

# Análise Exploratória de Dados de Tráfego de Rede Maliciosa no Dataset IoT-23: CTU-IoT-Malware-Capture-34-1 (Mirai)

Yago P. de Aquino<sup>1</sup>

<sup>1</sup>Instituto de de Ciência Exatas – Universidade Federal de Minas Gerais (UFMG)  
Caixa Postal 31.270-901 – Belo Horizonte – MG – Brazil

**Resumo.** *Este relatório apresenta uma análise exploratória de dados do dataset IoT-23, especificamente do cenário CTU-IoT-Malware-Capture-34-1 contendo 23.145 fluxos de tráfego de rede do malware Mirai. O objetivo é identificar padrões que distinguem tráfego malicioso de benigno através de análises estatísticas e visualizações usando Python e bibliotecas como Pandas, Matplotlib e Seaborn. O trabalho documenta o tratamento dos dados e análises sobre distribuição de ataques, portas visadas, estados de conexão, comportamento temporal e correlações entre variáveis. Os resultados revelaram 91,7% de tráfego malicioso, com predominância de DDoS (67,8%) e C&C (31,6%), concentração de ataques na porta 80, uso da porta 6667 para IRC, e correlações temporais entre atividades maliciosas, fornecendo insights sobre o comportamento do Mirai em ambientes IoT.*

## 1. Introdução

Com o crescimento dos dispositivos IoT, questões de segurança se tornaram críticas devido às vulnerabilidades frequentemente presentes nesses dispositivos. O malware Mirai, descoberto em 2016, exemplifica essa problemática ao ter infectado milhões de dispositivos IoT para formar uma das maiores botnets já registradas, realizando ataques DDoS massivos. Este trabalho apresenta uma análise exploratória de dados (EDA) do dataset IoT-23, especificamente do cenário CTU-IoT-Malware-Capture-34-1, que contém 23.145 fluxos de tráfego de rede capturados durante uma infecção do Mirai. O objetivo é identificar padrões e características distintivas entre tráfego malicioso e benigno, analisando a distribuição dos tipos de ataques (62,19% DDoS, 28,97% C&C, 0,53% PortScan, 8,31% benigno), protocolos utilizados, comportamento temporal e correlações entre variáveis, visando contribuir com insights para detecção de atividades maliciosas em redes IoT.

## 2. Apresentação dos Dados

O conjunto de dados utilizado nesta análise foi obtido da plataforma Kaggle, um repositório online para ciência de dados e aprendizado de máquina. A base pode ser acessada em <https://www.kaggle.com/datasets/agungpambudi/network-malware-detection-connection-analysis/data?select=CTU-IoT-Malware-Capture-34-1conn.log.labeled.csv>.

Ressalta-se que os dados correspondem a um único cenário de captura, sendo que outros cenários estão disponíveis na mesma fonte.

O dataset é composto por 23.145 registros e 23 colunas. A Tabela 1 apresenta os metadados do conjunto de dados.

**Tabela 1. Metadados do dataset CTU-IoT-Malware-Capture-34-1.**

Nome do Campo	Descrição (Tradução)	Tipo de Dado
<b>ts</b>	O timestamp do evento de conexão.	float64
<b>uid</b>	Um identificador único para a conexão.	object
<b>id.orig_h</b>	O endereço IP de origem.	object
<b>id.orig_p</b>	A porta de origem.	int64
<b>id.resp_h</b>	O endereço IP de destino.	object
<b>id.resp_p</b>	A porta de destino.	int64
<b>proto</b>	O protocolo de rede utilizado (ex: 'tcp').	object
<b>service</b>	O serviço associado à conexão.	object
<b>duration</b>	A duração da conexão.	object
<b>orig_bytes</b>	O número de bytes enviados da origem para o destino.	object
<b>resp_bytes</b>	O número de bytes enviados do destino para a origem.	object
<b>conn_state</b>	O estado da conexão.	object
<b>local_orig</b>	Indica se a conexão é considerada local ou não.	object
<b>local_resp</b>	Indica se a conexão é considerada local ou não.	object
<b>missed_bytes</b>	O número de bytes perdidos na conexão.	int64
<b>history</b>	Um histórico dos estados da conexão.	object
<b>orig_pkts</b>	O número de pacotes enviados da origem para o destino.	int64
<b>orig_ip_bytes</b>	O número de bytes IP enviados da origem para o destino.	int64
<b>resp_pkts</b>	O número de pacotes enviados do destino para a origem.	int64
<b>resp_ip_bytes</b>	O número de bytes IP enviados do destino para a origem.	int64
<b>tunnel_parents</b>	Indica se esta conexão faz parte de um túnel.	object
<b>label</b>	Um rótulo associado à conexão (ex: 'Malicious' ou 'Benign').	object
<b>detailed-label</b>	Uma descrição ou rótulo mais detalhado para a conexão.	object

### 3. Tratamento dos Dados

A etapa de pré-processamento, que antecede a Análise Exploratória de Dados (EDA), consistiu primeiramente na substituição do caractere '-' por valores nulos ('NaN'). Este tratamento é fundamental para que os dados ausentes sejam corretamente interpretados por bibliotecas de análise, como o Pandas, viabilizando as operações de filtragem e os cálculos subsequentes.

Posteriormente, os tipos de dados das colunas **duration**, **orig\_bytes** e **resp\_bytes** foram convertidos para formato numérico, adequando-os para a análise quantitativa.

### 3.1. Sanity check dos dados

Nesta etapa, foram realizadas validações para aferir a qualidade dos dados, por meio da análise de valores nulos, duplicatas, cardinalidade, estatísticas descritivas e distribuição das variáveis numéricas.

A Tabela 2 apresenta as colunas que possuem valores ausentes, onde se observa que **local\_orig**, **local\_resp**, **tunnel\_parents** e **detailed-label** são inteiramente compostas por valores nulos.

**Tabela 2. Quantitativo de valores nulos por coluna no dataset.**

Nome da Coluna	Qtd. Nulos	% Nulos
duration	17824	77.01
orig_bytes	17824	77.01
resp_bytes	17824	77.01
service	21298	92.02
local_orig	23145	100.00
local_resp	23145	100.00
tunnel_parents	23145	100.00
detailed-label	23145	100.00

Observa-se que o tráfego malicioso concentra a maior proporção de valores nulos, conforme detalha a Tabela 3. Optou-se pelo uso de percentuais para mitigar a influência do desbalanceamento dos dados nesta análise.

