

Employing machine learning in a security environment: A data science-driven approach

Table of contents

3 The evolving need for machine learning, AI, and data science

The hype and the reality

Applying machine learning to security

Applying machine learning to user and entity behaviour analytics

7 UEBA use cases for machine learning

Account compromise

Insider threats

Privilege account abuse

Data exfiltration

Potential pitfalls of machine learning

9 Conclusion: Machine learning enables better, smarter, and faster security

10 About LogRhythm

11 Glossary



No matter where you look in the security world today, you'll see the terms machine learning and artificial intelligence (AI). There's been a great deal of interest in machine learning and AI as security vendors and their customers look for better ways to improve their security posture and fight against advancing cyberattacks. Machine learning and AI offer breakthroughs in solving problems in many other areas of our lives, so it's only natural to try to use them to make similar breakthroughs in the field of security.

Unfortunately, there's a lot of hype and misinformation surrounding what machine learning and AI can do to improve security. In this paper, you will discover the most critical things you need to know about applying machine learning and AI in your security environment. You will also learn to recognise the most significant opportunities and challenges for using machine learning and AI to improve your security team's ability to swiftly detect and respond to cyberthreats.

The evolving need for machine learning, AI, and data science

Machine learning, artificial intelligence, and data science are terms with shifting definitions. For the purposes of this paper, we've defined the following terms:

- **Artificial Intelligence (AI):** The science of enabling a computer to automate something a human would normally do that requires intelligence, analysis, and decision making.
- **Machine Learning:** The science of enabling computers to learn without being explicitly programmed to do so. Machine learning applies statistics and algorithms at scale on large amounts of data. One of the goals for machine learning is to achieve artificial intelligence.
- **Data Science:** The discipline of extracting information from data. Data science is a broad field that includes machine learning.



Machine learning has been around for decades, but until recently, it wasn't feasible for most organisations for two reasons:

1. Machine learning needs an incredible amount of computational power in order to apply its algorithms to data and get reliable results quickly.
2. Machine learning requires vast stores of data to mine.

Supervised and unsupervised learning

Machine learning algorithms can evolve through unsupervised or supervised learning. In unsupervised learning, the tuneless algorithm has all the information and context it needs to fully understand the training data provided to it, so it can learn on its own. In supervised learning, the algorithm benefits from additional information and organisational context, either within the training data or provided separately, in order for the machine to get smarter.

Both unsupervised and supervised learning have roles to play in machine learning. Supervised learning is often necessary for data sets with benign anomalies, especially if the intent of using machine learning is to predict future anomalies that aren't benign. For example, in the security field, high-quality training data is difficult to obtain because of the numerous occurrences of false positives and negatives. Humans in the security operations centre (SOC) need to review training data and provide adjustments, such as indicating certain sets of events represent valid security threats, while others are benign. It's expected that humans will always be necessary for supervised learning for security.

Employing machine learning in a security environment: A data science-driven approach

As the cost of storage has gone down, it has become increasingly accessible for data storage needs (e.g., building repositories such as data warehouses or data lakes). Processing power continues to nearly double each year. All of these advances in technology have come together to make machine learning practical and accessible. In fact:

Forrester estimated in a recent report that investments in AI-based technologies, which are driven by machine learning, would triple from 2016 to 2017.¹

Machine learning and AI have become frequent buzzwords in the security space. Security teams have an urgent need for more automated methods for detecting threats and malicious user behaviour—and this need is driving increased interest in these topics. Automation is vital for overwhelmed security teams. This is because prevention measures are not infallible, and many of today's detection methods rely on manual investigation and decision making to find advanced threats, malicious user behaviour, and other serious issues.

