



# Cybersecurity Data Science (CSDS)

Practitioner Methods and Best Practices

## PART 1: FRAME

Scott Allen Mongeau  
Cybersecurity Data Scientist  
SAS Institute  
[scott.mongeau@sas.com](mailto:scott.mongeau@sas.com)

# Cybersecurity Data Science (CSDS)

TOPIC
1. FRAME
2. DATA
3. DISCOVER
4. DETECT
5. DEPLOY

# Introductions, Expectations, Agenda



# Scott Allen Mongeau

Cybersecurity Data Scientist

MA GD MA MBA PhD (ABD)



+31 (0)68 370 3097



scott@sark7.com

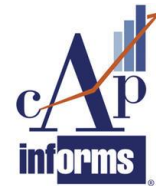


Scott Allen Mongeau



SARK7

## Cybersecurity Data Science (CSDS)



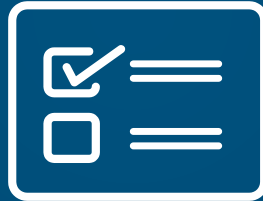
- **SAS (2015 – present)**
  - F&SI Cyber – Global (EU/Netherlands based)
  - Fraud & Financial Crime Analytics (London)
- **SARK7 (2008 – present)**
  - Data analytics consulting (Leiden)
- **Deloitte (2013 – 15)**
  - Fraud, Financial Crime, Cyber (Amsterdam)
- **Genentech Inc. (2000 – 2008)**
  - Data analytics (San Francisco, CA)

Scott Mongeau

# Cybersecurity Data Science: Perspectives on an Emerging Profession

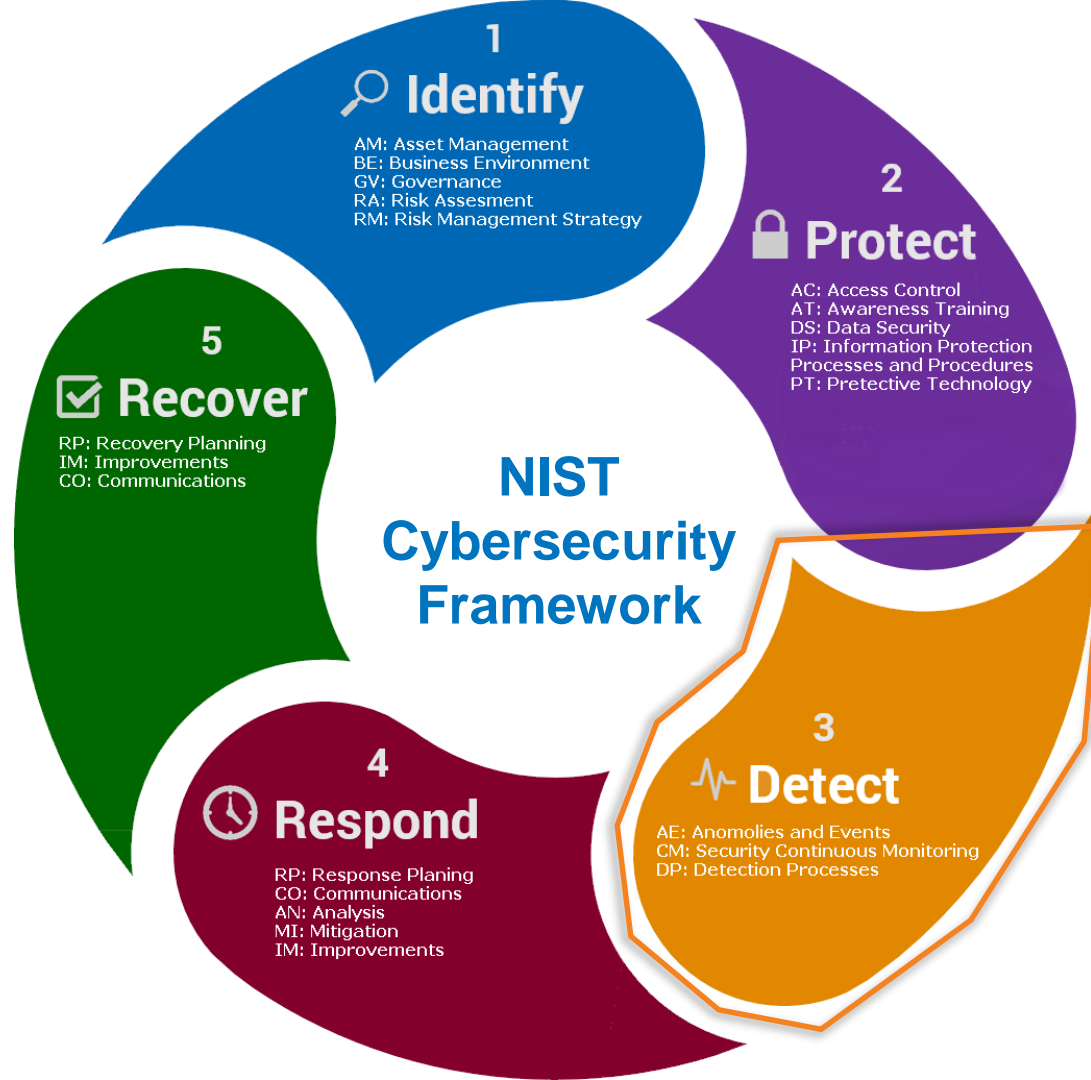
- Defining key CSDS challenges
- Outline best practices
- Guidance to managers and practitioners in implementation of CSDS programs

