOS-ICMP_FRAGMENTATION	
10	DDOS-ACK_FRAGMENTATION	
11	DDOS-HTTP_FLOOD	
12	DDOS-SLOWLORIS	
13	DOS-UDP_FLOOD	DOS
14	DOS-SYN_FLOOD	
15	DOS-TCP_FLOOD	
16	DOS-HTTP_FLOOD	
17	MIRAI-GREETH_FLOOD	MIRAI
18	MIRAI-GREIP_FLOOD	
19	MIRAI-UDPPLAIN	
20	DNS_SPOOFING	SPOOFING
21	MITM-ARPSPOOFING	
22	RECON-PINGSWEEP	RECON
23	RECON-OSSCAN	
24	RECON-PORTSCAN	
25	VULNERABILITYSCAN	
26	RECON-HOSTDISCOVERY	
27	BROWSERHIJACKING	
28	BACKDOOR_MALWARE	
29	XSS	
30	UPLOADING_ATTACK	
31	SQLINJECTION	
32	COMMANDINJECTION	
33	DICTIONARYBRUTEFORCE	BRUTE FORCE

BENIGNTRAFFIC	BENIGN
---------------	--------

Fuente: Elaboración propia.

Luego se cambian las categorías a valores numéricos, cero (0) para tráfico benigno y uno (1) para tráfico malicioso:

Cuadro 5: Sustitución de valores por tráfico

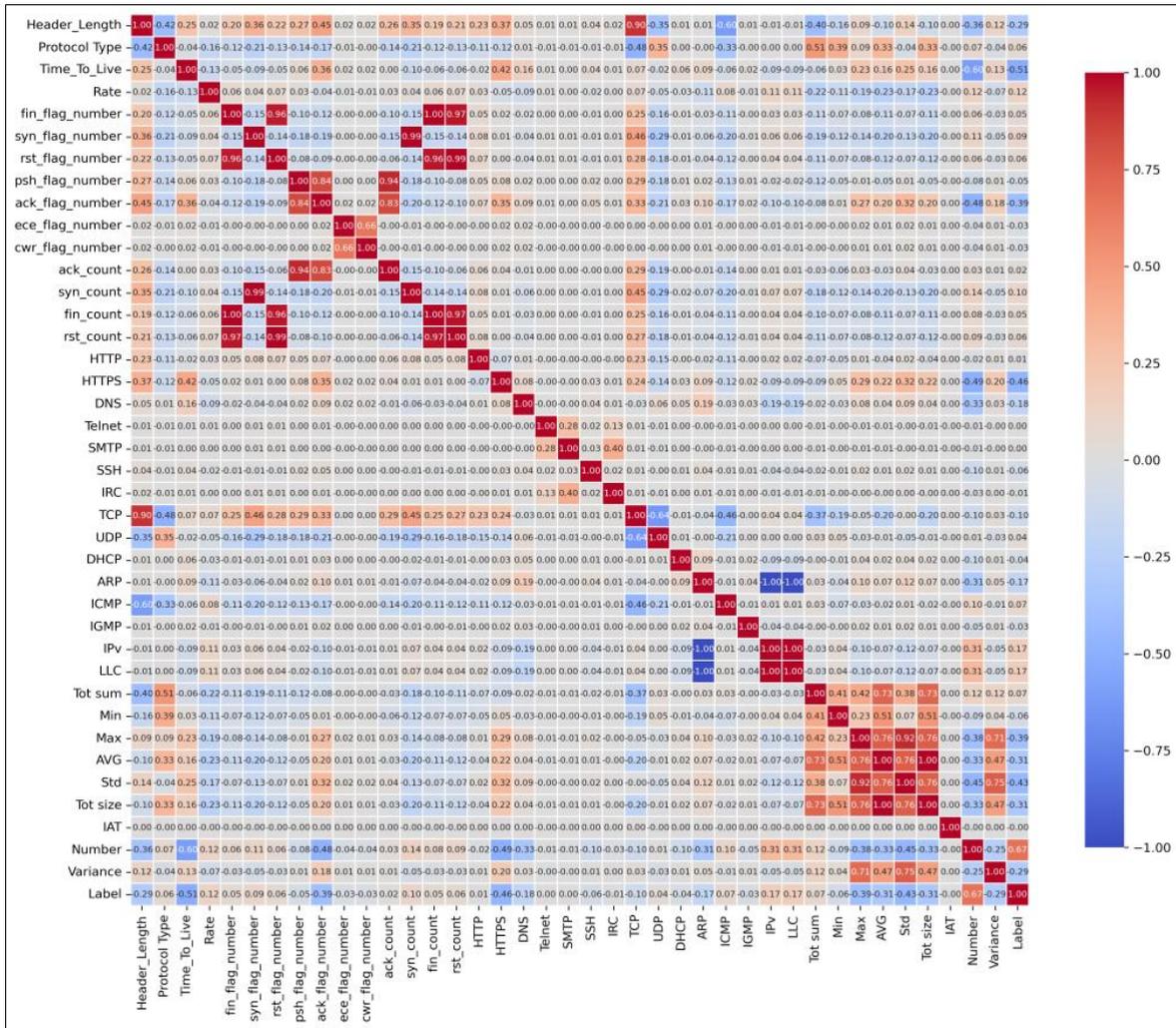
Nº	CATEGORÍA DE ATAQUE	CATEGORÍA BINARIA	CANTIDAD
1	DDOS	1	1078975
2	DOS	1	324651
3	MIRAI	1	145131
4	RECON	1	38039
5	SPOOFING	1	25928
6	WEB	1	1392
7	BRUTE FORCE	1	729
8	BENIGN	0	60875
TOTAL			1675720

Fuente: Elaboración propia.

Selección de características más relevantes

Proceso que tiene como objetivo principal identificar y seleccionar un subconjunto de características (columnas) que son más significativas o relevantes para los modelos de aprendizaje automático, dejando de lado aquellas que son redundantes o irrelevantes, con la actividad previa ya se puede generar la matriz de correlación de variables, ver la figura 6, donde se aprecia que la columna más relacionada con la variable objetivo "Label", es "Number" (Número total de paquetes de flujo) con 67% de correlación.

Figura 6: Matriz de correlación de dataframe con sus 40 características.



Fuente: Elaboración propia.

Ya en la matriz de correlación anterior se aprecia a simple vista, que algunas columnas no tienen relación significativa con “Label”, pero ello no es suficiente para dejarlas de lado, por tal motivo se procede a realizar la selección de las 15 características más relevantes del “DF Limpio”, dicho proceso permitirá posteriormente, optimizar el rendimiento de los modelos de clasificación, mejorar la eficiencia computacional y facilitar la interpretación de los resultados, esta selección de características se realiza con las siguientes herramientas:

1. SelectKBest, selecciona las mejores características utilizando ANOVA (Análisis de Varianza).
2. Correlación de Pearson, calcula la correlación con la variable objetivo y se seleccionan las mejores.

3. Chi-cuadrado, mide la relación con la variable objetivo y se seleccionan las más altas.

SelectKBest y Correlación de Pearson (ver cuadro 6), así como la matriz de correlación previa, identifican como característica más relevante a la columna “Number”, en cambio para Chi-cuadrado es la columna “Variance”, a la vez que también considera a columna “Number” pero no como la más relevante”.

Cuadro 6: Selección de características con SelectKBest, Correlación de Pearson y Chi2

Nº	Característica	SelectKBest	Nº	Característica	Correlación de Pearson	Nº	Característica	Chi2 Score
1	Number	1117199.79	1	Number	0.67	1	Variance	4.90E+11
2	Time_To_Live	477640.59	2	Time_To_Live	0.51	2	Rate	9.06E+08
3	HTTPS	351117.44	3	HTTPS	0.46	3	Max	3.22E+08
4	Std	310941.62	4	Std	0.43	4	Tot sum	2.06E+08
5	Max	236775.24	5	Max	0.39	5	Std	1.92E+08
6	ack_flag_number	233441.65	6	ack_flag_number	0.39	6	AVG	5.95E+07
7	AVG	145450.18	7	Tot size	0.31	7	Tot size	5.95E+07
8	Tot size	145450.18	8	AVG	0.31	8	Number	3.75E+06
9	Header_Length	122853.5	9	Header_Length	0.29	9	Time_To_Live	1.64E+06
10	Variance	122053.83	10	Variance	0.29	10	Min	1.06E+06
11	DNS	46130.5	11	DNS	0.18	11	syn_count	1.01E+06
12	ARP	39558.47	12	LLC	0.17	12	Header_Length	6.23E+05
13	IPv	39558.46	13	IPv	0.17	13	rst_count	4.14E+05
14	LLC	39558.46	14	ARP	0.17	14	fin_count	3.61E+05
15	Rate	21149.08	15	Rate	0.12	15	HTTPS	2.18E+05

Características Seleccionadas	Comunes en las 3 selecciones
	Comunes con diferente ubicación
	No comunes con las otras 2 selecciones

Fuente: Elaboración propia.

Para confirmar la selección anterior se procedió a aplicar Regresión Lasso y Random Forest, generando una selección adicional. Regresión Lasso se ejecutó muy rápido y seleccionó sólo a “Number” como característica única relevante, la ejecución de Random Forest tomó más tiempo, pero también identificó a “Number” como la característica más relevante (ver cuadro 7).

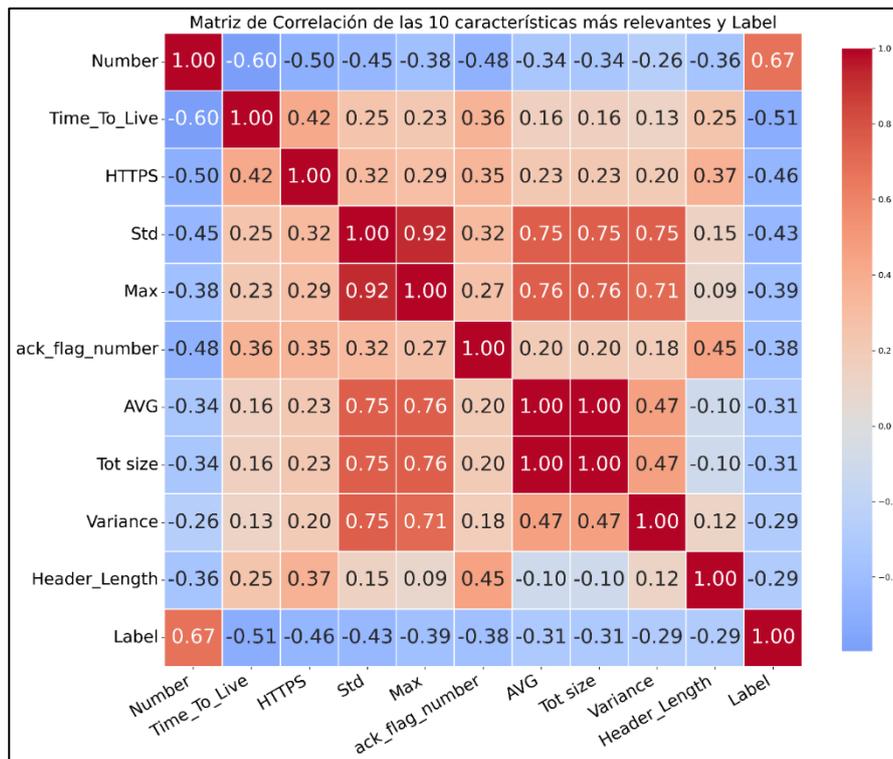
Cuadro 7: Selección de 10 características con Random Forest

Nº	Característica	Importancia
1	Number	0.099
2	Header_Length	0.093
3	Rate	0.084
4	Time_To_Live	0.083
5	IAT	0.081
6	HTTPS	0.076
7	Max	0.068
8	Tot sum	0.060
9	Std	0.045
10	Tot size	0.043

Precisión del modelo: 0.988

Fuente: Elaboración propia.

Figura 7: Matriz de correlación 1 de las diez características seleccionadas.



Fuente: Elaboración propia.

Por los procesos y resultados previos de selección de variables con las herramientas ya vistas, se decide utilizar sólo a la Correlación de Pearson en la selección final de las diez variables para el dataframe de trabajo, por detectar a las mismas variables más relevantes, ser confiable por las comparaciones anteriores y su rápido tiempo de ejecución.

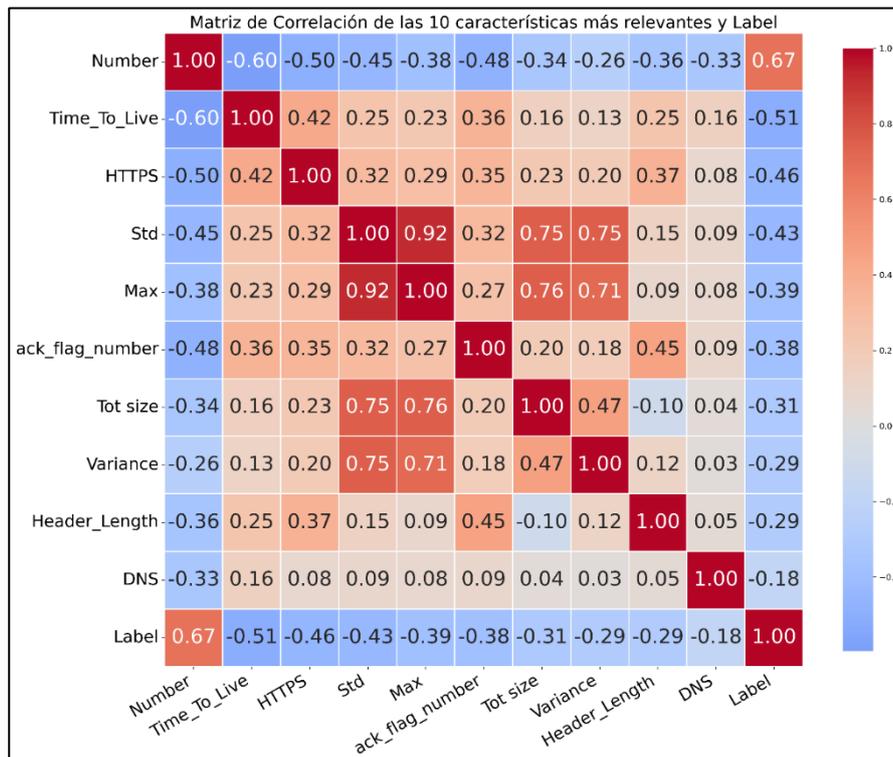
Es importante mencionar que las características “AVG” y “Tot size” tienen la misma correlación (ver cuadro 6 y figura 7), es decir que tiene valores idénticos, son columnas con valores duplicados, donde: “AVG” es la longitud promedio del paquete de flujo de datos y “Tot size” es el tamaño del paquete, por lo cual se decide considerar a “Tot size” y evitar el uso de más de una característica con los mismos datos (redundantes), en consecuencia se considera a la siguiente característica con correlación más alta, la característica “DNS” (ver cuadro 6).

Cuadro 8: Diez características seleccionadas con la Correlación de Pearson

Nº	Característica	Correlación de Pearson
1	Number	0.67
2	Time_To_Live	0.51
3	HTTPS	0.46
4	Std	0.43
5	Max	0.39
6	ack_flag_number	0.39
7	Tot size	0.31
8	Header_Length	0.29
9	Variance	0.29
10	DNS	0.18

Fuente: Elaboración propia.

Figura 8: Matriz de correlación 2 de las diez características seleccionadas



Fuente: Elaboración propia.

La matriz de correlación anterior muestra valores negativos en nueve de las diez características seleccionadas, lo que indica que existe en ellas una correlación negativa, con respecto a la variable objetivo "Label", lo que significa que una de las características aumenta y otra disminuye, por lo cual para la presente investigación, por los modelos de clasificación que se implementarán y su eficiente modo de trabajo, es más importante la fuerza de correlación (cercanía a -1 y +1) y no su dirección (signo positivo o negativo).