**Tabela 3. Comparativo de valores nulos entre tráfego Maligno e Benigno.**

Nome da Coluna	Qtd. Nulos	% Nulos	Tipo
duration	16904	79.65	Maligno
orig_bytes	16904	79.65	Maligno
resp_bytes	16904	79.65	Maligno
service	19581	92.27	Maligno
local_orig	21222	100.00	Maligno
local_resp	21222	100.00	Maligno
tunnel_parents	21222	100.00	Maligno
detailed-label	21222	100.00	Maligno
duration	920	47.84	Benigno
orig_bytes	920	47.84	Benigno
resp_bytes	920	47.84	Benigno
service	1717	89.29	Benigno
local_orig	1923	100.00	Benigno
local_resp	1923	100.00	Benigno
tunnel_parents	1923	100.00	Benigno
detailed-label	1923	100.00	Benigno

A Tabela 4 apresenta a cardinalidade de cada variável. Visto que 4 colunas foram previamente removidas, a análise subsequente considera um total de 19 colunas.

**Tabela 4. Cardinalidade das variáveis (valores distintos).**

Nome da Coluna	Qtd. Distintos	% Distintos
ts	23145	100.00
uid	23145	100.00
duration	4653	20.10
id.orig_p	4383	18.94
orig_ip_bytes	108	0.47
resp_ip_bytes	62	0.27
orig_pkts	53	0.23
id.resp_h	49	0.21
resp_bytes	43	0.19
resp_pkts	28	0.12
orig_bytes	28	0.12
history	26	0.11
id.resp_p	10	0.04
conn_state	6	0.03
label	4	0.02
service	4	0.02
id.orig_h	2	0.01
proto	2	0.01
missed_bytes	3	0.01

A Tabela 5 apresenta as estatísticas descritivas das variáveis numéricas.

**Tabela 5. Estatística descritiva das variáveis numéricas.**

Coluna	count	mean	std	min	25%	50%	75%	max
duration	5321.0	22.8065	722.5223	0.0005	2.0758	3.1110	3.1537	48976.8191
orig_bytes	5321.0	14788.6846	1036441.4659	0.0	0.0	0.0	62.0	75546624.0
resp_bytes	5321.0	350.4294	5378.2628	0.0	0.0	0.0	243.0	164266.0
missed_bytes	23145.0	2.1271	102.4908	0.0	0.0	0.0	0.0	5792.0
orig_pkts	23145.0	6.3752	178.5487	0.0	0.0	0.0	1.0	18444.0
orig_ip_bytes	23145.0	3664.3116	500376.1653	0.0	0.0	0.0	76.0	76063056.0
resp_pkts	23145.0	0.6110	8.3059	0.0	0.0	0.0	0.0	1070.0
resp_ip_bytes	23145.0	111.2190	2713.0828	0.0	0.0	0.0	0.0	168910.0

Foi realizada também a verificação das categorias presentes nas variáveis qualitativas. Devido à sua grande quantidade, o que tornaria o relatório exaustivo, optou-se por apresentar apenas a Tabela 6 como exemplo. A análise completa de todas as variáveis pode ser acessada no notebook do projeto em referências.

Adicionalmente, foram criadas duas colunas para rotular os dados e facilitar as análises: **class** e **class\_name**. A coluna **class** é numérica, utilizando 0 para tráfego benigno e 1 para malicioso. A coluna **class\_name** segue a mesma lógica, porém com os rótulos categóricos "Benigno" e "Maligno".

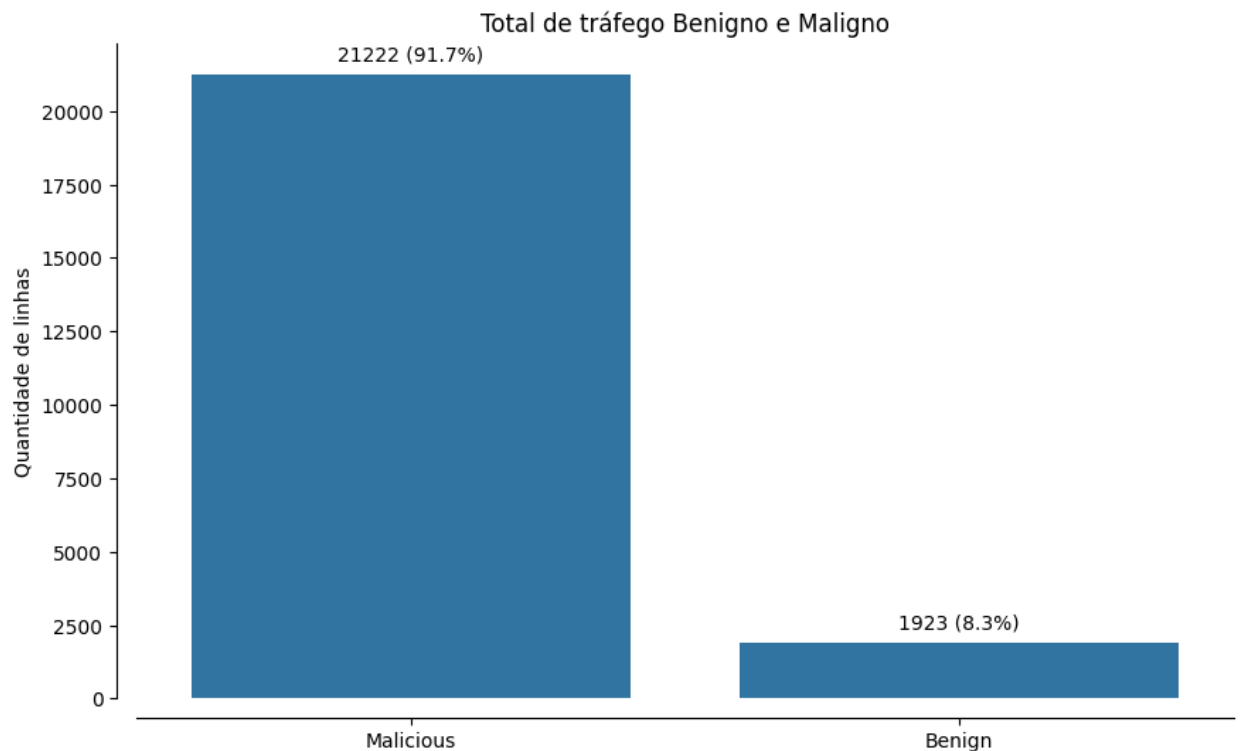
**Tabela 6. Distribuição de frequência da variável `detailed-label`.**