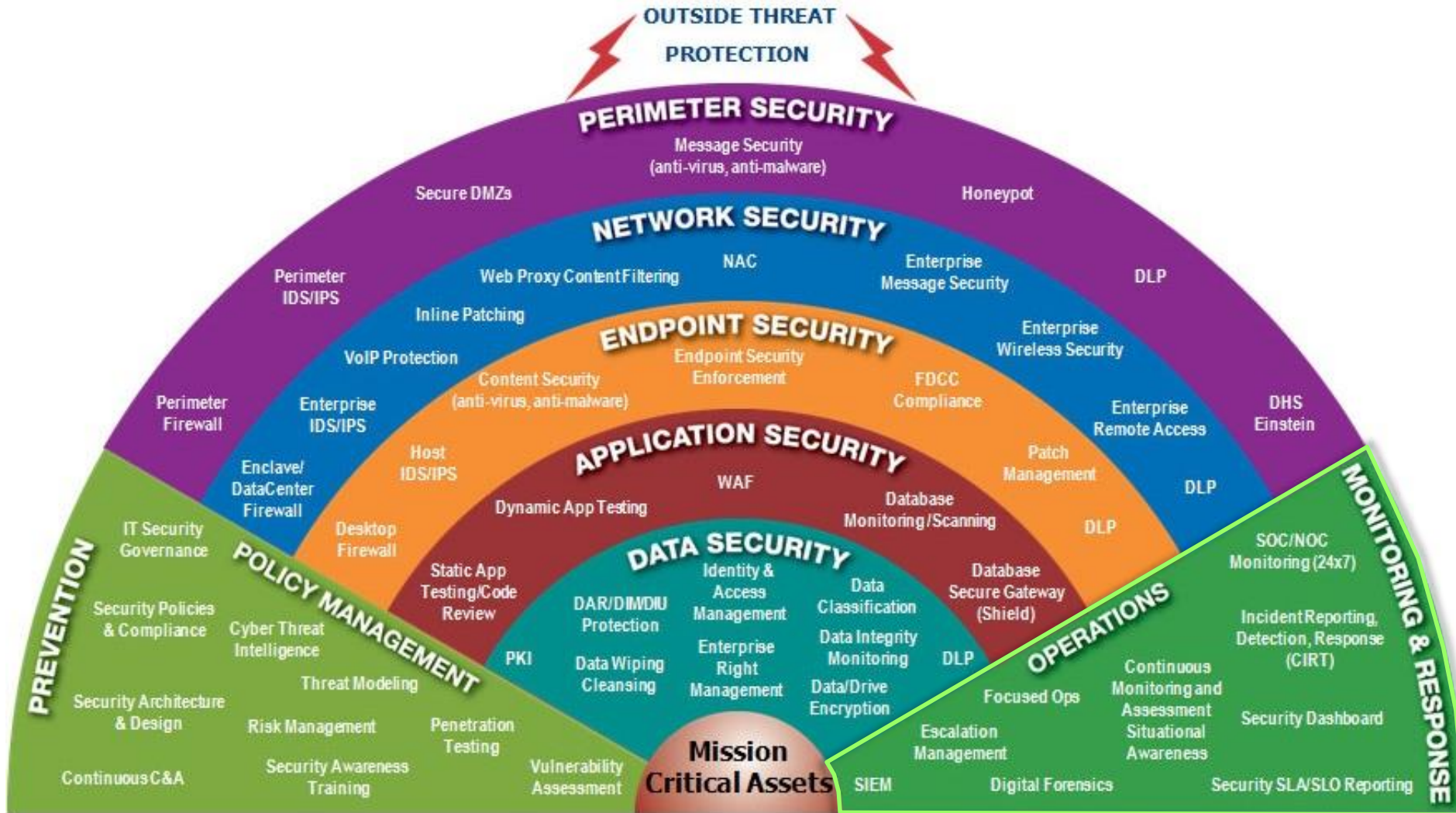
# Key Learning Objectives



# Cybersecurity Services and Solutions: Overview of Major Areas







# Cybersecurity Data Science as a Process

## Data Engineering



## Advanced Analytics

Diagnostics & patterns    Establishing baselines    Predictive modelling    Anomaly detection    Behavioral insights



