

A DDOS general-purpose classifier

Progetto di Manutenzione Preventiva

Davide De Zuane, Rahmi El Mechri, Roberto Mustillo

May 2024

1 Introduzione

Il progetto è stato avviato a seguito delle riflessioni presentate nell'articolo [4], che mette in luce la mancanza di generalità nei classificatori attuali e l'importanza della lunghezza delle sequenze nell'analisi degli attacchi DDoS.

L'obiettivo principale di questo lavoro è dunque sviluppare un classificatore per gli attacchi DDoS che sia il più generale possibile rispetto al dataset utilizzato. Per cercare di raggiungere tale obiettivo abbiamo cercato di escludere quelle che sono le caratteristiche specifiche della rete e seguendo le indicazioni che emergono dalla letteratura scientifica su come raggiungere tale scopo. Nella relazione che segue, descriveremo i vari passaggi del nostro percorso e analizzeremo i risultati ottenuti. Il codice sorgente relativo a questo lavoro è disponibile in una repository pubblica [2].

Il lavoro svolto è stato suddiviso nelle seguenti fasi:

- **Dataset analysis:** in questa fase analizziamo i dataset forniti per comprendere la loro struttura e contenuto.
- **Initial Feature Selection:** selezioniamo inizialmente un sottoinsieme di caratteristiche che riteniamo potenzialmente utili per il classificatore. Questa selezione si basa su considerazioni preliminari, come la pertinenza delle feature rispetto al problema da risolvere, l'intuizione o la conoscenza preliminare del dominio.
- **Feature Generation:** in questa fase, creiamo nuove feature a partire da quelle esistenti con l'obiettivo di migliorare la rappresentazione dei dati per facilitare l'apprendimento del modello.
- **Data preprocessing:** comprende varie attività per rendere i dati pronti per l'analisi, come per esempio la feature encoding e la gestione dei valori nulli.
- **Feature Selection:** eseguiamo una selezione delle feature più rilevanti tramite l'utilizzo di tecniche (ANOVA), l'importance nel caso del random forest. L'obiettivo è ridurre il numero di feature eliminando quelle meno utili, migliorando così l'efficienza e le prestazioni del modello.
- **Training :** addestriamo il classificatore utilizzando i dati selezionati e utilizziamo diverse strategie per mitigare il problema dell'overfitting e a fornire una stima più accurata delle prestazioni del modello.
- **Testing e Confronto:** andiamo a valutare le prestazioni reali. Questa fase include il confronto dei risultati ottenuti con diversi modelli o con diverse configurazioni dello stesso modello.

2 Dataset Analysis

Tutte le informazioni riguardo il dataset utilizzato sono disponibili al seguente link, ne riportiamo le caratteristiche principali.

2.1 Generazione

Questo dataset è stato generato in uno scenario di automazione di processo su piccola scala utilizzando apparecchiature MODBUS/TCP, per la ricerca sull'applicazione delle tecniche di machine learning alla sicurezza informatica nei Sistemi di Controllo Industriale. Il banco di prova emula un processo CPS (Cyber-Physical System) controllato da un sistema SCADA che utilizza il protocollo MODBUS/TCP. Esso consiste in una pompa di liquido simulata da un motore elettrico controllato da un inverter di frequenza variabile (che permette diverse velocità del rotore), a sua volta controllato da un Controllore Logico Programmabile (PLC). La velocità del motore è determinata da una serie di soglie di temperatura del liquido predefinite, la cui misurazione è fornita da un dispositivo MODBUS Remote Terminal Unit (RTU) che fornisce un misuratore di temperatura, simulato da un potenziometro collegato a un Arduino. Il PLC comunica orizzontalmente con l'RTU, fornendo informazioni preziose su come questo tipo di comunicazioni possa influenzare l'intero sistema. Il PLC comunica anche con l'interfaccia uomo-macchina (HMI) che controlla il sistema.

2.2 Structure

La struttura del dataset è mostrata nel seguente albero e all'interno di ogni cartella abbiamo dei file pcap che adottano la seguente convenzione per il naming:

```
1 <capture interface>dump-<attack>-<attack subtype>-<attack
  duration>-<capture duration>
```

```
captures1/
|-- clean/
|-- mitm/
|-- modbusQuery1/
|-- modbusQuery2/
|-- pingFloodDDoS/
'-- tcpSYNFloodDDoS/

captures2/ & captures3/
|-- modbusQueryFlooding/
|-- pingFloodDDoS/
'-- tcpSYNFloodDDoS/
```

2.3 Osservazioni

Intuitivamente i risultati migliori si ottengono andando a considerare dataset in cui la presenza di pacchetti malevoli è maggiore in relazione alla cattura totale. Questa osservazione è confermata dai risultati riportati nel paper di riferimento [4], che evidenzia come le migliori performance dei classificatori siano associate a dataset contenenti attacchi della durata di circa 15 minuti rispetto ai 30 minuti di cattura.

