# Artificial Intelligence and Machine Learning Case Studies

Aaron Turner, IANS Faculty

# Agenda

- Why Artificial Intelligence?

- Machine Learning and Artificial Intelligence

- Analytics and Machine Learning in Security

- Security Analytics Tools

- Recommendations

**IANS**

# Why Artificial Intelligence?

# Audience Questions

- How do you define Machine Learning (ML) / Artificial Intelligence (AI)?

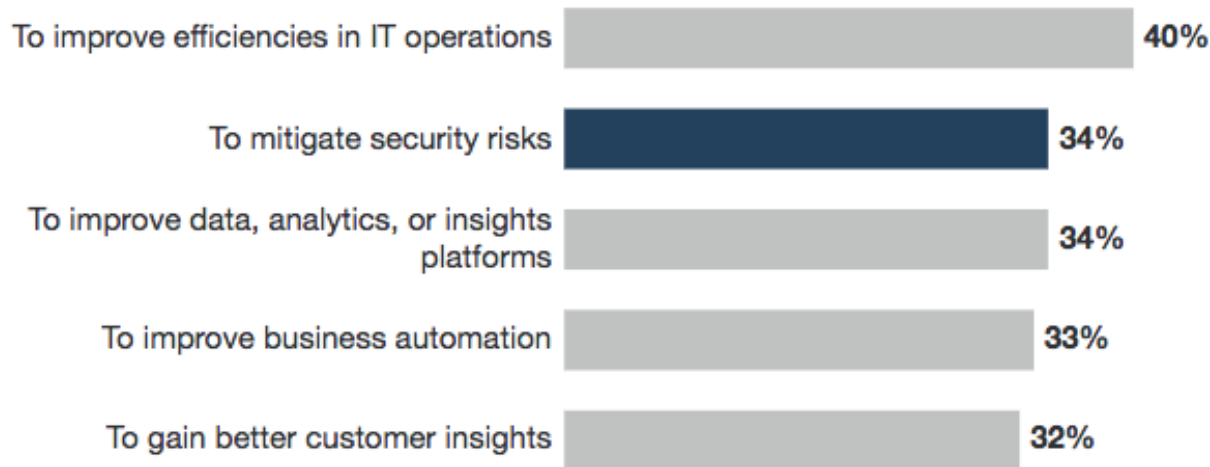- Who has a ML project / budget?

- What are you doing with ML / AI?

# The Promise of Artificial Intelligence (AI)

- Analyze massive amounts of data

- Address the skills gap - giving your workforce better tools

- Constantly adapt to a changing threat landscape and attack patterns

- Learn from user feedback and actions (remediations, triage, etc.)

- Make analysts more effective leading to better and faster decisions

  - Automation where possible

  - Detection systems with less false positives

  - Tools that provide more context, better visualizations, etc.

**IANS**

5

# What Are Companies Doing With AI?



FIGURE 1 Firms Plan To Use AI To Mitigate Security Risks

Top five use cases/application scenarios firms are planning to use or are currently using artificial intelligence technologies for

| Use case | Percentage |
|---|---|
| To improve efficiencies in IT operations | 40% |
| To mitigate security risks | 34% |
| To improve data, analytics, or insights platforms | 34% |
| To improve business automation | 33% |
| To gain better customer insights | 32% |

Source: Forrester Data Global Business Technographics® Data And Analytics Survey, 2017

IANS

6

# Machine Learning and Artificial Intelligence

*"Everyone calls their stuff 'machine learning' or even better 'artificial intelligence' - It's not cool to use statistics!"*

*"Companies are throwing algorithms on the wall to see what sticks - see security analytics market"*

# Machine Learning & Artificial Intelligence

- **Machine Learning (ML)**
  - Algorithms ways to 'describe' data
  - **Supervised**
    - We are giving the system a lot of training data and it learns from that
  - **Unsupervised**
    - We give the system some kind of optimization to solve (clustering, dimensionality reduction)

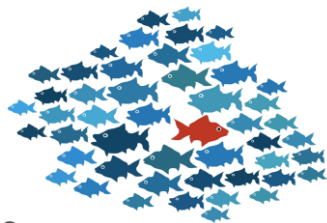- **Anomaly Detection (Outlier Detection)**
  - Can be done with ML but simple statistics often work much better
  - Statistical outliers are hardly ever security relevant
  - 2 decades of anomaly detection research in security!