En los cuadros 9 y 10 se presenta información de las características seleccionadas.

Cuadro 9: Características con valores originales (vista parcial)

Nº	Number	Time_To_Live	HTTPS	Std	Max	ack_flag_number	Tot size	Variance	Header_Length	DNS
1	100.00	64.64	0.00	5.50	115.00	0.00	60.55	30.25	8.00	0.01
2	100.00	65.91	0.00	0.00	60.00	1.00	60.00	0.00	20.00	0.00
3	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	0.00	0.00
4	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	8.00	0.00
5	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	8.00	0.00
6	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	0.00	0.00
7	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	0.00	0.00
8	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	8.00	0.00
9	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	20.00	0.00
10	100.00	63.36	0.00	0.00	60.00	0.99	60.00	0.00	19.80	0.00
11	100.00	65.70	0.01	76.86	554.00	0.01	540.68	5907.65	8.24	0.01
12	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	20.00	0.00
13	10.00	96.60	1.00	949.45	2882.00	1.00	879.60	901450.49	32.00	0.00
14	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	0.00	0.00
15	100.00	65.70	0.00	1.40	74.00	0.01	60.14	1.96	20.20	0.00
16	10.00	168.20	0.40	62.54	218.00	0.50	113.50	3910.72	20.00	0.00
17	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	0.00	0.00
18	100.00	65.69	0.01	0.60	66.00	0.01	60.06	0.36	0.32	0.00
19	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	20.00	0.00
20	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	0.00	0.00
...

Fuente: Elaboración propia.

Cuadro 10: Estadísticas de las características seleccionadas con valores originales

	Number	Time_To_Live	HTTPS	Std	Max	ack_flag_number	Tot size	Variance	Header_Length	DNS
count	1675720	1675720	1675720	1675720	1675720	1675720	1675720	1675720	1675720	1675720
mean	93.09	67.94	0.09	64.34	314.15	0.15	170.65	53082.74	14.05	0.00
std	23.92	17.82	0.26	221.23	710.69	0.33	278.32	479929.94	8.82	0.03
min	2.00	0.00	0.00	0.00	46.00	0.00	46.00	0.00	0.00	0.00
25%	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	8.00	0.00
50%	100.00	64.00	0.00	0.00	60.00	0.00	60.00	0.00	19.88	0.00
75%	100.00	64.54	0.00	8.87	214.00	0.01	64.45	78.63	20.00	0.00
max	100.00	255.00	1.00	9730.18	33370.00	1.00	9430.30	94676391.51	60.00	1.00

Fuente: Elaboración propia.

Escalamiento de características

Es un proceso de ajuste o transformación de todos los valores de cada característica en una misma escala, la escala de cada variable depende de la amplitud del rango de sus respectivos valores y la cual debe tener como media cero o un valor cercano y desviación estándar uno, para que los datos estén centrados, uniformizados y sean más fácilmente comparables, este proceso se realiza con el objetivo de mejorar la eficiencia y efectividad de los modelos de clasificación que se ejecutarán, para realizar el escalamiento se utiliza el método “StandardScaler”, el cuadro 11 muestra parte del resultado del escalamiento y el cuadro 12 sus estadísticas descriptivas.

Cuadro 11: Características escaladas (Vista parcial)

Nº	Number	Time_To_Live	HTTPS	Std	Max	ack_flag_number	Tot size	Variance	Header_Length	DNS
1	0.289	-0.185	-0.330	-0.266	-0.280	-0.458	-0.396	-0.111	-0.686	0.215
2	0.289	-0.114	-0.330	-0.291	-0.358	2.578	-0.398	-0.111	0.674	-0.149
3	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	-1.594	-0.149
4	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	-0.686	-0.149
5	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	-0.686	-0.149
6	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	-1.594	-0.149
7	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	-1.594	-0.149
8	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	-0.686	-0.149
9	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	0.674	-0.149
10	0.289	-0.257	-0.330	-0.291	-0.358	2.548	-0.398	-0.111	0.652	-0.149
11	0.289	-0.126	-0.291	0.057	0.337	-0.428	1.330	-0.098	-0.659	0.215
12	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	0.674	-0.149
13	-3.474	1.608	3.536	4.001	3.613	2.578	2.547	1.768	2.035	-0.149
14	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	-1.594	-0.149
15	0.289	-0.126	-0.330	-0.284	-0.338	-0.428	-0.397	-0.111	0.697	-0.149
16	-3.474	5.625	1.217	-0.008	-0.135	1.060	-0.205	-0.102	0.674	-0.149
17	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	-1.594	-0.149
18	0.289	-0.126	-0.291	-0.288	-0.349	-0.428	-0.397	-0.111	-1.557	-0.149
19	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	0.674	-0.149
20	0.289	-0.221	-0.330	-0.291	-0.358	-0.458	-0.398	-0.111	-1.594	-0.149
...

Fuente: Elaboración propia.

Cuadro 12: Estadísticas de las características seleccionadas escaladas

	Number	Time To Live	HTTPS	Std	Max	ack_flag number	Tot size	Variance	Header Length	DNS
count	1675720	1675720	1675720	1675720	1675720	1675720	1675720	1675720	1675720	1675720
mean	0	0	0	0	-0	-0	0	-0	0	-0
std	1	1	1	1	1	1	1	1	1	1
min	-3.81	-3.81	-0.33	-0.29	-0.38	-0.46	-0.45	-0.11	-1.59	-0.15
25%	0.29	-0.22	-0.33	-0.29	-0.36	-0.46	-0.4	-0.11	-0.69	-0.15
50%	0.29	-0.22	-0.33	-0.29	-0.36	-0.46	-0.4	-0.11	0.66	-0.15
75%	0.29	-0.19	-0.33	-0.25	-0.14	-0.43	-0.38	-0.11	0.67	-0.15
max	0.29	10.49	3.54	43.69	46.51	2.58	33.27	197.16	5.21	36.21

Fuente: Elaboración propia.

En el cuadro 12 se aprecian las características seleccionadas ya escaladas, donde la media de cada una, es uno y la desviación estándar es cero o al menos un valor cercano al mismo, en las características “Max”, “ack_flag_number”, “Variance” y “DNS” el valor de las desviaciones es menos cero (-0), lo cual indica que los datos están cercanos al cero, pero con un predominio de valores negativos.

Al dataframe con las características seleccionadas y escaladas, se le agrega la columna “Label”, quedando ya preparado (“DF Preparado”) para ser utilizado con los modelos de aprendizaje automático, para la presente investigación, modelos de clasificación.

Como resultado de los procesos ya mencionados se presenta el cuadro 13 con la comparación del “DS Original” y el “DF Preparado”

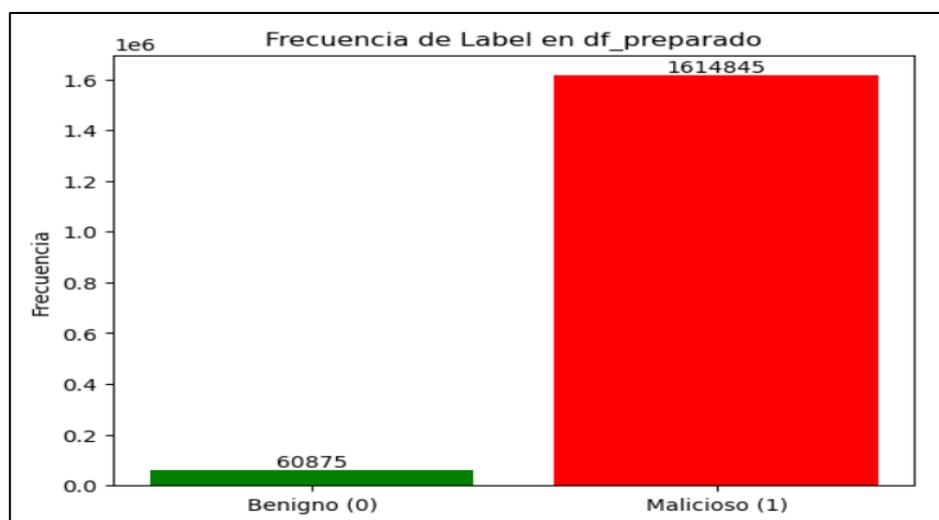
Cuadro 13: DS Original y DF Preparado

CONJUNTO DE DATOS	Nº DE FILAS	Nº DE COLUMNAS	TAMAÑO (MB)
DS Original	7 829 178	40	788
DF Preparado	1 675 720	11	140
Diferencia (Eliminadas por limpieza)	6 153 458	29	648

Fuente: Elaboración propia.

El “DF Preparado” conserva 1 675 720 filas y once características, diez seleccionadas, las más relevantes y una más, “Label”, que indica si el tráfico de la fila es benigno o malicioso. El “DF Preparado” contiene 60 875 filas de tráfico benigno y de tráfico malicioso 1 614 845 filas, ver la figura 9.

Figura 9: Gráfico de barras de Label en DF Preparado



Fuente: Elaboración propia.

Capítulo V: Resultados

5.1. Análisis de resultados

5.1.1. Análisis exploratorio de datos.

En el análisis se abordan las diez características más relevantes del dataset seleccionadas en el capítulo anterior, para ello son necesarias las siguientes aclaraciones:

Paquete de datos

Es la unidad básica de información que se transmite entre dispositivos o desde dispositivos a aplicaciones en la nube. Estos paquetes pueden contener datos de sensores, actuadores, instrucciones para dispositivos, o cualquier otra información necesaria para comunicaciones en entornos IoT.

Flujo de datos

Es el movimiento continuo de estos paquetes de datos entre dispositivos, aplicaciones y sistemas. Este flujo es crucial para el funcionamiento eficiente de aplicaciones IoT.

Flag

En el contexto de protocolos de comunicación como TCP (Protocolo de Control de Transmisión), un flag es un bit en la cabecera del paquete de datos que indica la acción específica que se debe ejecutar y/o su estado. Los flags son esenciales para el manejo de la comunicación.

Flag ACK (Acknowledgment: Reconocimiento)

Es un tipo específico de flag que se utiliza para confirmar la recepción de paquetes de datos, cuando un dispositivo recibe un paquete, responde con un paquete que tiene el flag ACK activado, indicando que los datos han sido recibidos correctamente, esto es parte del flujo de datos ya que garantiza que la comunicación sea confiable.

Según las aclaraciones mencionadas, los paquetes de datos son transmitidos como parte del flujo de datos, cada paquete de datos puede contener flags en su cabecera para indicar acciones específicas a realizar. El flag ACK es un tipo de flag

que se utiliza para confirmar la recepción de paquetes, asegurando que el flujo de datos sea confiable. Por lo tanto, los paquetes de datos se transmiten en un flujo continuo y dentro de estos paquetes, los flags juegan un papel crucial en el manejo y confirmación de la comunicación.

A continuación, se describen las características seleccionadas.

1. Number.

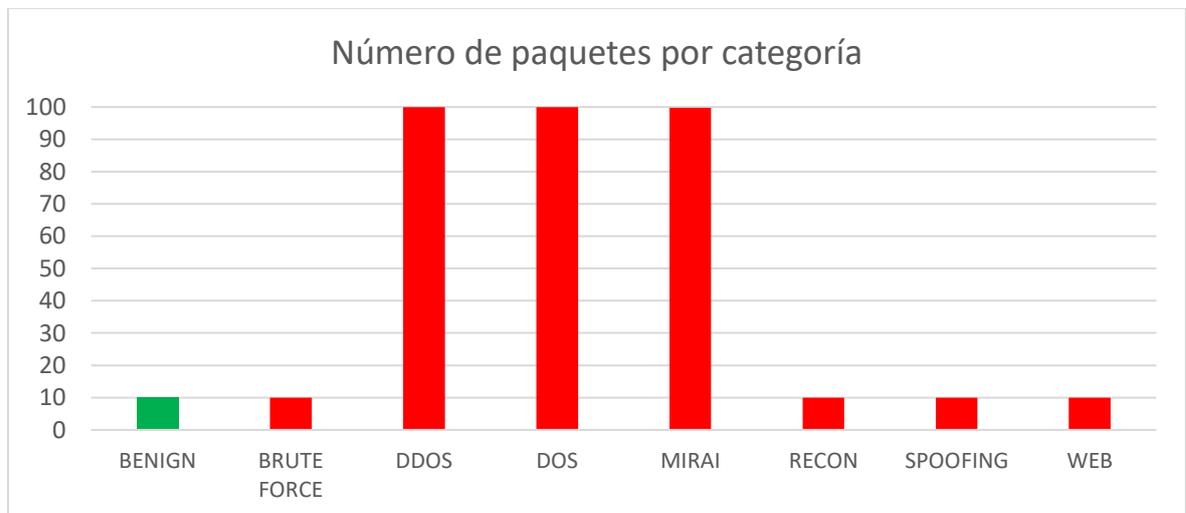
Representa el número total de paquetes de cada flujo de datos capturado, en redes IoT los dispositivos generan múltiples flujos de datos simultáneamente y esta métrica puede ayudar a identificar y rastrear cada flujo individual. Permite segmentar y organizar grandes volúmenes de tráfico generado por dispositivos IoT, esenciales para análisis y auditorías. Ver su información en la tabla 1 y la figura 10.

Tabla 1: Estadísticas descriptivas de la característica Number

LABEL	COUNT	MEAN	STD	MIN	25%	50%	75%	MAX
BENIGN	60875	10.00	0.12	2.00	10.00	10.00	10.00	10.00
BRUTE FORCE	729	10.00	0.00	10.00	10.00	10.00	10.00	10.00
DDOS	1078975	99.91	2.37	2.00	100.00	100.00	100.00	100.00
DOS	324651	99.94	2.00	2.00	100.00	100.00	100.00	100.00
MIRAI	145131	99.71	4.40	2.00	100.00	100.00	100.00	100.00
RECON	38039	10.00	0.06	3.00	10.00	10.00	10.00	10.00
SPOOFING	25928	10.00	0.13	2.00	10.00	10.00	10.00	10.00
WEB	1392	9.99	0.21	2.00	10.00	10.00	10.00	10.00

Fuente: Elaboración propia.

Figura 10: Gráfico de barras de la característica Number



Fuente: Elaboración propia.

Patrón Observado: La mayor parte del tráfico benigno tiene un número constante de paquetes (10), mientras que los ataques DDoS, DoS y Mirai presentan un número significativamente mayor de paquetes, con promedios muy cercanos a 100. Además, por el número de paquetes los ataques restantes podrían ser tomados como tráfico benigno.

Interpretación: Esto indica que los ataques DDoS, DoS y Mirai generan un flujo masivo de paquetes, posiblemente para saturar la red o el servicio objetivo.

2. Time_To_Live (TTL).

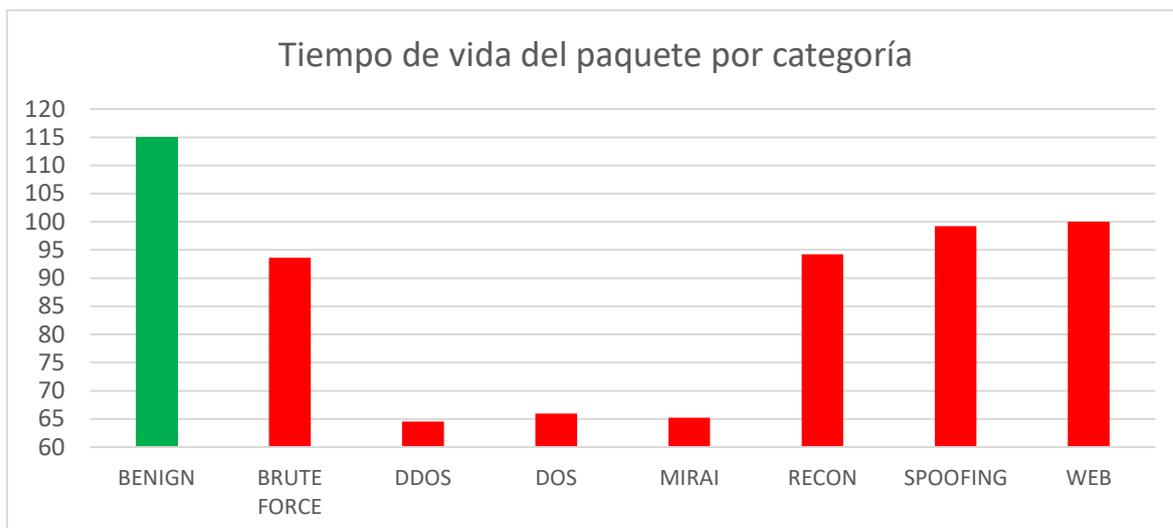
Tiempo de vida del paquete en el flujo de datos, normalmente en segundos, según la configuración de los dispositivos, representa cuánto tiempo puede permanecer activo un paquete durante la comunicación requerida, antes de ser descartado, en IoT los paquetes pueden tener TTLs bajos si la comunicación es local (por ejemplo, dentro de una red doméstica) o altos si los datos se envían a la nube. TTL anómalo podría indicar problemas de configuración, retrasos en la red o diferentes intentos de ataque, esta característica ayuda a entender el alcance de la comunicación en entornos IoT (local o global). Ver su información en la tabla 2 y la figura 11.