Security analysts encounter huge numbers of false positives and negatives. The threat surface has increased exponentially due to the expansion of mobile devices, cloud storage, and the Internet of Things—all of which only increase the number of false positives. Security teams are buried in alarm fatigue. They can't work quickly enough to keep up with the activity to be analysed, or they simply can't identify emerging threats. "Unfortunately, more security doesn't necessarily mean better security. In fact, the current strategy of most organisations—layering on many different technologies—is not only proving ineffective, it is overly complex and expensive. The status quo is not sustainable," says Keith Weiss, head of U.S. software coverage for Morgan Stanley. "Even as companies spend more on security, losses related to cybercrime have nearly doubled in the last five years."² Improving detection means improving accuracy and efficiency, and that requires figuring out how to make detection technologies smarter. That's where AI and machine learning come in.

Machine learning offers far better capabilities than humans can deliver in recognising and predicting certain types of

patterns. Security technologies can use machine learning to identify patterns in their data, enabling them to make decisions and to help humans make decisions faster and more accurately. With machine learning, security technologies can also move beyond rule-based approaches that require prior knowledge of known patterns. For example, security technologies using machine learning can learn the typical patterns of activity within a networking environment to recognise pattern deviations. These departures are possibly indicative of threats and identify these threats earlier in the Cyber Attack Lifecycle. This could prevent many incidents and reduce the impact of others by stopping them sooner.

A recent report from the National Science and Technology Council (NSTC) reiterated this, stating that:

"Using AI may help maintain the rapid response required to detect and react to the landscape of ever-evolving cyberthreats. There are many opportunities for AI, and specifically machine learning systems, to help cope with the sheer complexity of cyberspace and support effective human decision making in response to cyberattacks."³

The effectiveness of machine learning relies on having access to large sets of high-quality, rich, structured data capturing network activities across numerous endpoints. The old phrase "garbage in/garbage out" perfectly explains this situation. If machine learning algorithms ingest data sets that aren't accurate, clear, well-organised, and comprehensive, they're not going to produce the desired results. In other words, just because there are machine learning algorithms in place doesn't necessarily mean what they learn is intelligent and useful. If you teach the algorithms the wrong lessons, they're going to deliver the wrong answers.

¹ https://go.forrester.com/wp-content/uploads/Forrester_Predictions_2017_Artificial_Intelligence_Will_Drive_The_Insights_Revolution.pdf

² <http://www.morganstanley.com/ideas/cybersecurity-needs-new-paradigm>

³ https://obamawhitehouse.archives.gov/sites/default/files/whitehouse_files/microsites/ostp/NSTC/preparing_for_the_future_of_ai.pdf

The hype and the reality

In a perfect world, machine learning would be the silver bullet for defeating your organisation's security challenges. It would enable full automation of security operations, eliminating the need for human involvement. It would learn what every user, system, and application does in incredible detail, enabling immediate identification and handling of user impersonation, malicious intent, and other issues.

Applying AI to security via machine learning is frequently presented as an easy solution. It's not. Contrary to many vendors' claims, no product can do this effectively today. And, it will likely take considerable time and advancement to achieve. Consider the similarities between a SOC that identifies and responds to security incidents and a fraud department that uses fraud analytics techniques to identify and respond to credit card misuse. Even though analysis and identification may be automated, humans are still required to respond and recover (e.g., deciding an issue is a false positive, communicating with the affected people, and coordinating actions with other organisations). Today's security products cannot fully automate the SOC and completely eliminate the need for security analysts, incident responders, and other SOC staff.

While it is not a silver bullet, there is a tremendous amount of value in applying machine learning to solving security challenges. Achieving AI would significantly reduce the mundane work performed by highly skilled and highly paid people. It would also make incident response much faster, effective, efficient, and accurate. However, instead of striving for the unrealistic goal of having AI today, we need to make incremental progress. For instance, applying machine learning pattern recognition to automatically link a threat model from six weeks ago to a similar one today is a realistic goal.

Today, machine learning is most helpful in threat detection by learning the patterns of normal activities and recognising anomalies: the introduction or prediction of a new pattern, a change in an existing pattern, or the removal of a pattern. Given the sheer volume of activities occurring in today's systems and applications, machine learning's pattern recognition and predictive capabilities have become incredibly important.



15–20% of threats are unknown. This is where machine learning comes in.