- **Deep learning**
  - Is just another ML algorithm - significantly improved results for classification problems
  - Basically eliminates the feature engineering step

- **Artificial Intelligence (AI)**
  - *"A program that doesn't simply classify or compute model parameters, but comes up with **novel knowledge** that a security analyst finds insightful."*

# Some Other Analytics Concepts

- **Predictive Analytics**
  - Make statements about future or unknown events using methods like machine learning, etc.
  - In security some simple approaches exist that try to predict future attacks (you know how hard security is!)
    - Generally identify patterns of suspicious behavior to indicate that something might soon go wrong
  - Can we look at the kill-chain and connect a threat actor across the different steps?

- **Natural Language Processing (NLP)**
  - How to process language data / speech
    - Understand words (syntax parsing)
    - Extracting meaning (semantics)
  - Applications: DGA detection, analyzing threat reports, analyzing emails (SPAM, phishes), source code analysis
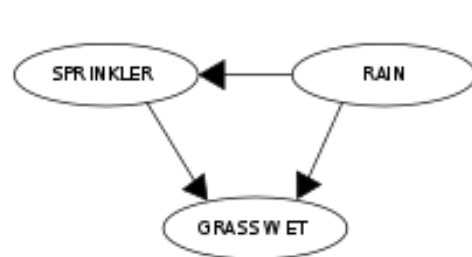
- **Expert System**
  - 'If-then' rules
  - Knowledge represented as facts and rules
  - Inference engine applies the rules to the known facts to deduce new facts.

- **Belief Networks** (extensions of expert systems)
  - Probabilistic graph model describing knowledge
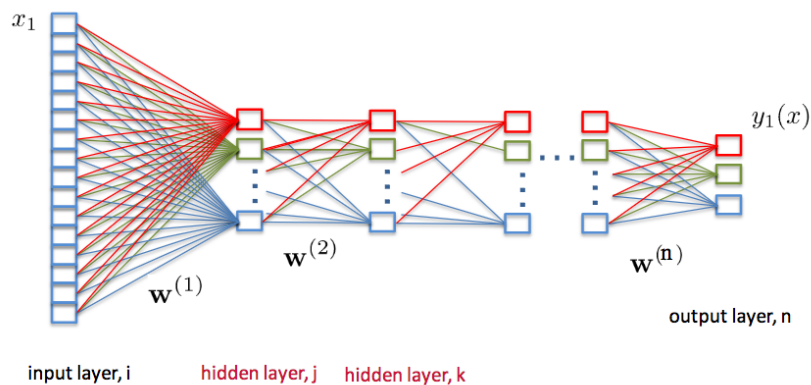  - Used to model export knowledge (e.g., tier-1 analyst automation)

|  | SPRINKLER | |
|---|---|---|
| RAIN | T | F |
| F | 0.4 | 0.6 |
| T | 0.01 | 0.99 |

|  | RAIN | |
|---|---|---|
| | T | F |
| | 0.2 | 0.8 |

A simple belief network

| | | GRASS WET | |
|---|---|---|---|
| SPRINKLER | RAIN | T | F |
| F | F | 0.0 | 1.0 |
| F | T | 0.8 | 0.2 |
| T | F | 0.9 | 0.1 |
| T | T | 0.99 | 0.01 |

# Deep Learning - Details

- Deep learning performs much better than traditional neural networks (see graph)

- Deep learning has more complex networks (many hidden layers, fully connected neurons)

  - Possible due to progress in hardware technology (GPUs, FPGAs, etc.)