Tabla 2: Estadísticas descriptivas de la característica Time_To_Live

LABEL	COUNT	MEAN	STD	MIN	25%	50%	75%	MAX
BENIGN	60875	115.00	51.22	0.00	74.00	98.90	147.60	248.00
BRUTE FORCE	729	93.64	34.94	36.30	64.00	83.10	112.00	234.90
DDOS	1078975	64.52	3.49	4.77	64.00	64.00	64.00	251.18
DOS	324651	65.93	8.90	36.91	64.00	64.00	64.64	255.00
MIRAI	145131	65.24	5.82	2.08	64.00	64.00	64.64	252.45
RECON	38039	94.23	38.06	0.00	64.10	85.70	113.40	255.00
SPOOFING	25928	99.19	50.99	0.00	63.10	81.70	118.90	248.70
WEB	1392	100.02	39.72	31.60	67.70	93.30	120.98	247.00

Fuente: Elaboración propia.

Figura 11: Gráfico de barras de la característica Time_To_Live



Fuente: Elaboración propia.

Patrón Observado: Los ataques DDoS, DoS y Mirai tienen un TTL promedio bajo, alrededor de 65, en comparación con el tráfico benigno que tiene un TTL de 115.

Interpretación: Un TTL más bajo puede indicar que los paquetes son generados por máquinas automatizadas en lugar de usuarios finales, lo que es común en ataques digitales.

3. HTTPS.

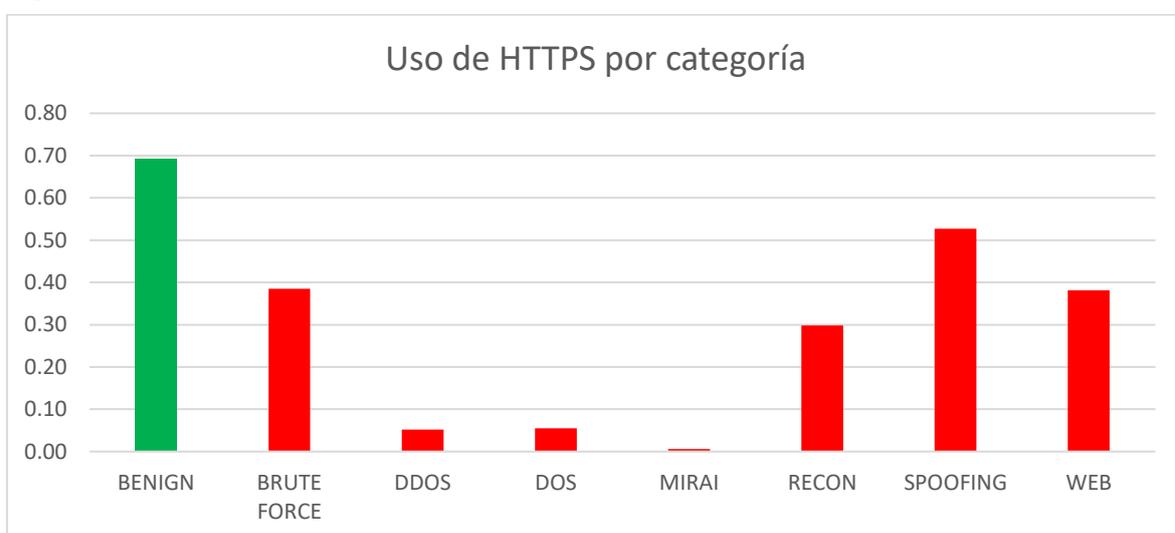
Indicador o medida del tráfico que utiliza el protocolo de navegación segura HTTPS (Protocolo de Transferencia de Hipertexto Seguro), que cifra los datos para protegerlos para evitar que terceros los intercepten y los lean durante su transmisión, en IoT es esencial porque muchos dispositivos transmiten información sensible, como datos personales, imágenes u otros; un bajo uso de HTTPS puede indicar configuraciones inseguras en dispositivos IoT, aumentando el riesgo de ciberataques. Monitorear esta métrica ayuda a garantizar la seguridad de las comunicaciones. Ver su información en la tabla 3 y la figura 12.

Tabla 3: Estadísticas descriptivas de la característica HTTPS

LABEL	COUNT	MEAN	STD	MIN	25%	50%	75%	MAX
BENIGN	60875	0.69	0.31	0.00	0.50	0.80	1.00	1.00
BRUTE FORCE	729	0.38	0.35	0.00	0.10	0.30	0.70	1.00
DDOS	1078975	0.05	0.22	0.00	0.00	0.00	0.00	1.00
DOS	324651	0.05	0.22	0.00	0.00	0.00	0.00	1.00
MIRAI	145131	0.01	0.04	0.00	0.00	0.00	0.00	1.00
RECON	38039	0.30	0.31	0.00	0.00	0.20	0.50	1.00
SPOOFING	25928	0.53	0.42	0.00	0.10	0.60	1.00	1.00
WEB	1392	0.38	0.35	0.00	0.10	0.30	0.70	1.00

Fuente: Elaboración propia.

Figura 12: Gráfico de barras de la característica HTTPS



Fuente: Elaboración propia.

Patrón Observado: La mayor parte del tráfico benigno, el 69% utiliza HTTPS, mientras que el tráfico malicioso DDoS y DoS utilizan el 5% y en Mirai su uso es casi nulo, el 1%.

Interpretación: Los atacantes pueden optar por no usar HTTPS para evitar la sobrecarga de cifrado o porque sus objetivos no requieren conexiones seguras.

4. Std (Desviación estándar de la longitud del paquete en el flujo de datos).

Mide la variabilidad del flujo de datos en el tráfico IoT, como el tamaño de paquetes o el tiempo entre ellos, en IoT la variabilidad puede reflejar patrones normales o anomalías en dispositivos que funcionan mal o están siendo atacados. Una alta variabilidad puede sugerir problemas de red o ataques de DDoS, una baja

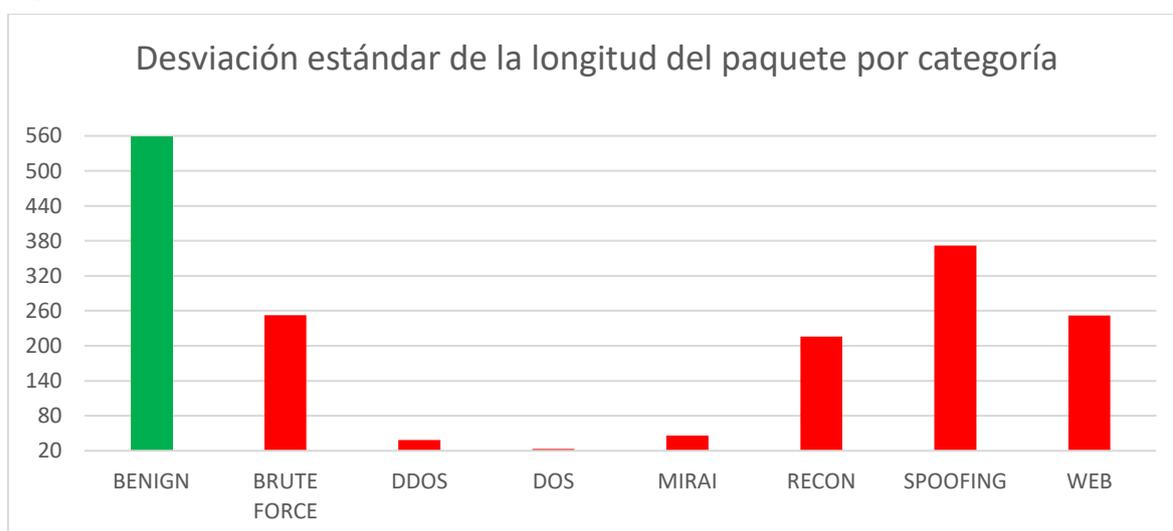
desviación estándar suele indicar un tráfico más predecible, típico en dispositivos configurados correctamente. Ver su información en la tabla 4 y la figura 13.

Tabla 4: Estadísticas descriptivas de la característica Std

LABEL	COUNT	MEAN	STD	MIN	25%	50%	75%	MAX
BENIGN	60875	559.50	677.23	0.00	56.32	391.59	896.85	9730.18
BRUTE FORCE	729	252.52	369.14	0.00	41.10	68.65	305.01	2137.10
DDOS	1078975	38.25	140.95	0.00	0.00	0.00	1.41	3428.32
DOS	324651	22.87	101.60	0.00	0.00	0.00	4.00	2200.64
MIRAI	145131	46.01	75.31	0.00	0.00	44.10	69.51	3744.64
RECON	38039	215.73	359.33	0.00	30.00	73.62	150.12	4666.37
SPOOFING	25928	371.97	371.31	0.00	53.25	329.18	609.69	5877.45
WEB	1392	252.24	408.38	0.00	52.83	81.19	228.57	4815.37

Fuente: Elaboración propia.

Figura 13: Gráfico de barras de la característica Std



Fuente: Elaboración propia.

Patrón Observado: La desviación estándar es considerablemente más alta, 559.50 en el tráfico benigno, en comparación con los ataques DDoS, DoS y Mirai es 46.01 o menos.

Interpretación: Esto indica que el tráfico benigno tiene una variedad más amplia en la longitud de los paquetes, mientras que los ataques tienden a generar paquetes de longitud similar.

5. Max (Longitud máxima del paquete en el flujo de datos).

Representa el valor máximo observado del tamaño de paquete dentro de un flujo de datos, en IoT podría indicar picos de tráfico debido a eventos específicos, como

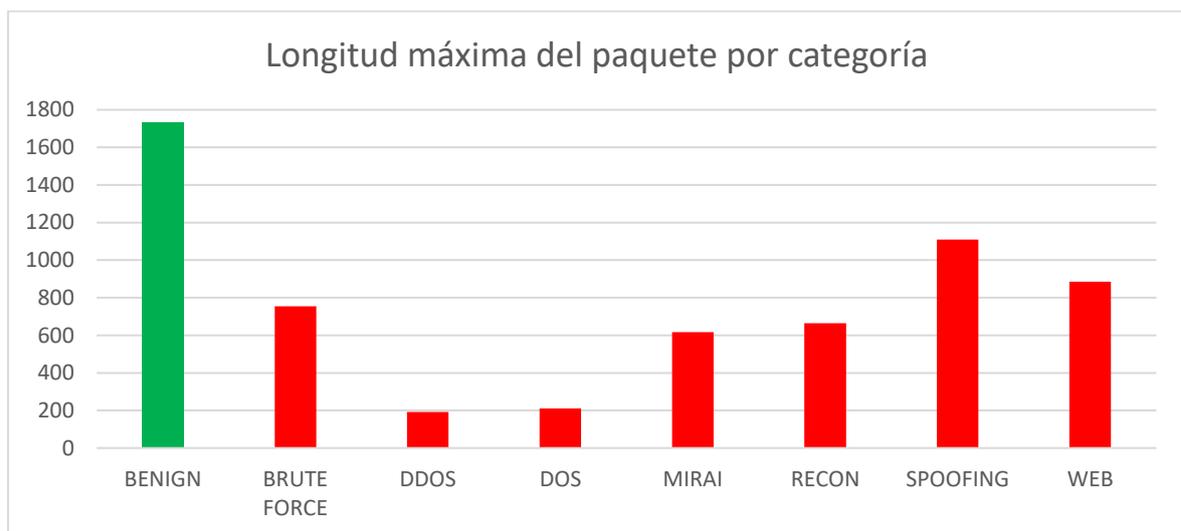
actualizaciones de firmware o intentos de saturación u otras anomalías en la red, ayuda a identificar comportamientos extremos, como dispositivos que envían paquetes inusualmente grandes. Ver su información en la tabla 5 y la figura 14.

Tabla 5: Estadísticas descriptivas de la característica Max

LABEL	COUNT	MEAN	STD	MIN	25%	50%	75%	MAX
BENIGN	60875	1730.70	2116.12	60.00	230.00	1469.00	2962.00	31922.00
BRUTE FORCE	729	753.70	1027.09	60.00	176.00	230.00	1042.00	5858.00
DDOS	1078975	192.23	465.82	60.00	60.00	60.00	74.00	33370.00
DOS	324651	210.38	491.75	60.00	60.00	66.00	162.00	14546.00
MIRAI	145131	617.05	345.28	78.00	554.00	578.00	592.00	33370.00
RECON	38039	663.48	1012.66	60.00	156.00	262.00	519.00	15994.00
SPOOFING	25928	1110.09	988.08	46.00	230.00	1208.00	1500.00	18890.00
WEB	1392	884.36	1227.75	60.00	230.00	345.00	1494.00	15774.00

Fuente: Elaboración propia.

Figura 14: Gráfico de barras de la característica Max



Fuente: Elaboración propia.

Patrón Observado: La longitud máxima del paquete es notablemente alta en el tráfico benigno (1730.70) frente a los ataques DDoS (192.23).

Interpretación: Esto indica que el tráfico benigno incluye datos más complejos o voluminosos, como los archivos multimedia, mientras que los ataques DDoS y DoS se centran en paquetes más pequeños, pero numerosos como se puede ver también en la figura 10, además es notorio en Mirai, que a pesar de ser un ataque más complejo, en magnitud es aproximadamente sólo un tercio del tráfico benigno.

6. ack_flag_number (Valor del flag Ack).

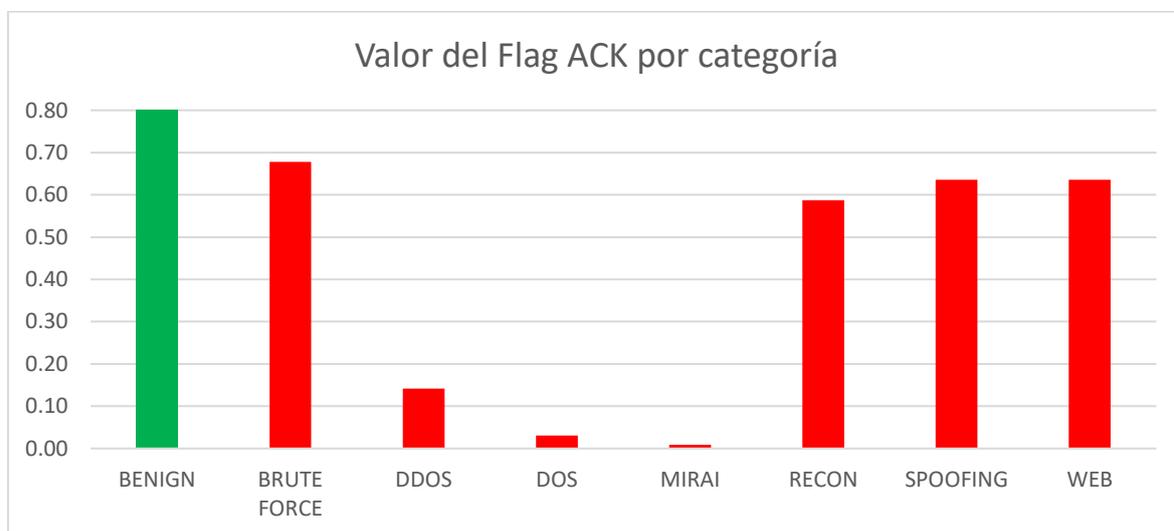
Es un flag que confirma la recepción de paquetes, en IoT es esencial para mantener conexiones confiables entre dispositivos, especialmente en sistemas críticos como dispositivos médicos, un número bajo puede reflejar problemas en las conexiones o intentos de ataque de tipo reset (RST: interrumpe comunicaciones de red), monitorear esta métrica es clave para evaluar la calidad, estabilidad y seguridad de las comunicaciones. Ver su información en la tabla 6 y la figura 15.