Label	Frequência Absoluta	Frequência Relativa (%)
DDoS	14394	62.19
C&C	6706	28.97
Benign	1923	8.31
PartOfAHorizontalPortScan	122	0.53

#### 4. Análises Exploratórias

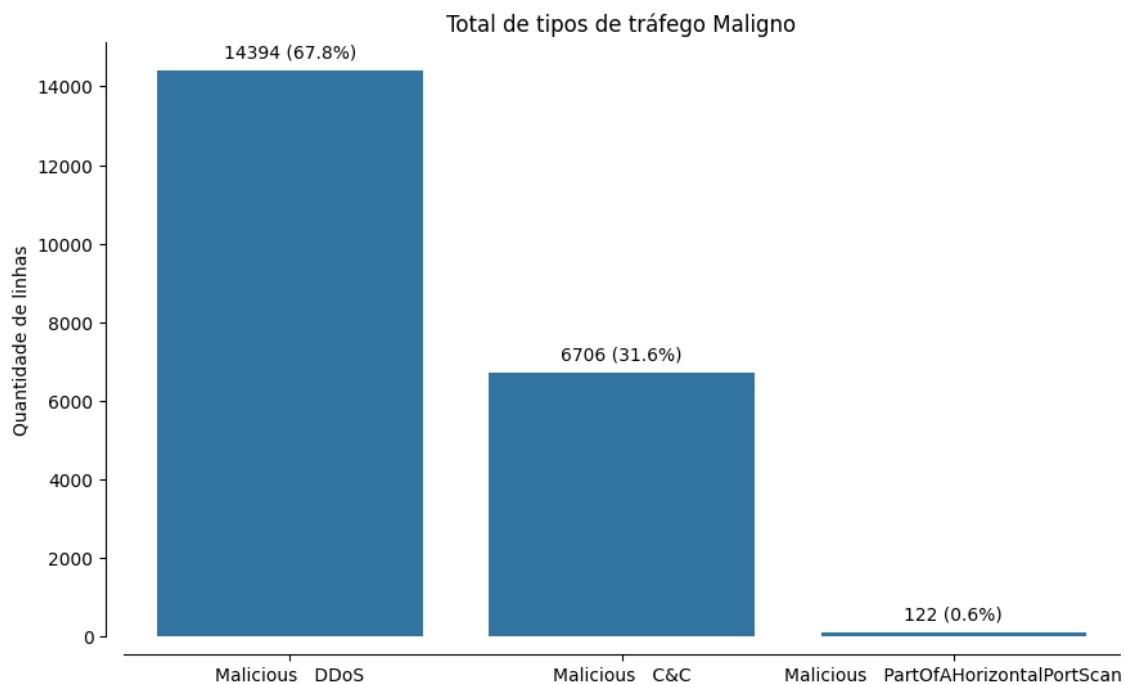
A análise exploratória iniciou-se com a verificação da distribuição dos dados entre as classes benigna e maligna, seguida pela análise da proporção de cada tipo de ataque na classe maligna. Adicionalmente, foi examinada a distribuição das variáveis numéricas para ambas as classes, a fim de identificar diferenças significativas entre elas.

A Figura 1 ilustra a contagem total de registros benignos e malignos no dataset.

