# Cybersecurity Data Science as a Process

Practitioner Insights and Best Practices

John Hopkins IACD cyber conference
Baltimore
May 2 – 3, 2019

https://www.iacdautomate.org/may-2019-integrated-cyber

DATA

DETECTION     DISCOVERY

§sas
THE POWER TO KNOW.

# Cybersecurity Context

DATA

DETECTION          DISCOVERY

§sas

# Evolving Threats



Internal Threats

Automated Attacks

Social Engineering

Fraud-Cyber Hybrids

State Actors

Ransomware & Cryptojacking

SAS

SURFACE WEB · 9.99%
Google
Bing
Wikipedia

DEEP WEB (full text not accessible via search engines)
Academic Databases
Medical Records
Multilingual Databases
Financial Records
Legal Documents
Academic Records
Subscription Information
Government Resources
Scientific Reports
Organizational Repositories
Competitor Websites
90%

DARK WEB (only partially searchable via Dark Web browsers)
Private Communication
Contraband Sales
Encrypted Sites
Illegal Information
.01%

8    UC San Diego Library    percentages from Greenberg (2014), Wired.com

# CYBERCRIME PRICE LIST

## ATTACK TOOLS

### MALWARE
| | |
|---|---|
| $200 | REMOTE ACCESS TROJAN |
| $50 | PASSWORD STEALER |

### RANSOMWARE
| | |
|---|---|
| $200 | SOPHISTICATED LICENSE FOR WIDESPREAD ATTACKS |
| $50 | UNSOPHISTICATED LICENSE FOR TARGETED ATTACKS |
| $1 | PC MALWARE INSTALLATION |
| $400 | 1 MILLION MALICIOUS SPAM |

### SOFTWARE
| | |
|---|---|
| $100 | REMOTE DESKTOP CONTROL TOOL |
| $700 | DISTRIBUTED DENIAL OF SERVICE ATTACK SOFTWARE |

### PAYMENT AND LOG-IN INFO
| | |
|---|---|
| $5 | CREDIT/DEBIT CARD FOR ONLINE USE |
| $10 | CREDIT/DEBIT CARD INFO THAT CAN BE CLONED ON PLASTIC |
| $5 | BANK ACCOUNT LOG-IN (USERNAME AND PASSWORD) |
| $25 | BANK ACCOUNT LOG-IN WITH ACCESS TO EMAIL, SECURITY ANSWERS, ETC. |
| $1 | EXISTING PAYPAL ACCOUNT |

## DATA

### PERSONAL INFORMATION
| | |
|---|---|
| $3 | SOCIAL SECURITY AND DATE OF BIRTH VERIFICATION |
| $150 | CREDIT REPORT 750+ CREDIT SCORE |

### DATABASE RECORDS
| | |
|---|---|
| $25 | 1 MILLION COMPROMISED EMAIL/PASSWORDS |

## SERVICES

### HACKING
| | |
|---|---|
| $100 | EMAIL ACCOUNT |
| $100 | SOCIAL MEDIA ACCOUNT |
| $300 | CMS WEBSITE (WORDPRESS, ETC.) |

### USER OBFUSCATION
| | |
|---|---|
| $150 | BULLETPROOF HOSTING IN LAX JURISDICTION (CHINA, EASTERN EUROPE, ETC.) |
| $20 | VIRTUAL PRIVATE NETWORK (VPN) |

### MALWARE
| | |
|---|---|
| $1 | PC MALWARE INSTALLATION |
| $25 | MALICIOUS FILE ENCRYPTION |

### SPAM
| | |
|---|---|
| $20 | 500 SMS (FLOODING) |
| $400 | 1 MILLION MALICIOUS SPAM |
| $20 | 500 PHONE CALLS (FLOODING) |
| $200 | 1 MILLION EMAIL SPAM (LEGAL) |