- Automatic feature engineering (the input to machine learning algorithms)

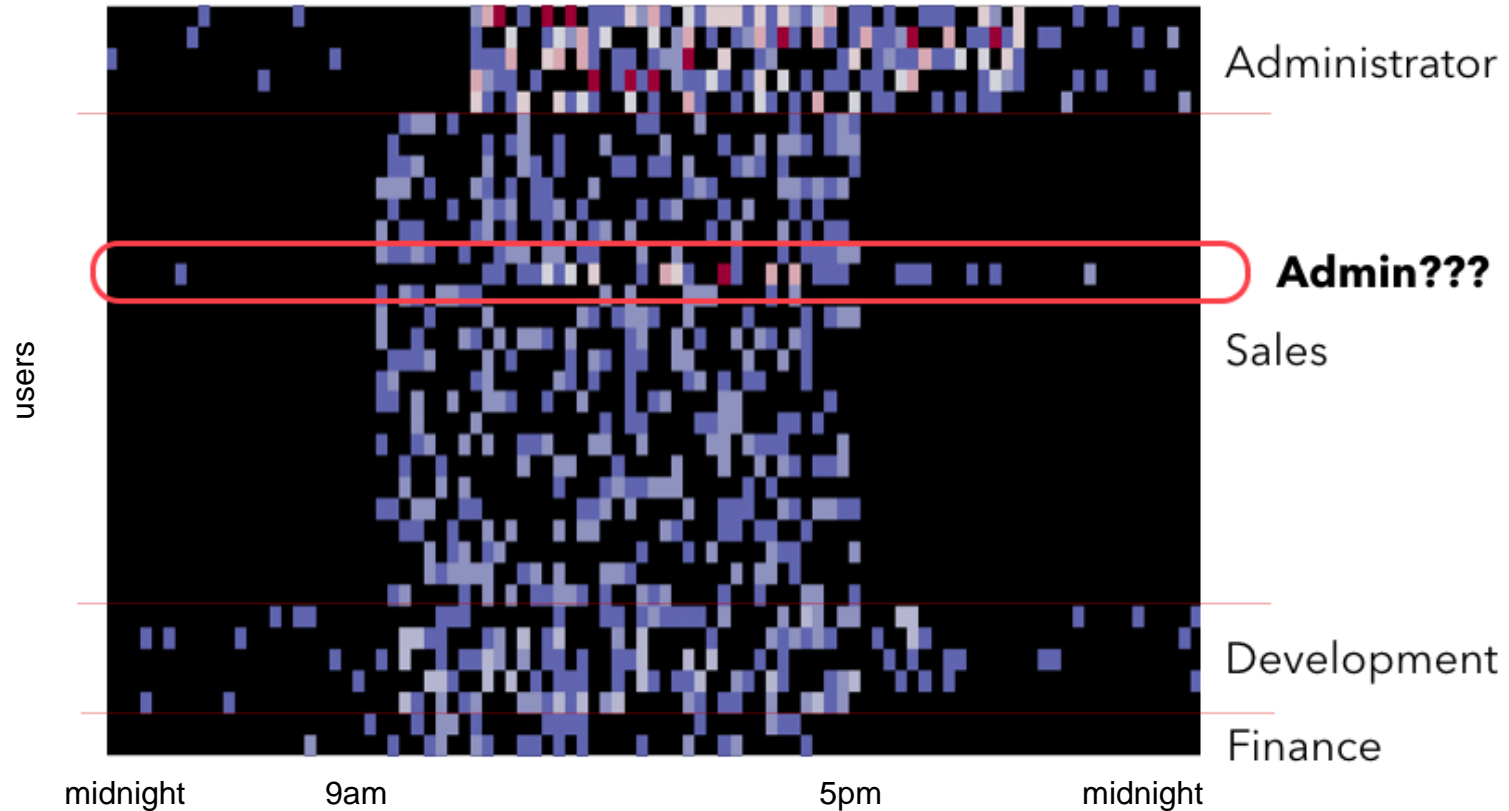  - Overall works well with large amounts of training data