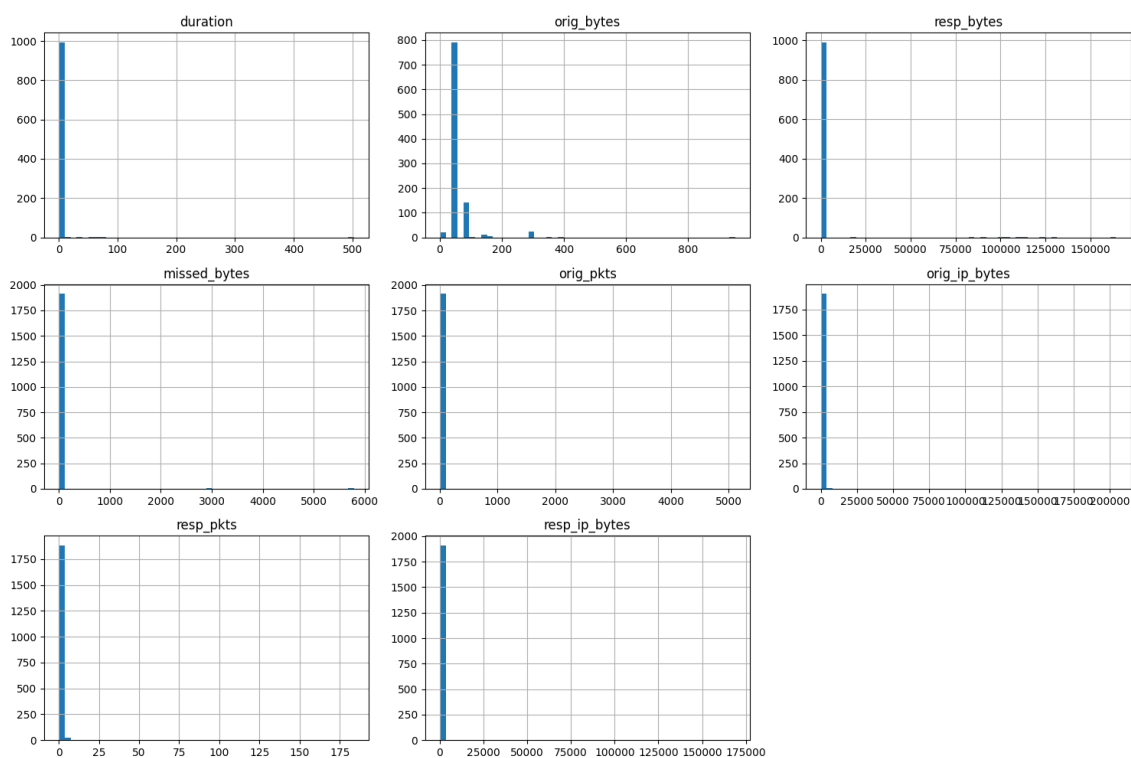


**Figura 1. Total de tráfego Benigno e Maligno**

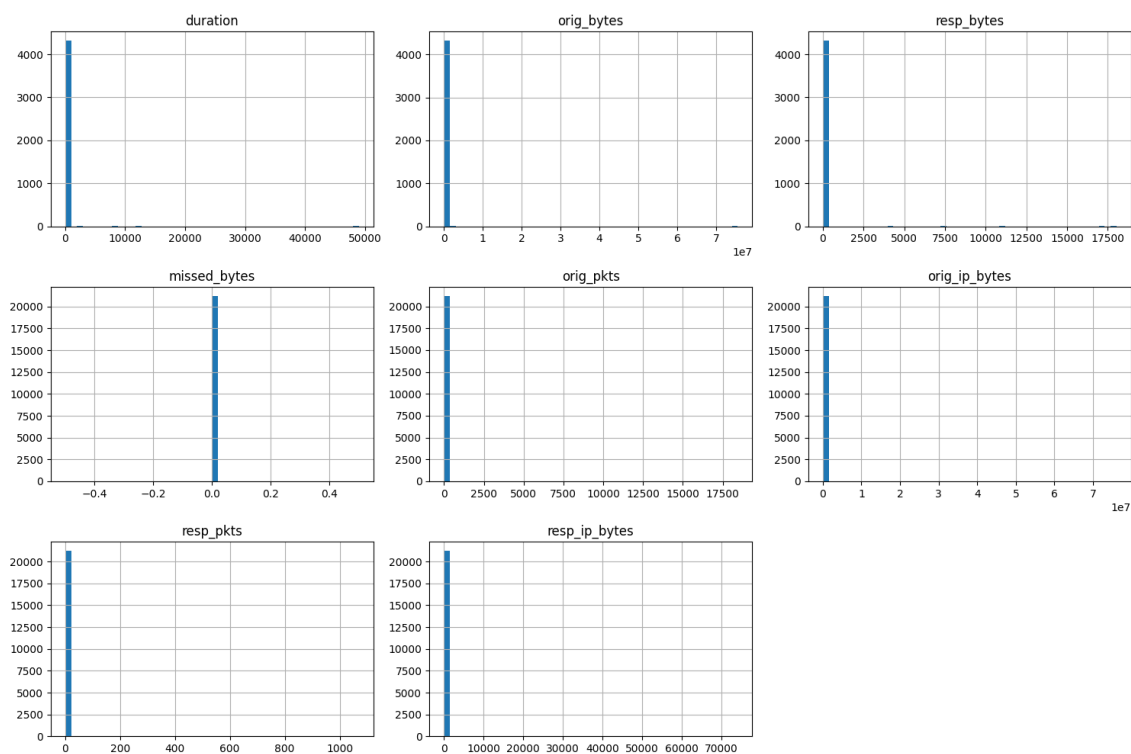
A Figura 2, por sua vez, apresenta a distribuição dos ataques dentro da classe maligna. As Figuras 3 e 4 exibem as distribuições das variáveis para o tráfego benigno e malicioso, respectivamente. Uma das diferenças mais notáveis ocorre na variável **orig.bytes**: o tráfego benigno registra concentrações de bytes em valores mais altos, enquanto o tráfego malicioso demonstra uma forte concentração próxima a zero.



**Figura 2. Total de cada tipo de tráfego maligno**



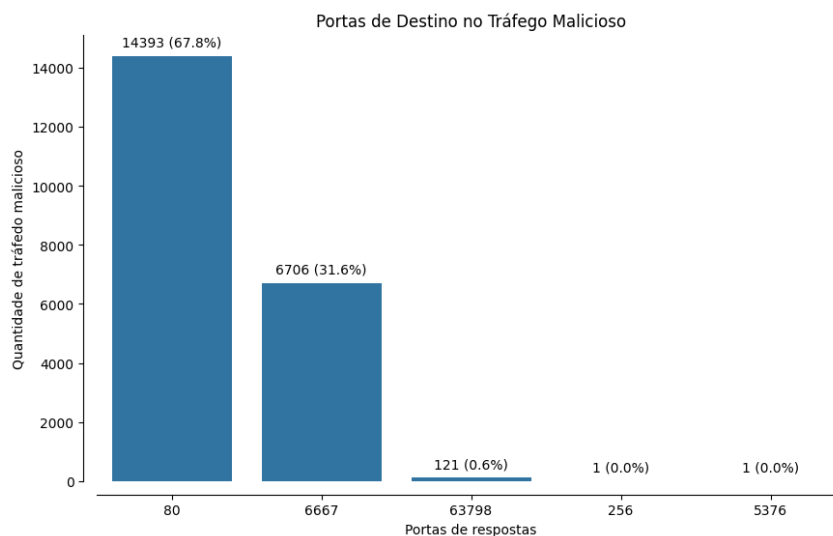
**Figura 3. Distribuição das colunas numéricas de dados benignos**



**Figura 4. Distribuição das colunas numéricas de dados malignos**

#### 4.1. Portas e serviços mais visadas no tráfego maligno e benigno

A análise da coluna **service** no tráfego malicioso revela a presença de duas categorias principais: valores ausentes e IRC. Essa ocorrência está diretamente relacionada às táticas de ataque: os fluxos de DDoS, por serem tentativas de conexão TCP/UDP sem dados de aplicação, resultam em valores ausentes, enquanto a comunicação de Comando e Controle (C&C) utiliza o protocolo IRC para se conectar aos seus servidores.



**Figura 5. Portas de destino no tráfego maligno**

A análise das portas de destino, apresentada na Figura 5, indica que a porta 80 é o principal alvo dos ataques. Essa preferência é justificada por dois fatores: a predominância de ataques DDoS no conjunto de dados e as características intrínsecas da porta 80, que a tornam amplamente exposta e raramente bloqueada. Adicionalmente, observa-se a relevância da porta 6667, que está associada aos ataques de C&C devido à sua utilização em serviços IRC. As características que tornam a porta 80 um alvo comum são:

- Utilização por praticamente todos os dispositivos e serviços web.
- Alta probabilidade de estar aberta na maioria das redes.
- Uso frequente para interfaces de gerenciamento web em dispositivos IoT, como câmeras, roteadores e smart TVs.

A Figura 6 corrobora essas observações, ao mostrar que os ataques estão concentrados (DDoS na porta 80 e C&C na porta 6667). A Tabela 7 complementa a análise, apresentando a distribuição percentual dos ataques para cada porta.