Tabla 6: Estadísticas descriptivas de la característica ack_flag_number

LABEL	COUNT	MEAN	STD	MIN	25%	50%	75%	MAX
BENIGN	60875	0.80	0.23	0.00	0.70	0.90	1.00	1.00
BRUTE FORCE	729	0.68	0.28	0.00	0.50	0.70	0.90	1.00
DDOS	1078975	0.14	0.33	0.00	0.00	0.00	0.01	1.00
DOS	324651	0.03	0.12	0.00	0.00	0.00	0.01	1.00
MIRAI	145131	0.01	0.06	0.00	0.00	0.00	0.00	1.00
RECON	38039	0.59	0.29	0.00	0.40	0.60	0.80	1.00
SPOOFING	25928	0.64	0.40	0.00	0.20	0.80	1.00	1.00
WEB	1392	0.64	0.30	0.00	0.40	0.70	0.90	1.00

Fuente: Elaboración propia.

Figura 15: Gráfico de barras de la característica ack_flag_number



Fuente: Elaboración propia.

Patrón Observado: El valor medio del flag ACK es alto en el tráfico benigno (0.80), pero considerablemente bajo en Mirai (0.01).

Interpretación: Un menor uso del flag ACK en ataques puede indicar que estos no buscan establecer conexiones confiables, sino simplemente inundar la red.

7. Tot size (Tamaño total del paquete).

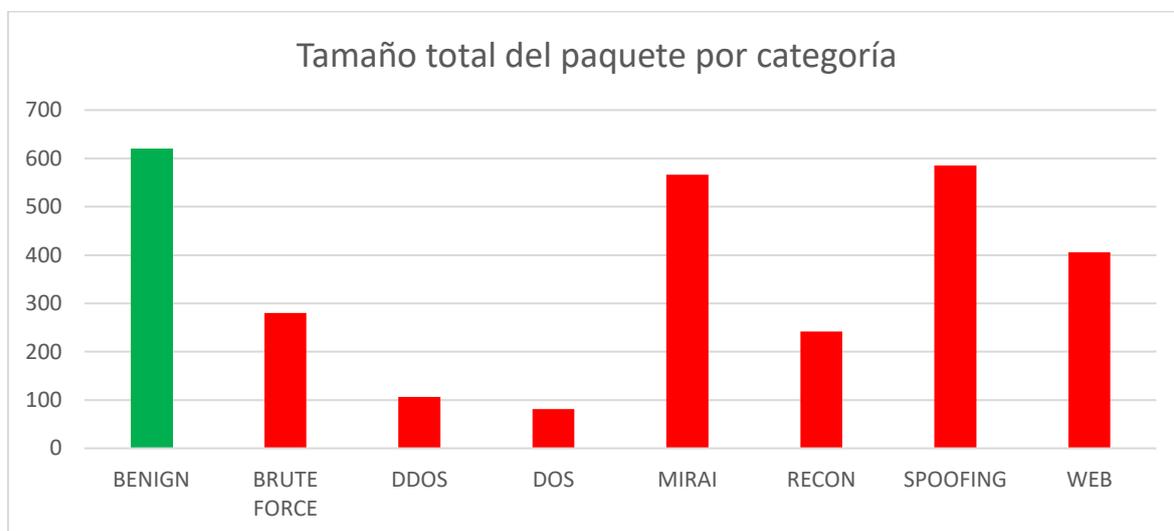
Es el tamaño total de los datos enviados o recibidos en un flujo (tráfico total), en IoT, dispositivos como cámaras IP generan grandes volúmenes de datos, mientras que sensores más simples generan volúmenes pequeños, esta característica permite identificar dispositivos que consumen más ancho de banda de lo esperado, lo que podría indicar mal uso o ataques. Ver su información en la tabla 7 y la figura 16.

Tabla 7: Estadísticas descriptivas de la característica Tot size

LABEL	COUNT	MEAN	STD	MIN	25%	50%	75%	MAX
BENIGN	60875	619.73	674.42	60.00	100.80	246.60	1057.20	9430.30
BRUTE FORCE	729	280.39	400.91	60.00	91.60	117.60	207.60	3120.00
DDOS	1078975	106.18	188.16	58.73	60.00	60.00	60.18	2948.72
DOS	324651	81.11	80.32	60.00	60.00	60.06	61.70	2935.58
MIRAI	145131	566.38	39.35	60.23	554.00	572.82	586.68	1879.61
RECON	38039	241.83	383.89	60.00	77.30	110.10	173.20	4410.00
SPOOFING	25928	585.33	576.85	46.00	115.80	361.00	926.20	4989.20
WEB	1392	405.57	545.99	60.00	105.30	141.25	414.23	3558.00

Fuente: Elaboración propia.

Figura 16: Gráfico de barras de la característica Tot size



Fuente: Elaboración propia.

Patrón Observado: El tamaño total promedio del paquete de datos, es mayor en el tráfico benigno (619.73) comparado con DoS (81.11).

Interpretación: Esto indica que el tráfico benigno tiene paquetes más grandes, posiblemente por la naturaleza de las aplicaciones utilizadas.

8. Variance (Varianza del tamaño de los paquetes).

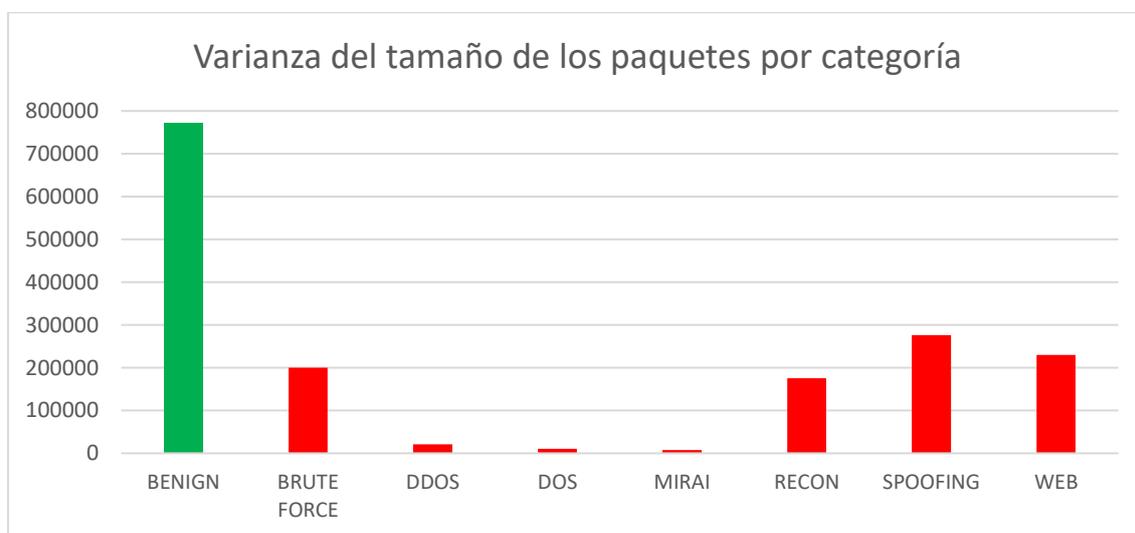
Es dispersión del tamaño de los paquetes del flujo entrante y saliente de datos, como la fluctuación en los tiempos de envío y llegada de paquetes, en IoT, una alta varianza puede ser señal de eventos inusuales, como dispositivos que enfrentan congestión de red, es una métrica clave para detectar eventos anómalos en el tráfico de red. Ver su información en la tabla 8 y la figura 17.

Tabla 8: Estadísticas descriptivas de la característica Variance

LABEL	COUNT	MEAN	STD	MIN	25%	50%	75%	MAX
BENIGN	60875	771683	2254840	0	3172	153341	804341	94676392
BRUTE FORCE	729	199839	535887	0	1689	4712	93028	4567190
DDOS	1078975	21330	98181	0	0	0	2	11753355
DOS	324651	10845	89576	0	0	0	16	4842827
MIRAI	145131	7788	77788	0	0	1945	4831	14022313
RECON	38039	175655	585876	0	900	5420	22536	21775054
SPOOFING	25928	276226	688753	0	2836	108357	371718	34544455
WEB	1392	230284	1105656	0	2791	6591	52243	23187747

Fuente: Elaboración propia.

Figura 17: Gráfico de barras de la característica Variance



Fuente: Elaboración propia.

Patrón Observado: La varianza es significativamente alta en el tráfico benigno (771 683) frente a los ataques DDoS, DoS y MIRAI (21 329 o menos).

Interpretación: Esto refuerza la idea de que el tráfico benigno es más diverso en términos de tamaño de paquete, mientras que los ataques son más homogéneos.

9. Header_Length (Longitud de cabecera del paquete de datos).

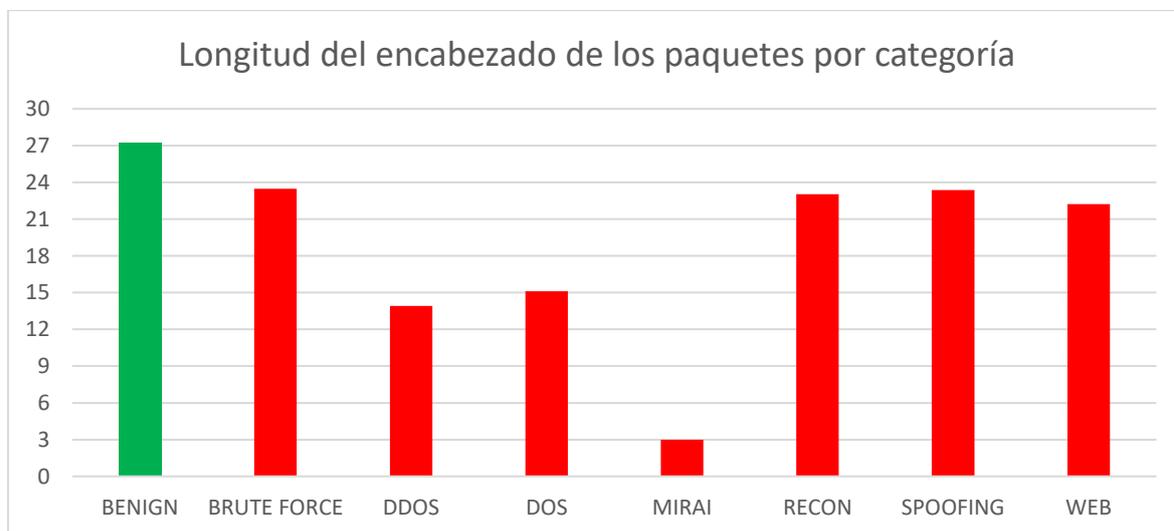
Es la longitud del encabezado de los paquetes, que contiene información sobre cómo se debe manejar el paquete, en IoT esto puede variar dependiendo del tipo de protocolo y configuración del dispositivo. Longitudes inusuales podrían sugerir configuraciones erróneas o tráfico malicioso (como encabezados personalizados en ataques). Ver su información en la tabla 9 y la figura 18.

Tabla 9: Estadísticas descriptivas de la característica Header_Length

LABEL	COUNT	MEAN	STD	MIN	25%	50%	75%	MAX
BENIGN	60875	27.21	6.40	0.00	24.00	29.60	32.00	60.00
BRUTE FORCE	729	23.48	7.22	4.80	18.40	24.40	29.60	38.40
DDOS	1078975	13.92	8.23	0.00	8.00	19.88	20.00	57.76
DOS	324651	15.11	6.29	0.00	8.00	19.88	20.00	53.56
MIRAI	145131	2.98	4.09	0.00	0.00	0.24	8.00	54.36
RECON	38039	23.03	6.92	0.00	19.20	22.80	27.60	60.00
SPOOFING	25928	23.37	11.24	0.00	12.80	24.80	32.00	60.00
WEB	1392	22.23	7.76	4.00	16.40	22.40	29.60	41.60

Fuente: Elaboración propia.

Figura 18: Gráfico de barras de la característica Header_Length



Fuente: Elaboración propia.

Patrón Observado: La longitud media de cabecera es mayor en el tráfico benigno (27.21) comparado con Mirai (2.92).

Interpretación: Esto puede reflejar la complejidad y la cantidad de información adicional necesaria para las conexiones seguras o protocolos utilizados en el tráfico benigno.

10. DNS (Sistema de nombres de dominio).

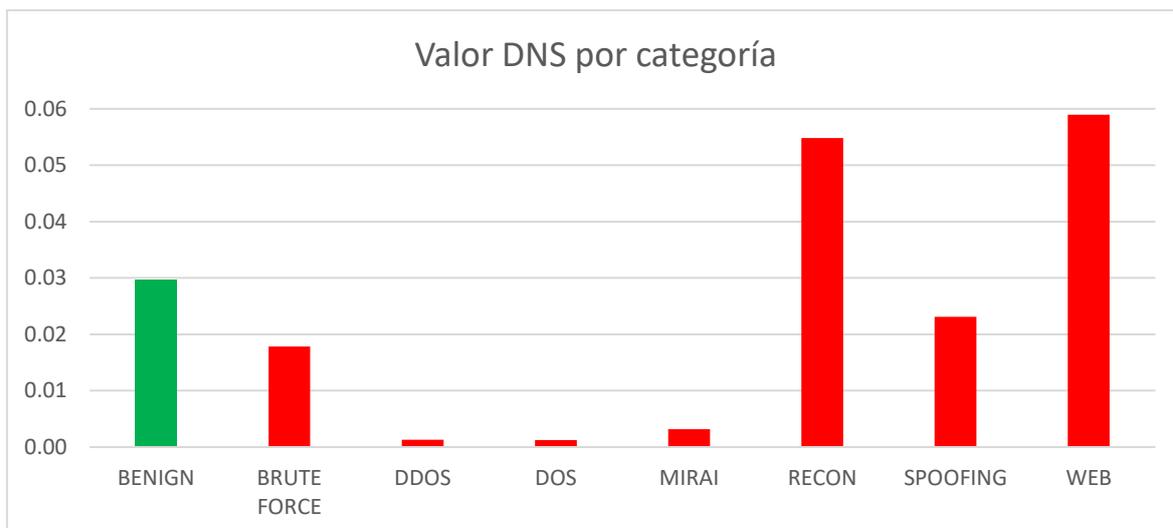
Es el valor que indica si el tráfico IoT utiliza DNS, este sistema es fundamental para la navegación en internet, ya que traduce nombres de dominio que entienden los usuarios, en direcciones IP, que las computadoras y otros dispositivos utilizan para identificarse en la red. El DNS es fundamental para la seguridad y el funcionamiento de los dispositivos IoT, permite detectar tráfico malicioso al identificar patrones anómalos que pueden indicar ataques como DDoS u otros, además, facilita el análisis del comportamiento del tráfico, ayuda a distinguir entre actividades normales y sospechosas, también proporciona información valiosa para diagnosticar problemas de conectividad, ya que la dificultad de resolver nombres de dominio pueden señalar fallos en la configuración del dispositivo o en la red. Ver su información en la tabla 10 y la figura 19.