### FAKE DOCUMENTS
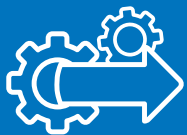| | |
|---|---|
| $25 | DIGITAL COPY OF FAKE CREDIT/DEBIT CARD |
| $25 | DIGITAL COPY OF FAKE DRIVER'S LICENSE OR PASSPORT |
| $15 | DIGITAL COPY OF FAKE UTILITY BILL OR SOCIAL SECURITY CARD |



**CRIMEWARE TOOLKITS**

Cyber Threat Professional

CRIMEWARE
HACKER TOOLKITS · FAKE ANTIVIRUS SOFTWARE
www.CyberThreatTraining.com
To View Full Video Course

# Security Operations Center (SOC)

# Emerging SOC Operational Drivers

**Big & fast streaming data** needs to be stitched into 'smart data'

Limitations of traditional signature and rules-based approaches, **requiring probabilistic and risk-focused models**

**Integrated situational awareness** of network, device, and user behavior while **reducing false alerts**

Need to build and validate efficacious **machine learning models**

**Automation of manual investigation** and remediation processes

§sas

# Cybersecurity Data Science (CSDS)

§sas

# Level of difficulty in reducing false alerts*



*Survey of 621 global IT security practitioners*

# CSDS: Cybersecurity Data Science

Replacing rules with **machine learning** to reduce false alerts

Moving to **real time detection** and decisioning

**Automation** of manual processes and routine decisions

**Data engineering** to structured and integrate distributed big data into '**smart data**'

Investigation tools that **visualize complexity** to improve investigator efficiency and decision making

§sas

# Cybersecurity Analytics as-a-Process

# Cybersecurity Data Science (CSDS) Lifecycle

DATA

DETECTION     DISCOVERY

# CSDS Process

§.sas

# CSDS Data

# The devil is in the data

# Hidden threats want to remain hidden (in the data)

# Anomaly Detection:  Simply Complex

## Identifying focused anomalies amongst an ocean of noise…



*SOURCE*  Aggarwal, Charu C. (2017). "Outlier Analysis: Second Edition". Springer International Publishing AG.

# Challenges preventing successful use of cybersecurity analytics*



When Seconds Count: How Security Analytics Improves Cybersecurity Defenses

Sponsored by SAS Institute
Independently conducted by Ponemon Institute LLC
Publication Date: January 2017

Ponemon Institute® Research Report

| Challenge | Percentage |
|---|---|
| Data challenges | 65% |
| Lack of in-house expertise | 58% |
| Insufficient technologies | 50% |
| Insufficient resources | 40% |
| Lack of clear leadership | 27% |
| Executives do not see cybersecurity as a significant risk | 24% |
| Lack of collaboration with other functions | 19% |
| No understanding how to protect against cyber attacks | 11% |
| Not a priority issue | 6% |

*Survey of 621 global IT security practitioners*

# Cybersecurity Analytics Maturity Curve

| Anomaly Detection | Data-aware Investigations | Predictive Detection | Risk Awareness / Resource Optimization |
|---|---|---|---|

**Anomaly Detection**
- Big data overload
- Flags, rules, and alerts

Chasing phantom patterns

?

?

?

§.sas

**What data sources are available within your organization, should a security analytics program happen?**

| Data source | Percentage |
|---|---|
| Log files | 60% |
| Network flow | 48% |
| Identity and access management systems | 43% |
| Physical security systems | 43% |
| Endpoint monitoring | 40% |
| Packet capture | 39% |
| SIEM | 19% |

Nearly two-third of respondents to a new SC Magazine survey said they were in the process of implementing security analytics.

