

Exploring statistical techniques for threat detection and AI security

Raffaele Castagna

Statistics Academic Year 2025-2026

Contents

1	Introduction	2
1.1	The Dual Role of Artificial Intelligence: Tool and Target	2
2	Statistical Anomaly Detection in Network Security	4
2.1	Gaussian Assumption and Univariate Analysis	4
2.2	Multivariate Analysis and Mahalanobis Distance	4
2.3	Payload Detection	5
2.4	Dimensionality Reduction and PCA	5
2.5	Industrial Control Systems and Operational Technology applications	6

1 Introduction

Cybersecurity has always been shaped by data analysis, ever since its inception. Even at the start when the *fortress model* was the optimal choice, sysadmins typically relied on the binary classification of traffic, files and behaviours that were either safe or were known malicious patterns, but to actually analyze a file most relied on static analysis of signatures, where the hash of a file was computed and then denoted as malicious, this could also be applied to byte sequences in network traffic or to generalize even more IP ranges could be blocked, either based on suspicious activity or on countries. Under the hood there was a statistical assumption: the distribution of malicious behaviour was disjointed from the distribution of normal behaviour, and therefore there was a boundary between the two that was immutable.

However with the continued expansion of the internet, and the users therein, this assumption became obsolete, the world has changed and with it the amount of data that is ingested, take for example Industrial Control systems or autonomous vehicles navigating public roads[3], in these kinds of environments the amount of data and heterogeneity of the devices make signature based maintenance an impossible task. Expanding on this, the convergence of IT and OT (Operational Technology) has exposed critical industrial machine to a wide area of attack, and therefore shifted what can be considered normal behaviour, as well as the attack surface.[8]

With these new surfaces, the entire landscape changed, particularly with the rise of APT (Advanced Persistent Threats), which can even be state-sponsored threats, that are particularly interested in IO, these APTs usually utilize a technique called "living off the land" which utilizes legitimate software and functions to carry out their attacks, e.g. utilizing group policy changes to gain persistence and access to restricted controls and information, another type of vulnerability that is continually exploited are Zero-Day attacks, which due to their nature are impossible to defend against with static hashing, so there's a need to move to a dynamic analysis of behaviour, which is where statistics come into play, so that we see a threat as something that *deviates from a probabilist baseline*, this is what in the modern cybersecurity field is called **statistical anomaly detection**[1].

Anomaly detection is rooted in the statistical hypothesis that malicious activity is rare and different from normal activity. It does not ask, "Is this specific packet known to be bad?" but rather, "What is the probability that this packet belongs to the distribution of normal traffic observed over the last month?" This shift requires security professionals to abandon binary certainty in favor of probabilistic reasoning. It demands a rigorous understanding of distributions, variance, and correlation. The security analyst of the future must be as fluent in covariance matrices and p-values as they are in firewalls and encryption protocols.

1.1 The Dual Role of Artificial Intelligence: Tool and Target

As statistical methods evolved, they gave rise to Machine Learning (ML) and Artificial Intelligence (AI), which are essentially statistical inference engines operating at scale. AI has become a powerful tool for defense, capable of ingesting terabytes of log data to identify subtle correlations that are impossible to spot for a human analysts[2]. AI-driven statistical anomaly detection has been shown to outperform traditional analytical techniques in mitigating cybersecurity risks in complex environments like wireless networks[10].

However, this reliance on AI introduces a recursive vulnerability: the statistical models themselves are now targets. We are witnessing the rise of Adversarial Machine Learning (AML), where attackers exploit the probabilistic nature of AI. By carefully crafting inputs that are statistically indistinguishable from benign data to a human but catastrophic to a model's mathematical logic, attackers can blind surveillance systems or cause autonomous vehicles to misinterpret stop signs[9][11].

