

# PROYECTO EDA

31-10-2024

“Estudio del Comportamiento del  
tráfico de datos en la red:  
Distinguiendo entre el tráfico  
normal del malicioso”



JON AMEYUGO  
ALUMNO DE THE BRIDGE

# Introducción

El descubrimiento de anomalías en el tráfico de la red, es un punto clave para la ciberseguridad, porque proporciona la posibilidad de identificar comportamientos sospechosos que podrían indicar intentos de intrusión, ataques de denegación de servicio (DDoS), propagación de malware, etc. En un entorno de red cada vez más complejo y en constante expansión, la eficacia para distinguir el tráfico benigno del malicioso puede suponer la capacidad de prevenir incidentes de seguridad importantes y mejorar la defensa de los sistemas informáticos.

En este análisis exploratorio de datos (EDA) se centra en el estudio de un conjunto de datos de tráfico de red para la identificación de comportamientos anómalos, especialmente para los ataques DDoS. Un ataque DDoS (Distributed Denial of Service o Ataque de Denegación de Servicio Distribuido) es un intento de hacer que un sitio web, servicio o red se vuelva inaccesible. Para lograrlo, el atacante envía una gran cantidad de tráfico (solicitudes o "peticiones") desde muchos dispositivos a la vez, sobrecargando el servidor o sistema objetivo hasta que no puede responder normalmente. El objetivo de este proyecto, es poder determinar el tráfico benigno y el tráfico anómalo en busca de ciertos patrones o de indicadores que nos permitan diferenciar ambos tipos de tráfico, proporcionándonos en el futuro una buena base para la detección de anomalías.

Nuestra **hipótesis principal** se trata de la siguiente:

*“Saber si el tráfico de red etiquetado como ataque presenta patrones y características significativamente distintos al tráfico benigno, lo que permite su identificación mediante el análisis de variables clave.”*

A través de este EDA, se busca analizar la relación que existe entre las diferentes variables del tráfico de red y su relación con comportamientos anómalos, utilizando herramientas gráficas y estadísticas para determinar diferencias relevantes. El estudio que se hará se dedicará a la exploración de distribuciones y en determinar la relación entre variables, a fin de verificar la hipótesis de que el tráfico malicioso presenta unas características diferenciables del tráfico normal.

## Información acerca del Dataset

Para el análisis, se emplea un dataset de tráfico de red sacado del “Instituto Canadiense en Ciberseguridad” creado como experimento de análisis, el cual consta de varias características asociadas con los flujos de red. La columna **Label** presenta la etiqueta **BENIGN** o **DDoS**. Esta variable (Label) facilita la clasificación del tráfico ya sea en tráfico normal (benigno) y tráfico anómalo (DDoS) sirviendo como base para la identificación de comportamientos indeseables.

El dataset consta de 79 columnas y de 225745 elementos lo cual lo hace ser válido y eficaz para el análisis.

Como variables **clave** del dataset se ha supuesto las siguientes:

- **“Flow Duration”**: Que se trata de la duración del flujo.
- **“Total Fwd Packets”** y **“Total Bwd Packets”**: Siendo el número total de paquetes enviados hacia adelante y hacia atrás en el flujo.

- **"Flow Bytes/s"**: Es la tasa de bytes por segundo en el flujo.
- **"Label"**: Variable objetivo que clasifica el tráfico como BENIGN o una etiqueta de ataque específica.
- **"Flow Packets/s"**: Se refiere a la cantidad de paquetes de flujo por segundo.

# Metodología del análisis aplicado

## 1. Introducción

Este análisis se centra en el estudio y diferenciación del comportamiento del tráfico de datos en la red, con el fin de identificar patrones característicos entre tráfico normal y tráfico malicioso, específicamente ataques de denegación de servicio distribuido (DDoS). La hipótesis principal establece que el tráfico identificado como ataque debería mostrar patrones y características diferenciados, lo cual permitiría una identificación precisa de estos eventos en contraste con el tráfico benigno.

## 2. Conjunto de datos

Para este análisis, se utilizó el conjunto de datos *"Friday-WorkingHours-Afternoon-DDos.pcap\_ISCX"* disponible en Kaggle. Este conjunto de datos contiene un total de 225.711 entradas de tráfico de red, cada una caracterizada por 79 variables que reflejan diferentes aspectos y métricas del tráfico. Estas variables incluyen tanto detalles de tiempo y volumen, como aspectos relacionados con los paquetes y flujos, lo que permite un análisis integral de las distintas interacciones de la red.

## 3. Importación de Librerías y Carga de Datos

Se importaron diversas librerías fundamentales para el análisis de datos y la visualización de patrones en el dataset:

- **numpy y pandas**: utilizado para la manipulación de datos, facilitando la carga, limpieza y gestión de las grandes cantidades de datos en el conjunto de datos.
- **matplotlib y seaborn**: empleadas para la visualización de datos, permitiendo una representación gráfica de la distribución y patrones en el tráfico.
- **scipy**: usado específicamente para realizar pruebas estadísticas que permitirán validar las hipótesis planteadas.

El conjunto de datos fue cargado en un DataFrame de pandas para su posterior análisis, y se llevó a cabo una limpieza inicial de los nombres de las columnas, eliminando espacios en blanco y otros caracteres innecesarios que podrían interferir en los procesos de análisis y visualización.