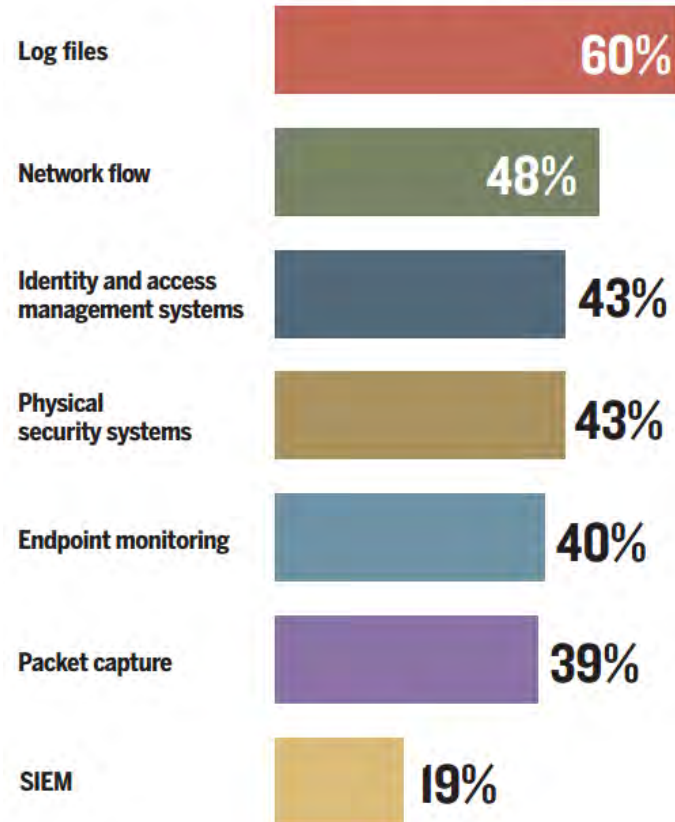**Analyze this!**

Sponsored by SAS

*SOURCE*
Security Brief Magazine. (2016). "Analyze This! Who's Implementing Security Analytics Now?" Available at
https://www.sas.com/en_th/whitepapers/analyze-this-108217.html

# IP address

# time stamp



# userid

# destination port

§sas

# Security Data Management Challenge: Speed and Volumes

Information overload, poor business context, lack of composite risk, mostly signature-based

Firewalls, End Point, Web Proxy, DNS, Vulnerabilities

Ad Hoc Query for organizational impact analysis

**?**
*Thousands*

**POINT SOLUTION ALERTS**
*Millions*

**FLOW**
*Billions*

**PCAP**
*Trillions*

Business Context Enriched, Composite Risk Ranked, Signature-based & Signature-less detection

Firewalls, End Point, Web Proxy, DNS, Vulnerabilities

Stream Processing and Unsupervised Machine Learning at Scale

§sas

# Many data sources... increasing data volumes



High false alerts... slow investigation processes

LACK OF CONTEXT

DISCONNECTED & FRAGMENTED

UNVALIDATED ALERTS

VOLUME & SPEED

?

MULTIPLE SYSTEMS

§sas

# Data Engineering:  Fusion, Quality and Delivery



**SOURCES**

**DESTINATIONS**

**DATA MANAGEMENT**

**DATA GOVERANCE**

**REAL TIME**

**DATA INTEGRATION**

**DATA ACCESS**

**DATA QUALITY & VALIDATION**

**DATA VIRTUALIZATION**

**MASTER DATA MGMT**

"Organizing data is a critical first step in figuring out what data means"

Larry Alton, Information Management Feb 14th, 2019

§.sas

# Cybersecurity Frameworks & Ontologies

## FRAMEWORKS

- MITRE Cyber Observable eXpression

- NIST Cybersecurity Framework

- Intrusion Kill Chain (Lockheed Martin)