There is a shortcoming to machine learning, however. Alone, it lacks the understanding of security context to recognise the importance or unimportance of each anomaly. Machine learning can identify that a user is acting in an atypical manner, but atypical behaviour is not necessarily good or bad. For example, a user connecting to a server for the first time might be an anomaly, but is it a malicious act?

In business analytics and other fields, machine learning works well on its own because it looks at anomaly-free data and needs no additional context to predict trends. In the security field, there are many benign anomalies, so the ability to identify anomalies, while important, can't possibly provide the whole explanation of what's happened and enable accurate predictions of what will happen.

It's best to think of machine learning's anomaly recognition capabilities as one of the tools in your toolbox. Imagine that 80-85% of threats are known or recognisable by your security information and event management (SIEM) platform. If this is true, 15-20% of threats are unknown, and therefore unrecognisable, by your SIEM. This is where machine learning comes in. You need the right tool for the right job.

To effectively detect threats, you need to employ the correct algorithm for that threat type. The rest of your tools provide the security context and relevancy. A SIEM solution can integrate and correlate information from many tools, such as human resources (HR) systems, identity management solutions, vulnerability scanners, and asset-management systems. When used together, machine learning and the other tools generate the risk information needed to prioritise human actions. Without prioritisation, there are so many anomalies that it's impossible to examine them all and find the truly important ones.



Benefits of machine learning

- Pattern learning
- Anomaly recognition
- Predictive capabilities

Applying machine learning to security

Machine learning can provide many potential opportunities to improve your security operations, such as:

Threat detection

- *Threat prediction and detection*: analysing anomalous activity in order to recognise emerging threats so they can be stopped before attackers achieve their intended results
- *Risk management*: monitoring and analysing user activity, asset contents and configurations, network connections, and other asset attributes to create and maintain dynamic risk profiles for all enterprise assets
- *Vulnerability information prioritisation*: using learned information about the organisation's assets and the vulnerabilities being actively exploited to prioritise their mitigation
- *Threat intelligence curation*: refining the information within threat intelligence feeds to improve quality

Threat response and recovery

- *Event and incident investigation and response*: reviewing and analysing information on events and incidents in order to identify next steps and organise the incident response processes and workflow, such as selecting and implementing the appropriate incident playbook
- *Forensics*: supplementing existing forensics information by identifying additional information likely to be related and potentially worth investigating
- *Deception and misdirection*: learning about the existing environment and developing intelligent techniques for deceiving and misdirecting attackers so they won't accomplish their goals



UEBA is a perfect application for machine learning as long as the necessary security context is available for understanding the significance of each anomaly.

Applying machine learning to user and entity behaviour analytics

Threat prediction and detection is a critical area of security that can benefit from machine learning. Consider the challenges in performing user and entity behaviour analytics (UEBA). Gartner defines UEBA as such:

User and entity behaviour analytics (UEBA)

Profiling and anomaly detection based on a range of analytics approaches, usually using a combination of basic analytics methods (e.g., rules that leverage signatures, pattern matching and simple statistics) and advanced analytics (e.g., supervised and unsupervised machine learning). Vendors use packaged analytics to evaluate the activity of users and other entities (hosts, applications, network traffic, and data repositories) to discover potential incidents.⁴

UEBA is a perfect application for machine learning as long as the necessary security context is available for understanding the significance of each anomaly. Machine learning can make UEBA considerably more effective for the following reasons:

- Machine learning can handle the volumes of data to be analysed and the environment to be understood. This includes being able to incorporate many types of data sets, from network traffic patterns and application data to records of user authentication attempts and user access to sensitive data.
- Machine learning-driven UEBA is well suited for identifying "qualified" threats—those that are legitimate and require action. Machine learning can take many more factors into consideration than humans can when looking at potential threats, and it can do so in near real time.
- Machine learning-driven UEBA is able to identify the threats that are hardest to find, such as insider threats, privileged account takeovers, and unknown threats by recognising shifts in behaviour. An organisation can use machine learning to dynamically identify asset risks. Machine learning-enhanced UEBA can leverage that risk information to identify new activity that conflicts with expected patterns, such as a low-risk user suddenly connecting to a high-risk system and transferring large amounts of data from it to a laptop.