## Triage / Validate



## Remediate



Data Manager



Data Scientist



Cyber Investigator



Infosec Response

# Cybersecurity Analytics Maturity Curve

## Anomaly Detection

- Big data management
- Flags, rules, and alerts

CHASING  
PHANTOM  
PATTERNS



## Data-aware Investigations



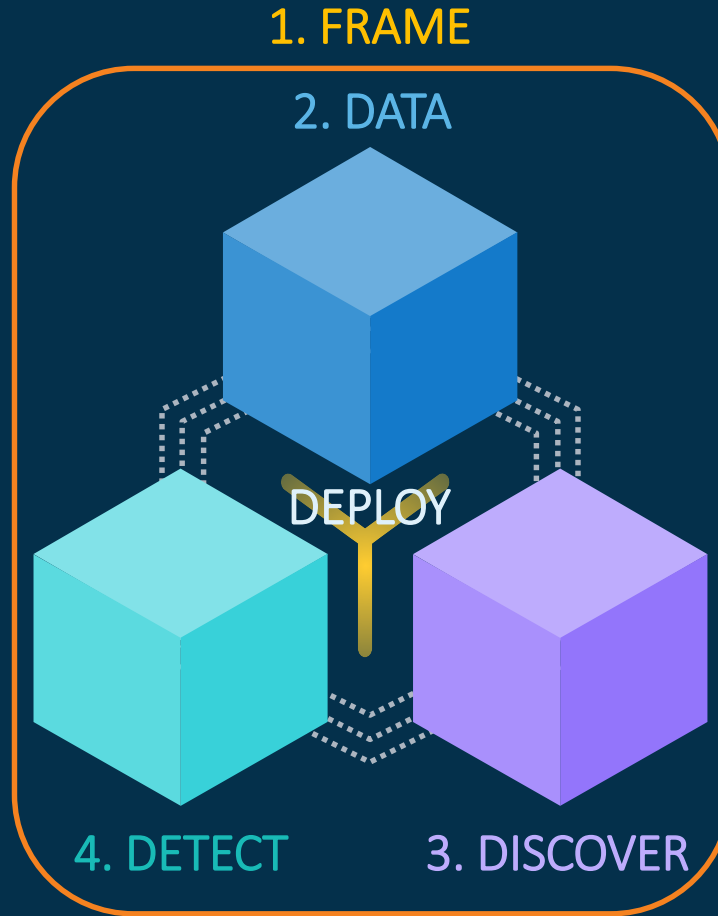
## Predictive Detection



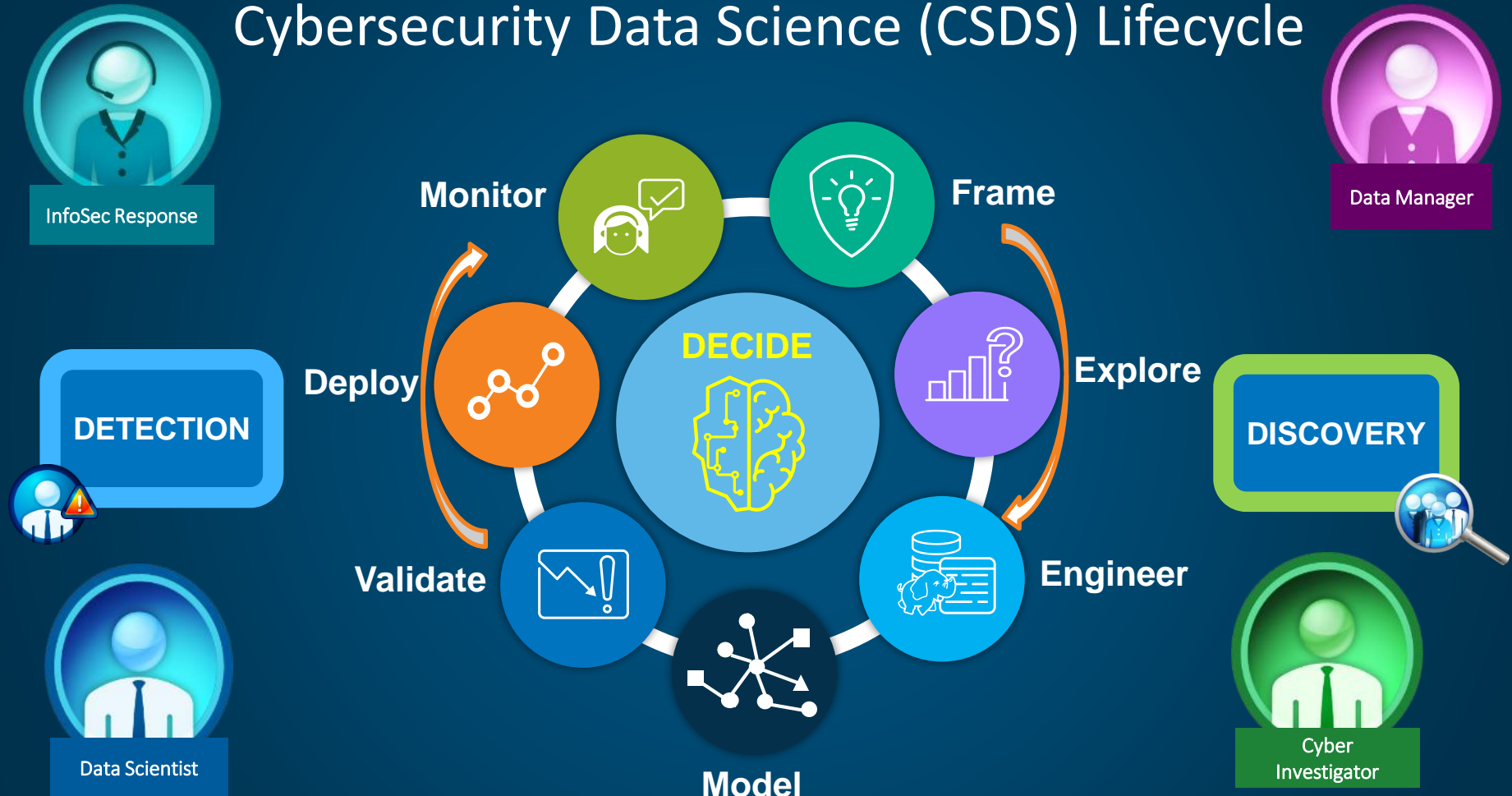
## Risk Optimization



# CSDS Process



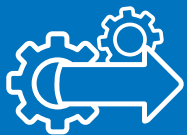
# Cybersecurity Data Science (CSDS) Lifecycle



# 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



# 1. FRAME

Cybersecurity, Data Science, Logs

# Cybersecurity Data Science (CSDS)

## TOPIC

**1. FRAME**

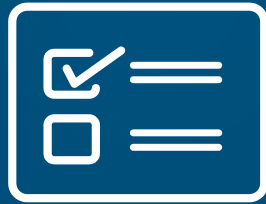
**2. DATA**

**3. DISCOVER**

**4. DETECT**

**5. DEPLOY**

# Learning Objectives



# Objectives of Context Setting

## Establishing a foundation

- Data challenges (opportunities) in cybersecurity
- Data science foundations
- Demonstration / exercises
  - Insights into the 'data deluge' (network analytics)
  - Log file analytics – from unstructured to structured