## ONTOLOGIES

- **DFAX** Digital Forensic Analysis eXpression

- **CVE** Cyber Intelligence Ontology
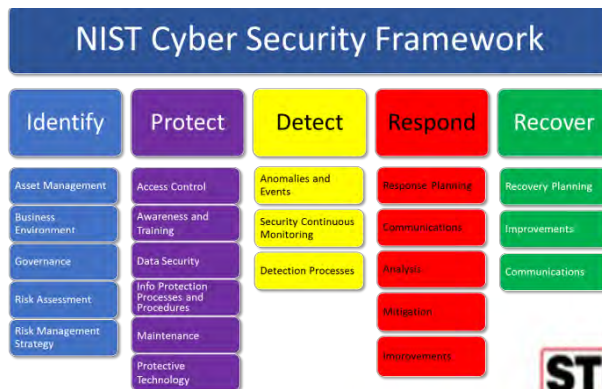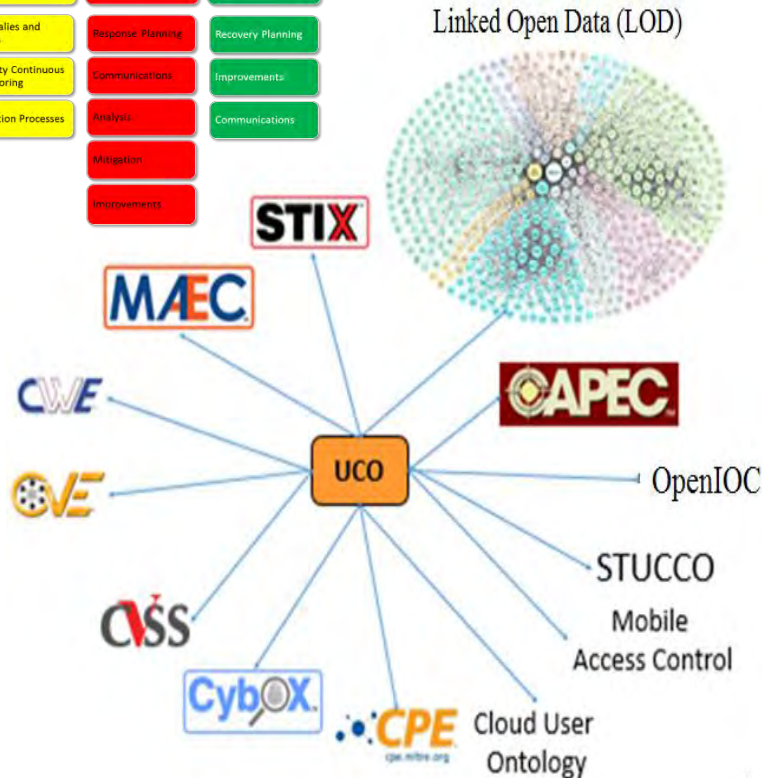
- **ICAS** Information Security (example)

- **UCO / UCO (OWL)**

  Unified Cybersecurity Ontology



NIST Cyber Security Framework

| Identify | Protect | Detect | Respond | Recover |
|----------|---------|--------|---------|---------|
| Asset Management | Access Control | Anomalies and Events | Response Planning | Recovery Planning |
| Business Environment | Awareness and Training | Security Continuous Monitoring | Communications | Improvements |
| Governance | Data Security | Detection Processes | Analysis | Communications |
| Risk Assessment | Info Protection Processes and Procedures | | Mitigation | |
| Risk Management Strategy | Maintenance | | Improvements | |
| | Protective Technology | | | |

Linked Open Data (LOD)

STIX
MAEC
CVE
CVSS
CybOX
CPE cpe.mitre.org
UCO
CAPEC
OpenIOC
STUCCO
Mobile Access Control
Cloud User Ontology

- Ingest
- Digest
- Expel

- Cleansing
- Integration
- Discovery

- Lineage
- Governance
- Security

Whitepaper: A Comprehensive Approach to Big Data Governance, Data Management and Analytics

# CSDS Discovery

DATA

DETECTION

DISCOVERY

§sas

# Cybersecurity Data Science (CSDS) Lifecycle

# CSDS Process
## Unified Orchestration

| DATA | | DISCOVERY | | | | DETECTION | | |
|------|------|------|------|------|------|------|------|------|
| Access | Prepare | Cleanse | Explore | Analyze | Model | Govern | Embed | Monitor |