⁴ <http://blogs.gartner.com/anton-chuvakin/2016/12/12/ueba-clearly-defined-again/>

UEBA use cases for machine learning

User and entity-based threats are a growing concern for security teams and, therefore, a growing need in the market. According to a recent Verizon Data Breach Incident Report, 63 percent of confirmed data breaches involved attackers posing as legitimate users (using stolen access credentials) or legitimate users maliciously exploiting their access.⁵ But to detect insider threats, your technology must first be able to understand and baseline user behaviour. And it must do so while reducing false-positive alarms in order to pinpoint a threat for effective detection.

This is where machine learning begins to provide real value. By establishing baseline behaviours and patterns, then detecting anomalies by combining statistical models, machine learning algorithms, and rules, a UEBA solution can compare incoming transactions with the existing baseline profile. It can then present analytic results to highlight patterns of unauthorised access and users, allowing the team to act upon the infractions—either through manual decisions or automated actions.

Account compromise

A UEBA solution should be able to easily detect if a hacker has accessed a network user's credentials, regardless of the attack vector or malware used. This includes the detection of attacks such as pass the hash, pass the token, and brute force attacks. For successful account compromise detection, the technology will need to recognise indicators of compromise across any asset the user touches (including endpoints and networks).



Potential indicators of account compromise

1. Unusual authentication patterns (e.g., dormant account access)
2. Lateral movement following an attack
3. Concurrent logins from multiple locations
4. Account activity from blacklisted locations

Attack terms defined:

Pass-the-hash or pass-the-token attacks

Password attacks, such as password guessing or password cracking, are time-consuming attacks. Tools that make use of the precomputed hashes reduce the time needed to obtain passwords greatly. In the pass-the-hash attack, the goal is to use the hash directly without cracking it. This makes time-consuming password attacks less needed.

Brute force attacks

Brute force, also known as brute force cracking, is a trial-and-error method used by application programs to decode encrypted data, such as passwords or data encryption. Standard (DES) keys, through exhaustive effort (using brute force) rather than employing intellectual strategies. A brute force cracking application proceeds through all possible combinations of legal characters in a sequence.

Insider threats

An insider threat is one of the primary UEBA drivers for many security teams because of the lack of confidence around accurately detecting when an insider threat is occurring. These threats include malicious insiders, compromised insiders, and negligent insiders, and they often result in destruction of data, threat of data, and so forth. This is an area UEBA solutions shine. By establishing baseline behaviour for your users, the solution should be able to detect and alarm on unusual, high-risk behaviour that falls out of that baseline profile based on several factors, including time, source host, and location.



Potential indicators of insider threats include

1. Deviation from peer group
2. New or unusual system access
3. Unusual login times
4. Disabled account logins
5. Unusual file access and modifications
6. Abnormal password activity
7. Excessive authentication failures
8. Multiple account lockouts

⁵<http://www.darkreading.com/endpoint/verizon-dbir-over-half-of-data-breaches-exploited-legitimate-passwords-in-2015/d/d-id/1325242>

Privileged account abuse

A UEBA solution should be able to identify specific attacks on privilege users who have access to sensitive information by detecting compromised credentials and lateral movement to the systems that contain this privileged data. Defining and maintaining a list of privileged users and groups can help your UEBA solution to validate permission changes and quickly disable accounts with observed privilege escalation.

In addition to privileged accounts, you'll also want your UEBA solution to monitor when sensitive, high-value assets are accessed. By identifying and assigning threat risk levels, your UEBA solution should be able to monitor high-profile or high-value assets to generate high-priority alarms for your security team.

Your UEBA solution should also be able to monitor for other indicators of risk, such as account lockouts, new account creation, account sharing, and accounts that have gone dormant.