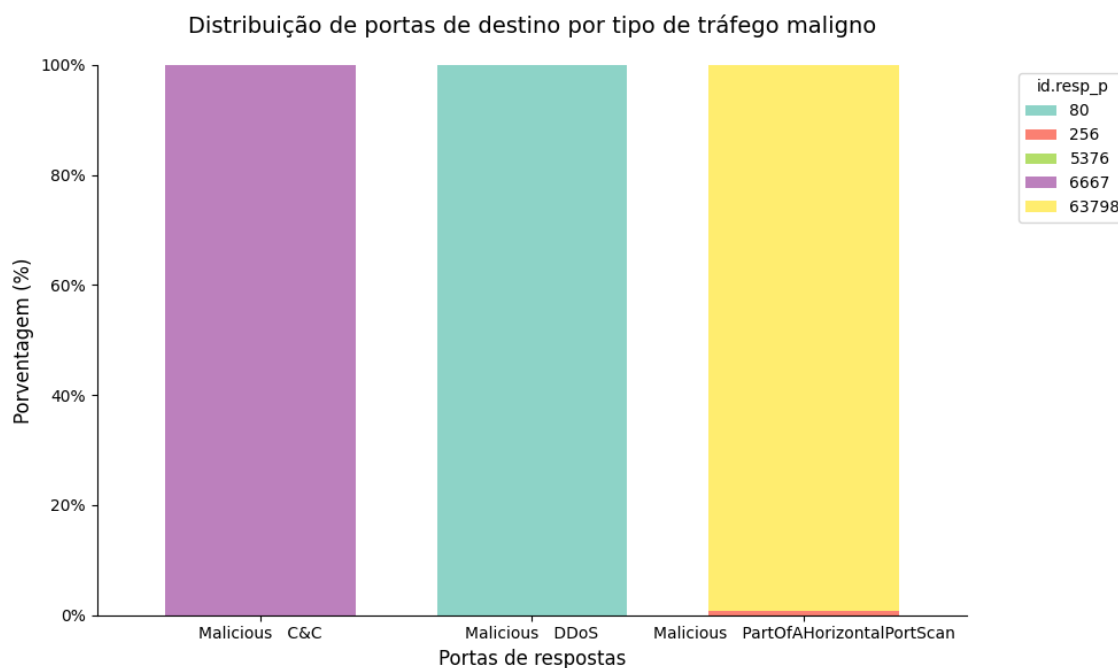


Figura 6. Distribuição de portas de destino por tipo de tráfego maligno

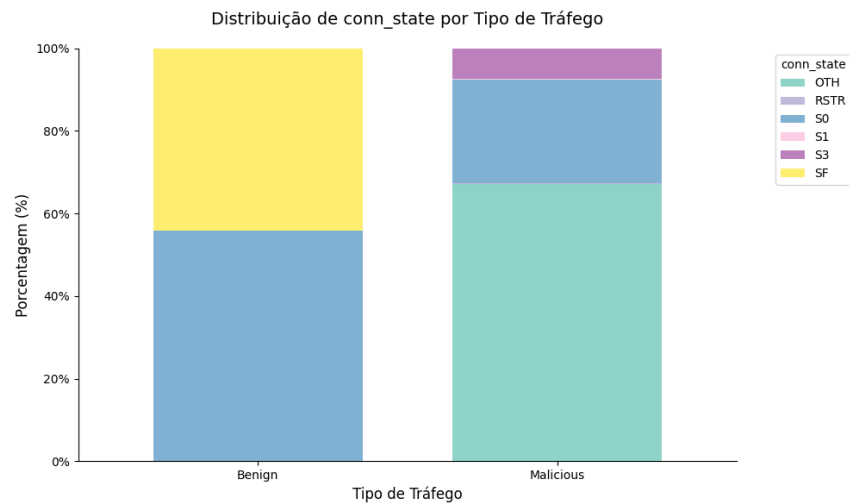
Tabela 7. Distribuição percentual dos labels por porta de destino (id.resp\_p).

Label	1	22	53	67	80	123	256	5376	6667	63798
Benign	0.2081	0.0520	12.1165	0.1040	3.3801	84.1394	0.0000	0.0000	0.0000	0.0000
Malicious C&C	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	100.0000	0.0000
Malicious DDoS	0.0000	0.0000	0.0000	0.0000	99.9931	0.0000	0.0000	0.0069	0.0000	0.0000
Malicious PartOfAHorizontalPortScan	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.8197	0.0000	0.0000	99.1803

#### 4.2. Análise dos estados da conexão

Ao analisar as conexões conseguimos obter padrões conforme a Figura 7 apresenta. O tráfego malicioso tem sua conexão concentrada em S0, OTH e S3. O Benigno está concentrado em S0 e SF, ter o S0 quando não se tem um ataque é esperado visto que são tentativas de conexão normais que não obtiveram resposta.





**Figura 7. Distribuição de conn\_state por tipo de tráfego**

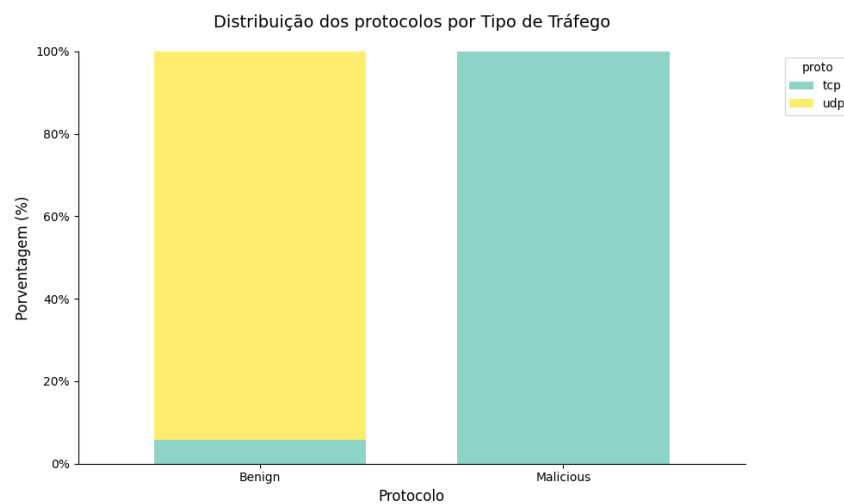
A Tabela 8 apresenta os percentual em cada estado de conexão.