11

# Statistics, Rules, and Models

- Threshold-based systems
  - Simple metric thresholds, moving averages, etc.
  - Doesn't deal with outliers
  - Doesn't adapt well to a changing baseline
  - Simple to implement - your SIEM can do this
- Rules
  - A way to model expert input or scenarios (if-then-*else, correlations, …)*
  - Cannot always express complexity of scenarios
  - Do not adapt well to changing scenarios and inputs
  - Supported by any SIEM product
- Models
  - *What a machine learning algorithm 'learns'*
  - *Automatically learns 'thresholds' or parameters and updates them over time*
  - For example authentication behavior - what time of the day are users active?
- Models from ground truth
  - Use ground truth (e.g., an incident) and learn what the factors are that make up the incident
  - For example, for a successful attack, learn what IDS alerts, what logs (activity) lead to the attack

**IANS**

# Rules versus Models - An Example

- User profiling - detect suspicious or malicious activity from users
- With **rules**:
  - *"Alert if 'sales' user active before 9am or after 5pm'*
  - Problem: way too many false positives!
- Build a **model for each user**:
  - *Learn what the normal time is for a user.*
  - Problem: Taking into account exceptions like travel, vacation, insomnia, etc.
  - How do we get 'clean' training data for users?
- Build a **model for a group of users** (e.g., sales):
  - *Learn what the normal times of activity are for users. Helps model some 'global' phenomena such as holidays.*
  - Problems: Not all user groups are homogeneous (e.g., sysadmins)
- Improvement:
  - *Add domain knowledge, such as vacation (from HR system), etc.*
- Can we do better?
  - ***Ensemble models*** *- use multiple models at the same time*
  - Each component contributes to an 'anomaly' score
  - Look at individual users, their peers, detected cohorts, etc.

**IANS**

13

# Visualization For Model Creation



users

Administrator

Admin???