In base a queste considerazioni, abbiamo deciso di focalizzare il nostro studio su questa specifica categoria di dataset. Pertanto, tutte le fasi successive della relazione fanno riferimento ai dataset strutturati con una durata di attacco di 15/30 minuti. Questo approccio ci ha permesso di garantire che le condizioni di test e valutazione fossero allineate con le configurazioni che hanno dimostrato di offrire le migliori performance nel contesto della rilevazione e classificazione degli attacchi.

3 Initial Feature Selection

In primo luogo abbiamo utilizzato la nostra conoscenza riguardo le reti per estrarre quelle che, secondo noi tra le informazioni contenute all'interno del pcap, sono le caratteristiche da considerare.

In particolare per il classificatore le feature più importanti che abbiamo deciso di considerare, e che poi utilizzeremo in fase di feature generation, sono quelle riportate in *Tabella 1*.

Feature	Description
Time	Il timestamp in cui il pacchetto è stato catturato.
Source	L'indirizzo IP del mittente del pacchetto.
Destination	L'indirizzo IP del destinatario del pacchetto.
Protocol	Il protocollo utilizzato nel pacchetto (ad es., TCP, UDP).
Length	La lunghezza del pacchetto.
S-Port	Il numero di porta sorgente utilizzato nel pacchetto.
D-Port	Il numero di porta di destinazione utilizzato nel pacchetto.
Delta-Time	La differenza di tempo tra pacchetti consecutivi.
ACK	Flag di riconoscimento che indica la ricezione di un pacchetto.
SYN	Flag di sincronizzazione usato per iniziare una connessione.

Table 1: Feature iniziali

4 Feature Generation

Basandoci sulla letteratura presente, riferendoci in particolare a [1], ed analizzando i problemi specifici che stiamo prendendo in considerazione, abbiamo generato le seguenti features:

- **Entropia source IP:** è una feature che misura la diversità degli indirizzi IP di origine all'interno di un determinato intervallo temporale. L'entropia, in questo contesto, quantifica il livello di dispersione o casualità delle sorgenti di traffico.
- **Entropia source Port:** è una misura della diversità delle porte di origine utilizzate nei pacchetti di rete all'interno di una determinata finestra temporale. L'entropia, in questo contesto, quantifica il livello di casualità o dispersione delle porte di origine.
- **Entropia destination IP:** misura la diversità degli indirizzi IP di destinazione all'interno di un periodo di tempo specifico. Questa metrica riflette quanto il traffico di rete sia distribuito verso diverse destinazioni.
- **Entropia destination Port:** misura la diversità delle porte di destinazione all'interno di una serie temporale. Questa metrica riflette quanto il traffico di rete sia distribuito verso diverse porte di destinazione. Un alto valore di destination port entropy suggerisce che il traffico è diretto verso una varietà di porte di destinazione, indicando la possibile presenza di vari servizi o una scansione di porte. Un basso valore, invece, indica che il traffico è concentrato su poche porte.
- **Packet rate:** misura il numero di pacchetti trasmessi in un'unità di tempo specifica all'interno della serie temporale. Questo valore fornisce un'indicazione del volume di traffico di rete, permettendo di quantificare l'intensità delle comunicazioni in un determinato intervallo temporale.
- **Flow packet number:** è una feature che rappresenta il numero totale di pacchetti associati a un flusso di rete specifico all'interno di un intervallo temporale definito. Un flusso di rete è comunemente identificato dall'insieme di parametri come l'indirizzo IP sorgente, l'indirizzo IP di destinazione, la porta sorgente, la porta di destinazione e il protocollo utilizzato. Questa feature quantifica il volume di pacchetti che costituiscono il flusso, offrendo una panoramica sul traffico generato da una particolare connessione o interazione.
- **Byte rate:** misura il volume totale di dati (in byte) trasmessi per unità di tempo all'interno di una serie temporale. Questa metrica indica la quantità di dati scambiati in un determinato intervallo temporale, fornendo una chiara indicazione del carico di traffico di rete.

- **Rapporto Syn/Ack:** misura il rapporto tra il numero di pacchetti SYN (synchronize) e i pacchetti ACK (acknowledge) osservati in una serie temporale. Un rapporto elevato di SYN rispetto agli ACK può indicare situazioni anomale come tentativi di attacchi SYN flood, mentre se è bilanciato rispecchia il normale funzionamento della rete.
- **Protocol distribution:** si tratta di una serie di feature che descrivono la distribuzione percentuale del traffico in termini di protocolli utilizzati all'interno di una serie temporale. Nello specifico, indica la proporzione del traffico che utilizza TCP, Modbus o Ping. Un improvviso aumento del traffico di un certo protocollo che potrebbe indicare un'attività malevola.

Alcune di queste feature sono calcolate tenendo conto di tutta la storia del dataset, come ad esempio il flow packet number, mentre altri sono calcolati in fase di creazione delle sequenze, come dimensione media pacchetti e packet rate. Questo tipo di dato è quindi dipendente dalla lunghezza della sequenza.

Il codice per l'estrazione delle feature si trova in appendice.