**Tabela 8. Distribuição percentual da variável conn\_state por classe.**

Classe	OTH	RSTR	S0	S1	S3	SF
Benign	0.0000	0.0000	55.8502	0.0000	0.0000	44.1498
Malicious	67.1944	0.2497	25.0730	0.0188	7.4640	0.0000

### 4.3. Análise dos protocolos

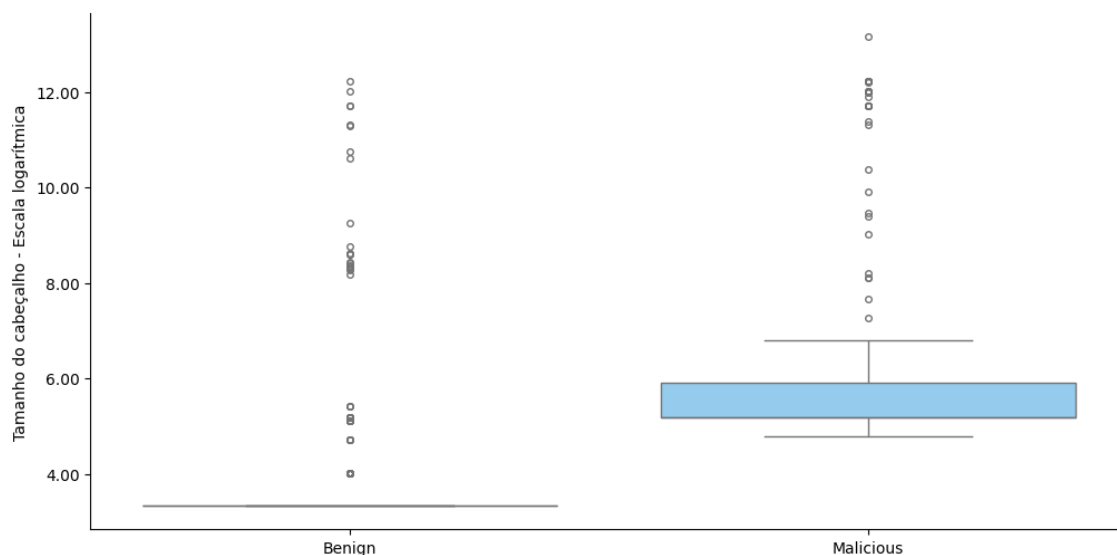
A análise dos protocolos de rede demonstra que o tráfego malicioso se concentra predominantemente em TCP. Tal comportamento é esperado, dada a natureza dos ataques de DDoS e C&C presentes no dataset, que majoritariamente utilizam este protocolo. A Figura 8 ilustra a distribuição de protocolos para o tráfego benigno e malicioso.



**Figura 8. Distribuição dos protocolos por tipo de tráfego**

#### 4.4. Análise do volume de dados

O tamanho do cabeçalho do pacote é uma métrica relevante para a detecção de ataques, que frequentemente manipulam essa estrutura. Nesta análise, o tamanho do cabeçalho foi calculado pela diferença entre **orig\_ip\_bytes** e **orig\_bytes**. O resultado confirmou a hipótese de que o tráfego malicioso apresenta cabeçalhos maiores, conforme ilustra o boxplot na Figura 9, que utiliza uma escala logarítmica para melhor visualização.



**Figura 9. Boxplot do tamanho do cabeçalho em escala logarítmica**

A Tabela 9 corrobora essa descoberta, apresentando em escala linear os valores reais da média e da mediana do tamanho do cabeçalho. Observa-se que ambas as métricas são superiores para o tráfego malicioso.

**Tabela 9. Comparativo da média e mediana do tamanho do cabeçalho**

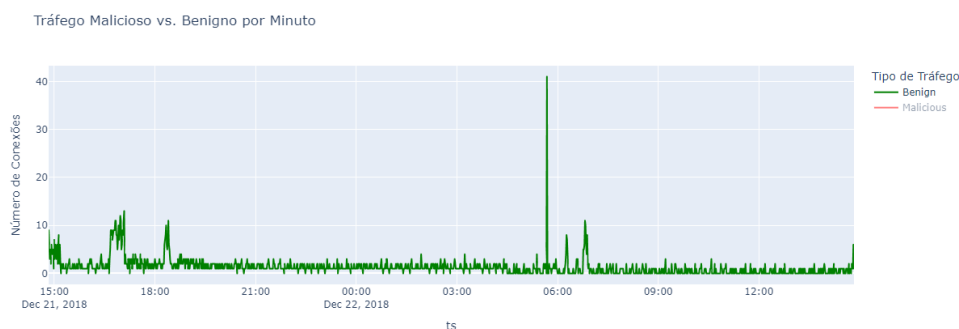
Tipo	Média do Cabeçalho	Mediana do Cabeçalho
Maligno	1141.62	180.0
Benigno	965.71	28.0

Por outro lado, a análise do volume de pacotes de origem e de resposta mostrou-se inconclusiva, devido à alta incidência de valores ausentes e à concentração de dados em torno de zero, o que impediu a extração de padrões claros.

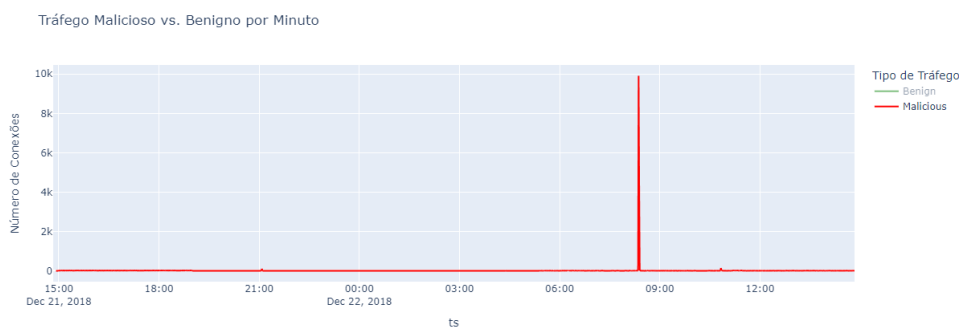
#### 4.5. Análise temporal

A análise temporal dos dados revela padrões de atividade distintos, com destaque para os picos de tráfego que caracterizam os ataques DDoS. As Figuras 10 e 11 apresentam o comportamento temporal dos fluxos de tráfego benignos e malignos, respectivamente.

A análise focada no tráfego malicioso demonstra que o pico principal de atividade corresponde a um ataque DDoS, enquanto o segundo maior é decorrente de operações de PortScan.