Sales

Development

Finance

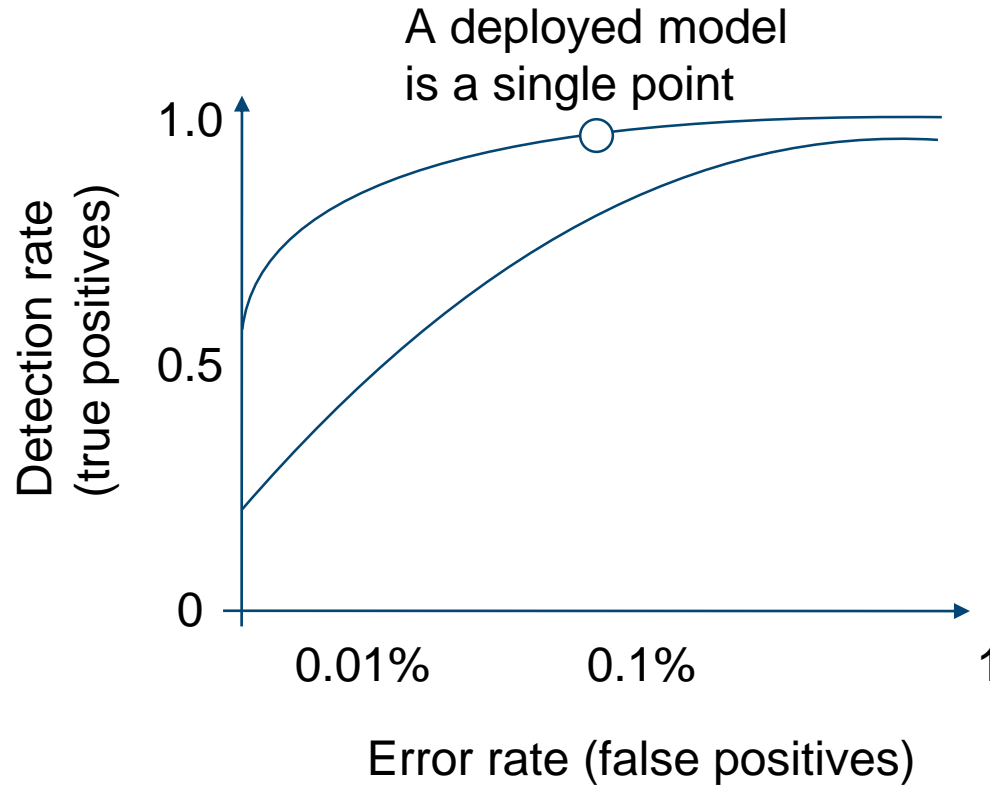midnight     9am     5pm     midnight

14

# ML & AI Challenges

- Explainability
  - What did the system just learn?
  - Especially with neural networks and deep learning
- Verifiability
  - Did the system just learn the correct thing?
  - How do you assess false positives and negatives of your anomaly detection tool?
- Data and Context
  - Need 'clean' data
  - Need contextual information for the data, such as machine roles, user groups, input from HR systems, etc.
  - Even deep learning needs well engineered features! (e.g., malware detection)
- "Technical" challenges
  - Parameter choices, such as distance functions
  - Algorithmic approach (drop outs, etc)

**IANS**

# Some Important Principles

- False negatives are very expensive

    - Could cause arbitrary damage to our environment by not detecting attacks

- False positives are expensive too

    - Analyst time is valuable

- Alerts should make **sense** to a human

    - False positives + inexplicable results → signal fatigue

**IANS**

# The ROC Curve

A deployed model
is a single point

- low detection,
  low false positives

1.0

Detection rate
(true positives)

0.5

better

0

0.01%          0.1%              1

Error rate (false positives)

# Analytics + ML in Cyber Security

# Frustrations That ~~AI~~ Analytics Should Solve?

- Better understanding of all security data (e.g. assist our hunters)

- Find security problems (anomaly detection)

- Prioritization of data (e.g. helping with alert triage)

- Reduce false positive (address alert fatigue)

- Improve analyst efficiency

- Increase retention for security analysts (automate boring tasks)

- Retain expert knowledge (document / capture tribal knowledge)

# What Industry Analysts Say (Forrester)

*[Vendor] conversations that begin with "We have the best data science" are not helpful.*

- Data science in security is as old as security itself

- … not a panacea for the prevention of all cyberattacks …

- Useful for recognizing

  - **Patterns** in large quantities of data

  - Informing decision making as a **supplement** to rules-based or signature-based detection

- What you should do

  - Ignore vendor claims about data science, and concentrate on **use cases**.

  - Ask for **referenceable customers** in your industry

  - Challenge them to prove the use-cases, preferably on your own data

**IANS**

# ML In Security

# ML for Malware Detection - Some History

- Starting with signature based approaches
  - Polymorphic malware becomes common
- AV responds by building decision trees
  - Malware authors respond by encrypting
- AV responds with software emulators
- By 2005 – AV companies with just 'check sums' (signatures) were dead
- By 2015 decision trees are beginning to fail

Seen on 100 network workstations

NO — Has suspicious headers
YES — Has suspicious headers

Has suspicious headers:
- NO → Legitimate
- YES → Malicious

Has suspicious headers:
- NO → Legitimate
- YES → Is digitally signed

Is digitally signed:
- NO → Contains Encrypted Data
- YES → Contains Encrypted Data

Contains Encrypted Data:
- NO → Legitimate
- YES → Malicious

Contains Encrypted Data:
- NO → Legitimate
- YES → Malicious

# Machine Learning in Malware Detection (~ 2012)

- Detection rate was pretty good

- But legitimate software gets identified as malware too often
  - False Positives

### Machine learning model evaluation

| Malware | Legitimate software |
|---|---|
| **94.3%**<br>True positive<br>Real malware detected | **8.1%**<br>False positive<br>Legitimate software that is detected as malware |
| **5.7%**<br>False negative<br>Undetected malware | **91.9%**<br>True negative<br>Legitimate software that is detected as good |