5 Data Preprocessing

In questa sezione analizziamo le tecniche di Data Preprocessing applicate per risolvere alcuni dei classici problemi presenti nel Machine Learning, così da garantire una corretta classificazione.

5.1 Normalizzazione

Per loro natura i classificatori KNN e SVC risultano sensibili alla scala delle feature, ciò ha richiesto l'utilizzo di una tecnica di normalizzazione: il *Min-Max Scaler*. Tale tecnica va a normalizza le feature trasformandole in un intervallo, tipicamente tra 0 e 1, assicurando che tutte contribuiscano equamente alla classificazione.

$$x' = \frac{x - x_{min}}{x_{max} - x_{min}}$$

- x è il valore originale della feature
- x_{min} e x_{max} rispettivamente il valore max e min della feature nel dataset
- x' è il valore scalato che sarà compreso in $[0, 1]$

Al contrario, Random Forest e Decision Tree non necessitano di normalizzazione poiché non sono influenzati dalla scala delle feature, basando le loro decisioni su soglie di valore.

5.2 Gestione valori nulli

Durante la fase di processamento dei dati, non è stata necessaria la gestione dei valori nulli, poiché abbiamo progettato la fase di feature generation in modo tale da prevenire la presenza di dati mancanti.

5.3 Creazione tensori

Per la classificazione di serie temporali con HYDRA, i dati devono essere trasformati in un formato adatto alle convoluzioni 1D. Prima, i dati vengono suddivisi in finestre temporali di lunghezza fissa, creando sequenze. Poi, queste sequenze sono convertite in tensori tridimensionali utilizzando PyTorch: Batch Size (numero di sequenze), Sequence Length (lunghezza di ogni sequenza) e Feature Size (numero di caratteristiche per passo temporale). HYDRA applica convoluzioni 1D su questi tensori per estrarre pattern temporali per permettere una classificazione ottimale.

5.4 Gestione class imbalance

Nel contesto del nostro lavoro, che si concentra sulla classificazione di attacchi DDoS, il problema del class imbalance si manifesta in modo significativo. Questo squilibrio deriva dalla natura del dataset utilizzato, in cui le varie classi, rappresentanti diversi tipi di attacchi e traffico normale, presentano una distribuzione disomogenea.

Tradizionalmente, per affrontare il problema del class imbalance, si possono adottare tecniche di riequilibrio dei dati, come l'oversampling delle classi minoritarie o il downsampling delle classi maggioritarie. Tuttavia, data la specificità del problema trattato, abbiamo scelto di non intervenire direttamente sui dati per modificare l'equilibrio delle classi. Questa decisione si basa sull'osservazione che l'utilizzo di sequenze basate su finestre di tempo offre una soluzione più naturale ed efficace per mitigare l'imbalance. Questo poiché catturando l'evoluzione e la dinamica degli attacchi nel tempo, forniscono informazioni contestuali che possono bilanciare in modo indiretto l'influenza delle classi meno rappresentate.

Per la natura degli attacchi, i dataset di partenza presentano un numero di sample diverso. In particolare il dataset contenente il traffico di rete in assenza di attacchi contiene molti meno pacchetti catturati. Pertanto è presente una importante class imbalance. Effettuando la classificazione su sequenze divise per finestre temporali questo problema è mitigato. Contrariamente, nel caso di sequenze divise per finestre di pacchetti, questo sbilanciamento è presente, ma crediamo che dato che abbiamo questo unico insieme di dataset, la soluzione migliore è quella di tenere i dati così come sono, poiché è in fase di classificazione abbiamo utilizzato un k-fold stratificato per mantenere fissa la ...

5.5 Scaling

Spiega lo scaling eseguito per Hydra (SparseScaler)

6 Feature Selection

Nella fase di feature selection, abbiamo adottato due distinti approcci per la generazione delle serie temporali, focalizzati rispettivamente sui pacchetti e sul tempo. Questo per avere due prospettive diverse rispetto agli attacchi e confrontarle.

- **Serie Temporali Basate sul Numero di Pacchetti:** Il primo approccio si è concentrato su serie temporali costruite a parità di numero di pacchetti.
- **Serie Temporali Basate sul Tempo:** Il secondo approccio prevede serie temporali costruite a parità di intervallo temporale.

6.1 Introduzione

Dopo la fase di generazione delle feature siamo andati ad utilizzare tecniche per determinare quelle che sono le feature più importanti per la classificazione per ogni approccio. Abbiamo utilizzato diverse tecniche in modo da confrontare i vari risultati. Questo al fine di diminuire il numero di feature da considerare per evitare curse of dimensionality.

Tra le tecniche utilizzate abbiamo:

- **Chi-Square:** è una misura statistica utilizzata per valutare l'importanza delle feature in problemi di classificazione. Si basa sul test del chi-quadrato, una tecnica statistica comunemente usata per determinare se esiste una relazione significativa tra due variabili categoriali.
- **Random Forest Importances:** è una tecnica utilizzata per valutare e interpretare l'importanza delle feature in un modello di classificazione o regressione basato su Random Forest. L'importanza delle feature fornita da una Random Forest aiuta a identificare quali variabili contribuiscono maggiormente alle decisioni del modello.
- **Recursive Feature Elimination (RFE):** è una tecnica di selezione delle feature utilizzata per migliorare la performance dei modelli di machine learning riducendo il numero di variabili (o feature) presenti nel dataset. L'idea principale di RFE è quella di costruire modelli iterativi e di rimuovere progressivamente le feature meno importanti per ottenere un set di variabili più rilevante e gestibile.