# CSDS Process

## Unified Orchestration



# Cybersecurity Context



# Evolving Threats

Internal  
Threats

Automated  
Attacks



State Actors

Fraud-Cyber  
Hybrids



Social  
Engineering

Ransomware  
&  
Cryptojacking

Hidden threats want to remain hidden (in the data)



# SURFACE WEB

9.99%

Google

Bing

Wikipedia

# DEEP WEB

(full text not accessible via search engines)

Academic  
Databases

Multilingual  
Databases

Medical Records

Financial Records

Legal  
Documents

Subscription  
Information

Scientific  
Reports

Competitor  
Websites

Academic  
Records

Government  
Resources

Organizational  
Repositories

90%

# DARK WEB

(only partially searchable via Dark Web browsers)

Private Communication

Contraband Sales

Encrypted Sites

Illegal Information

.01%

## 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

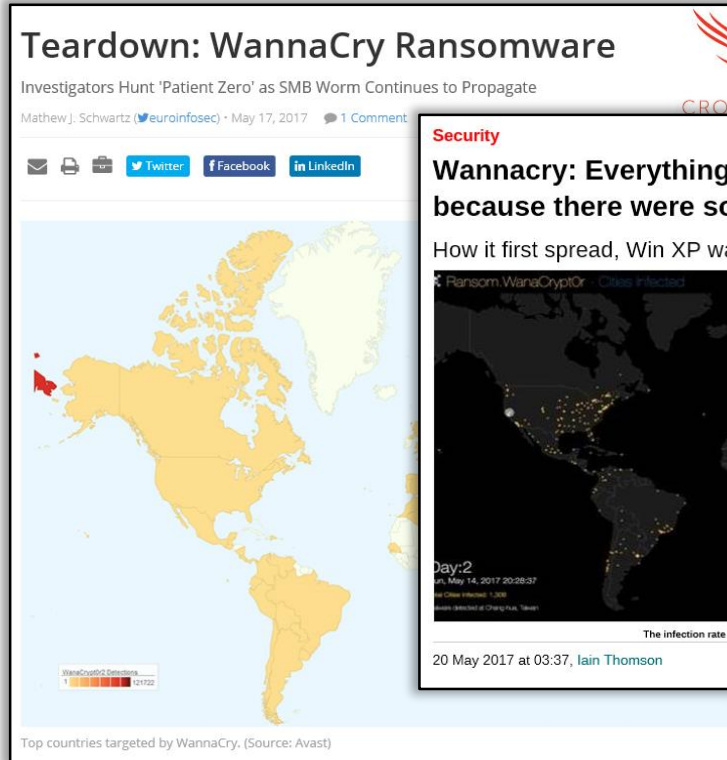
Cyber  
Threat  
Professional



Source: Recorded Future via Fortune Magazine 'A Hacker's Tool Kit'  
<http://fortune.com/2017/10/25/cybercrime-spyware-marketplace/>

# FUD Fear, Uncertainty, Doubt

Increasing FREQUENCY, sophistication, and speed of attacks



# FEAR IS THE MIND KILLER

## LIST OF COGNITIVE BIASES

[https://en.wikipedia.org/wiki/  
List of cognitive biases](https://en.wikipedia.org/wiki/List_of_cognitive_biases)

Ambiguity effect	Distinction bias	Information bias	Reactive devaluation
Anchoring or focalism	Dread aversion	Insensitivity to sample size	Recency illusion
Anthropocentric thinking	Dunning–Kruger effect	Interoceptive bias	Regressive bias
Anthropomorphism or personification	Duration neglect	Irrational escalation	Restraint bias
Attentional bias	Empathy gap	Law of the instrument	Rhyme as reason effect
Automation bias	Endowment effect	Less-is-better effect	Risk compensation / Peltzman effect
Availability heuristic	Exaggerated expectation	Look-elsewhere effect	Selection bias
Availability cascade	Experimenter's or expectation bias	Loss aversion	Selective perception
Backfire effect	Focusing effect	Mere exposure effect	Semmelweis reflex
Bandwagon effect	Forer effect or Barnum effect	Money illusion	Social comparison bias
Base rate fallacy or Base rate neglect	Form function attribution bias	Moral credential effect	Social desirability bias
Belief bias	Framing effect	Negativity bias or Negativity effect	Status quo bias
Ben Franklin effect	Frequency illusion or Baader–Meinhof effect	Neglect of probability	Stereotyping
Berkson's paradox	Functional fixedness	Normalcy bias	Subadditivity effect
Bias blind spot	Gambler's fallacy	Not invented here	Subjective validation
Choice-supportive bias	Groupthink	Observer-expectancy effect	Surrogation
Clustering illusion	Hard–easy effect	Omission bias	Survivorship bias
Confirmation bias	Hindsight bias	Optimism bias	Time-saving bias
Congruence bias	Hostile attribution bias	Ostrich effect	Third-person effect
Conjunction fallacy	Hot-hand fallacy	Outcome bias	Parkinson's law of triviality
Conservatism (belief revision)	Hyperbolic discounting	Overconfidence effect	Unit bias
Continued influence effect	Identifiable victim effect	Pareidolia	Weber–Fechner law
Contrast effect	IKEA effect	Pessimism bias	Well travelled road effect
Courtesy bias	Illicit transference	Planning fallacy	Zero-risk bias
Curse of knowledge	Illusion of control	Post-purchase rationalization	
Declinism	Illusion of validity	Present bias	
Decoy effect	Illusory correlation	Pro-innovation bias	
Default effect	Illusory truth effect	Projection bias	
Denomination effect	Impact bias	Pseudocertainty effect	
Disposition effect	Implicit association	Reactance	