Source – 'Analysis of Machine Learning techniques used in behavior-based malware detection' 2009
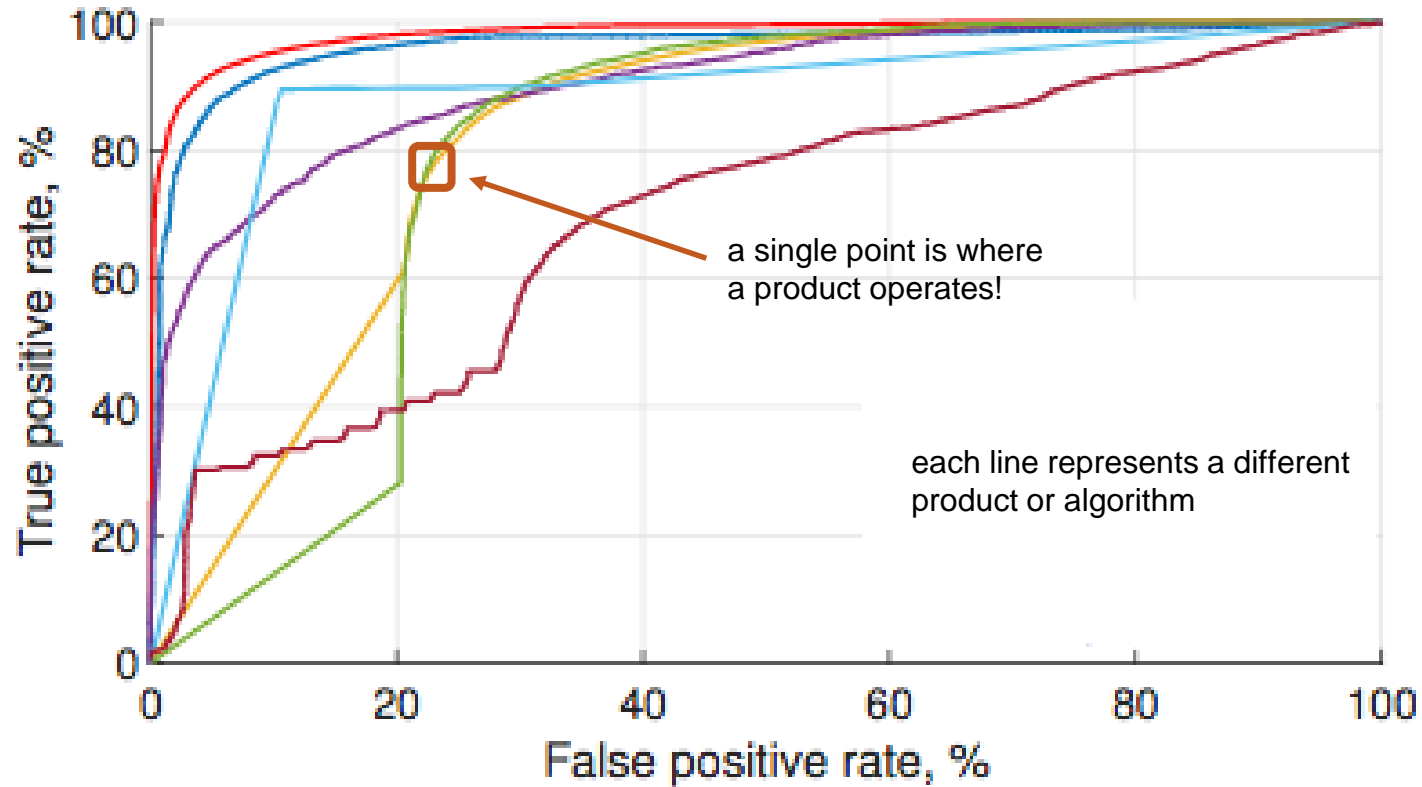
IANS

# Deep Learning For Malware Detection

- Self learning
  - No feature definition necessary

- Intrinsically scales
  - Hundreds of millions of malware samples used in training
  - Adding over 400K per day
  - Detects and stops threats within 20-100 milliseconds
  - Models are about 10-20 MB (Traditional ML models can get huge 500 MB-10 GB)

- Unparalleled accuracy
  - Proven ability to detect never before seen malware without signatures