Tabla 10: Estadísticas descriptivas de la característica DNS

LABEL	COUNT	MEAN	STD	MIN	25%	50%	75%	MAX
BENIGN	60875	0.03	0.08	0.00	0.00	0.00	0.00	1.00
BRUTE FORCE	729	0.02	0.06	0.00	0.00	0.00	0.00	0.40
DDOS	1078975	0.00	0.01	0.00	0.00	0.00	0.00	0.98
DOS	324651	0.00	0.01	0.00	0.00	0.00	0.00	1.00
MIRAI	145131	0.00	0.02	0.00	0.00	0.00	0.00	0.99
RECON	38039	0.05	0.10	0.00	0.00	0.00	0.10	0.90
SPOOFING	25928	0.02	0.07	0.00	0.00	0.00	0.00	0.80
WEB	1392	0.06	0.12	0.00	0.00	0.00	0.10	1.00

Fuente: Elaboración propia.

Figura 19: Gráfico de barras de la característica DNS



Fuente: Elaboración propia.

Patrón Observado: El uso del protocolo DNS es bastante bajo en todos los tipos de tráfico, especialmente en DDoS, DoS y Mirai donde es casi inexistente, en contraste Recon y Web tienen valores un poco más altos, casi el doble del benigno.

Interpretación: Esto indica que DDoS, DoS y Mirai no requieren resoluciones DNS (nombre de dominio) o están diseñados para operar sin interacciones DNS típicas (como anónimos), en contraste Recon y Web requieren más de lo necesario, alrededor del doble, para recopilar información y explorar vulnerabilidades de la red para futuros y diferentes ataques.

Interpretaciones generales del dataframe de trabajo.

Los datos muestran patrones claros entre el tráfico benigno y malicioso:

- El tráfico malicioso tiende a ser más uniforme y menos complejo, con longitudes de paquetes más cortas y menos variabilidad.
- Los ataques como DDoS, DoS y Mirai se caracterizan por su alto volumen, pero baja diversidad (son numerosos y compactos), lo que revela un objetivo estratégico enfocado en saturar recursos de la red.
- En contraste, el tráfico benigno presenta una variedad rica en longitudes y protocolos utilizados, reflejando la naturaleza diversa del uso legítimo de la red.
- Estos patrones son esenciales para desarrollar sistemas efectivos de detección y mitigación ante amenazas y ataques en redes IoT.

5.1.2. Modelos de clasificación para la detección de anomalías.

Las anomalías son de dos clases, tráfico benigno: 0 y tráfico malicioso: 1; y los modelos de clasificación de aprendizaje automático que se entrenaron para clasificar el tráfico de la red IoT en la presente investigación se presentan a continuación.

División de Datos

Previo al entrenamiento de los cinco modelos, se divide el conjunto de datos en dos conjuntos, para el entrenamiento y la prueba, esto se realiza utilizando la función "train_test_split", que separa los datos en un 80% para entrenamiento y un 20% para prueba, esta división permite evaluar el rendimiento del modelo en datos no vistos (conjunto de prueba), asegurando que las métricas obtenidas sean representativas al desempeño de cada modelo. Los resultados de los cinco modelos están trabajados con los parámetros por defecto (configuración interna y propia de cada modelo), en cuanto a los hiperparámetros se han ajustados para obtener un mejor rendimiento general de los modelos, cada caso se indica en la descripción posterior de cada modelo.

A. Regresión Logística.

Es un modelo de aprendizaje automático de clasificación, utilizado para predecir la probabilidad de que una instancia pertenezca a una clase específica, siendo especialmente útil en problemas de clasificación binaria, en el contexto de la detección de tráfico en redes IoT, para la presente investigación, este modelo se aplica para distinguir entre el tráfico benigno (0) y tráfico malicioso (1), la regresión logística calcula la respectiva probabilidad para el tráfico benigno o malicioso, según corresponda, utilizando una función logística que transforma la salida en un rango entre 0 y 1, menor a 0.5 benigno, mayor a 0.5 malicioso y 0.5 es el punto de corte o de decisión neutral, pero que también es configurable. En consecuencia, de este modo el modelo toma la decisión, si el tráfico debe ser clasificado como benigno o malicioso. La regresión logística es particularmente efectiva cuando las relaciones entre las características y la clase objetivo son lineales, lo que podría evidenciar patrones en tráfico de redes IoT.

A continuación, se presentan sus respectivos componentes y funcionamiento:

“Logistic Regression”

Es un algoritmo de clasificación utilizado para crear el modelo clasificación, que predice la probabilidad de que una determinada entrada pertenezca a una clase específica, la regresión logística se utiliza cuando la variable objetivo es categórica, en nuestro caso benigno o maligno.

“max_iter”

Es un hiperparámetro que especifica el número máximo de iteraciones que el algoritmo va a realizar durante el proceso de optimización, un valor de max_iter=1000 (valor por defecto) significa que el modelo intentará ajustar los parámetros hasta 1000 veces, para la presente investigación se trabajó con el valor de 60 iteraciones, pues valores más altos no incrementaron ni disminuyeron el rendimiento del modelo, solo hicieron que la ejecución tomara más tiempo y valores bajos (50 o menos) generaban error o no se completaba la ejecución, por lo mencionado se comprueba que la manipulación de hiperparámetros tiene un efecto en el rendimiento general del modelo, dicha manipulación ya depende de la aplicación donde se desee implementar el modelo, para nuestro caso se conserva el valor mencionado.

La evaluación de la exactitud (accuracy) de la Regresión Logística y de los demás modelos posteriores, se realiza mediante la métrica “accuracy_score” para comparar las predicciones con las etiquetas reales del conjunto de prueba, lo que da como resultado el rendimiento general del modelo. También se genera un reporte de clasificación con “classification_report” que incluye la exactitud, precisión, recall y F1-score, para cada clase, ayudando a entender cómo clasifica el tráfico benigno y malicioso. Además, con “confusion_matrix” se crea la matriz de confusión que muestra los verdaderos positivos, verdaderos negativos, falsos positivos y falsos negativos, brindando información general sobre las clasificaciones correctas e incorrectas de cada modelo. Los resultados obtenidos se presentan en la tabla 11 y el cuadro 14.

Tabla 11: Reporte de clasificación de la Regresión Logística

	PRECISION	RECALL	F1-SCORE	SUPPORT
BENIGNO	67.64%	67.54%	67.59%	12189
MALICIOSO	98.78%	98.78%	98.78%	322955
ACCURACY		97.64%		335144

Fuente: Elaboración propia.

Cuadro 14: Matriz de confusión de la Regresión Logística

	PREDICE BENIGNO	PREDICE MALICIOSO
ES BENIGNO	8233	3956
ES MALICIOSO	3939	319016

Fuente: Elaboración propia.

Interpretación:

- **Precision:** La precisión para la clase maliciosa es del 98.78%, lo que indica que cuando el modelo predice que un caso es malicioso, es correcto en un 98.78% de los casos, sin embargo, la precisión para la clase benigna es significativamente más baja (67.64%), lo que indica que el modelo tiene dificultades para identificar correctamente los casos benignos.
- **Recall:** El modelo detecta alrededor del 67.54% de todas las instancias realmente benignas.
- **F1-Score:** El 67.59% de esta métrica indica un equilibrio entre precisión y recall para la clase benigna, pero también indica que hay un margen importante por mejorar.
- **Exactitud (Accuracy):** El modelo tiene un alto rendimiento general del 97.64% al clasificar correctamente tanto el tráfico benigno como el malicioso. La Regresión Logística es efectiva en general, pero su rendimiento en la clase Benigno es deficiente.

B. Árbol de Decisión.

Es un modelo de aprendizaje automático utilizado para clasificación, que toma decisiones basadas en las características de entrada y las divide en subconjuntos según las características más relevantes, este modelo se representa como un árbol, donde cada nodo interno representa una prueba sobre una característica y cada rama representa el resultado de esa prueba, las hojas del árbol corresponden

a las predicciones de la clase, se aplica el árbol de decisión para clasificar el tráfico benigno y malicioso.

El modelo construye una estructura jerárquica que facilita la interpretación del proceso de clasificación, puesto que cada división en el árbol es una decisión sobre una característica del tráfico, ayudando a identificar patrones en los datos. Los árboles de decisión son útiles cuando las relaciones entre las características y la clase objetivo son relativamente simples y fácilmente separables.

A continuación, se presentan sus respectivos componentes y funcionamiento:

“DecisionTreeClassifier”

Es un algoritmo de aprendizaje supervisado que construye un modelo de predicción mediante una estructura de árbol jerárquico, donde cada nodo representa una decisión basada en características específicas de los datos, divide los datos en función de las características de entrada, tomando decisiones en cada nodo del árbol, para este caso, benigno o malicioso. En cuanto a la profundidad del árbol (número de decisiones para realizar la clasificación) se trabajó con el valor 54, que es un valor por defecto y funciona bien, se ejecutó el modelo variando su profundidad con valores de 1 a 500 y el rendimiento fue el mismo.

“random_state”

Es un parámetro que establece la semilla (valor inicial) para el generador de números aleatorios, lo que garantiza que los resultados sean reproducibles, es decir que siempre se generen los mismos valores aleatorios, al fijar un valor de “random_state = 42” (valor por defecto), se asegura que el modelo produzca los mismos resultados cada vez que se ejecute con los mismos datos, esto es útil para comparar el rendimiento de distintos modelos o para asegurar que los resultados obtenidos sean consistentes a lo largo de diferentes ejecuciones del código. Se modificó el valor por defecto, pero los resultados variaban muy poco, incrementaba o disminuían sin relación notoria, se conservó como estaba. Los resultados obtenidos se presentan en la tabla 12 y el cuadro 15.

Tabla 12: Reporte de clasificación del Árbol de decisión

	PRECISION	RECALL	F1-SCORE	SUPPORT
BENIGNO	75.82%	77.41%	76.60%	12189
MALICIOSO	99.15%	99.07%	99.11%	322955
ACCURACY		98.28%		335144

Fuente: Elaboración propia.

Cuadro 15: Matriz de confusión del Árbol de decisión

	PREDICE BENIGNO	PREDICE MALICIOSO
ES BENIGNO	9435	2754
ES MALICIOSO	3009	319946

Fuente: Elaboración propia.

Interpretación:

- **Precision:** La precisión para la clase maliciosa es de 99.15%, lo que indica un excelente desempeño en la identificación de dicha clase, la precisión para la clase benigna es también mejor (75.82%) en comparación con la regresión logística.
- **Recall:** El modelo detecta al 77.41% la clase benigna, el modelo identifica mejor las instancias benignas en comparación con la regresión logística.
- **F1-Score:** El 76.60% en esta métrica indica un buen equilibrio entre precisión y recall para la clase benigna, también por mejorar en esta clase.
- **Exactitud (Accuracy):** El modelo tiene un mejor rendimiento general que la Regresión Logística, 98.28% al clasificar correctamente tanto el tráfico benigno como el malicioso. El árbol de decisión supera a la regresión logística en la clasificación de instancias benignas y maliciosas, mostrando un buen balance entre ambas clases.

C. Random Forest (Bosque Aleatorio).

Es un modelo de aprendizaje automático que utiliza un conjunto de árboles de decisión para realizar la clasificación de datos, se basa en "ensembles" (conjuntos) de modelos, donde múltiples árboles de decisión independientes son entrenados y sus resultados son combinados para obtener una predicción final, esto ayuda a mejorar la precisión del modelo y reducir el riesgo de sobreajuste en comparación con un solo árbol de decisión, en el contexto de redes IoT, Random Forest puede

clasificar el tráfico como benigno o malicioso combinando las decisiones de múltiples árboles, ofreciendo una clasificación más robusta.

Cada árbol en el bosque se entrena con un subconjunto diferente de los datos, utilizando un proceso llamado "bagging" (proceso que busca mejorar la estabilidad y precisión de algoritmos, reduciendo la varianza y previniendo el sobreajuste), el que permite, que cada árbol sea ligeramente diferente, luego las predicciones de todos los árboles se agregan para decidir la clase final, este modelo es eficaz cuando hay una gran cantidad de características y permite identificar patrones complejos en los datos.

A continuación, se presentan sus respectivos componentes funcionamiento:

“RandomForestClassifier”

Es un algoritmo de aprendizaje supervisado que utiliza un conjunto de árboles de decisión para hacer la clasificación de datos, este modelo ayuda a reducir el sobreajuste que podría ocurrir al usar un solo árbol, ya que toma decisiones basadas en múltiples árboles y sus combinaciones, en este caso, se usa para clasificar el tráfico como benigno o malicioso. Se inició el entrenamiento con una profundidad de los árboles de 5 y tuvo un rendimiento general 98.27% un poco menos que el árbol de decisión, con 10 lo superó con 98.50% y 30 llegó a 98.63%, valores superiores a 30 no mejoran el rendimiento, se mantiene igual.

“random_state”

Es el mismo parámetro que también utiliza el árbol de decisión que establece la semilla (valor inicial) para el generador de números aleatorios, garantizando que los resultados sean consistentes cada vez que se ejecute el modelo. Se modificó el valor por defecto 42 por 30 y se obtuvo mejor rendimiento del modelo en consecuencia se conservó dicho valor. Los resultados obtenidos se presentan en la tabla 13 y el cuadro 16.

Tabla 13: Reporte de clasificación de Random Forest

	PRECISION	RECALL	F1-SCORE	SUPPORT
BENIGNO	81.03%	81.24%	81.13%	12189
MALICIOSO	99.29%	99.28%	99.29%	322955
ACCURACY		98.63%		335144

Fuente: Elaboración propia.

Cuadro 16: Matriz de confusión de Random Forest

	PREDICE BENIGNO	PREDICE MALICIOSO
ES BENIGNO	9902	2287
ES MALICIOSO	2318	320637

Fuente: Elaboración propia.

Interpretación:

- **Precision:** La precisión para la clase maliciosa es del 99.29% y para la clase benigna es del 81.03%, esto indica que el modelo es muy efectivo para detectar ambas clases y también mejora el reconocimiento de instancias benignas, superando a los modelos anteriores.
- **Recall:** El modelo detecta al 81.24% la clase benigna, el modelo identifica mejor dichas instancias que los dos modelos anteriores.
- **F1-Score:** El 81.13% en esta métrica indica tiene mayor equilibrio entre precisión y recall que los modelos anteriores para la clase benigna, lo cual refuerza la efectividad del modelo para dicha clase.
- **Exactitud (Accuracy):** El modelo tiene un buen rendimiento general del 98.63%, Random Forest tiene un rendimiento superior a la Regresión Logística y al Árbol de Decisión, clasificando con éxito tanto a las instancias benignas como a las maliciosas.

D. AdaBoost (Adaptive Boosting: Potenciación Adaptativa).

Es un algoritmo de aprendizaje automático de clasificación que combina varios modelos débiles (por lo general, árboles de decisión simples) para crear un modelo fuerte, en lugar de entrenar un solo modelo, AdaBoost ajusta iterativamente los modelos débiles en función de los errores cometidos en las iteraciones anteriores, lo que permite que el modelo final sea más preciso, en el contexto de la detección de tráfico en redes IoT, AdaBoost puede usarse para mejorar la clasificación de tráfico benigno o malicioso, combinando los resultados de varios árboles de decisión, AdaBoost es muy eficaz cuando se necesita mejorar un modelo débil y puede generar un modelo fuerte a partir de modelos individuales con baja precisión, cada modelo sucesivo se ajusta para corregir los errores cometidos por los modelos previos.

A continuación, se presentan sus respectivos componentes y funcionamiento:

“AdaBoostClassifier”

Es un algoritmo de aprendizaje automático clasificador que se utiliza para combinar un conjunto de modelos débiles (generalmente árboles de decisión de profundidad limitada) y ajusta sus pesos para mejorar la precisión general, en el presente trabajo, se usa para clasificar el tráfico de red como benigno o malicioso.