ORCHESTRATION

# Cybersecurity Events
## Irregular and Complex Events

time's arrow…

**behavioral profile**

DATA

SORTED

ARRANGED

PRESENTED VISUALLY
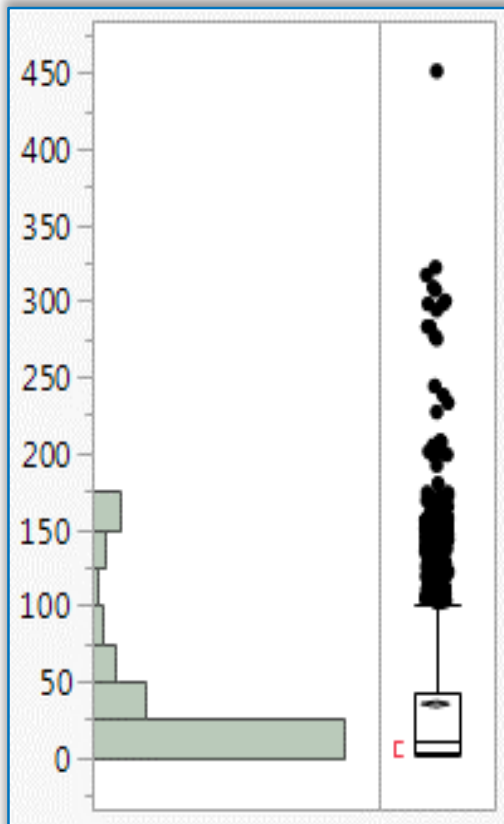
# Self-Service Visual Analytics