6.2 Serie basate su pacchetti

I risultati per questa tipologia di serie temporali, con le tecniche descritte, sono riportati in *Tabella 3*. La serie temporale su cui sono stati calcolati è caratterizzata da una window size pari a 20 e ad una packet size pari a 20.

Feature	Chi-2 Score	Importances	RFE
entropia_d_ip	500.85	0.010535	✗
entropia_d_port	6,785.74	0.080121	✓
entropia_s_ip	21,353.25	0.158119	✓
entropia_s_port	8,969.34	0.091637	✓
packet_rate	665,887,900	0.165575	✓
byte_rate	46,294,110,000	0.102822	✓
ping_rate	40,144.66	0.091592	✓
modbus_rate	3,545.83	0.018442	✓
tcp_rate	7,407.87	0.076059	✓
other_rate	940.75	0.004223	✗
synack_ratio	10,536.84	0.040974	✓
avg_deltatime	613.66	0.159903	✓

Table 2: Feature Selection Time Series Pacchetti

6.3 Serie basate sul tempo

Abbiamo ottenuto i seguenti risultati considerando una size di 20 e una window di 20:

Feature	Chi-2 Score	Importances	RFE
entropia_d_ip	4.252658	0.015600	✗
entropia_d_port	7.133164	0.034389	✗
entropia_s_ip	296.3207	0.093706	✓
entropia_s_port	260.3345	0.064228	✓
packet_rate	159,998.8	0.048697	✓
byte_rate	8,956,747	0.053681	✓
ping_rate	519.5008	0.022624	✓
modbus_rate	187.3043	0.187848	✓
tcp_rate	68.37837	0.100644	✓
other_rate	4.555297	0.023983	✓
synack_ratio	22,497.23	0.140747	✓
avg_deltatime	3.874184	0.213853	✓

Table 3: Feature Selection Time Series Tempo

6.4 Osservazioni

Per quanto riguarda RFE le feature selezionate sono praticamente le stesse, quello che risulta più evidente è la differenza nei Chi-2-Score. Infatti dal confronto dei valori di Chi-Square Score, risulta evidente che le feature basate sui pacchetti mostrano punteggi molto più elevati rispetto a quelle basate sul tempo. Questo suggerisce che le caratteristiche relative ai pacchetti (come `byte_rate`, `packet_rate`, e `synack_ratio`) sono molto più indicative nel distinguere tra traffico normale e attacchi rispetto a quelle basate sui tempi di osservazione.

L'approccio basato sui pacchetti, dal nostro punto di vista che lavoriamo sulle reti, risulta più realistico e potenzialmente più preciso nella pratica. Questi risultati indicano che, per una classificazione accurata del traffico in esame, è preferibile utilizzare feature che riflettono il comportamento dei pacchetti piuttosto che solo aspetti temporali. Questo approccio fornisce una base più solida per la rilevazione e la classificazione degli attacchi, migliorando così l'efficacia del sistema di monitoraggio e difesa.

7 Training

Durante la fase di training del nostro modello, abbiamo adottato due tecniche fondamentali per garantire una valutazione accurata e affidabile delle sue prestazioni: ten-fold cross-validation e stratificazione.

La *ten-fold cross-validation* è stata utilizzata per valutare la performance del nostro modello in modo robusto e generalizzabile. Questa tecnica prevede la suddivisione del dataset in dieci parti (fold) di dimensioni approssimativamente uguali. In combinazione con la ten-fold cross-validation, abbiamo applicato la *stratificazione* per garantire una rappresentazione proporzionale delle classi all'interno di ogni fold.

L'entità dell'hold-out utilizzato durante questa fase è

8 Testing e Confronto

In questa fase, abbiamo seguito due approcci distinti per valutare le time series generate. Nel primo approccio, le serie temporali sono state direttamente date in pasto ad un classificatore tradizionale, che ha analizzato i dati per eseguire la classificazione. Nel secondo approccio, invece, le time series sono state elaborate utilizzando Hydra [3], una rete neurale convoluzionale.

I classificatori che abbiamo utilizzato sono gli stessi del paper di riferimento:

- *Random Forest*
- *Decision Tree*
- *K-Nearest Neighbors (KNN)*
- *Support Vector Classifier (SVC)*

8.1 Hyperparameter tuning

È stato necessario fare hyperparameter tuning unicamente per il KNN, per poter impostare:

- **n**: numero di punti che si vuole considerare come "neighbours" in fase di classificazione.
- **weights**: parametro che specifica come i "neighbours" influiscono nella votazione. In particolare se si vuole che il voto contribuisca in modo inversamente proporzionale alla distanza (valore impostato a "distance"), o in modo uniform (valore impostato a "uniform").