Thus, the domain of AI security is not merely about securing the software code but about secur-

ing the statistical inference process. It involves using advanced statistical tests such as Kernel Density Estimation (KDE), Maximum Mean Discrepancy (MMD), and Benford's Law to detect when a distribution is being manipulated. The defender must now monitor the monitor, applying statistical tools to ensure that the AI systems protecting the network are not themselves compromised by statistical illusions.

2 Statistical Anomaly Detection in Network Security

Network intrusion detection is the practice of monitoring network traffic for suspicious activity and unauthorized access. While rule-based systems (e.g. Snort) look for known byte sequences, statistical anomaly detection techniques build a model of what is considered "normal" traffic for a network and flags certain deviations. This approach has been widely used since the 2000s, and is currently being amplified by the use of AI.

2.1 Gaussian Assumption and Univariate Analysis

To be able to distinguish normal from anomalous traffic we must assume that normality exists. The Gaussian distribution (or what we called the normal distribution) is essential for anomaly detection because it's a mathematically convenient way to define what we consider normal. In a normal distribution, around 68% of data points fall within one standard deviation (SD, σ) or the mean (μ), 95% within 2 SDs and 99.7% within 3 SDs. Therefore we can statistically define anomalies as those data points that are statistically improbable, e.g. $p < 0.0003$.

In the simplest scenario a security analyst may apply **univariate analysis** to examine a single variable in isolation, take for example the CPU load of a database server, to do so we have multiple ways:

- **Grubb's test** The statistical test is used to detect a single outlier in an univariate dataset that follows a normal distribution[1]. It calculates a G statistic based on the difference between the maximum deviation and the mean.
- **Z-score** Modern systems utilizing the Z-score, which normalizes a data point x by subtracting the mean and dividing by the standard deviation $Z = \frac{x-\mu}{\sigma}$. If the Z-score is greater than 3 this indicates that the data point is an outlier.

But, as we learned, univariate analysis is limited, and for modern cybersecurity applications it is insufficient, if we were to analyze data points like "high outbound network traffic", this might be benign, for example we could be doing an update or a backup, but if an attacker were to be using a low speed or low bandwidth exfiltration technique, they would stay within normal bounds ($< 2\sigma$), and even if we direct it to an unusual IP address at an unusual time, we are not analyzing this variable and therefore we ignore that. So multivariate analysis is something that isn't needed but required for cybersecurity applications.[6]

2.2 Multivariate Analysis and Mahalanobis Distance

In network analysis, variables are often correlated, take, for example, the number of packets, it is typically correlated with the number of bytes transferred, if we were to rely on euclidean distance (which assumes independence), we might not see anomalies that arise when the relationship is broken.

For a more practical example, consider a server that typically sends small packets (high packet count but small number of bytes for each packet) or large file transfers (high count and high packet count). If an attacker were to try and perform a buffer overflow attack, he would send a small number of extremely large packets, if we were to map this into an Euclidean space, this point might still be close to the center of a data cloud, but it has almost no **correlation** to usual traffic.

To actually solve this we use **Mahalanobis Distance**, which measures the distance between a point P and a distribution D , effectively it measures how many standard deviations away P is from the mean of D , while also accounting for the covariance among the variables.[6][7]

The Mahalanobis distance D_M of an observation vector $x = (x_1, x_2, \dots, x_n)^T$ from a set of observations with mean vector $\mu = (\mu_1, \mu_2, \dots, \mu_n)^T$ and covariance matrix Σ is defined as:

$$D_M(\mathbf{x}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}$$

where:

- \mathbf{x} is the feature vector of the incoming traffic (packet length, inter-arrival time, port number).
- $\boldsymbol{\mu}$ is the mean vector calculated from the normal training data

The term Σ^{-1} either whitens or decorrelates data, if two variables typically vary together (so we have an high covariance) the matrix Σ captures this. By multiplying the inverse we penalize deviations that go against the known correlation structures.

For example, if Variable A and Variable B usually increase together, a data point where A is high and B is low will result in a large Mahalanobis distance, identifying it as an anomaly even if the individual values of A and B are within their respective normal ranges.