# Network Graph Analytics

- Centrality
- Eigenvector
- Density
- Reach
- Strength
- Recopricity



System Manager

System Host
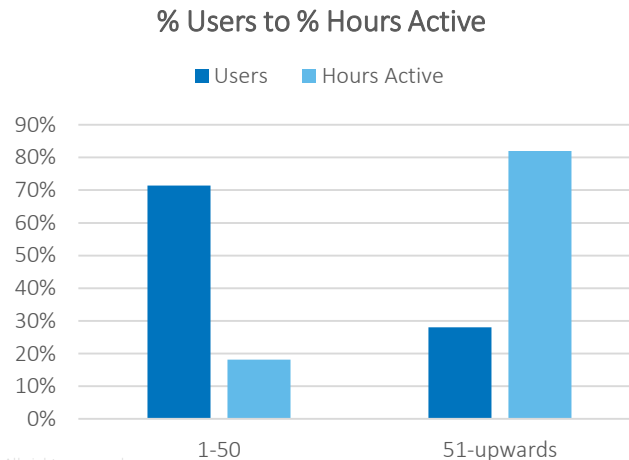
System-User

Interface-User

Human Users

§sas

# Feature Selection / Extraction
## Understanding Network Behavioral Patterns



## Pareto Principle

- **80/20%** pattern in network-usage
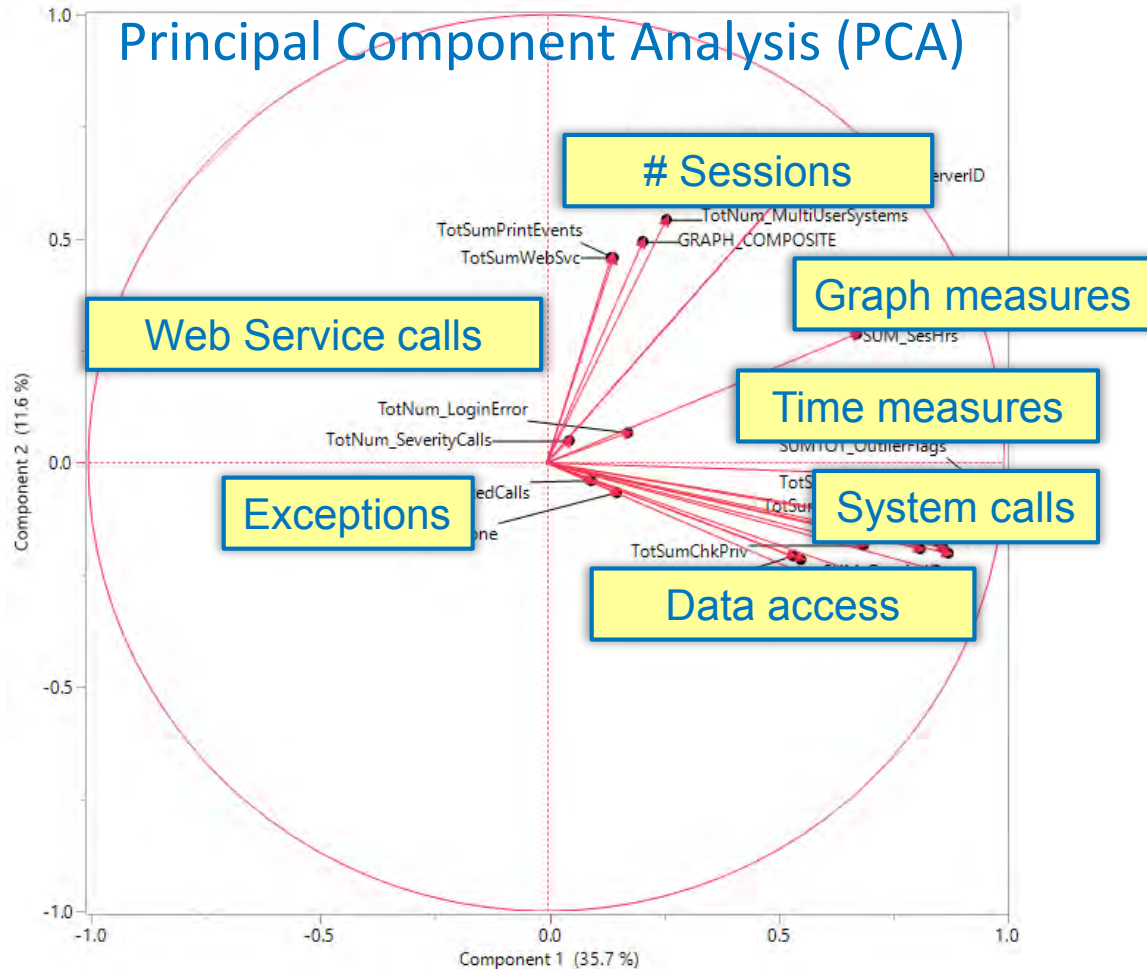- *Outliers*: multiple devices 24 hours online
- High correlation: hrs online and breadth of activities
- Pattern observed across multiple networks

**% Users to % Hours Active**

# Dimensionality Reduction



Principal Component Analysis (PCA)