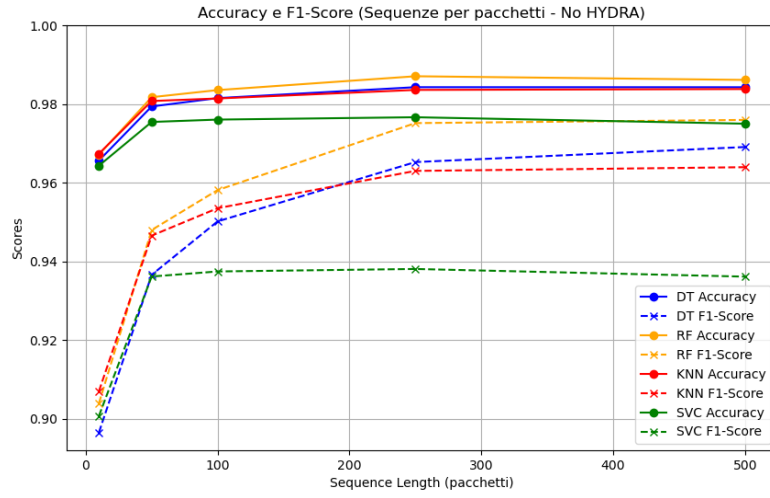
Come tecnica abbiamo utilizzato il Grid Search con 5-fold cross validation, che consiste nel testare tutte le possibili combinazioni di parametri specificando il range di prova. In particolare nel nostro caso abbiamo testato valori di n interi compresi tra 1 e 30, testando con peso uniforme e proporzionale alla distanza. In fase di classificazione abbiamo utilizzato i parametri ottenuti in questo modo.

8.2 Serie Basate su Pacchetti

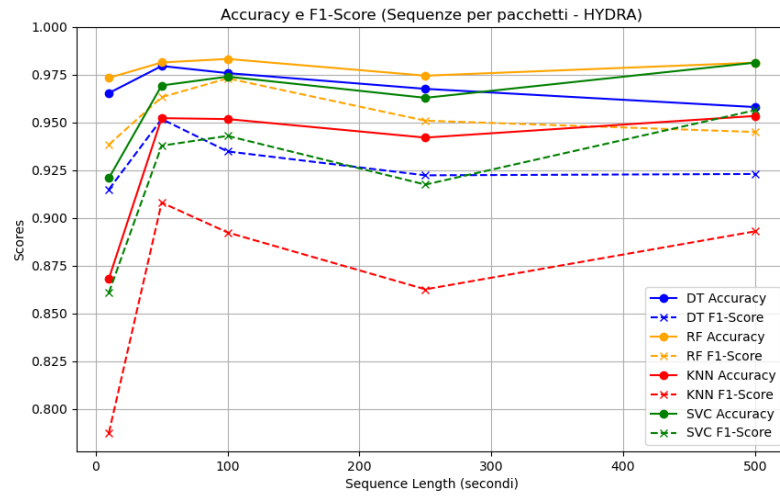
In *Figura 1* riportiamo il confronto tra i risultati dei vari classificatori utilizzati, riportandone l'accuracy e l'F1 (espressi in percentuale). Abbiamo anche considerato come variano tali valori al variare della dimensione della serie temporale considerata.

8.3 Serie Basate su Tempo

Un approccio analogo si è utilizzato nel considerare le serie temporali basate sul tempo, in cui siamo andati a valutare l'accuracy al variare della dimensione dell'intervallo di tempo considerato. I risultati sono riportati in *Figura 2*.

