# Anomaly Detection with Extreme Value Theory

A. Siffer, P-A Fouque, A. Termier and C. Largouet

May 30, 2017

# Contents

# Context

—o Massive usage of the Internet

—o Massive usage of the Internet
  · More and more vulnerabilities

**1 Tbps DDoS Attack**

Powered By 150,000 Hacked IoT Devices

**WannaCry ransomware used in widespread attacks all over the world**

—o Massive usage of the Internet
- More and more vulnerabilities
- More and more threats

**1 Tbps DDoS Attack**

Powered By 150,000 Hacked IoT Devices

**WannaCry ransomware used in widespread attacks all over the world**

—o Massive usage of the Internet
  · More and more vulnerabilities
  · More and more threats

—o Awareness of the sensitive data and infrastructures

**1 Tbps DDoS Attack**

Powered By 150,000 Hacked IoT Devices

**WannaCry ransomware used in widespread attacks all over the world**

- Massive usage of the Internet
  - More and more vulnerabilities
  - More and more threats

- Awareness of the sensitive data and infrastructures

**1 Tbps DDoS Attack**

Powered By 150,000 Hacked IoT Devices

⇒ Network security :
a major concern

--o IDS (Intrusion Detection System)
  - Monitor traffic
  - Detect attacks

—∘ IDS (Intrusion Detection System)
- Monitor traffic
- Detect attacks

—∘ Current methods : rule-based
- Work fine on common and well-known attacks
- Cannot detect new attacks

- —o IDS (Intrusion Detection System)
  - · Monitor traffic
  - · Detect attacks
- —o Current methods : rule-based
  - · Work fine on common and well-known attacks
  - · Cannot detect new attacks



- —o Emerging methods : anomaly-based
  - · Use the network data to estimate a normal behavior
  - · Apply algorithms to detect abnormal events ($\rightarrow$ attacks)

—o Basic scheme

$$\text{data} \longrightarrow \boxed{\text{ALGORITHM}} \longrightarrow \text{alerts}$$

—o Basic scheme

$$\text{data} \longrightarrow \boxed{\text{ALGORITHM}} \longrightarrow \text{alerts}$$

—o Many "standard" algorithms have been tested

—○ Basic scheme



data ⟶ ALGORITHM ⟶ alerts

—○ Many "standard" algorithms have been tested

—○ Complex pipelines are emerging (ensemble/hybrid techniques)



data ⟶ ⟶ alerts

—◦ Algorithms are not magic
  · They give some information about data (scores)

—○ Algorithms are not magic
- They give some information about data (scores)
- But the decision often rely on a human choice

```
if score>threshold then trigger alert
```

—◦ Algorithms are not magic
- They give some information about data (scores)
- But the decision often rely on a human choice

```
if score>threshold then trigger alert
```

—◦ The thresholds are often hard-set
- Expertise
- Fine-tuning
- Distribution assumption

—◦ Algorithms are not magic
  · They give some information about data (scores)
  · But the decision often rely on a human choice

    ```
    if score>threshold then trigger alert
    ```

—◦ The thresholds are often hard-set
  · Expertise
  · Fine-tuning
  · Distribution assumption

—◦ **Our idea**: provide dynamic threshold with a probabilistic meaning

# Providing better thresholds

—○ How to set $z_q$ such that $\mathbb{P}(X€ > z_q) < q$ ?

—o Drawbacks: stuck in the interval, poor resolution

—○ Drawbacks: manual step, distribution assumption

—o Different clients and/or temporal drift

| PROPERTIES | Empirical quantile | Standard model |
|---|---|---|
| *statistical guarantees* | Yes | Yes |
| *easy to adapt* | Yes | No |
| *high resolution* | No | Yes |

⊸ Main result (Fisher-Tippett-Gnedenko, 1928)

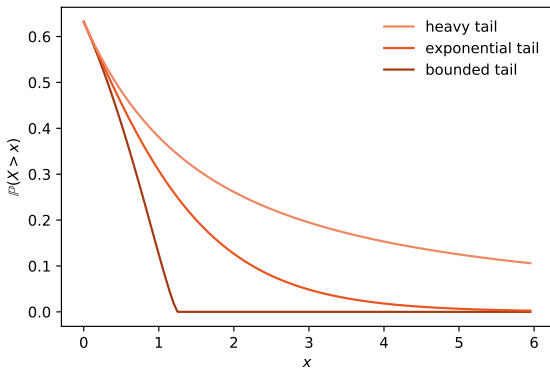*The extreme values of any distribution have nearly the same distribution (called Extreme Value Distribution)*

—∘ Main result (Fisher-Tippett-Gnedenko, 1928)

*The extreme values of any distribution have nearly the same distribution (called Extreme Value Distribution)*

—○ Let $X_1, X_2, \ldots X_n$ a sequence of i.i.d. random variables with

$$S_n = \sum_{i=1}^{n} X_i \qquad M_n = \max_{1 \leq i \leq n}(X_i)$$

⊸ Let $X_1, X_2, \ldots X_n$ a sequence of i.i.d. random variables with

$$S_n = \sum_{i=1}^{n} X_i \qquad M_n = \max_{1 \leq i \leq n}(X_i)$$

⊸ Central Limit Theorem

$$\frac{S_n - n\mu}{\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, \sigma^2)$$