# Pattern Extrapolation Machine Learning (Unsupervised)



Cluster analysis finds 20 statistically similar user groups

Dendrogram

sas

# CSDS: Diverse Analytics Toolkit



UNKNOWN PATTERNS

COMPLEX PATTERNS

UNSTRUCTURED PATTERNS

Process Analytics

Pattern Extraction

Predictive Modeling

Network Analysis

Text Analytics

Anomaly Detection

KNOWN ATTACKS

Multivariate Extrapolation

Cyber Analytics

Diagnostic Baselining

Unsupervised Machine Learning

KNOWN PATTERNS

UNEXPLAINED PATTERNS

§sas

# CSDS Detection

DATA

DETECTION

DISCOVERY

§sas

# CSDS Process
## Unified Orchestration

**DATA**

| Access | Prepare | Cleanse |

**DISCOVERY**

| Explore | Analyze |

**DETECTION**

| Model | Govern | Embed | Monitor |

ORCHESTRATION

§sas

Cybersecurity Data Science (CSDS) Lifecycle

# CSDS: Diverse Analytics Toolkit



UNKNOWN PATTERNS

COMPLEX PATTERNS

UNSTRUCTURED PATTERNS

Process Analytics

Pattern Extraction

Predictive Modeling

Network Analysis

Text Analytics

Anomaly Detection

KNOWN ATTACKS

Multivariate Extrapolation

Cyber Analytics

Diagnostic Baselining

Unsupervised Machine Learning

KNOWN PATTERNS

UNEXPLAINED PATTERNS

§sas

# Data Science: Continuum of Analytics Methods



Understanding Patterns

Data visualization

Optimizing Systems

PRESCRIPTIVE

SEMANTIC

Network Context & Meaning

Identifying Factors & Causes

DIAGNOSTIC

PREDICTIVE

Forecasting & Probabilities

DESCRIPTIVE

DATA ENGINEEERING

Business Intelligence

HIGH

VALUE

LOW

OVERHEAD

HIGH

§sas

# Predictive Machine Learning (Supervised)

NORMAL

UNUSUAL

Secure profile

Known devices

Average amount

Known location

Known destination

Non-hacker

Hacker

At risk profile

Device

Time of day

Source location

Destination location

Amount

Threat intelligence

IP

Training set

Validation set

§sas

Decision trees

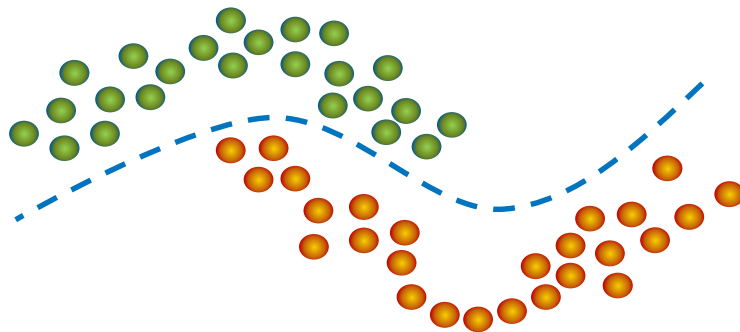Gradient boosting

Random forests

Naïve Bayes

SVM

Gaussian processes

# Supervised Learning

- Trained on labeled examples. We have a target we are predicting.

- Map inputs to desired output.

- Suitable for classification and prediction.

## Considerations

- Obtaining labeled data for rare events can be a challenge

- Suspicion is not a cyber incident!

- Data is skewed 99-1

§sas

# Machine Learning Model = Active Data Vehicle

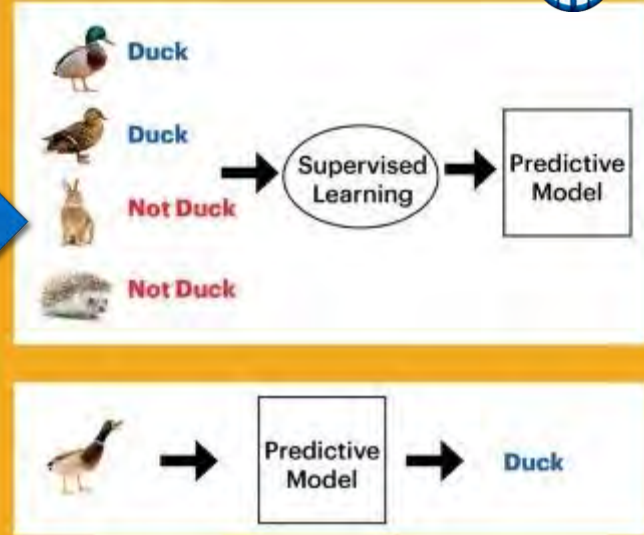# Machine Learning Segmentation and Classification

# Model Building Process ⬌ Analytics Life Cycle

**Raw Data**  **Feature Selection**  **Features**  **Feature Engineering**  **Modeling**  **Insights**



SAS: 'Managing the Analytics Life Cycle for Decisions at Scale'

# Bootstrapping Machine Learning Facilitated Cyber Detection



**1. TRADITIONAL RULE-BASED PATH**

SIEM / Big Data Repository → InfoSec Hunters → Rule-based Alerts → InfoSec Investigation

*Human-in-the-Loop*

**2. UNSUPERVISED DISCOVERY**

Unsupervised Machine Learning → Unsupervised Outliers & Anomalies

Labeled Instances (TP or FP)

**3. SUPERVISED LEARNING**

Network Discovery / Data Quality — Feature Engineering ← Supervised Machine Learning

