

## Feature Summary

Feature	Formula / Definition	Scale?	Benefit
<b>rate_ratio</b>	$(src\_rate + \epsilon) / (dst\_rate + \epsilon)$	Yes	Highlights asymmetric traffic patterns
<b>syn_to_ack</b>	$(syn\_packets + 1) / (ack\_packets + 1)$	Yes	Flags incomplete handshakes (scans)
<b>rst_to_fin</b>	$(rst\_packets + 1) / (fin\_packets + 1)$	Yes	Distinguishes abortive vs. clean teardowns
<b>avg_pkt_size</b>	$(overall\_rate \times flow\_time) / total\_packets$	Yes	Captures average payload size per packet
<b>mean_interpkt</b>	Mean of inter-packet arrival times	Yes	Detects uniform vs. bursty traffic
<b>std_interpkt</b>	Std. dev. of inter-packet arrival times	Yes	Measures consistency of packet spacing
<b>p90_interpkt</b>	90th percentile of inter-packet times	Yes	Spots occasional long gaps vs. floods
<b>burstiness</b>	$max(interpkt) / mean(interpkt)$	Yes	High in volumetric flood attacks
<b>payload_entropy</b>	Shannon entropy of payload byte histogram	Yes	Differentiates encrypted/compressed vs. repetitive payloads
<b>value_range</b>	$max\_value - min\_value$	Yes	Captures payload variability
<b>flows_last_10s</b>	Rolling count of flows by source IP in 10 s	Yes	Detects fast-scanning or botnet behavior
<b>unique_dsts_last_10s</b>	Rolling count of unique dst IPs/ports in 10 s	Yes	Quantifies scanning breadth
<b>hour_sin</b>	$\sin(2\pi \cdot hour\_of\_day/24)$	Yes	Encodes cyclical time-of-day patterns
<b>hour_cos</b>	$\cos(2\pi \cdot hour\_of\_day/24)$	Yes	Encodes cyclical time-of-day patterns
<b>handshake_complete</b>	1 if full SYN→SYN-ACK→ACK seen, else 0	No	Flags half-open (SYN flood) attempts
<b>abrupt_reset</b>	1 if $rst\_flags=1$ AND $fin\_flags=0$ , else 0	No	Captures abrupt, reset-only connection drops

Feature	Formula / Definition	Scale?	Benefit
tcp_syn_ratio	syn_flags * protocol_tcp	No	Emphasizes SYN usage specifically in TCP flows
udp_psh	psh_flags * protocol_udp	No	Highlights PSH behavior in UDP traffic

## Detailed Descriptions

### Continuous Features (Scale)

- **rate\_ratio**
  - **Formula:**  $(\text{src\_rate} + \epsilon) / (\text{dst\_rate} + \epsilon)$  /  $(\text{src\_rate} + \epsilon) / (\text{dst\_rate} + \epsilon)$
  - **Scaling:** Yes
  - **Benefit:** Amplifies asymmetric flows (e.g., reflection attacks) by comparing source vs. destination throughput.
- **syn\_to\_ack**
  - **Formula:**  $(\text{syn\_packets} + 1) / (\text{ack\_packets} + 1)$  /  $(\text{syn\_packets} + 1) / (\text{ack\_packets} + 1)$
  - **Scaling:** Yes
  - **Benefit:** Flags flows with many SYNs but few ACKs—typical of port scans or half-open connection attempts.
- **rst\_to\_fin**
  - **Formula:**  $(\text{rst\_packets} + 1) / (\text{fin\_packets} + 1)$  /  $(\text{rst\_packets} + 1) / (\text{fin\_packets} + 1)$
  - **Scaling:** Yes
  - **Benefit:** Distinguishes abortive connections (high RST) from clean teardowns (high FIN).
- **avg\_pkt\_size**
  - **Formula:**  $(\text{overall\_rate} \times \text{flow\_time}) / \text{total\_packets}$  /  $(\text{overall\_rate} \times \text{flow\_time}) / \text{total\_packets}$
  - **Scaling:** Yes
  - **Benefit:** Captures typical payload size per packet—different for volumetric floods vs. interactive traffic.
- **mean\_interpkt, std\_interpkt, p90\_interpkt**

- **Definition:** Summary statistics (mean, standard deviation, 90th percentile) of the time gaps between successive packets.
- **Scaling:** Yes
- **Benefit:** Helps differentiate uniform packet streams (e.g., botnet traffic) from bursty, human-driven flows.
- **burstiness**
  - **Formula:**  $\max(\text{interpkt}) / \overline{\text{interpkt}}$
  - **Scaling:** Yes
  - **Benefit:** Very high in DDoS floods (sharp bursts), lower in steady or interactive sessions.
- **payload\_entropy**
  - **Definition:** Shannon entropy of the histogram of payload byte values.
  - **Scaling:** Yes
  - **Benefit:** High entropy indicates encrypted/compressed content; low entropy suggests repetitive flood packets.
- **value\_range**
  - **Formula:**  $\max\_value - \min\_value$
  - **Scaling:** Yes
  - **Benefit:** Measures spread of payload bytes—small in ping floods, larger in application-layer attacks.
- **flows\_last\_10s, unique\_dsts\_last\_10s**
  - **Definition:** Rolling counts over a 10-second window per source IP—total flows and unique destinations/ports.
  - **Scaling:** Yes
  - **Benefit:** Detects rapid scanning or botnet behavior by volume and breadth of connection attempts.
- **hour\_sin, hour\_cos**
  - **Formula:**  $\sin(2\pi \times \text{hour}/24)$ ,  $\cos(2\pi \times \text{hour}/24)$
  - **Scaling:** Yes
  - **Benefit:** Encodes cyclical daily patterns without artificial boundaries between midnight and 0h.

## Binary Features (No Scale)

- **handshake\_complete**
  - **Definition:** 1 if a full SYN→SYN-ACK→ACK handshake was observed, else 0.
  - **Scaling:** No

- **Benefit:** Flags half-open (SYN flood) flows lacking complete handshakes.
  - **abrupt\_reset**
    - **Definition:** 1 if `rst_flags == 1` **and** `fin_flags == 0`, else 0.
    - **Scaling:** No
    - **Benefit:** Captures abrupt, reset-only connection terminations common in scans or attacks.
  - **tcp\_syn\_ratio**
    - **Definition:** `syn_flags * protocol_tcp` (product of two binary indicators).
    - **Scaling:** No
    - **Benefit:** Emphasizes SYN usage specifically in TCP traffic, ignoring other protocols.
  - **udp\_psh**
    - **Definition:** `psh_flags * protocol_udp`.
    - **Scaling:** No
    - **Benefit:** Highlights PSH-like behavior in UDP flows, which is atypical and can signal anomalies.
- 

**Pipeline Note:** Compute all above features **after** raw-data cleaning and outlier handling, but **before** any scaling or normalization steps. Then fit your scaler(s) on training data (including these continuous features) and apply to both train and test sets.