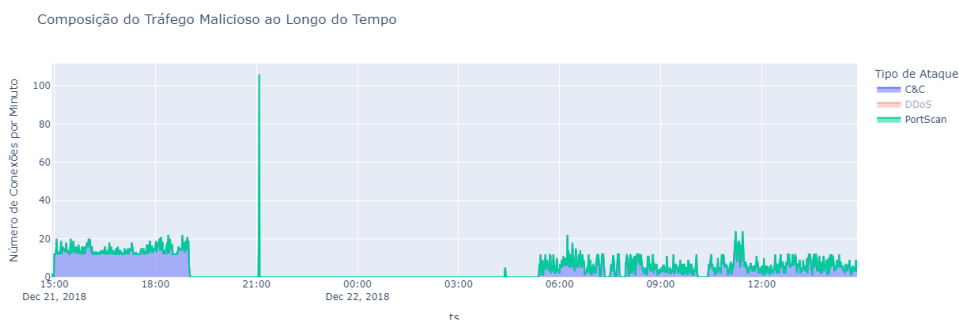
**Figura 10. Tráfego malicioso vs benigno por minuto, filtro benigno**



**Figura 11. Tráfego malicioso vs benigno por minuto, filtro maligno**

Observa-se uma forte correlação temporal entre o tráfego de C&C e as atividades de PortScan, que apresentam padrões sincronizados ao longo do tempo. Essa sincronia é um reflexo do ciclo operacional do Mirai: os bots recebem comandos via C&C para executar varreduras de portas, com o objetivo de identificar novos dispositivos vulneráveis e reportar os resultados aos servidores de comando e controle.

A Figura 12 evidencia esta correlação, especialmente no período entre 15:00 e 18:00, onde se observa atividade simultânea de ambos os tipos de tráfego, com um pico às 21:00 que sugere um comando massivo de varredura. Em contraste, o tráfego DDoS apresenta um padrão isolado de ataques intensos e de curta duração. Este comportamento é característico da fase de execução dos objetivos maliciosos, diferindo das fases de manutenção e expansão da botnet, representadas pelas atividades de C&C e PortScan.



**Figura 12. Tráfego malicioso por minuto - filtro C&C e PortScan**

A Figura 13, por sua vez, exibe a atividade temporal do tráfego malicioso quando

filtrado exclusivamente para ataques do tipo DDoS.

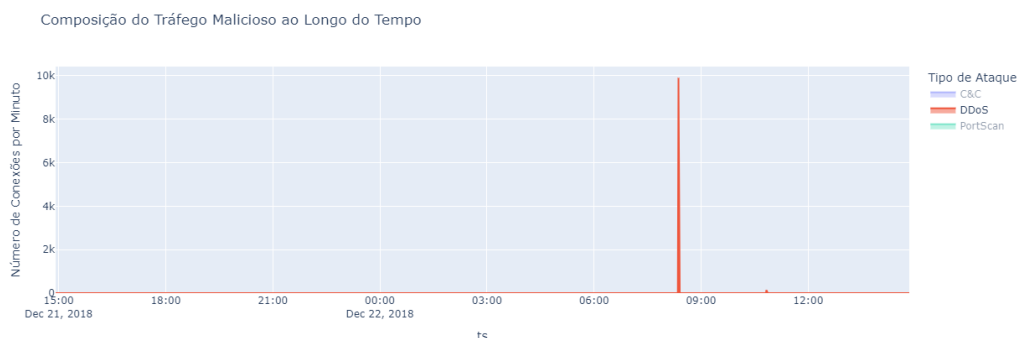


Figura 13. Tráfego malicioso por minuto - filtro DDoS

#### 4.6. Análise da correlação das variáveis numéricas

Adicionalmente, foi realizada uma análise de correlação de Spearman para investigar a relação entre as variáveis numéricas. A escolha deste método justifica-se por sua robustez a outliers e por não pressupor uma relação linear entre os dados. A matriz de correlação resultante é apresentada na Figura 14. A análise da matriz evidencia uma forte correlação da variável alvo **class** com as variáveis **header\_size** e **duration**.

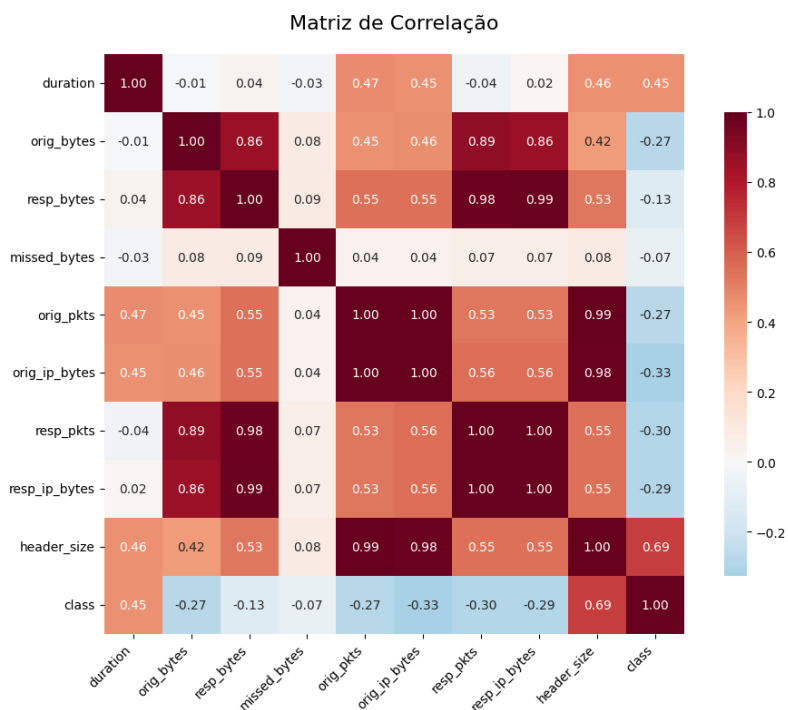


Figura 14. Matriz de correlação entre as variáveis numéricas

A diferença no comportamento da variável **duration** entre os tipos de tráfego, que valida a correlação previamente identificada, é detalhada na Figura 15 e na Tabela 10. Ressalta-se que a figura utiliza uma escala logarítmica para facilitar a visualização, enquanto a tabela apresenta os dados em escala linear para exibir os valores reais das métricas.

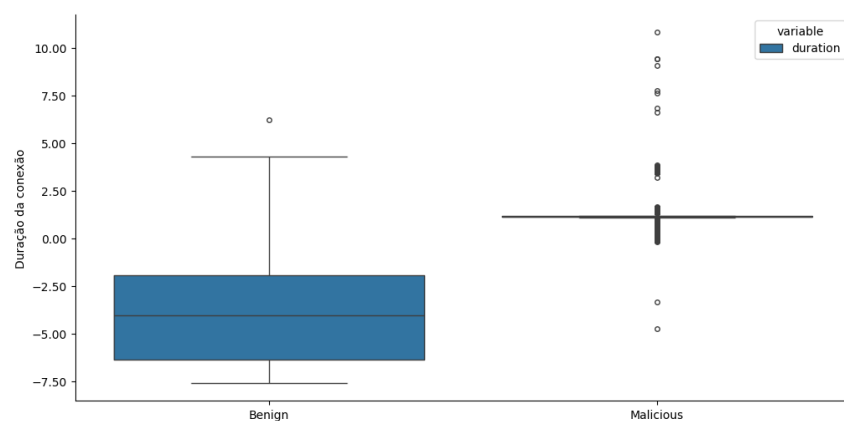


Figura 15. Boxplot da duração em escala logarítmica

Tabela 10. Comparativo da média e mediana da duração da conexão por classe.