2.3 Payload Detection

Mahalanobis distance has been used to detect anomalies in packet payloads, systems like PAYL compute the byte frequency distribution of incoming payloads and compare them to a historical profile. It is then given a score (referred to as the "distance") and if for example the payload contains a shellcode or exploit, this would alter the natural frequency of bytes found in ASCII text or common protocols, and therefore gives us a higher distance.

The decision is typically governed by a threshold τ :

$$\text{Status}(\mathbf{x}) = \begin{cases} \text{Anomaly} & \text{if } D_M(\mathbf{x}) > \tau \\ \text{Normal} & \text{if } D_M(\mathbf{x}) \leq \tau \end{cases}$$

Because the squared mahalanobis distance of a Gaussian vector still follows a Chi-squared distribution with n degrees of freedom where n is the number of dimensions, the threshold τ can be rigorously derived from Chi-squared tables for whatever significance level is desired.[6]

2.4 Dimensionality Reduction and PCA

A problem that arises when using any high dimensionality dataset, be for statistics or machine learning is the **curse of dimensionality**, as the number of features (dimensions) increases, the computational cost of inverting the covariance matrix (Σ^{-1}) becomes prohibitive ($O(n^3)$), and the data becomes sparse, so distance measuring becomes less impactful.

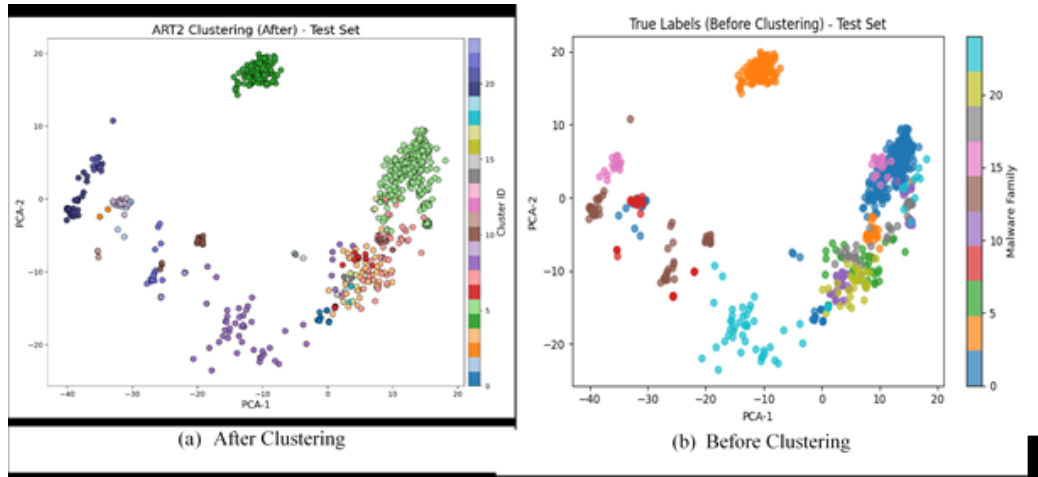
To counteract this we use **Principal Component Analysis** (PCA), which uses an orthogonal transformation to convert a set of data points that are possibly correlated with each other into a set of values of linearly uncorrelated variables which are then called *principal components*. The first principal component has the largest possible variance (so it accounts for as much of the variability in the data as possible), and each succeeding component has the highest variance possible under the constraint that it is orthogonal to the preceding components.

One of the main ways PCA is used in to visualize the clustering of malware families, most malware analysis often involves extracting thousands of features, this would be impossible manually, but even with powerful computers analyzing this raw data is expensive and difficult to visualize without the help of PCA and other dimensionality reduction techniques.

Using PCA we can project high dimensionality data into lower dimensionality spaces (2D or 3D) while preserving the maximum variance. This reveals the underlying clusters of data present in malware families[4], to be more precise:

- **Clustering:** Different Malware families (e.g. ransomware, trojans, worms) often exhibit distinct behavioral patterns. When projected into a lower-dimensional space using PCA, these patterns can form separate clusters, making it easier to identify and categorize malware samples.
- **Anomaly Detection:** PCA can also help identify outliers or novel malware variants that do not fit into existing clusters. These anomalies may represent new threats that require further investigation.