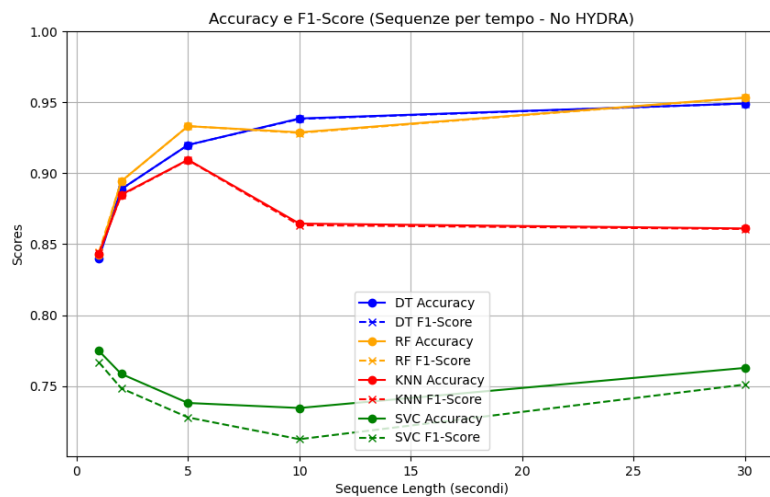


(a) Utilizzo di Dati Aggregati

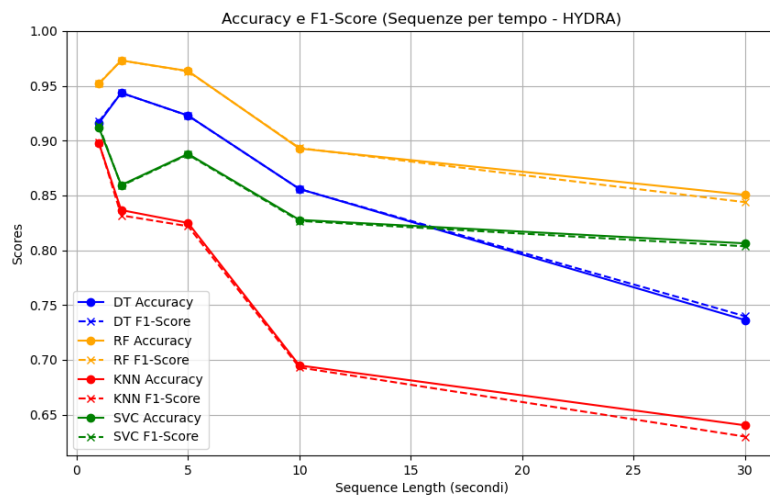


(b) Utilizzo di Hydra

Figure 1: Serie Temporali su Pacchetti



(a) Utilizzo di Dati Aggregati



(b) Utilizzo di Hydra

Figure 2: Serie Temporalı su Tempo

8.4 Osservazioni

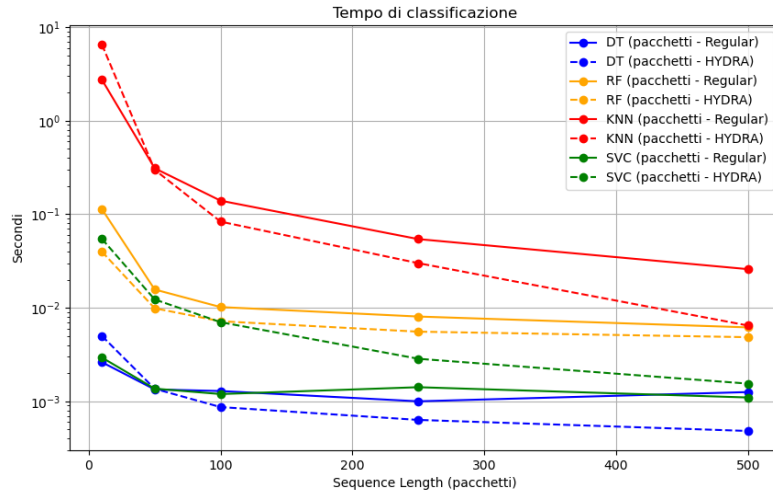
Possiamo notare che l'utilizzo di Hydra presenta, dal nostro punto di vista, le seguenti problematiche:

Questo confronto ci aiuta ad identificare quale combinazione di classificatore e lunghezza della serie temporale offrisse i migliori risultati in termini di accuratezza nella classificazione degli attacchi.

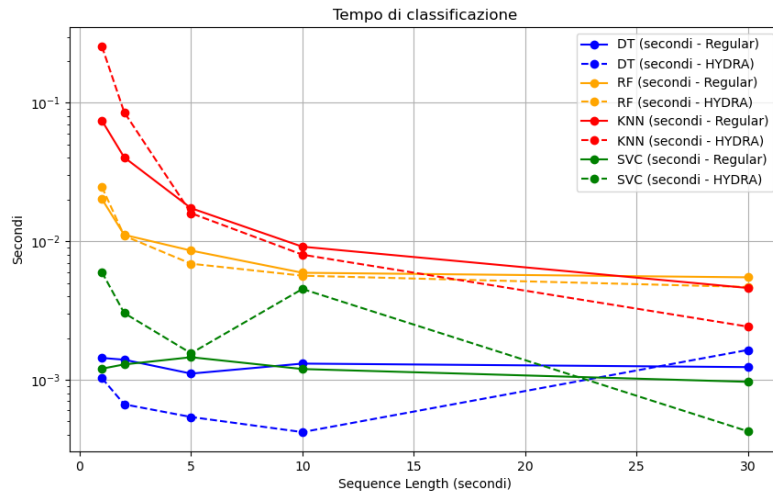
In *Figura ??* sono riportati i tempi di classificazione impiegati dai classificatori considerando le due tipologie di serie temporali. La scala dei tempi è logaritmica per evidenziare la distribuzione non uniforme dei tempi quando si considerano serie di piccole dimensioni.

Nella relazione, non riportiamo direttamente i parametri di valutazione come l'accuratezza media, matrici di confusione e di correlazione per ciascun modello. Tuttavia, questi dati sono disponibili nel codice allegato, che contiene tutte le analisi dettagliate e i risultati ottenuti durante le nostre sperimentazioni.