Machine learning model evaluation (Deep learning)

| Malware | Legitimate software |
|---|---|
| >99%<br>True positive<br>Real malware detected | <1%<br>False positive<br>Legitimate software that is detected as malware |
| <1%<br>False negative<br>Undetected Malware | >99%<br>True negative<br>Legitimate software that is detected as good |

Source – Sophos Labs testing of new executable malware, 2017

**IANS**

# The ROC Curve



a single point is where
a product operates!

each line represents a different
product or algorithm

IANS

# Other ML Uses in Security

- MLSec
  - Looking at firewall "block" data across a large number of networks to find likely attackers early
- DNS Analytics
  - Co-occurance, domain name classification
  - DNS lookup analysis (frequency)
- Threat Intelligence Feed Analysis
  - IOC prioritization, de-duplication, campaign association, removing false positives
- URL Analytics
  - Identify malicious URLs
  - Turns out, you have to analyze the content of the Web site behind the URL as well
- Security Analytics Solutions (see later)
  - Risk scoring



www.mlsecproject.org

# Don't Use Machine Learning and Deep Learning If …

- Not enough or no quality **labeled data**
  - Don't use for network traffic analysis - you don't have labeled data - really, you don't!
- No well trained **domain experts** and **data scientists** to oversee the implementation
  - Not enough domain expertise to engineer good **features**
- Need to understand what ML actually learned (**explainability**)

Also remember:

- Data cleanliness issues (timestamps, normalization across fields, etc.)
- Operational challenges  (scalability and adaptability) of implementing machine learning models in practice

**IANS**

# What Do Security Analytics Tools Do?

# Security Analytics - A Set of Products

Attackers are using 'allowed' channels and mask in benign looking activity that traditional security tools cannot detect.

## User and Entity Behavior Analytics (UEBA)

- Identify anomalies based on user and/or machine behavior.
- Most vendors don't use real machine learning, don't fall for snake oil – ask for real-world proof
- Two groups of products: based on logs or based on network traffic

## Automation & Orchestration

- Sit on top of SIEM (and some other data) to close the loop of a) prioritizing important attacks and b) automating response.

## Hunting

- Enable senior security analysts to explore data within a SIEM or big data store to find environment specific attacks and breaches.

All these products are really features of a larger platform:

- They should all be under one single product
- If they are sold as individual products, make sure they interoperate well. Where is the data stored? etc.

**IANS**

# User and Entity Behavior Analytics (UEBA)

- Risk scoring of entities (devices, users)
  - List of top suspicious entities
- Anomaly detection for entities

- Mutli-vector approach for risk scoring
  - Never seen before
  - Cohort behavior
  - Group behavior
  - Hard-coded known bad (countries, etc.)
  - …
- Bayesian Belief Networks (BBN)

Not all anomalies are security problems or attacks

- Former employee requests an authorization token
  - Account revocation bug? Attack?
  - Nope: username typo
- Actor fails authentication 20K times
  - Brute-force attack?
  - Nope: actor changed password, forgot to update script
- Email address in RPC to location service
  - Privacy violation?
  - Nope: address is "test@123.com"

**IANS**

# Recommendations

# Practical Considerations for Analytics / ML / AI Projects

- What were the use-cases you wanted to cover?
  - Lateral Movement detection, Exfiltration detection, C2 detection, DAG detection, etc
- Do you have the **right data**?
  - Logs
  - Access to taps / SPAN ports to intercept network traffic
- Do you have **context** for data and do tools and processes incorporate it?
  - Do yo have a **dynamic asset inventory** that can be integrated? Solve this problem first!
  - What other contextual data feeds do you have and would be useful?
- What is the process to deal with alerts?
  - Manual Automation / orchestration capability?
- How do you **capture** expert **knowledge**?
  - Manual entry of rules? How do you verify?
  - Collaboration with others?
- Figure out how to **share your models**
  - STIX technically supports that, but nobody is doing it