Potential indicators of privilege account abuse include:

1. Suspicious temporary account activity
2. Abnormal account administration
3. Unusual privilege escalation

Data exfiltration

A UEBA solution should also be able to monitor and alarm your team on indicators that data exfiltration appears to be happening. This should occur in as real time as possible so that your team can investigate and stop the exfiltration before damage occurs. This is where automated responses can be extremely valuable in lowering your team's mean time to respond—ultimately protecting your organisation from a high-profile data breach.



Potential indicators of data exfiltration include:

1. Suspicious data transfers
2. Malicious payload drops
3. Abnormal traffic patterns
4. Blacklisted communication

Potential pitfalls of machine learning

Using machine learning to improve your organisation's security is a great idea, but there are some potential pitfalls of which you should be aware. Every product that implements machine learning does it in different ways—some great and some not so great. As a prospective consumer of security technologies, you should be aware of potential pitfalls so you can ask vendors the right questions and make an educated decision as to which products are using machine learning effectively.

Some potential pitfalls to keep in mind include:

- 1. Poor data quality.** Machine learning can't magically transform bad inputs into good outputs. The security solution should be able to convert a wide breadth of security data into a standard format for forensic visibility. Even a single activity monitored by multiple systems (in some cases systems from the same vendor) could be recorded and detailed quite differently by each of them. This is why it's important to have sufficient detail and consistency in the data that's been prepared for analysis. Machine learning algorithms need rich, structured data that's well organised. Otherwise, results will likely be unpredictable and AI will not be achieved.
- 2. Lack of context.** Using machine learning to look at security events alone, without any additional contextual information, will produce much less accurate and less useful results. Taking contextual information into account changes the understanding of risk, providing greater clarity to the nature and severity of activities. Simple examples of this context are understanding:
 - a. The purpose of a targeted server application
 - b. The identity of a user conducting suspicious activity
 - c. The reputation of an external IP address being connected to one of the organisation's laptops

Strong machine learning can ingest many forms of context and incorporate that information into its analysis and decision-making processes. Visualisation is a critical, customisable layer of machine learning. Without visualisations that accurately represent the threat, the analyst lacks the tools and insight they need for incident response and investigation.

3. Insufficient training data. Machine learning algorithms need high-quality training data provided via supervised learning. Such data is key for achieving precise threat detection. Without it, algorithms won't necessarily know what's good or bad, or which anomalies are more important to look at than others. Machine learning training data must also take contextual information into account.

4. Failure to produce useful information. Another important factor to consider is how well the solution learns on an ongoing basis. Will the machine continue to learn? Will it adapt to changes in the environment over time and continue to output high-quality, useful, security-relevant information? Just because a product uses machine learning doesn't mean it's achieving security value and improving SOC effectiveness.

5. Lack of focus on machine learning value. Many vendors of machine learning-based security products talk extensively about their machine learning algorithms. Although this information is often fascinating, it's not necessarily an indicator of how well the technology will perform. Algorithms are actually a relatively easy component of machine learning. Prospective consumers of security products should focus on the value that the technology can provide instead of the details of how the technology functions. The value of machine learning should be measured by how much an organisation's detection and response times improve once it's in use.

When used effectively, machine learning could help your team:



- Detect hidden threats and false positives
- Accelerate incident response
- Streamline SOC operations to reduce mean time to detect and respond to threats



Conclusion: Machine learning enables better, smarter, and faster security

Machine learning offers a great deal of promise in improving security by greatly reducing human effort and lowering the time to detect, respond to, and recover from incidents. Challenges in successfully using machine learning for security include the following:

- Giving machine learning real-time access to large sets of high-quality, rich structured data encompassing all security-related events from throughout the enterprise
- Providing machine learning with the contextual information necessary to understand the meaning and importance of each observed activity and detected anomaly
- Performing supervised learning with extensive sets of high-quality training data to educate the machine on which activities are good and which are bad

If your organisation is struggling to stay ahead of cyberthreats due to a shortage of resources and the costs of inefficient and manual workflows, machine learning could be a helpful technology for you. Machine learning could allow your analysts to focus on the problems that require intuition and creativity. It could also help your security operations scale as threats continue to evolve.