Tipo	Média	Mediana
Benigno	1.94	0.02
Maligno	27.65	3.12

#### 4.7. Análise da correlação das variáveis categóricas

Por fim, foi realizado um teste de dependência para avaliar a associação entre as variáveis categóricas. A variável alvo para esta análise foi a coluna **class\_name**, criada para agrupar todos os ataques sob o rótulo **Maligno** em contraste com o tráfego **Benigno**. Os resultados deste teste são apresentados na Figura 16.

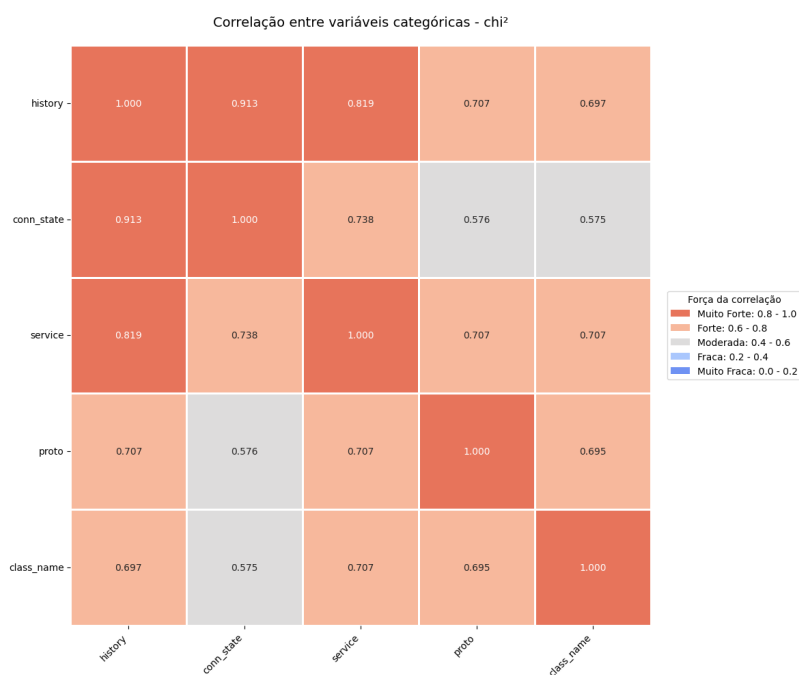


Figura 16. Matriz de teste de dependência de variáveis categóricas

A análise de correlação entre as variáveis categóricas revela que aquelas com maior associação à classe alvo seguem dois padrões distintos. O primeiro padrão é exemplificado pelas variáveis **service** e **proto**, cujas correlações elevadas derivam de características específicas: a variável **service** apresenta predominantemente valores nulos ou o serviço IRC no tráfego malicioso, enquanto a variável **proto** se concentra exclusivamente no protocolo TCP para os ataques.

A variável **history**, em contrapartida, exibe um comportamento distinto e mais informativo. Sua correlação não decorre de valores nulos ou de concentração categórica, mas sim do comprimento da string que registra os eventos TCP. Conexões benignas tendem a apresentar strings mais longas, enquanto ataques maliciosos resultam em strings curtas, que refletem tentativas de conexão não concluídas.

## 5. Conclusão

Esta análise exploratória do dataset IoT-23, especificamente do cenário CTU-IoT-Malware-Capture-34-1, revelou padrões para a compreensão do comportamento de tráfego malicioso em redes IoT. A investigação de 23.145 fluxos de rede permitiu identificar características que diferenciam claramente o tráfego benigno (8,3%) do malicioso (91,7%), com predominância de ataques DDoS (67,8%) e comunicações C&C (31,6%).

Os padrões identificados demonstram o modus operandi característico do Mirai: concentração massiva de ataques na porta 80 devido à sua "onipresença" em dispositivos IoT, uso da porta 6667 para comunicações C&C via IRC, e correlação temporal entre atividades de comando e controle com operações de PortScan, revelando o ciclo de propagação da botnet. A análise revelou ainda que o tráfego malicioso apresenta características quantificáveis distintas, como maior tamanho de cabeçalho, conexões com durações mais longas, e estados de conexão predominantemente incompletos (S0, OTH), refletindo tentativas massivas de conexão não estabelecidas.

A análise temporal apresenta o comportamento coordenado dos ataques, com picos sincronizados de C&C e PortScan, particularmente entre 15:00 e 18:00, seguidos de ataques DDoS concentrados, demonstrando as fases distintas de operação: reconhecimento, infecção e ataque. As correlações identificadas, especialmente entre as variáveis **duration**, **header\_size** e a classificação do tráfego, fornecem indicadores robustos para sistemas de detecção. Este estudo, embora limitado a um único cenário do dataset, contribui com insights valiosos para o desenvolvimento de mecanismos de detecção de malware IoT, destacando a importância de monitorar não apenas assinaturas de ataques conhecidos, mas também padrões comportamentais anômalos na rede.

## Referências

- Ben, I. and Dima (2016). Breaking Down Mirai: An IoT DDoS Botnet Analysis. <https://www.imperiva.com/blog/malware-analysis-mirai-ddos-botnet/>. Acessado em: 31-08-2025.
- de Aquino, Y. P. (2025). Repositório do trabalho . [https://github.com/Yago-Pacheco/ufmg/blob/main/especializacao\\_ciencia\\_de\\_dados/topicos\\_em\\_ciencia\\_de\\_dados/tp\\_1/TP\\_01\\_Ciberseguranca.ipynb](https://github.com/Yago-Pacheco/ufmg/blob/main/especializacao_ciencia_de_dados/topicos_em_ciencia_de_dados/tp_1/TP_01_Ciberseguranca.ipynb). Acessado em: 31-08-2025.

Garcia, S., Parmisano, A., and Erquiaga, M. J. (2020). IoT-23: A labeled dataset with malicious and benign IoT network traffic (Version 1.0.0).