Lack of statistical and analytics  
approaches to establish alert validity

High volumes of  
disconnected and  
disparate data sources

Data Silos

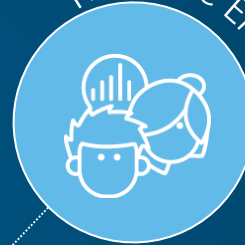


Alert Context



Too many alerts combined  
with slow and resource  
constrained investigations

Resource Efficiency



Disconnected log  
sources of variable  
quality

Complex Data Environment



Proliferation of point  
analytics solutions  
impedes holistic risk-  
based approach

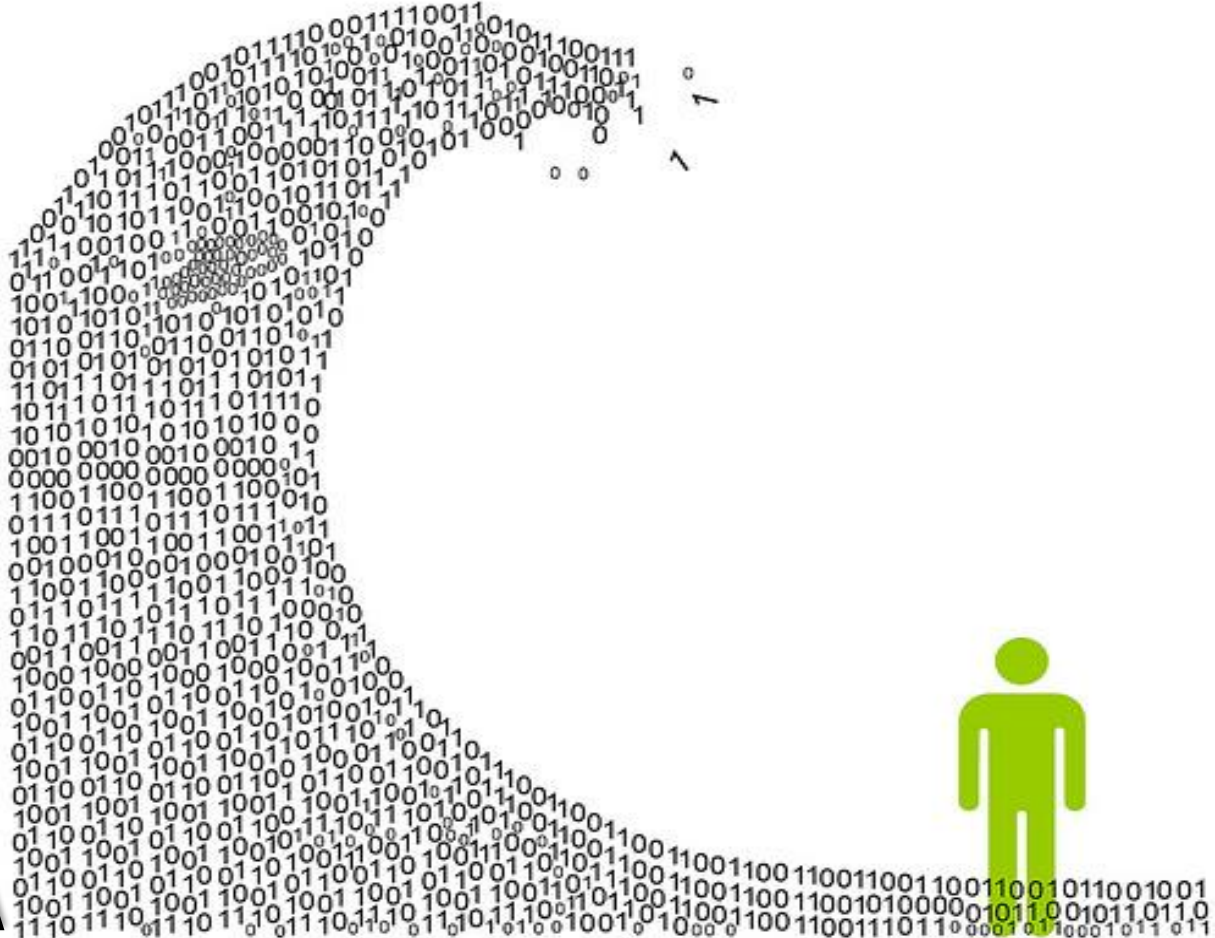
Comprehensive Approach



Difficulty  
Actioning  
Indicators

# ‘Drowning’ in Data Lakes

BIG DATA





# Network Traffic

Tracking network traffic on a single device

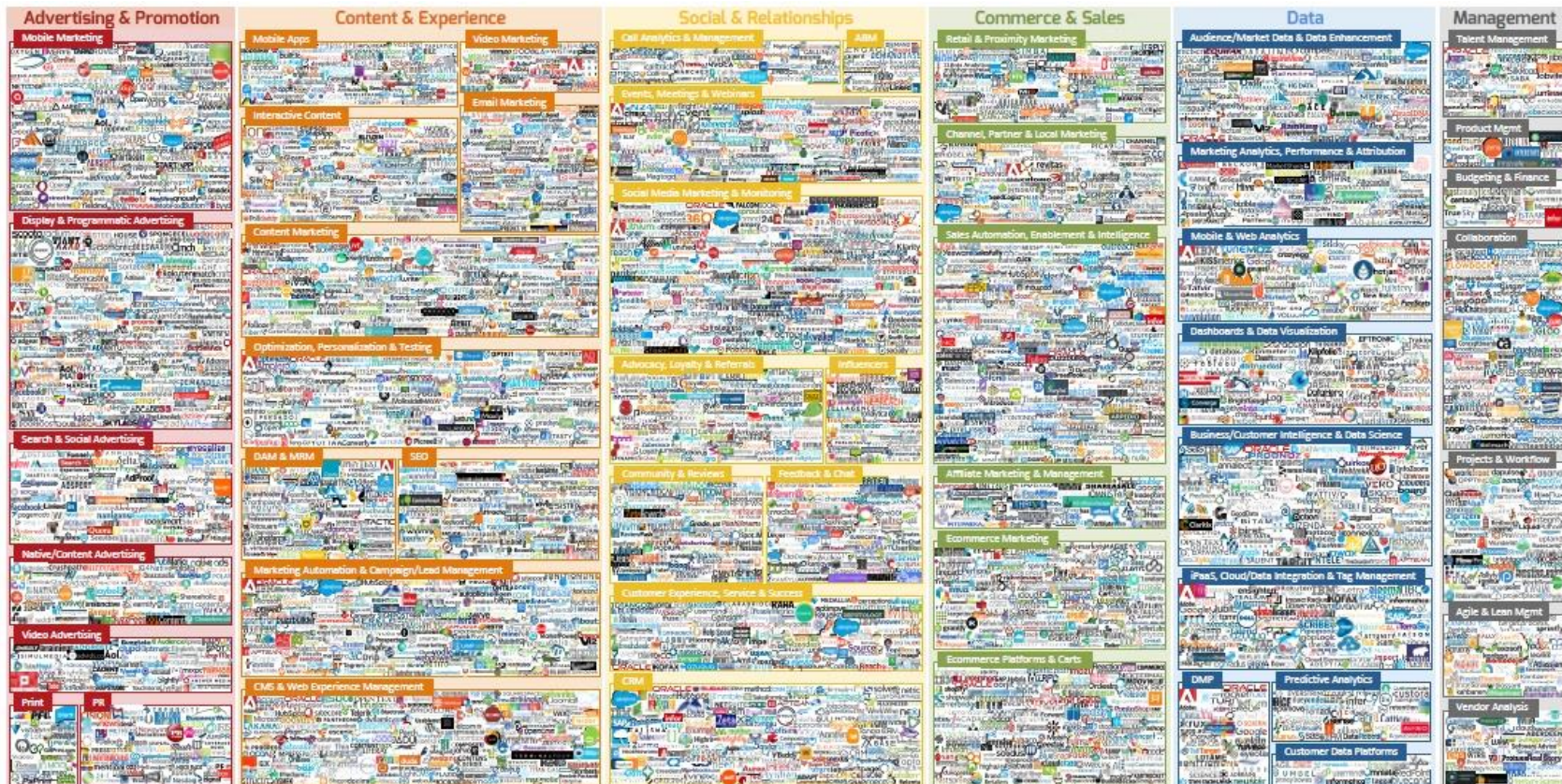
# Exponential Growth in Network Traffic



Using cookies and other tracking techniques, many “normal” web sites aggressively access and store our:

- Device details
- IP address
- Present geolocation
- Changing physical locations
- Browser history
- Online purchases
- Profile information
- Buying behavior
- Personal finances
- Religious beliefs
- Political affiliations
- Health concerns and problems
- Camera / photos / microphone

# Marketing Analytics Ecosystem (>4000 companies)



Sources: CabinetM (<http://cabinetm.com>), Captera, G2 Crowd, Google, Growthverse, LUMA Partners, Siftify, TrustRadius, VBProfiles — see <http://chiefmartec.com/2016/03/marketing-technology-supergaphic-2016/> for details.

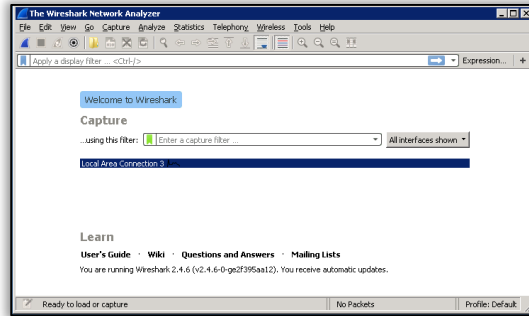
Created by Scott Brinker (@chiefmartec).

# Demonstration / Exercise

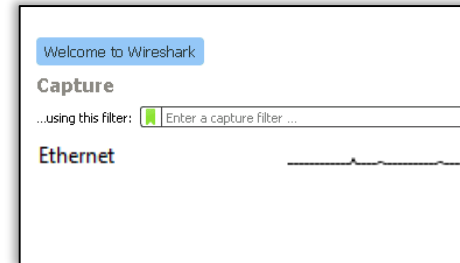
## Network Traffic from Web Browsing



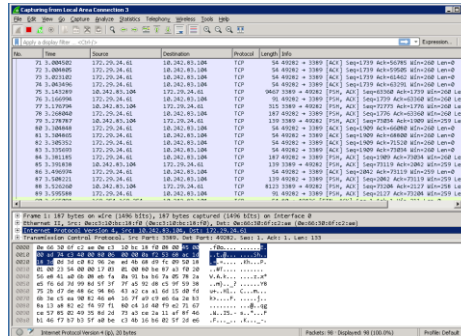
1. Open 'Wireshark'