─○ Let $X_1, X_2, \ldots X_n$ a sequence of i.i.d. random variables with

$$S_n = \sum_{i=1}^{n} X_i \qquad M_n = \max_{1 \leq i \leq n}(X_i)$$

─○ Central Limit Theorem

$$\frac{S_n - n\mu}{\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, \sigma^2)$$

─○ FTG Theorem

$$\frac{M_n - a_n}{b_n} \xrightarrow{d} \mathrm{EVD}(\gamma)$$

—◦ Second theorem of EVT (Pickands-Balkema-de Haan, 1974)

*The excesses over a high threshold follow a Generalized Pareto Distribution (with parameters $\gamma, \sigma$)*

─○ Second theorem of EVT (Pickands-Balkema-de Haan, 1974)

*The excesses over a high threshold follow a Generalized Pareto Distribution (with parameters $\gamma, \sigma$)*

─○ What does it imply ?
   · we have a model for extreme events
   · we can compute $z_q$ for $q$ as small as desired

- Get some data $X_1, X_2 \ldots X_n$
- Set a high threshold $t$ and retrieve the excesses $Y_j = X_{k_j} - t$ when $X_{k_j} > t$
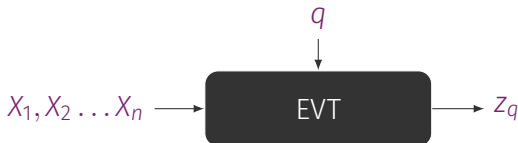
- ⊸ Get some data $X_1, X_2 \ldots X_n$
- ⊸ Set a high threshold $t$ and retrieve the excesses $Y_j = X_{k_j} - t$ when $X_{k_j} > t$
- ⊸ Fit a GPD to the $Y_j$ ($\rightarrow$ find parameters $\gamma, \sigma$)

- Get some data $X_1, X_2 \ldots X_n$
- Set a high threshold $t$ and retrieve the excesses $Y_j = X_{k_j} - t$ when $X_{k_j} > t$
- Fit a GPD to the $Y_j$ ($\to$ find parameters $\gamma, \sigma$)
- Compute $z_q$ such as $\mathbb{P}(X > z_q) < q$

- Get some data $X_1, X_2 \ldots X_n$
- Set a high threshold $t$ and retrieve the excesses $Y_j = X_{k_j} - t$ when $X_{k_j} > t$
- Fit a GPD to the $Y_j$ ($\to$ find parameters $\gamma, \sigma$)
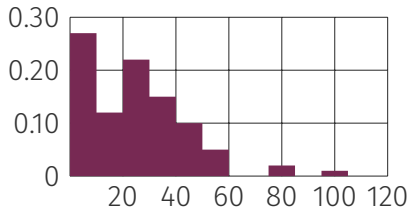- Compute $z_q$ such as $\mathbb{P}(X > z_q) < q$

$$X_1, X_2 \ldots X_n \longrightarrow \boxed{\text{EVT}} \longrightarrow z_q$$

with $q$ as input to EVT.

# Finding anomalies in streams

(initial batch)

$X_1, X_2 \ldots X_n$

(initial batch)

$X_1, X_2 \ldots X_n \longrightarrow$ **Calibration** $\longleftarrow q$

(initial batch)

$X_1, X_2 \ldots X_n$ → Calibration → 

$q$

(initial batch)

$X_1, X_2 \ldots X_n$ → Calibration →

$q$

$t$

(initial batch)

$X_1, X_2 \ldots X_n$ $\longrightarrow$ CALIBRATION

$q$

16

(initial batch)

$X_1, X_2 \ldots X_n$ → Calibration

$q$

(stream)

$X_{i>n}$

(initial batch)

$X_1, X_2 \ldots X_n \longrightarrow$ | CALIBRATION |

$q$

(stream)

$X_{i>n} \longrightarrow$ | $X_i > z_q$ |

(initial batch)

$X_1, X_2 \ldots X_n \longrightarrow$ CALIBRATION

$q$

$t$    $z_q$

0.30
0.20
0.10
0
20  40  60  80  100  120

TRIGGER ALARM

YES

(stream)

$X_{i>n} \longrightarrow$ $X_i > z_q$

16

—◦ An example with ground truth : a Gaussian White Noise
- · 40 streams with 200 000 iid variables drawn from $\mathcal{N}(0, 1)$
- · $q = 10^{-3} \Rightarrow$ theoretical threshold $z_{th} \simeq 3.09$

⊸ An example with ground truth : a Gaussian White Noise
  · 40 streams with 200 000 iid variables drawn from $\mathcal{N}(0,1)$
  · $q = 10^{-3} \Rightarrow$ theoretical threshold $z_{th} \simeq 3.09$
⊸ Averaged relative error