Nowadays, we combine PCA with other techniques and clustering algorithms like K-MEANS and ART2 (Adaptive Resonance Theory) to visualize malware families. The spatial alignment of color-coded clusters in PCA space corresponds to true malware families, proving that the statistical variance capture by PCA corresponds to the function differences in malware code.[5]



As shown in the image above, PCA plots often show lines or trajectories within clusters, these can represent the evolution of a malware family over time (e.g. Racoon Stealer), this allows for the tracking of the genealogy of threats.[4]

2.5 Industrial Control Systems and Operational Technology applications

These statistical methods are particularly useful in ICS and OT, which both prioritize availability over confidentiality, in these kind of networks we can identify 3 major anomalies:

- **Cyber Anomalies:** Old-fashioned intrusions.
- **Operation Anomalies:** System malfunctions or misconfigurations.
- **Service Disruptions:** Availability issues.

The convergence of IT and OT means that industrial networks are now exposed to a wider range of threats. However, OT traffic is often more regular and deterministic than IT traffic (e.g., a sensor polling a valve every 500ms). This regularity actually makes statistical anomaly detection highly effective. Deviations from the strict, deterministic patterns of industrial protocols (like Modbus or DNP3) are easier to detect using statistical baselines than the chaotic traffic of a corporate web network.[8]

References

- [1] VARUN CHANDOLA. “Anomaly Detection: A Survey”. In: (2009).
- [2] SONG WANG JUAN FERNANDO BALAREZO SITHAMPARANATHAN KANDEEPAN AKRAM AL-HOURANI KARINA GOMEZ CHAVEZ and BENJAMIN RUBINSTEIN. *Machine Learning in Network Anomaly Detection: A Survey*. Tech. rep. IEEE, 2021.
- [3] Meera Sridhar Danial Abshari. “A Survey of Anomaly Detection in Cyber-Physical Systems”. In: *ArXiv* 1.1 (18 Feb 2025).
- [4] Walid El-Shafai. “Visualized Malware Multi-Classification Framework Using Fine-Tuned CNN-Based Transfer Learning Models”. In: *Appl. Sci* 14.11 (2021).
- [5] David george. “Static Malware Family Clustering via Structural and Functional Characteristics”. In: *SMU Data Science Review* 7.2 (2023).
- [6] Hamid Ghorbani. “MAHALANOBIS DISTANCE AND ITS APPLICATION FOR DETECTING MULTIVARIATE OUTLIERS”. In: *Ser. Math. Inform.* 32.3 (2019).
- [7] Salvatore J. Stolfo Ke Wang. “Anomalous Payload-based Network Intrusion Detection”. In: ().
- [8] Karel Kuchar and Radek Fujdiak. “Anomaly Detection in Industrial Networks: Current State, Classification, and Key Challenges”. In: *IEEE SENSORS JOURNAL* 25.3 (2025).
- [9] Kathrin Grosse†. Praveen Manoharan†. Nicolas Papernot‡. Michael Backes†. Patrick McDaniel‡. “On the Statistical Detection of Adversarial Examples”. In: *CISPA, Saarland Informatics Campus†; Penn State University‡; MPI SWS* (2024).
- [10] Ali Mohanad Faris Mohammed Q. Mohammed Mohammed G. S. Al-Safi. *Statistical Anomaly Detection for Enhanced Cybersecurity Using AI-Based Wireless Networks*. Tech. rep. Ingénierie des Systèmes d’Information, 2024.
- [11] Palo Alto Networks. “What Is an Adversarial AI Attack?” In: (2024). <https://www.paloaltonetworks.com/cyb-are-adversarial-attacks-on-AI-Machine-Learning>.