When used effectively, machine learning could help your team:

- Detect hidden threats and minimise false positives
- Accelerate incident response
- Streamline SOC operations to reduce mean time to detect and respond to threats

Machine learning can enable your team and technology to be better, smarter, and faster by having advanced analytics at its fingertips to solve real problems—like detecting user-based threats such as UEBA—quickly.

But machine learning is not a silver bullet, and no application of this technology will solve all your security problems. There are limitations and pitfalls you must be aware of when researching a potential machine learning solution to apply to your security environment.

To learn more about how LogRhythm can make your security smarter with artificial intelligence, visit us at [LogRhythm.com/CloudAI](https://www.logrhythm.com/CloudAI)

About LogRhythm

LogRhythm is the pioneer in Threat Lifecycle Management™ (TLM) technology, empowering organisations on six continents to rapidly detect, respond to and neutralise damaging cyberthreats. LogRhythm's TLM platform unifies leading-edge data lake technology, artificial intelligence, security analytics and security automation and orchestration in a single end-to-end solution. LogRhythm serves as the foundation for the AI-enabled security operations centre, helping customers secure their cloud, physical and virtual infrastructures for both IT and OT environments. Among other [accolades](#), LogRhythm is positioned as a Leader in Gartner's SIEM Magic Quadrant.



Glossary

Alarm Fatigue:

Security teams encounter high numbers of false positives and negatives. They can't work quickly enough to keep up with the activity to be analysed, making it difficult to identify emerging threats.

Artificial Intelligence (AI):

The science of enabling a computer to automate something a human would normally do that requires intelligence, analysis, and decision making.

Brute Force Attacks:

This trial-and-error method is used by application programs to decode encrypted data through an exhaustive effort of an application proceeding through all possible combinations of legal characters in a sequence.

Cyber Attack Lifecycle:

When a threat breaches a network, it undergoes a process beginning with initial intrusion and ending with a final attack execution. A threat must be detected and remediated in the early stages to minimise impact. This process is known as the Cyber Attack Lifecycle (also referred to as the Cyber Kill Chain). Phases of the Cyber Attack Lifecycle include Reconnaissance, Initial Compromise, Command & Control, Lateral Movement, Target Attainment, and Exfiltration/Corruption.

Data Science:

The discipline of extracting information from data. Data science is a broad field that includes machine learning.

Insider Threat:

This type of threat comes from the people within an organisation who have access to data and computer systems or information of the organisation's security practices, whether intentionally malicious or not.

Machine Learning:

The science of enabling computers to act without being explicitly programmed to do so. Machine learning applies statistics and algorithms at scale on large amounts of data. One of the goals for machine learning is to achieve artificial intelligence.

Mean Time to Detect (MTTD):

The average time it takes to recognise a threat requiring further analysis and response efforts.

Mean Time to Respond (MTTR):

The average time it takes to respond and ultimately resolve an incident.

Pass-the-Hash (Also Known as Pass-the-Token) Attacks:

In a pass-the-hash attack, the attacker uses a hash without cracking it to obtain passwords faster than the usual time-consuming efforts.

Security Operations Center (SOC):

A centralised unit that deals with security issues on an organisational and technical level. A SOC within a building or facility is a central location from where the staff supervises the site, using data processing technology.

Structured Data:

Structured data refers to information with a high degree of organisation, such that inclusion in a relational database is seamless and readily searchable by simple, straightforward search engine algorithms or other search operations.

Supervised Learning:

In supervised learning, the machine learning algorithm requires additional information and organisational context—either within the training data or provided separately—in order for the machine to get smarter.