“base_estimator”

Es el modelo base sobre el que se construye el conjunto, en este caso se utiliza un árbol de decisión con profundidad limitada (`max_depth=1`) para asegurarse de que cada árbol sea débil, AdaBoost ajustará los pesos de los ejemplos mal clasificados para mejorar el modelo en cada iteración. Se cambió la profundidad de los árboles y a partir de 2, se generó un ligero incremento en el rendimiento general del modelo, con la profundidad igual a 5 se alcanzó en mejor rendimiento, valores mayores no presentaron mejora e incluso disminuye, además toma más tiempo la ejecución del modelo.

“n_estimators”

Es el número de modelos clasificadores que se van a combinar para mejorar el rendimiento general del modelo, en este caso se especifica 50 estimadores, lo que significa que el modelo final será el resultado de 50 iteraciones del algoritmo AdaBoost, el estimador es como un aprendiz que aprende a partir de sus errores.

“random_state”

Como en los otros modelos, “random_state” se utiliza para garantizar que los resultados sean reproducibles, es decir controlar la aleatoriedad, en la selección de muestras, la inicialización de pesos o la división de datos, esto significa que cada vez que se ejecute el mismo código, se generen los mismos valores aleatorios y se obtengan los mismos resultados (controlando así la naturaleza aleatoria del proceso). También se trabajó con 42 como valor inicial, no hay cambio notorio al modificarlo. Los resultados obtenidos se presentan en la tabla 14 y el cuadro 17.

Tabla 14: Reporte de clasificación de AdaBoost

	PRECISION	RECALL	F1-SCORE	SUPPORT
BENIGNO	79.42%	81.22%	80.31%	12189
MALICIOSO	99.29%	99.21%	99.25%	322955
ACCURACY		98.55%		335144

Fuente: Elaboración propia.

Cuadro 17: Matriz de confusión de AdaBoost

	PREDICE BENIGNO	PREDICE MALICIOSO
ES BENIGNO	9900	2289
ES MALICIOSO	2566	320389

Fuente: Elaboración propia.

Interpretación:

- **Precision:** La precisión para la clase maliciosa es del 99.29%, mientras que para la clase benigna es del 79.42%, esto muestra que tiene la capacidad para detectar dicha clase, pero aún con un margen de mejora en la identificación de instancias benignas, aunque no supera a Random Forest, presenta una mejora respecto a la Regresión Logística y al Árbol de Decisión.
- **Recall:** El modelo detecta al 81.22% la clase benigna, superando en esta métrica también a la Regresión Logística y al Árbol de Decisión.
- **F1-Score:** El 80.31% en esta métrica para la clase benigna muestra un buen equilibrio entre precisión y recall.
- **Exactitud (Accuracy):** El modelo tiene un buen rendimiento general, el 98.55% indica que AdaBoost ofrece resultados competitivos y es efectivo en la clasificación de ambas clases, aunque no supera a Random Forest en términos de precisión para la clase Benigno.

E. Perceptrón

Es un algoritmo de aprendizaje automático supervisado que sirve principalmente para tareas de clasificación binaria, es uno de los modelos más simples de red neuronal, con una sola capa de neuronas artificiales que simulan el comportamiento de las neuronas biológicas, el perceptrón toma una serie de entradas (características), las multiplica por un conjunto de pesos, las suma y a continuación, aplica una función de activación para determinar la salida, en el caso del perceptrón, la función de activación es típicamente una función escalón, que

produce una salida binaria, en el contexto de clasificación de tráfico IoT el perceptrón se puede utilizar para distinguir entre dos clases (benigna y maliciosa), a medida que se entrena el modelo, ajusta los pesos de las entradas en función de los errores cometidos, permitiendo así que el modelo mejore su capacidad de clasificación con cada iteración.

A continuación, se presentan sus respectivos componentes y funcionamiento:

“Step Function” (Función Escalón)

Es la función de activación que utiliza el perceptrón para decidir si una neurona se activa o no, recibe y evalúa varias señales de entrada que se multiplican por determinados pesos (valores numéricos) para determinar cuanta influencia o importancia tiene cada entrada en la salida final y luego se suman, si esta suma total o ponderada es mayor o igual a cero, la neurona se activa y la salida será 1 (tráfico malicioso), si la suma ponderada es menor que cero la neurona no se activa y la salida será 0 (tráfico benigno).

"random_state"

Al igual que modelos anteriores, el parámetro "random_state" se utiliza para asegurar la reproducibilidad del modelo, al establecer un valor específico (42 en este caso), se puede controlar la aleatoriedad, esto asegura que cada vez que se ejecute el mismo código, el modelo obtenga los mismos resultados y se mantenga la confiabilidad de los resultados obtenidos en diferentes ejecuciones, también se modificó el valor 42 y sus resultados fueron más errados, con 1 presentó un ligero incremento en el rendimiento general y con 100 o más disminuyó su rendimiento. Los resultados obtenidos se presentan en la tabla 15 y el cuadro 18.

Tabla 15: Reporte de clasificación del Perceptrón

	PRECISION	RECALL	F1-SCORE	SUPPORT
BENIGNO	64.27%	73.58%	68.61%	12189
MALICIOSO	99.00%	98.46%	98.73%	322955
ACCURACY		97.55%		335144

Fuente: Elaboración propia.

Cuadro 18: Matriz de confusión del Perceptrón

	PREDICE BENIGNO	PREDICE MALICIOSO
ES BENIGNO	8969	3220
ES MALICIOSO	4987	317968

Fuente: Elaboración propia.

Interpretación:

- **Precision:** La precisión para la clase maliciosa es del 99.00%, pero la precisión para la clase benigna es solo del 64.27%, esto indica problemas significativos al identificar correctamente a los casos benignos, es la precisión más baja de los cinco modelos.
- **Recall:** El modelo detecta al 73.58% a la clase benigna, superando sólo a la Regresión Logística.
- **F1-Score:** El 68.61% en esta métrica, muestra las dificultades del modelo para equilibrar la precisión y recall.
- **Exactitud (Accuracy):** El modelo tiene un rendimiento general considerable del 97.55% es alto, pero puede ser engañoso debido al desbalance de precisión y recall. El perceptrón tiene la exactitud más baja en comparación con los otros modelos, especialmente en la detección de instancias benignas, lo que puede ser causa de bloqueo de tráfico benigno, importante y necesario para la comunicación.

Los cinco modelos ya analizados muestran un buen rendimiento general en términos de exactitud (ver cuadro 19 y figura 20), pero varían significativamente en su capacidad para detectar correctamente las instancias benignas y maliciosas:

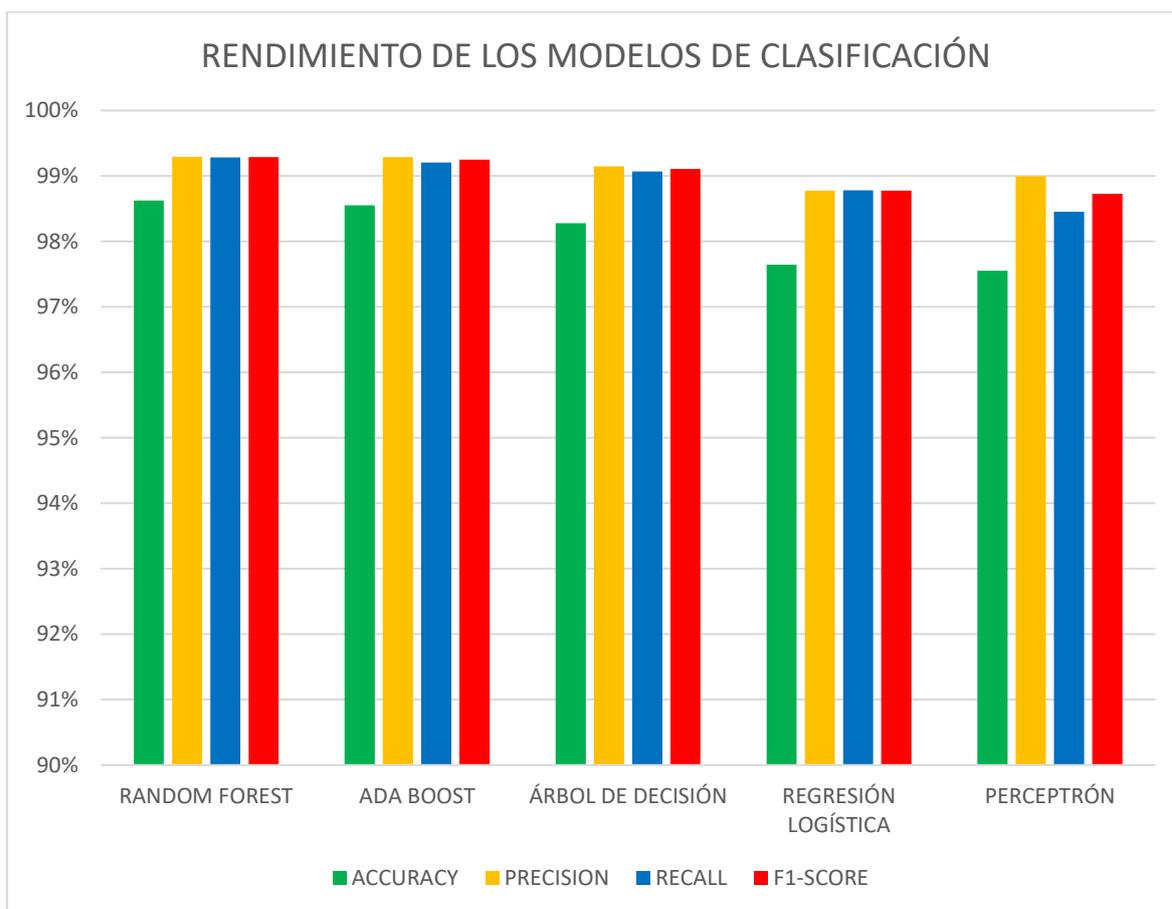
- El Random Forest se destaca como el modelo más equilibrado y robusto, ofreciendo una alta exactitud así como una alta precisión en ambas clases.
- El Árbol de Decisión y Ada Boost también presentan buenos resultados, aunque con ligeras diferencias en sus capacidades de detección.
- La Regresión Logística tiene una alta exactitud, pero presenta debilidades significativas en la identificación correcta de instancias benignas.
- El Perceptrón muestra el menor rendimiento general y presenta desafíos importantes en su capacidad para clasificar correctamente los casos benignos.

Cuadro 19: Rendimiento de los modelos de clasificación

MODELO DE CLASIFICACIÓN	ACCURACY	PRECISION	RECALL	F1-SCORE
RANDOM FOREST	98.63%	99.29%	99.28%	99.29%
ADA BOOST	98.55%	99.29%	99.21%	99.25%
ÁRBOL DE DECISIÓN	98.28%	99.15%	99.07%	99.11%
REGRESIÓN LOGÍSTICA	97.64%	98.78%	98.78%	98.78%
PERCEPTRÓN	97.55%	99.00%	98.46%	98.73%

Fuente: Elaboración propia.

Figura 20: Rendimiento de los modelos de clasificación



Fuente: Elaboración propia.

Interpretación:

En el cuadro 19 y la figura 20 se aprecia que Random Forest es el modelo más confiable, con un rendimiento notablemente alto y equilibrado en cuanto a exactitud y precisión en la clasificación de tráfico malicioso y benigno, esto indica que es el modelo de clasificación más eficiente para aplicaciones críticas donde la detección correcta de amenazas es esencial.

5.2. Discusión de resultados

En esta sección se presenta un análisis de los resultados obtenidos en la investigación y su relación con los objetivos y las hipótesis planteadas.

A. Resultados relevantes.

Los resultados de la evaluación de los cinco modelos de clasificación (Regresión Logística, Árbol de Decisión, Random Forest, Ada Boost y Perceptrón) indican que se lograron niveles significativos en el rendimiento general de los modelos en la detección de anomalías en redes IoT utilizando el dataset CICIoT2023, a continuación, se presentan los hallazgos más relevantes.

- Rendimiento General (Exactitud: Accuracy):
Todos los modelos alcanzaron una exactitud superior al 90%, cumpliendo así con la hipótesis general que planteaba, que el modelo sería capaz de detectar anomalías con una precisión superior al 90%. El modelo Random Forest obtuvo el mejor rendimiento con una exactitud del 98.63%, seguido por el Árbol de Decisión y Ada Boost, lo que demuestra también el cumplimiento de la hipótesis específica relacionada a la efectividad los modelos en la detección de anomalías en redes IoT.
- Rendimiento Parcial (Precisión):
La precisión para la clase maliciosa fue alta en todos los modelos, especialmente en Random Forest (99.29%) y Ada Boost (99.29%), lo que valida la hipótesis específica que sugiere que el diseño de un modelo de clasificación permitirá identificar patrones anómalos en el tráfico de red con una tasa de detección significativamente alta, sin embargo, se observó una variabilidad en la precisión para la clase benigna, siendo más baja en la Regresión Logística (67.64%).
- Impacto del análisis exploratorio y procesamiento de datos:
La limpieza y preparación de los datos, que incluyeron la eliminación de filas con valores faltantes y duplicados, así como la selección de características más relevantes mediante la Correlación de Pearson, mejoraron la calidad del dataset, esto respalda la hipótesis específica que indica que el análisis exploratorio y el procesamiento de datos mejorarán la calidad del dataset, generando un rendimiento superior del modelo.

- **Identificación de Patrones:**

A través del análisis exploratorio inicial, se generaron gráficos de barras y cuadros de con descripciones estadísticas para cada característica, que ayudaron a visualizar patrones relevantes en el tráfico de red, aunque algunos patrones fueron difíciles de discernir debido a la magnitud y dispersión de los datos, esta etapa fue crucial para entender mejor las características del dataset antes del entrenamiento del modelo.

B. Validación de los objetivos.

Objetivo General

El objetivo general de desarrollar un modelo de clasificación para la detección de anomalías en redes IoT, ha sido cumplido satisfactoriamente, ya que se implementaron cinco modelos diferentes y se evaluaron sus rendimientos utilizando métricas clave como "Accuracy" y "Precision".

Objetivos Específicos

- **Análisis exploratorio y procesamiento de datos.** se realizó un análisis exhaustivo del dataset CICIoT2023, lo que incluyó la limpieza de datos y la selección adecuada de características, esto ha permitido mejorar la calidad del dataset, se disminuyó su tamaño eliminando información duplicada o irrelevante.
- **Diseño del modelo,** se diseñaron e implementaron cinco modelos diferentes para clasificar el tráfico benigno y malicioso (Regresión Logística, Árbol de Decisión, Random Forest, Ada Boost y Perceptrón).
- **Evaluación de los modelos,** se evaluaron los cinco modelos en cuanto a su rendimiento general de exactitud (accuracy) es decir su capacidad de clasificación de la clase benigna y maliciosa, también se evaluó su rendimiento parcial de precisión, es decir su capacidad de clasificación solo para clase benigna o solo para la clase maliciosa.

C. Validación de las hipótesis.

La evaluación los modelos ejecutados, determinó que todos alcanzaron un rendimiento superior al 98% en términos de exactitud (accuracy), validando así su efectividad en la detección de anomalías.

Hipótesis general:

La hipótesis general fue validada al demostrar que el modelo puede detectar anomalías con una exactitud superior al 90%.

Las hipótesis específicas también fueron respaldadas por los resultados obtenidos:

- El análisis exploratorio y procesamiento mejoraron efectivamente la calidad del dataset.
- El diseño del modelo permitió identificar patrones anómalos con altas tasas de detección.
- Las métricas de rendimiento mostraron resultados positivos, confirmando que el modelo es efectivo para detectar anomalías.