# Application to intrusion detection

—o Lack of relevant public datasets to test the algorithms …

—◦ Lack of relevant public datasets to test the algorithms …

—◦ KDD99 ? See [McHugh 2000] and [Mahoney & Chan 2003]

—o Lack of relevant public datasets to test the algorithms …

—o KDD99 ? See [McHugh 2000] and [Mahoney & Chan 2003]

—o We rather use MAWI

· 15 min a day of real traffic (.pcap file)
· Anomaly patterns given by the MAWILab [Fontugne *et al.* 2010] with taxonomy [Mazel et al. 2014]

—o Lack of relevant public datasets to test the algorithms …

—o KDD99 ? See [McHugh 2000] and [Mahoney & Chan 2003]

—o We rather use MAWI
- 15 min a day of real traffic (.pcap file)
- Anomaly patterns given by the MAWILab [Fontugne *et al.* 2010] with taxonomy [Mazel et al. 2014]

—o Preprocessing step : raw .pcap → NetFlow format (only metadata)

- The ratio of SYN packets : relevant feature to detect network scan [Fernandes & Owezarski 2009]

—o The ratio of SYN packets : relevant feature to detect network scan [Fernandes & Owezarski 2009]

—o The ratio of SYN packets : relevant feature to detect network scan [Fernandes & Owezarski 2009]



—o Goal: find peaks

—∘ Parameters : $q = 10^{-4}, n = 2000$ (from the previous day record)

—o Parameters : $q = 10^{-4}, n = 2000$ (from the previous day record)

—◦ The main parameter $q$: a False Positive regulator

—o The main parameter $q$: a False Positive regulator

—○ The main parameter $q$: a False Positive regulator



—○ 86% of scan flows detected with less than 4% of FP

# A more general framework

- A single main parameter $q$
  - With a probabilistic meaning $\rightarrow \mathbb{P}(X > z_q) < q$
  - False Positive regulator

- A single main parameter $q$
  - With a probabilistic meaning $\rightarrow \mathbb{P}(X > z_q) < q$
  - False Positive regulator
- Stream capable
  - Incremental learning
  - Fast ($\sim 1000$ values/s)
  - Low memory usage (only the excesses)

—o SPOT
  · performs dynamic thresholding without distribution assumption
  · uses it to detect network anomalies

–○ SPOT
  · performs dynamic thresholding without distribution assumption
  · uses it to detect network anomalies
–○ But it could be adapted to

—◦ SPOT
- performs dynamic thresholding without distribution assumption
- uses it to detect network anomalies

—◦ But it could be adapted to
- compute upper and lower thresholds

—○ SPOT
  · performs dynamic thresholding without distribution assumption
  · uses it to detect network anomalies
—○ But it could be adapted to
  · compute upper and lower thresholds
  · other fields

—○ SPOT
  · performs dynamic thresholding without distribution assumption
  · uses it to detect network anomalies

—○ But it could be adapted to
  · compute upper and lower thresholds
  · other fields
  · drifting contexts (with an additional parameter) → DSPOT

—o Thursday the 9th of February 2017

—○ Thursday the 9th of February 2017
- **9h** : explosion at Flamanville nuclear plant

—o Thursday the 9th of February 2017
- **9h** : explosion at Flamanville nuclear plant
- **11h** : official declaration of the incident by EDF

—○ Thursday the 9th of February 2017
- **9h** : explosion at Flamanville nuclear plant
- **11h** : official declaration of the incident by EDF

—○ What about the EDF stock prices ?

— Context: A great deal of work has been done to develop anomaly detection algorithms

—∘ <u>Context</u>: A great deal of work has been done to develop anomaly detection algorithms

—∘ <u>Problem</u>: Decision thresholds rely on either distribution assumption or expertise

- —∘ <u>Context</u>: A great deal of work has been done to develop anomaly detection algorithms
- —∘ <u>Problem</u>: Decision thresholds rely on either distribution assumption or expertise
- —∘ <u>Our solution</u>: Building dynamic threshold with a probabilistic meaning

—○ <u>Context</u>: A great deal of work has been done to develop anomaly detection algorithms

—○ <u>Problem</u>: Decision thresholds rely on either distribution assumption or expertise

—○ <u>Our solution</u>: Building dynamic threshold with a probabilistic meaning
  - Application to detect network anomalies

- ○ <u>Context</u>: A great deal of work has been done to develop anomaly detection algorithms
- ○ <u>Problem</u>: Decision thresholds rely on either distribution assumption or expertise
- ○ <u>Our solution</u>: Building dynamic threshold with a probabilistic meaning
  - · Application to detect network anomalies
  - · But a general tool to monitor online time series in a blind way

–○ <u>Context</u>: A great deal of work has been done to develop anomaly detection algorithms

–○ <u>Problem</u>: Decision thresholds rely on either distribution assumption or expertise

–○ <u>Our solution</u>: Building dynamic threshold with a probabilistic meaning
  · Application to detect network anomalies
  · But a general tool to monitor online time series in a blind way

–○ <u>Future</u>: Adapt the method to higher dimensions