Unstructured Data:

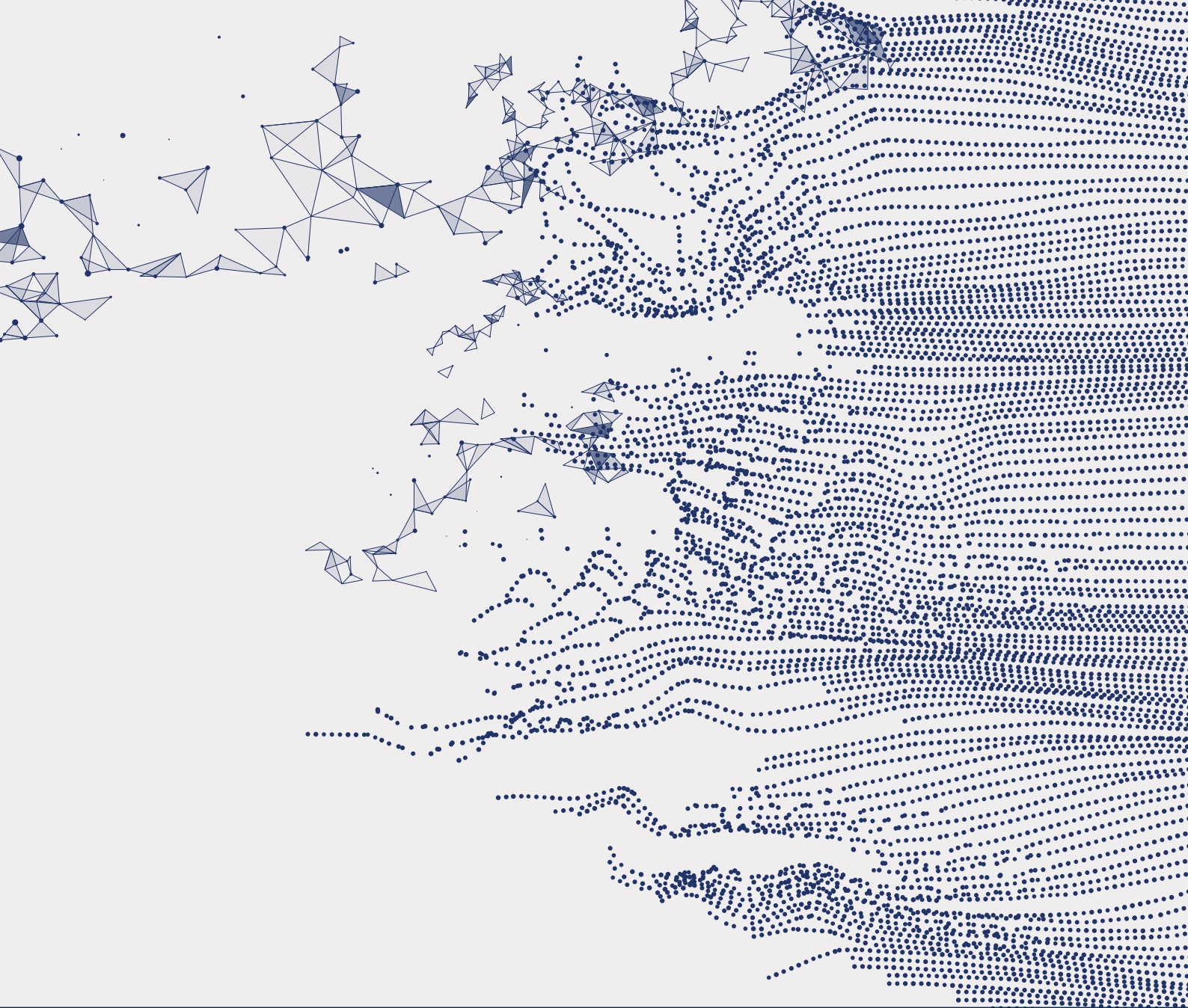
Information that either does not have a predefined data model or is not organised in a predefined manner.

Unsupervised Learning:

In unsupervised learning, the machine learning algorithm is tuneless. It has all the information and context it needs to fully understand the training data provided to it, so it can learn on its own.

User and Entity Behaviour Analytics (UEBA):

UEBA combines basic analytics methods and advanced analytics to detect anomalous behaviour of users and entities within a network.



+44 1628 918300
europe@logrhythm.com
Clarion House, Norreys Drive, Maidenhead, SL6 4FL
United Kingdom

 **LogRhythm**[®]
The Security Intelligence Company