## 4. Limpieza de Datos

Para asegurar la calidad y consistencia de los datos, se eliminaron valores infinitos y valores nulos (NaN) presentes en el conjunto de datos. Esta limpieza es esencial para

evitar sesgos en el análisis y garantizar que las métricas calculadas sean representativas de la muestra total.

## 5. Exploración Inicial de los Datos

Se llevó a cabo una exploración preliminar del conjunto de datos, la cual incluyó:

- **Dimensiones y Estructura de las Columnas:** Se revisó la estructura general del conjunto de datos, confirmando el número de filas y columnas, así como los nombres y tipos de las variables.
- **Clasificación de Tráfico:** Se creó una nueva columna binaria denominada “es\_ataque”, que asigna un valor de 0 para tráfico normal y 1 para tráfico identificado como ataque, facilitando el análisis diferenciado.
- **Distribución del Tráfico:** Se analizó la proporción entre tráfico normal y tráfico anómalo, revelando que el 56.7% de los datos corresponden a ataques DDoS y el restante al tráfico benigno, lo cual es importante para contextualizar el análisis y entender el sesgo potencial en la muestra.

## 6. Análisis de Características Principales

Para profundizar en la diferenciación de tráfico, se seleccionan las características más relevantes para el análisis:

- **Selección de Características Clave:** Entre las variables evaluadas, se eligieron como claves *Flow Duration*, *Flow Packets/s*, *Flow Bytes/s*, *Total Fwd Packets* y *Total Backward Packets*, debido a su capacidad para reflejar patrones en el tráfico.
- **Estadísticas Descriptivas:** Se calcularon estadísticas descriptivas (media, mediana, desviación estándar, entre otras) tanto para el tráfico normal como para el tráfico anómalo, proporcionando una visión comparativa preliminar.
- **Visualización Gráfica:** Mediante gráficos de caja (boxplots) e histogramas, se visualizó la distribución de las características seleccionadas, lo que permitió identificar diferencias clave entre los dos tipos de tráfico.

## 7. Análisis de Correlaciones

Para profundizar en las relaciones entre las variables seleccionadas, se realizó un análisis de correlaciones específicas para tráfico normal y anómalo:

- **Tráfico Normal:** Se calculó la matriz de variación para el tráfico normal, mostrando una fuerte relación positiva entre *Total Fwd Packets* y *Total Backward Packets*, lo cual es coherente con un tráfico regular. Indica una comunicación bidireccional que está balanceada como si se tratase de un tráfico web normal.
- **Tráfico Anómalo:** En el caso del tráfico anómalo, en relación hacia las variables *Total Fwd Packets* y *Total Backward Packets*, estas no presentan una relación dependiente. Son independientes la una con la otra. Se suele producir un desbalance en la dirección de los paquetes.

## 8. Análisis de Patrones Específicos

Para identificar patrones distintivos, se realizaron gráficos de dispersión (diagramas de dispersión) que muestran las relaciones entre las principales características seleccionadas. Esto permitió identificar visualmente diferencias en la distribución y patrones de comportamiento entre tráfico normal y tráfico de ataque, facilitando una comprensión intuitiva de los rasgos que distinguen a los ataques.

## 9. Pruebas Estadísticas

Con el fin de validar estadísticamente las diferencias observadas, se realizaron pruebas de Mann-Whitney para cada una de las características seleccionadas. Estas pruebas me permitieron verificar si las diferencias entre las características del tráfico normal y el tráfico anómalo son estadísticamente significativas. Un valor p menor a 0.05 me indicó que las diferencias son relevantes, confirmando la hipótesis de que las características del tráfico de ataque son diferentes a las del tráfico benigno.

## 10. Aplicación de Resultados

Para explorar en mayor profundidad las relaciones específicas, se analizó la relación entre la duración del flujo (*Flow Duration*) y la clasificación entre *Total Fwd Packets* y *Total Backward Packets* en diferentes intervalos de duración. Los resultados de este análisis se representaron gráficamente, proporcionando una visualización clara de las diferencias de comportamiento entre tráfico normal y tráfico anómalo.

Esta metodología constituye un marco detallado y replicable para el análisis de tráfico en redes, permitiendo la identificación de patrones y características que pueden mejorar significativamente los sistemas de detección de ataques en entornos de ciberseguridad.

# Conclusiones

Se confirma la hipótesis principal: *“Saber si el tráfico de red etiquetado como ataque presenta patrones y características significativamente distintos al tráfico benigno, lo que permite su identificación mediante el análisis de variables clave”*. Esta confirmación se sustenta en múltiples evidencias estadísticas y patrones observados.

### Hallazgos clave detectados:

- **Flow Duration:** Muestra patrones significativamente diferentes entre tráfico normal y ataques
- **Flow Packets/s:** Muestra comportamientos distintivos y predecibles en ataques.
- **Flow Bytes/s:** Presenta correlaciones fuertes con patrones de ataque.
- **Total Forward/Backward Packets:** Revela desequilibrios característicos en ataques.

Este análisis exploratorio de datos (EA) ha cumplido su objetivo principal (el poder resolver la hipótesis principal), proporcionando una comprensión profunda de las diferencias entre el tráfico normal y el tráfico de ataque (DDoS), y estableciendo una base sólida para futuras implementaciones y mejoras en la seguridad de red.