La classificazione tramite hydra richiede più tempo, rendendo l'utilizzo di classificatori di questo tipo non consigliabili per applicazioni che richiedono rilevamento real time. Questo si può notare dal confronto tra i grafici presenti in *Figura ??*

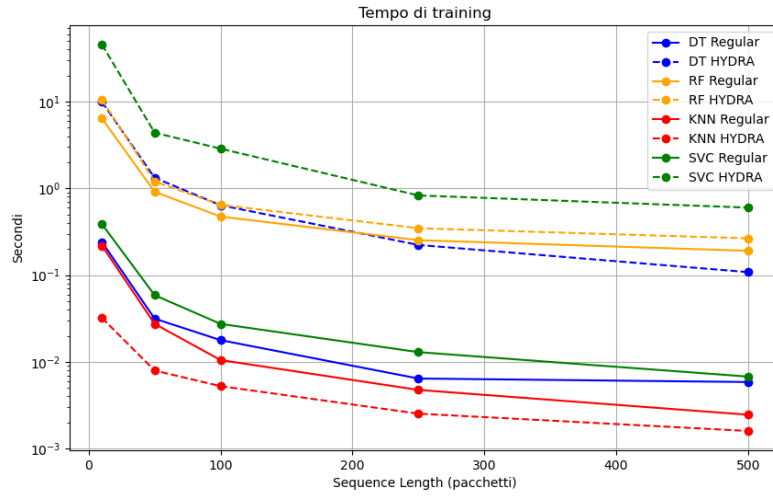


(a) Sequenze per pacchetti

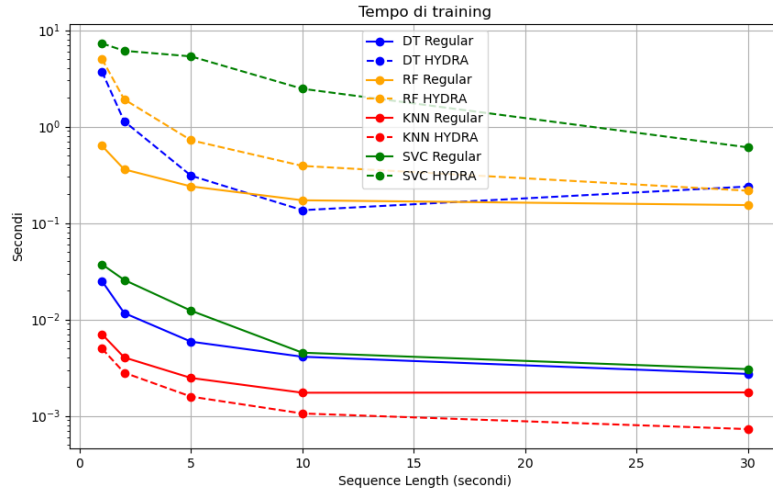


(b) Sequenze per tempo

Figure 3: Tempi di classificazione



(a) Sequenze per pacchetti



(b) Sequenze per tempo

Figure 4: Tempi di training

9 Conclusioni

- Considerazioni sui risultati ottenuti e sulla validità dei classificatori
- Considerazioni sulla generalità e sui parametri scelti
- Considerazioni su come futuri lavori necessitano l'utilizzo di più dataset generati da reti diverse per validare il valore delle nostre feature

References

- [1] Muhammad Aamir and Syed Mustafa Ali Zaidi. “DDoS attack detection with feature engineering and machine learning: the framework and performance evaluation”. In: *International Journal of Information Security* 18 (2019), pp. 761–785.
- [2] Roberto Mustillo Davide De Zuane Rahmi El Mechri. URL: <https://github.com/rahmec/general-ddos-classification>.
- [3] Angus Dempster, Daniel F Schmidt, and Geoffrey I Webb. URL: <https://github.com/angus924/hydra>.
- [4] Ivo Frazão et al. “Denial of Service Attacks: Detecting the Frailties of Machine Learning Algorithms in the Classification Process.” in: Dec. 2018, pp. 230–235. ISBN: 978-3-030-05848-7. DOI: 10.1007/978-3-030-05849-4_19.