La investigación ha demostrado que es posible desarrollar un modelo efectivo para la detección de anomalías en redes IoT utilizando el dataset CICIoT2023, los resultados obtenidos no solo cumplen con las expectativas planteadas al inicio de la investigación, sino que también resaltan la importancia del preprocesamiento adecuado y análisis exploratorio para maximizar el rendimiento del modelo. Los modelos más robustos como Random Forest y Ada Boost ofrecen una gran capacidad para detectar tráfico malicioso con una alta exactitud y precisión, lo cual es crucial para mejorar las medidas de seguridad en entornos IoT. Estos hallazgos son valiosos para investigadores y profesionales en ciberseguridad interesados en implementar soluciones efectivas para proteger redes IoT.

Contraste con los antecedentes de la investigación

A continuación, se realiza una comparación detallada de los resultados obtenidos en la presente investigación con los antecedentes más relacionados:

1. Investigación de Tseng et al. (2024) – Taiwán.

- Exactitud:
 - Tseng et al., alcanzaron una exactitud del 99.40% con un modelo de aprendizaje profundo basado en Transformer.
 - En la presente investigación, la mayor exactitud alcanzada fue del 98.63% utilizando el modelo Random Forest.
- Similitud o Relación:
 - Ambos estudios trabajaron con el dataset CICIoT2023 y también se enfocaron en la detección de anomalías en redes IoT, aunque Tseng et al. realizaron una clasificación multiclase, mientras que la presente investigación se centró en clasificación binaria.
 - Ambos estudios trabajaron con diferentes modelos de clasificación de los cuales ya mostró se exactitud más alta.
- Divergencia o Diferencia:
 - La exactitud de Tseng et al. es superior debido al uso de modelos avanzados como Transformer, diseñado para manejar grandes cantidades de datos y relaciones complejas.
 - La presente investigación utilizó modelos más simples como Random Forest, lo que puede explicar la ligera diferencia en el rendimiento.
 - Tseng et al., trabajaron con una computadora de muy alto rendimiento, en la presente investigación se trabajó con una computadora básica.
- Aporte Novedoso:
 - La presente investigación aporta un enfoque más accesible y menos exigente en cuanto capacidad de cómputo, facilitando su implementación en entornos con recursos limitados y obteniendo también resultados satisfactorios.

2. Investigación de Pinto et al. (2023) – Canadá.

- Exactitud:
 - En ambas investigaciones se trabajó con Random Forest, la presente investigación alcanzó una exactitud del 98.63%, mientras que Pinto et al. lograron 99.68%.
- Similitud o Relación:
 - Ambos estudios utilizaron el dataset CICIoT2023, lo que permite cierta comparación en la efectividad de los modelos.
- Divergencia o Diferencia:
 - La investigación de Pinto et al. mostró un rendimiento superior, posiblemente por haber trabajado con todo el dataset y también porque indican que para el manejo del dataset completo se necesita de bastante capacidad de cómputo, por lo que se deduce que utilizaron una o varias computadoras como la que mencionan.
 - La presente investigación obtuvo un rendimiento inferior y las causas probables pueden ser que se utilizó sólo una muestra del total del dataset (6 de los 63 archivos csv) y también porque la capacidad de cómputo es limitada, pues se trabajó con una computadora básica.
- Aporte Novedoso:
 - La presente investigación valida la aplicación y el rendimiento satisfactorio de Random Forest en un contexto diferente, lo que puede ser útil para futuras investigaciones con recursos de cómputo limitados.
 - Un aporte notable es que la exactitud del Árbol de Decisión ejecutado en la presente investigación (98.28%) es superior a la del Perceptrón ejecutado por Pinto et al. (98.18%), evidenciando que métodos más simples pueden ser también efectivos en ciertos contextos.

3. Investigación de Ragab et al. (2023) – Arabia Saudita.

- Exactitud:
 - La presente investigación utilizó Random Forest con el cual alcanzó 98.63% de exactitud con siete categorías de ataque.

- En la investigación de Ragab et al. lograron 99.20% con un clasificador de aprendizaje profundo específico, sólo para DDoS.
- Similitud o Relación:
 - Ambos estudios abordan la detección de anomalías en redes IoT, aunque Ragab et al. se enfocaron específicamente en DDoS y en la presente investigación se trabajó con siete categorías de ataque (DDoS, DoS, Reconocimiento, Basados en la web, Fuerza bruta, Suplantación y Mirai).
 - Ambos estudios utilizaron datasets reales.
- Divergencia o Diferencia:
 - Ragab et al. utilizaron un enfoque especializado para ataques DDoS, lo que podría explicar su mayor exactitud en ese contexto específico, sólo para ese tipo de ataque y trabajaron con el dataset Bot-IoT.
 - En la presente investigación se trabajó con siete categorías de ataque, incluido DDoS con el dataset CICIoT2023.
- Aporte Novedoso:
 - La presente investigación amplía el enfoque a la detección general de anomalías, no limitándose solo a un tipo específico de ataque.

4. Investigación de Vigoya (2023) – España.

- Exactitud:
 - La presente investigación con Random Forest alcanzó 98.63% de exactitud, mientras que Vigoya logró 99.99% con el mismo modelo.
- Similitud o Relación:
 - Ambos estudios utilizaron técnicas de aprendizaje automático para detectar anomalías en el tráfico IoT.
- Divergencia o Diferencia:
 - La diferencia en exactitud puede deberse a la utilización de escenarios virtuales y diferentes configuraciones experimentales por parte de Vigoya.

- Vigoya trabajó con datasets generados artificialmente en escenarios virtuales, mientras que en la presente investigación se utilizó el dataset CICIoT2023 generado en condiciones reales.
- Aporte Novedoso:
 - La presente investigación ofrece una comparación práctica y realista sobre la efectividad del modelo Random Forest en condiciones distintas, reales y simuladas.

5. Investigación de Plata (2022) – Colombia.

- Exactitud:
 - La presente investigación con Random Forest alcanzó 98.63% de exactitud, mientras que el trabajo de Plata obtuvo 99.53% con KNN.
- Similitud o Relación:
 - Ambos estudios se centran en la detección de anomalías en redes IoT utilizando diferentes algoritmos y con datasets reales.
- Divergencia o Diferencia:
 - Plata utilizó el modelo KNN con el total de su propio dataset, en la presente investigación se trabajó con Random Forest y cuatro modelos más, con una muestra del dataset CICIoT2023.
 - Plata se centró exclusivamente en anomalías físicas relacionadas con fallas a nivel de sensores, mientras que la presente investigación abarcó tráfico malicioso y benigno.
- Aporte Novedoso:
 - La presente investigación aborda un problema más amplio relacionado con ciberseguridad IoT, ampliando el rango de aplicaciones posibles.

6. Investigación de Becerra et al. (2024) – Perú.

- Exactitud:
 - La presente investigación con Random Forest alcanzó 98.63% de exactitud, mientras que Becerra et al. lograron 99.40% con CNN (Red Neuronal Convolucional).

- Similitud o Relación:
 - Ambos estudios evaluaron los modelos con el dataset CICIoT2023 y también en ambos trabajos ejecutó la clasificación binaria.
- Divergencia o Diferencia:
 - Becerra et al. utilizaron modelos más avanzados como CNN, que generalmente es más efectivo para datos estructurados como imágenes, lo que podría explicar su mayor precisión, en la presente investigación se trabajó con modelos más simples, pero también eficientes.
 - Becerra et al. trabajaron con el 1% de cada categoría, seleccionando dicha muestra de forma aleatoria, en la presente investigación se trabajó con una muestra representativa de 7 829 178 registros (6 de los 63 archivos csv).
- Aporte Novedoso:
 - La presente investigación valida el uso de Random Forest como una alternativa viable y menos compleja para la detección de anomalías, cuyos resultados también son confiables.

7. Investigación de Becerra et al. (2024) – Perú.

- Exactitud:
 - La presente investigación con Random Forest alcanzó 98.63% de exactitud, mientras que Becerra et al., en una segunda investigación, también con Random Forest lograron 99.97% de exactitud, utilizando otro dataset para ataques DDoS.
- Similitud o Relación:
 - Ambos estudios abordan la detección de ataques cibernéticos, ejecutando la clasificación binaria.
- Divergencia o Diferencia:
 - La diferencia significativa en exactitudes puede deberse a las características específicas del dataset CICDDoS2019 utilizado por Becerra et al., exclusivo de ataques DDoS.
 - La presente investigación utilizó el dataset CICIoT2023 con siete categorías de ataque.

- Aporte Novedoso:
 - Ambas investigaciones demuestran que los modelos de clasificación como Random Forest pueden ser aplicados a diversas amenazas.

8. Investigación de León et al. (2024) – Perú.

- Exactitud:
 - La presente investigación con Random Forest alcanzó 98.63% de exactitud y León et al., alcanzaron niveles de efectividad de según el tamaño de la muestra, desde 65.50% hasta 100.00% con Redes Neuronales, para cuatros dispositivos diferentes.
- Similitud o Relación:
 - Ambos estudios se enfocan en la detección de anomalías, pero León et al. se centran en logs IoT específicos para cuatro dispositivos.
- Divergencia o Diferencia:
 - La variación en la efectividad reportada por León et al. depende del volumen y calidad de los logs analizados, utilizan ocho muestras de cuatro datasets (2 muestras de cada uno), la muestra más pequeña es de 610 datos y la más grande de 138 193 datos, cada dataset pertenece a cada uno de los cuatro dispositivos IoT, estudiados en su investigación.
 - En la presente investigación se trabaja sólo con una muestra de 7 829 178 datos del dataset CICIoT2023.
- Aporte Novedoso:
 - La presente investigación confirma que un modelo clasificación como Random Forest, más simple que una Red Neuronal, obtiene un mejor rendimiento cuando se lo entrena y evalúa con más datos.

9. Investigación de Bazan et al. (2024) – Perú.

- Exactitud:
 - La presente investigación con Random Forest alcanzó 98.63% de exactitud y Bazan et al., alcanzaron 99.60% también con el mismo modelo, pero para ransomware.

- Similitud o Relación:
 - Ambos estudios utilizan Random Forest para clasificar amenazas específicas dentro del ámbito IoT.
- Divergencia o Diferencia:
 - Bazan et al., trabajaron con el dataset RandomFore_ransomware_detection_and_classification y se enfocaron exclusivamente en ataques ransomware, lo cual puede haber influido positivamente en su exactitud debido a características específicas del dataset utilizado.
 - En la presente investigación se trabajó con el dataset CICIoT2023 y siete categorías de ataque.
- Aporte Novedoso:
 - La presente investigación resalta cómo el modelo de Random Forest, puede adaptarse y ser efectivo con diferentes tipos de amenazas más allá del ransomware.

10. Investigación de Nolasco (2023) – Perú.

- Exactitud:
 - Con este trabajo la comparación de la exactitud no se aplica, porque Nolasco centra su investigación en pronósticos meteorológicos, utilizando aprendizaje automático sin relación directa con la detección de anomalías IoT.
- Similitud o Relación:
 - Aunque ambos trabajos utilizan aprendizaje automático, los contextos de aplicación son completamente diferentes.
- Divergencia o Diferencia:
 - No hay relación directa entre ambos trabajos, la presente investigación está centrada exclusivamente en seguridad IoT y detección de anomalías, en cambio Nolasco trabaja con imágenes de radar y pronósticos meteorológicos.

- Aporte Novedoso:
 - Ambas investigaciones son aplicadas a distintas áreas, lo cual sirve para demostrar que los modelos de aprendizaje automático son adaptables y eficientes en diferentes contextos, según donde se los requiera aprovechar.

Por el análisis ya presentado se concluye que, la presente investigación aplicó el modelo Random Forest y alcanzó una exactitud del 98.63% con una muestra del dataset CICIoT2023, se demuestra así que es viable y accesible para la detección de anomalías en redes IoT, especialmente en entornos con recursos de cómputo limitados, si bien las otras investigaciones han logrado mayores niveles de exactitud utilizando modelos más avanzados, algunas con diferente dataset e incluso con mejores recursos de cómputo, este estudio alcanza un equilibrio entre el rendimiento y la eficiencia, validando que modelos simples pueden ser efectivos en ciertos contextos. La comparación con otros trabajos resalta la importancia de considerar el tipo de dataset, la capacidad de cómputo disponible y el enfoque específico de la detección de anomalías al seleccionar un modelo de aprendizaje automático, por lo tanto esta investigación contribuye al campo de la ciberseguridad al ofrecer una alternativa práctica, realista, accesible y confiable para la detección de anomalías en entornos IoT, adaptable incluso a diferentes áreas, necesidades y recursos.

5.3. Conclusión general

Los resultados obtenidos en esta investigación revelan un notable desempeño de los modelos de clasificación aplicados al dataset CICIoT2023, donde la capacidad de detección de anomalías en redes IoT se ejecutó utilizando una muestra representativa del 10% de archivos de la población original, este enfoque, aunque limitado por las capacidades computacionales, permitió cumplir adecuadamente con los objetivos planteados.

La comparación de los cinco modelos de clasificación, Regresión Logística, Árbol de Decisión, Random Forest, Ada Boost y Perceptrón, ha arrojado hallazgos valiosos, los cuales muestran que el modelo Random Forest se destaca con una

precisión general del 98.63%, evidenciando su capacidad para detectar tanto tráfico benigno como malicioso de manera eficaz, con el resultado obtenido se tiene presente a la hipótesis general, donde se esperaba que el modelo alcanzara una precisión superior al 90%, por lo ya demostrado se valida la efectividad del enfoque adoptado.

Con la matriz de confusión de Random Forest, no solo se demostró que es el modelo más preciso, sino que también mostró una reducción significativa en el número de falsos positivos y negativos (clasificaciones incorrectas) en comparación con los demás modelos, este aspecto es crucial, dado que se busca garantizar la confiabilidad del modelo en aplicaciones reales de seguridad IoT.

El reporte de clasificación del Perceptrón, revela que dicho modelo alcanzó una alta precisión en la detección de tráfico malicioso, pero su bajo rendimiento en la identificación de tráfico benigno sugiere la necesidad de un balance en la capacidad de detección, esta observación es especialmente relevante, dado que las aplicaciones en entornos críticos reales requieren minimizar los falsos positivos y negativos.

Los resultados de esta investigación demuestran la efectividad del enfoque de clasificación abordado en la detección de anomalías en redes IoT, mientras que también se reconoce la necesidad de mayores recursos computacionales y conjuntos de datos completos para optimizar los resultados.

De forma global, por lo presentado y desarrollado en la presente investigación y el contexto abordado, se afirma que los resultados obtenidos pueden beneficiar principalmente a:

- ✓ Países en desarrollo con recursos de cómputo e infraestructura tecnológica restringidos.
- ✓ Organizaciones gubernamentales responsables de proteger infraestructuras críticas y servicios públicos con redes IoT
- ✓ Pequeñas y medianas empresas (PyMEs) con presupuestos limitados en ciberseguridad.

- ✓ Y también a usuarios domésticos de IoT que buscan salvaguardar sus datos personales, sus dispositivos y demás implicados.

Finalmente es importante mencionar que con el presente trabajo también se podrían beneficiar a:

- ✓ Fabricantes de dispositivos IoT interesados en diseñar productos más seguros.
- ✓ Consultoras de ciberseguridad que apoyan y asesoran en asuntos de protección de redes IoT.
- ✓ Presentes y futuros proveedores de servicios de seguridad IoT, que buscan ofrecer soluciones de detección de anomalías.
- ✓ Investigadores y académicos que buscan antecedentes para futuras investigaciones.
- ✓ Desarrolladores independientes de hardware y software de automatización y domótica.
- ✓ Así como también estudiantes y personas interesadas en la seguridad IoT.

Conclusiones

Conclusión 1: Efectividad del Modelo Random Forest.

Se determinó que el modelo Random Forest es el más eficaz para la detección de anomalías en redes IoT, logrando una exactitud del 98.63%, lo que valida la hipótesis general de la investigación.

Conclusión 2: Impacto del Preprocesamiento de Datos.