**IANS**

# Practical Considerations Buying a Product

- Does the solution really detect **behavioral anomalies**?
- Does the solution provide policy **enforcement** features?
- Does the solution **integrate** with the rest of your infrastructure (e.g. SIEM)?
- How does the solution affect employee **experience**? Does the product learn from user input / feedback
- Does the product deal with **containers, VMs,** and the **cloud** ?
- **How long** does it take to begin recognizing suspicious patterns? How long does it take to establish a baseline?
- How does the solution adapt to completely **novel attacks**?
- Ask for **results** that have been seen in actual customer environments
  - What data does the solution work on best and have you used the tool in companies of my industry?
  - Do you have metrics on the improvement in capabilities (e.g., detection, analysis, prioritization, investigations, response)?
- Do a **PoC on your network** to learn
  - How hard it is to **install** the product and how much time does it take to **tune**
  - How much time it will take on **ongoing** maintenance
  - What does it actually **detect** in your environment?
    - Are all the detections trivial? Or could they be modeled in your SIEM?
- For log-based SA tools, **authentication** logs are most useful; then **proxy** logs.
  - Others are harder to collect and not that useful

**IANS**

# Analytics - Do It Yourself

- Do you have enough **expertise** to tackle some of the use-cases in house?
- **People**
  - Data scientists to build models
  - Data scientists that understand security
  - Security engineers that can help build and validate the models and provide security expertise
- **Infrastructure**
  - Necessary data is centralized and easily accessible
  - Backend is in place that allows for running rules, models, etc. on all the data necessary
  - Make sure you can run these things on Splunk / your SIEM!
- **Algorithms**
  - Simple works better (for example monitoring counts over time)
  - Don't start with choosing an algorithm - EVER
  - Always identify the use-cases, the data, and then figure out what algorithm helps most

**IANS**

# Finally

# Action Plan

- Define your **use-cases** first - understand where you want and should use ML
  - Define a **holistic** approach (NIST framework? visibility, …)
  - Make sure you can retain your **export's knowledge** in case they should ever leave
  - **Collaborate** with your peers on use-cases and solutions?!
  - How and where does ML support your **other security efforts** (e.g., continuous risk attestation and enforcement)
- Make sure you have the right **data and context**
  - Beware the **over-collection** of data - capture the same data many times Asset inventory - up to date!
- **Understand your environment** inside out!
  - Invest in **data exploration** capabilities?
- Buy **products** for you most pressing problems. Make sure they solve them cost effectively!
- Don't ever have an "AI project"



**NIST Cyber Security Framework**

**1 🔍 Identify**
AM: Asset Management
BE: Business Environment
GV: Governance
RA: Risk Assessment
RM: Risk Management Strategy

**2 🔒 Protect**
AC: Access Control
AT: Awareness Training
DS: Data Security
IP: Information Protection Processes and Procedures
PT: Protective Technology

**3 🎧 Detect**
AE: Anomalies and Events
CM: Security Continuous Monitoring
DP: Detection Processes

**4 🕐 Respond**
RP: Response Planning
CO: Communications
AN: Analysis
MI: Mitigation
IM: Improvements

**5 ✅ Recover**
RP: Recovery Planning
IM: Improvements
CO: Communications

**IANS**

# Resources

- [Cut Through the AI/ML Hype](), IANS Faculty Aaron Turner and John Strand

- **Artificial Security Will Revolutionize Cybersecurity** - But Security Leaders Must View All Vendor Claims With Skepticism by Chase Cunningham, Joseph Blankenship, and Mike Gualtieri -September 2017

- **Apache SPOT** - machine learning routines: https://github.com/apache/incubator-spot/tree/master/spot-ml

- Many **ML resources**: https://github.com/wtsxDev/Machine-Learning-for-Cyber-Security

- **Even more**: https://github.com/RandomAdversary/Awesome-AI-Security

**IANS**

# BlackHat Workshop



**Applied Machine Learning**
for
Identity and Access Management

ML | AI | IAM

**August 4,5 & August 6,7 - Las Vegas, USA**

**http://secviz.org**

**IANS**

# Questions?

info@iansresearch.com