2. Click 'Ethernet'



3. Observe results



# Demonstration / Exercise

## Network Traffic from Simple Web Browsing

4. Open 'Internet Explorer' (next to Wireshark)

5. Go to a popular news website

6. Monitor Wireshark – what do you see?

7. Sort on destination

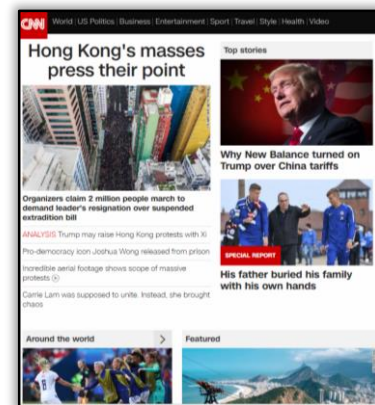
7. View 'red' & 'black' events

*See codes View => Coloring Rules*

*Note 'Destination' IP*

8. Go to <http://ip-lookup.net>

*Where and what are these IPs?*

A screenshot of the Wireshark packet list pane. The table shows various network packets with columns for No., Time, Source, Destination, Protocol, Length, and Info. The packets are sorted by destination IP address. Several packets are highlighted in red and black, corresponding to the 'red' and 'black' events mentioned in the text. The highlighted packets show destinations like 192.31.22.2 and 192.31.22.1.

Lookup an IP address :

192.31.22.2 [x] [list]

IP addresses can be entered using IPv4 or IPv6 address format.

# Wireshark Coloring Rules (default)



Wireshark · Coloring Rules Default



Name	Filter
<input checked="" type="checkbox"/> Bad TCP	tcp.analysis.flags && !tcp.analysis.window_update
<input checked="" type="checkbox"/> HSRP State Change	hsrp.state != 8 && hsrp.state != 16
<input checked="" type="checkbox"/> Spanning Tree Topology Change	stp.type == 0x80
<input checked="" type="checkbox"/> OSPF State Change	ospf.msg != 1
<input checked="" type="checkbox"/> ICMP errors	icmp.type eq 3    icmp.type eq 4    icmp.type eq 5    icmp.type eq 11    icmpv6.type eq 1    icmpv6.type eq 2    icmpv6.type eq 3    icmpv6.type eq 4    icmpv6.type eq 5    icmpv6.type eq 11
<input checked="" type="checkbox"/> ARP	arp
<input checked="" type="checkbox"/> ICMP	icmp    icmpv6
<input checked="" type="checkbox"/> TCP RST	tcp.flags.reset eq 1
<input checked="" type="checkbox"/> SCTP ABORT	sctp.chunk_type eq ABORT
<input checked="" type="checkbox"/> TTL low or unexpected	( ! ip.dst == 224.0.0.0/4 && ip.ttl < 5 && !pim && !ospf )    ( ip.dst == 224.0.0.0/24 && ip.dst != 224.0.0.251 && ip.ttl != 1 && !(vrrp    carp) )
<input checked="" type="checkbox"/> Checksum Errors	eth.fcs.status == "Bad"    ip.checksum.status == "Bad"    tcp.checksum.status == "Bad"    udp.checksum.status == "Bad"    sctp.checksum.status == "Bad"
<input checked="" type="checkbox"/> SMB	smb    nbss    nbns    netbios
<input checked="" type="checkbox"/> HTTP	http    tcp.port == 80    http2
<input checked="" type="checkbox"/> DCERPC	dcerpc
<input checked="" type="checkbox"/> Routing	hsrp    eigrp    ospf    bgp    cdp    vrrp    carp    gvrp    igmp    ismp
<input checked="" type="checkbox"/> TCP SYN/FIN	tcp.flags & 0x02    tcp.flags.fin == 1
<input checked="" type="checkbox"/> TCP	tcp
<input checked="" type="checkbox"/> UDP	udp
<input checked="" type="checkbox"/> Broadcast	eth[0] & 1
<input checked="" type="checkbox"/> System Event	systemd_journal    sysdig