La correcta ejecución del análisis exploratorio y el procesamiento de datos contribuyeron a la mejora en el rendimiento del modelo clasificador, confirmando la hipótesis específica relacionada con la calidad del dataset.

Conclusión 3: Desbalance en la Detección.

Se observó que el Perceptrón, es efectivo al detectar la clase maliciosa, pero falla en la clase benigna, evidenciando la importancia del balance en la detección de anomalías.

Conclusión 4: Comparativa con Investigaciones Previas.

A pesar de los buenos resultados, la presente investigación presentó rendimientos inferiores en comparación con estudios que utilizaron el dataset completo, modelos de aprendizaje automático más avanzados y mejores recursos computacionales, lo que sugiere que los aspectos anteriores, son claves en los resultados finales.

Conclusión 5: Relevancia del Dataset CICIoT2023.

El uso del dataset CICIoT2023 se reafirma como un pilar pertinente para el desarrollo de futuras investigaciones en el área de seguridad IoT, resaltando la importancia de su uso correcto y análisis para la mejora continua en la detección de anomalías.

Recomendaciones

Recomendación 1: Ampliar el Tamaño de la Muestra.

Se recomienda realizar investigaciones futuras utilizando el dataset CICIoT2023 completo para obtener resultados más robustos y representativos, mejorando así la generalización del modelo.

Recomendación 2: Explorar Modelos Avanzados.

Considerar la implementación de técnicas de aprendizaje profundo y modelos híbridos que puedan complementar los modelos clásicos, para aumentar potencialmente la exactitud y la precisión, por consiguiente la fiabilidad en la detección.

Recomendación 3: Optimizar el Preprocesamiento.

Invertir tiempo y recursos en procesos de preprocesamiento más exhaustivos y eficientes, que incluyan una limpieza y selección de características más refinadas del dataset, para mejorar la calidad de la información utilizada.

Recomendación 4: Realizar Pruebas en Entornos Reales.

Implementar pruebas en entornos reales y bajo condiciones controladas, lo que permitirá evaluar la efectividad del modelo en situaciones prácticas y ajustar los parámetros según las necesidades específicas de cada aplicación en IoT.

Recomendación 5: Evaluar el Impacto de Diferentes Clases de Ataques.

Investigar cómo el modelo responde a diferentes tipos de ataques y comportamientos anómalos, con diferentes datasets, adaptando la clasificación a un enfoque más flexible que contemple una variedad más extensa de anomalías.

Referencias

1. Alpaydin, E. (2020). *Introduction to Machine Learning*. 4ª ed. The MIT Press. <https://mitpress.ubliish.com/ebook/introduction-to-machine-learning--4-preview/9935/iv>
2. Amazon Web Services. (29 de abril de 2022). *¿Qué es el Internet de las cosas (IoT)?* <https://aws.amazon.com/es/what-is/iot/>
3. Ashton, K. (2009). *That 'Internet of Things' Thing*. RFID Journal. <https://www.rfidjournal.com/articles/view?4986>
4. Asociación Peruana de Telecomunicaciones (10 de julio de 2024). *Análisis del mercado de telecomunicaciones en Perú: Perspectivas 2024-2032*. <https://aptcp Peru.org/2024/07/10/analisis-del-mercado-de-telecomunicaciones-en-peru-perspectivas-2024-2032/>
5. Ávila, F. y Moreno, L. (2023). *Internet de las Cosas (IoT) Retos para las Empresas en la era de la Industria 4.0*. Pádi Boletín Científico de Ciencias Básicas e Ingenierías del ICBI. <https://repository.uaeh.edu.mx/revistas/index.php/icbi/article/view/9516>
6. Baccelli, E. (2021). *Internet de las cosas (IoT) Retos sociales y campos de investigación científica en relación con la IoT*. 1ª ed. Inria.
7. Bazan, A. y Perez, R. (2024). *Método de clasificación de ataques ransomware utilizando algoritmos a través de machine learning*. Universidad Señor de Sipán. <https://repositorio.uss.edu.pe/handle/20.500.12802/12634>
8. Becerra, F., Fernández, I. y Forero, M. (25 de abril de 2024). *Improvement of Distributed Denial of Service Attack Detection through Machine Learning and Data Processing*. Mathematics, 12(9), 1294. <https://www.mdpi.com/2227-7390/12/9/1294>
9. Becerra, F., Tuesta, V., Mejia, H., y Arcila, J. (17 de mayo de 2024). *Performance Evaluation of Deep Learning Models for Classifying Cybersecurity Attacks in IoT Networks*. Informatics. <https://doi.org/10.3390/informatics11020032>
10. Centro Nacional de Planeamiento Estratégico (agosto de 2024). *Desarrollo de sistemas de defensa digital resilientes usando blockchain, IA e IoT*. https://observatorio.ceplan.gob.pe/ficha/o25_2024
11. Datascientest. (20 de abril de 2022). *Machine Learning: definición, funcionamiento, usos*. <https://datascientest.com/es/machine-learning-definicion-funcionamiento-usos>
12. Diccionario sobre inteligencia artificial. (2024). 1 ed. TN Editorial.
13. Glosario de términos de ciberseguridad. (2020). España: INCIBE.
14. Glosario de términos TIC. (2021). 1 ed. UNAM.
15. González, L., Laguía, D., Gesto, E., & Hallar, K. (2020). *Internet del Futuro - Estudio de tecnologías IoT*. Future Internet. <https://dialnet.unirioja.es/descarga/articulo/7756122.pdf>
16. Google Cloud. (2024). *¿Qué es la inteligencia artificial?* <https://cloud.google.com/learn/what-is-artificial-intelligence?hl=es-+%22419+%28%22+site%3Ayoutube.com%29>

17. IBM (31 de enero de 2023) *¿Qué es el análisis exploratorio de datos (EDA)?*
<https://www.ibm.com/mx-es/topics/exploratory-data-analysis>
18. IBM. (2023). *¿Qué es el Internet de las cosas (IoT)?*
<https://www.ibm.com/mx-es/topics/internet-of-things>
19. IBM. (22 de febrero de 2024). *¿Qué es la detección de anomalías?*
<https://www.ibm.com/mx-es/topics/anomaly-detection>
20. ISO. (2024). Aprendizaje automático (AA): todo lo que hay que saber
<https://www.iso.org/es/inteligencia-artificial/aprendizaje-automático>
21. Kaspersky (29 de febrero de 2024). More than half of companies use AI and IoT in their business processes.
<https://www.kaspersky.com/about/press-releases/more-than-half-of-companies-use-ai-and-iot-in-their-business-processes>
22. León, E., Chaffo, F. y Chavez, J. (15 de abril de 2024). *Sistema que alerta la presencia de anomalías en logs de equipos IoT para prevenir el acceso no autorizado*. Universidad Peruana de Ciencias Aplicadas.
<http://hdl.handle.net/10757/675786>
23. Moabits (30 de mayo de 2024). *El Potencial en ventas del mercado de IoT en Latinoamérica: Perspectivas y oportunidad*. Portal especializado en soluciones de conectividad IoT.
<https://www.moabits.com/post/el-potencial-en-ventas-del-mercado-de-iot-en-latinoamérica-perspectivas-y-oportunidad>
24. Muñoz, L. A. (2023). *Sistema de detección de anomalías para Infraestructuras IoT*.
<https://rua.ua.es/dspace/handle/10045/135258>
25. Navarro, S. (17 de abril de 2024). *¿Qué son los datasets?* KeepCoding.
<https://keepcoding.io/blog/que-son-datasets/>
26. Nolasco, P. (marzo de 2023). *Aplicación de Machine Learning para pronóstico de desplazamiento de lluvias usando imágenes del radar de lluvias de UDEP*. Universidad de Piura.
<https://pirhua.udep.edu.pe/item/e9e965e9-3727-4994-a875-61e4b0881a5a>
27. Osio, J., Salvatore, J., Salina, M., Montezanti, D., Denon, N., Doti, S., Olivera, L., Busum Fradera, M., Alonso, D., Cappelletti, M., Encinas, D., & Morales, M. (2021). *Tecnologías de IoT y aprendizaje automático para la solución de problemas en el medio productivo y el cuidado del medioambiente*. <https://sedici.unlp.edu.ar/bitstream/handle/10915/120043/Ponencia.pdf-PDFA.pdf?isAllowed=y&sequence=1>
28. Ouaisa M., Ouaisa M., Boulouard Z., Kumar A., Sharma V., Kaushik K. (2025). *Artificial Intelligence for Blockchain and Cybersecurity Powered IoT Applications*. 1ª ed. CRC Press.
29. Pinto E., Dadkhah, S., Ferreira, R., Zohourian, A., Lu, R. y Ghorbani, A. (26 de junio de 2023). *CICIoT2023: A real-time dataset and benchmark for large-scale attacks in IoT environment*. Sensors. <https://doi.org/10.3390/s23135941>
30. Plata, N. (18 de enero de 2023). *Detección de anomalías físicas en redes IoT empleando técnicas de machine learning*. Universidad de los Andes.
<https://repositorio.uniandes.edu.co/entities/publication/0dacd688-db7f-46e7-a588-ee9fe20525f7>
31. Pypro. (2024). *Modelos de Clasificación*.
<https://www.pypro.mx/app/curso/machine-learning-con-python/modelos-de-clasificacion>

32. Ragab, M., Alshammari, S. , Maghrabi, L., Als Salman, D., Althaqafi, T. y Al-Malaise A. (27 de octubre de 2023). *Robust DDoS attack detection using piecewise Harris Hawks optimizer with deep learning for a secure Internet of Things environment*. Mathematics. <https://doi.org/10.3390/math11214448>
33. Sujay, L. (11 de setiembre de 2024). *Number of Internet of Things (IoT) connections worldwide from 2022 to 2023, with forecasts from 2024 to 2033*. Plataforma global de datos e inteligencia empresarial Statista. <https://www.statista.com/statistics/1183457/iot-connected-devices-worldwide/>
34. Tseng, S., Wang, Y. y Wang Y. (8 de agosto de 2024). *Multi-class intrusion detection based on transformer for IoT networks using CIC-IoT-2023 dataset*. Future Internet. <https://doi.org/10.3390/fi16080284>
35. UIT (7 de enero de 2007). *Internet Reports 2005: The Internet of Things*. Unión Internacional de Telecomunicaciones. <https://www.itu.int/pub/S-POL-IR.IT-2005>
36. Vigoya, L. (enero de 2023). *Aplicación de algoritmos de aprendizaje automático para la detección de anomalías de tráfico en entornos IoT*. Universidade da Coruña. <https://ruc.udc.es/dspace/handle/2183/33306>
37. Villón, M. (16 de mayo de 2024). *Seguridad Nacional, Relaciones Internacionales y Bienestar Social en la Era Digital*. Centro de Estudios Estratégicos del Ejército del Perú. <https://ceeeep.mil.pe/2024/05/16/seguridad-nacional-relaciones-internacionales-y-bienestar-social-en-la-era-digital/#post-25008-endnote-27>

Anexos

Matriz de consistencia

Pregunta general	Preguntas específicas	Objetivo general	Objetivos específicos	Variables	Dimensiones	Enfoque, tipo y diseño	Población y muestra	Técnicas e instrumentos
¿Cómo desarrollar un modelo de clasificación para la detección de anomalías en redes IoT utilizando el dataset CICIoT2023?	¿Cómo analizar y procesar los datos del dataset CICIoT2023?	Desarrollar un modelo de clasificación para la detección de anomalías en redes IoT utilizando el dataset CICIoT2023.	Realizar el análisis exploratorio y el procesamiento de datos del dataset CICIoT2023.	Independiente: Características del dataset	Número de Características	Enfoque Cuantitativo	Población: Todos 63 archivos "csv" del dataset CICIoT2023, que contienen 45 019 243 de filas y 40 columnas, que incluye tráfico normal y anómalo generado por dispositivos IoT. Muestra: Seis Archivos "csv" del dataset, que contienen 7 829 178 filas y 40 columnas para diseñar y evaluar el modelo de aprendizaje automático.	Descarga del Dataset del sitio web del Instituto Canadiense de Ciberseguridad
	¿Cómo diseñar un modelo de clasificación para la detección de anomalías en el dataset CICIoT2023?		Diseñar un modelo de clasificación para la detección de anomalías en el dataset CICIoT2023.		Métricas de tráfico			
			¿Cómo evaluar el modelo de clasificación para la detección de anomalías en el dataset CICIoT2023?		Evaluar el modelo de clasificación para la detección de anomalías en el dataset CICIoT2023.			
		Flags						
				Dependiente: Detección de Anomalías	Clasificación de ataques	Tipo de Investigación Aplicada		
					Tráfico: Malicioso o Benigno.	Alcance Explicativo		
					Métodos de detección: Regresión Logística, Árbol de Decisión, Random Forest, Ada Boost y Perceptron.			

Fuente: Elaboración propia.

Comparación del CICIoT2023 con otros datasets de seguridad IoT existentes

Attack	IoTHIDS	N-BalIoT	Kitsune	IoTNIDS	IoT-SH	BoT-IoT	MedBlot	IoT-23 (2020)	IoTIDS	MQTT	MQTT-IoT-IDS	X-IoTID	WUSTL-IoT	Edge-IoTSet	CICIoT2023
ACK	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
Fragmentation	-	-	-	✓	-	✓	-	-	-	-	-	-	-	✓	✓
UDP Flood	-	✓	-	-	-	-	-	-	-	-	-	-	-	-	✓
SlowLoris	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
ICMP Flood	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
RSTFIN Flood	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
PSHACK Flood	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
HTTP Flood	-	✓	-	✓	-	✓	-	-	-	-	-	-	-	✓	✓
UDP	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
Fragmentation	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
ICMP	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
Fragmentation	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
TCP Flood	-	✓	-	-	-	-	-	-	-	-	-	-	✓	-	✓
SYN Flood	-	✓	-	-	-	-	-	-	-	-	-	-	-	✓	✓
SynonymousIP Flood	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
TCP Flood	-	✓	-	-	✓	✓	-	-	-	-	-	-	-	-	✓
HTTP Flood	-	✓	-	-	-	✓	-	-	✓	-	-	-	-	-	✓
SYN Flood	-	✓	-	✓	-	✓	-	-	✓	-	-	-	-	-	✓
UDP Flood	-	✓	-	-	✓	✓	-	-	✓	-	-	-	-	-	✓
Ping Sweep	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
OS Scan	-	-	✓	✓	✓	✓	-	-	✓	-	✓	✓	-	✓	✓
Vulnerability Scan	-	✓	-	-	-	-	-	-	-	-	✓	✓	-	✓	✓
Port Scan	-	✓	-	✓	✓	✓	-	-	✓	-	✓	✓	-	✓	✓
Host Discovery	-	-	-	✓	-	-	-	-	-	-	-	✓	-	-	✓
Sql Injection	-	-	-	-	-	-	-	-	-	-	-	-	-	✓	✓
Command Injection	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
Backdoor Malware	-	-	-	-	-	-	-	-	-	-	-	✓	-	-	✓
Uploading Attack	-	-	-	-	-	-	-	-	-	-	-	-	-	✓	✓
XSS	-	-	-	-	-	-	-	-	-	-	-	-	-	✓	✓
Browser Hijacking	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
Dictionary	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
Brute Force	-	-	-	✓	-	-	-	-	-	✓	✓	✓	-	✓	✓
Arp Spoofing	-	-	✓	✓	✓	-	-	-	✓	-	-	✓	-	✓	✓
DNS Spoofing	-	-	-	-	-	-	-	-	-	-	-	-	-	-	✓
GREIP Flood	✓	✓	✓	✓	-	-	✓	✓	✓	-	-	-	-	-	✓
Greeth Flood	✓	✓	✓	✓	-	-	✓	✓	✓	-	-	-	-	-	✓
UDP/Plain	✓	✓	✓	✓	-	-	✓	✓	✓	-	-	-	-	-	✓

Fuente: Instituto Canadiense de Ciberseguridad.