§sas

# Architecture: Exploratory & Detection Platforms*
## Functional Architectural Segmentation

**Exploratory 'big data' repository**

**Operationally focused detection**

Feature engineering
i.e. selection, refinement, binning, correlations

Canonical ontology / schemas

Analytical models
- Descriptive
- Unsupervised

Analytical models
- Semi-supervised
- Human-in-the-loop
- Reinforcement

Analytical models
- Statistical
- Supervised

*Runs counter to the vendor stance of store 'all-the-data-all-the-time'*

§.sas

# Summary

# Accelerate the Analytics Lifecycle
## A Unified Process Approach

| DATA | | | DISCOVERY | | | DETECTION | | |
|------|------|------|-----------|------|------|-----------|------|------|
| Access | Prepare | Cleanse | Explore | Analyze | Model | Govern | Embed | Monitor |

§sas

Cybersecurity Data Science (CSDS) Lifecycle

# Cybersecurity Analytics Maturity Model

| Anomaly Detection | Data-aware Investigations | Predictive Detection | Risk Awareness / Resource Optimization |
|---|---|---|---|

**Anomaly Detection**
- Big data overload
- Flags, rules, and alerts

**Chasing phantom patterns**

**Understanding**
- Feature engineering
- *Unsupervised ML*
- Labeling
- Diagnostics

**Learning**
- Human-in-the-loop *reinforcement learning*
- *Semi- and Supervised ML*

**Risk Optimal**
- Champion-challenger model management
- Automating alert triage
- Resource optimization

§.sas

# Want to Know More?

SAS whitepaper '*Data Management for Artificial Intelligence*'

SAS Cybersecurity Solution (SCS)

www.sas.com/en_us/software/cybersecurity.html

Scott Allen Mongeau
Data Scientist - Cybersecurity

scott.mongeau@sas.com

Scott Mongeau

www.sas.com/en_us/whitepapers/data-management-artificial-intelligence-109860.html

§sas

# REFERENCES

# REFERENCES

- Aggarwal, C. (2013). "Outlier Analysis." Springer. http://www.springer.com/la/book/9781461463955

- Harris, H., Murphy, S., and Vaisman, M. (2013). "Analyzing the Analyzers." O'Reilly Media. Available at https://www.oreilly.com/data/free/analyzing-the-analyzers.csp

- Kirchhoff, C., Upton, D., and Winnefeld, Jr., Admiral J. A. (2015 October 7). "Defending Your Networks: Lessons from the Pentagon." Harvard Business Review. Available at https://hbr.org/webinar/2015/10/defending-your-networks-lessons-from-the-pentagon

- Mongeau, S. (2018). "Cybersecurity Data Science (CSDS)." SCTR7.com. https://sctr7.com/2018/12/03/cybersecurity-data-science-csds-how-not-to-drown-in-your-cyber-data-lake/

- Mongeau, S. (2017). "Cybersecurity Big Data Overload?" SCTR7.com. https://sctr7.com/2017/10/22/cybersecurity-big-data-overload/

- Ponemon Institute. (2017). "When Seconds Count: How Security Analytics Improves Cybersecurity Defenses." Available at https://www.sas.com/en_us/whitepapers/ponemon-how-security-analytics-improves-cybersecurity-defenses-108679.html

- SANS Institute. (2015). "2015 Analytics and Intelligence Survey." Available at https://www.sas.com/en_us/whitepapers/sans-analytics-intelligence-survey-108031.html

- SANS Institute. (2016). "Using Analytics to Predict Future Attacks and Breaches." Available at https://www.sas.com/en_us/whitepapers/sans-using-analytics-to-predict-future-attacks-breaches-108130.html

- SAS Institute. (2016). "Managing the Analytical Life Cycle for Decisions at Scale." Available at https://www.sas.com/content/dam/SAS/en_us/doc/whitepaper1/manage-analytical-life-cycle-continuous-innovation-106179.pdf

- SAS Institute. (2017). "SAS Cybersecurity: Counter cyberattacks with your information advantage." Available at https://www.sas.com/en_us/software/fraud-security-intelligence/cybersecurity-solutions.html

- UBM. (2016). "Dark Reading: Close the Detection Deficit with Security Analytics." Available at https://www.sas.com/en_us/whitepapers/close-detection-deficit-with-security-analytics-108280.html

§.sas