A Feature generation

```
1 import pandas as pd
2 import threading
3
4 def add_flow_number_column(pcap_df):
5     pcap_df['Flow'] = pcap_df['Source'] + '->' + pcap_df['Destination']
6     pcap_df['FlowNumber'] = 0
7     flows = {flow_key: 0 for flow_key in pcap_df['Flow'].unique()}
8     for index, row in pcap_df.iterrows():
9         pcap_df.at[index, 'FlowNumber'] = flows[row['Flow']]
10        flows[row['Flow']] += 1
11
12 def aggregate_sequences(df, seq_length, window_size, mean):
13     aggregated_data = []
14
15     df_range = []
16     if mean == "packets":
17         df_range = range(0, len(df) - seq_length + 1, window_size)
18     elif mean == "time":
19         df_range = range(0, 60*30-seq_length, window_size)
20
21     for start in df_range:
22
23         if mean == "packets":
24             sequence = df.iloc[start:start + seq_length]
25             seq_seconds = sequence['Time'].iloc[seq_length-1] -
26                 sequence['Time'].iloc[0]
27             seq_packets = seq_length
28         elif mean == "time":
29             sequence = df[df['Time'] >= start]
30             sequence = sequence[sequence['Time'] < start +
31                 seq_length]
32             seq_seconds = seq_length
33             seq_packets = len(sequence)
```

```

32
33         if seq_packets == 0:
34             continue
35
36         length_sum = sequence['Length'].sum()
37         byte_rate = length_sum/seq_seconds
38         packet_rate = seq_packets/seq_seconds
39         avg_flow_number = sequence['FlowNumber'].mean()
40         avg_deltatime = sequence['Delta-Time'].mean()
41         num_modbus = (sequence['Protocol'] == 'Modbus/TCP').sum()
42         num_tcp = (sequence['Protocol'] == 'TCP').sum()
43         num_ping = (sequence['Protocol'] == 'ICMP').sum()
44         num_other = seq_packets - num_modbus - num_tcp - num_ping
45         num_syn = (sequence['SYN'] == 'Set').sum()
46         num_ack = (sequence['ACK'] == 'Set').sum()
47         modbus_rate = num_modbus/seq_length
48         tcp_rate = num_tcp/seq_packets
49         ping_rate = num_ping/seq_packets
50         other_rate = num_other/seq_packets
51         if num_syn == 0:
52             synack_ratio = 0
53         elif num_ack == 0:
54             synack_ratio = 1
55         else:
56             synack_ratio = num_syn/num_ack
57         source_entropy = sequence['Source'].nunique()/seq_packets
58         destination_entropy = sequence['Destination'].nunique()/
59             seq_packets
60         s_port_entropy = sequence['S-Port'].nunique()/seq_packets
61         d_port_entropy = sequence['D-Port'].nunique()/seq_packets
62         label = sequence['Label'].iloc[0]
63
64         aggregated_row = {
65             'byte_rate' : byte_rate,
66             'packet_rate' : packet_rate,
67             'avg_flow_number' : avg_flow_number,
68             'avg_deltatime' : avg_deltatime,
69             'source_entropy' : source_entropy,
70             'destination_entropy' : destination_entropy,
71             's-port_entropy' : s_port_entropy,
72             'd-port_entropy' : d_port_entropy,
73             'synack_ratio' : synack_ratio,
74             'modbus_rate' : modbus_rate,
75             'tcp_rate' : tcp_rate,
76             'ping_rate' : ping_rate,
77             'other_rate' : other_rate,
78             'label' : label
79         }
80
81         aggregated_data.append(aggregated_row)
82
83         aggregated_df = pd.DataFrame(aggregated_data)
84
85         return aggregated_df
86
87 def aggregate_df(seq_length, mean):
88     window_size = seq_length

```

```

88     aggregated_dfs = []
89     for df in dfs:
90         agg_df = aggregate_sequences(df, seq_length, window_size,
91                                     mean)
92         aggregated_dfs.append(agg_df)
93
94     aggregated_df = pd.concat(aggregated_dfs)
95     if mean == "packets":
96         print(f'CREATED csv/aggregated_{seq_length}packets_df.csv')
97         aggregated_df.to_csv(f'csv/aggregated_{seq_length}
98                             packets_df.csv', index=False)
99     elif mean == "time":
100         print(f'CREATED csv/aggregated_time_{seq_length}s_df.csv')
101         aggregated_df.to_csv(f'csv/aggregated_time_{seq_length}s_df
102                             .csv', index=False)
103
104 clean_df = pd.read_csv('csv/clean15-30.csv')
105 tcpsynflood_df = pd.read_csv('csv/tcpsynflood15-30.csv')
106 pingflood_df = pd.read_csv('csv/pingflood15-30.csv')
107 modbusflood_df = pd.read_csv('csv/modbusflood15-30.csv')
108
109 clean_df['Label'] = 0
110 tcpsynflood_df['Label'] = 1
111 pingflood_df['Label'] = 2
112 modbusflood_df['Label'] = 3
113
114 dfs = [clean_df, tcpsynflood_df, pingflood_df, modbusflood_df]
115
116 packet_seq_lengths = [10,20,50,100,250,500]
117 time_seq_lengths = [1,2,5,10,30]
118 threads = []
119
120 print("Generating flow numbers")
121 threads = []
122 for df in dfs:
123     threads.append(threading.Thread(target=add_flow_number_column,
124                                     args=(df,)))
125 for thread in threads:
126     thread.start()
127 for thread in threads:
128     thread.join()
129 print("Generating threads")
130 threads = []
131 for seq_length in packet_seq_lengths:
132     threads.append(threading.Thread(target=aggregate_df, args=(
133         seq_length, "packets")))
134 print(f"Generated threads for packet sequences with length {
135     packet_seq_lengths}")
136 for seq_length in time_seq_lengths:
137     threads.append(threading.Thread(target=aggregate_df, args=(
138         seq_length, "time")))
139 print(f"Generated threads for time sequences with length {
140     time_seq_lengths}")
141 print("Starting threads")
142 for thread in threads:
143     thread.start()
144 print("Waiting until all threads are done")

```

```
137 for thread in threads:  
138     thread.join()  
139 print("CSV created!")
```

B Classificazione