# Data Science Foundations





## When Seconds Count: How Security Analytics Improves Cybersecurity Defenses

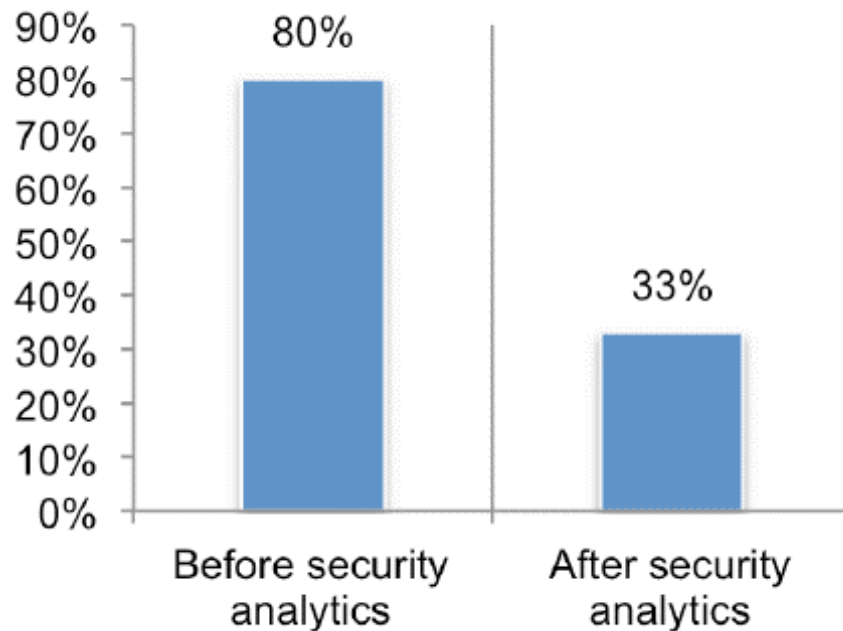
**Sponsored by SAS Institute**

Independently conducted by Ponemon Institute LLC

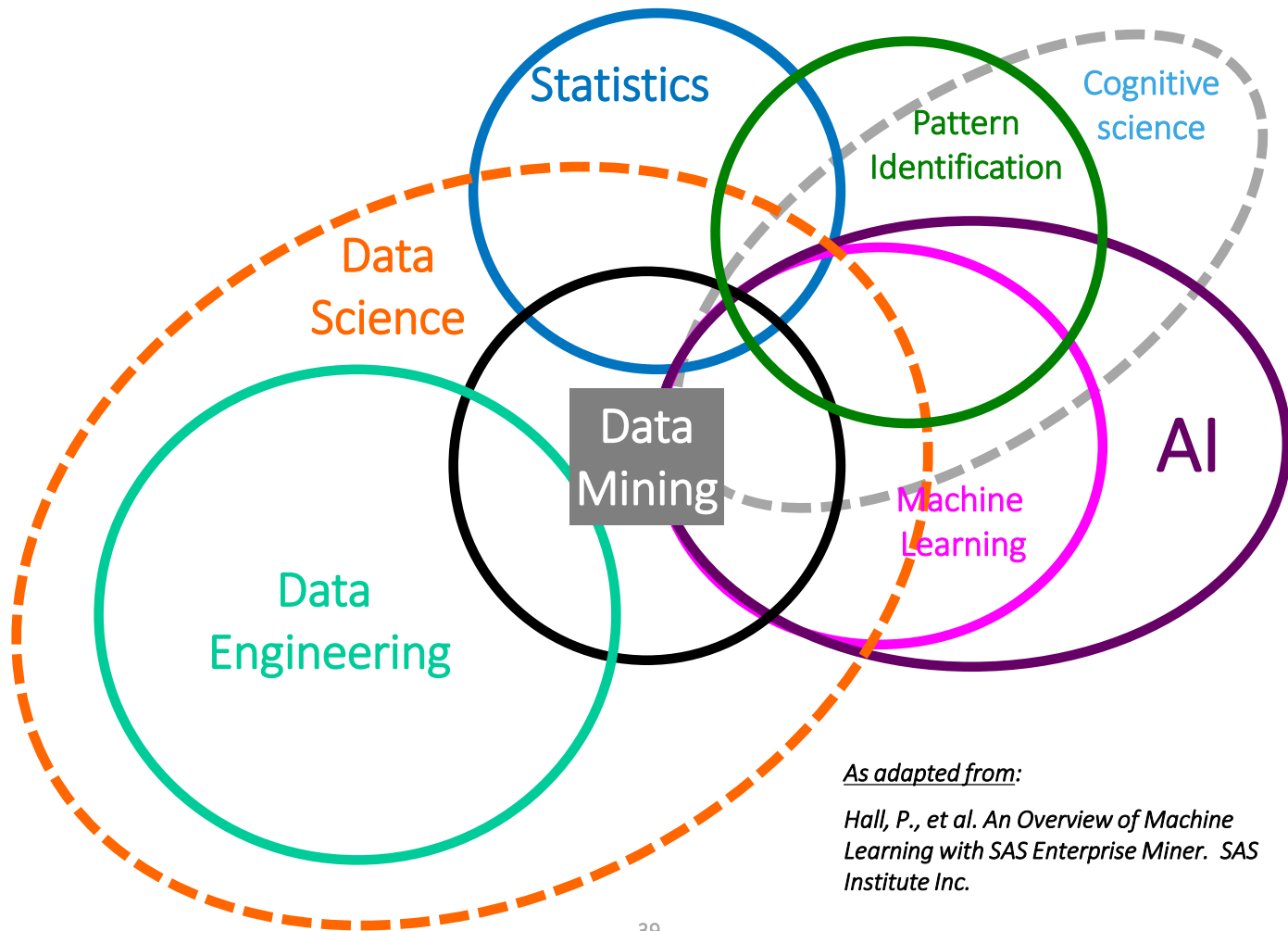
Publication Date: January 2017

Ponemon Institute® Research Report

## Level of difficulty in reducing false alerts\*



\* Survey of 621 global IT security practitioners



*As adapted from:*

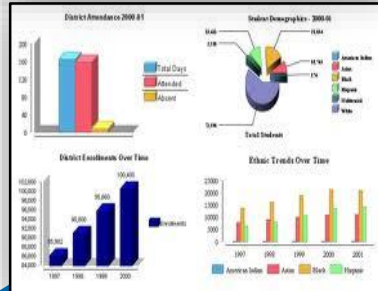
*Hall, P., et al. An Overview of Machine Learning with SAS Enterprise Miner. SAS Institute Inc.*

# DATA ANALYTICS

SOPHISTICATION

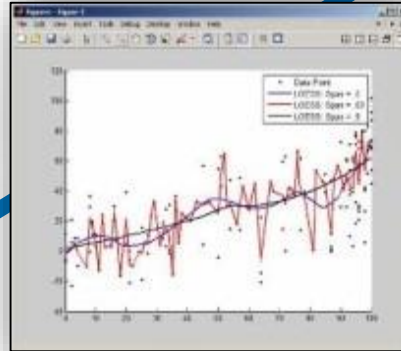
What happened?

DESCRIPTIVE



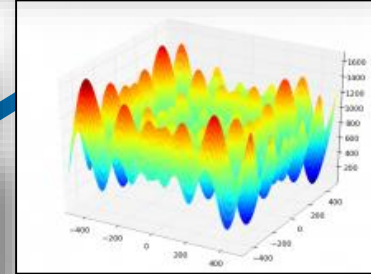
What are trends?

PREDICTIVE



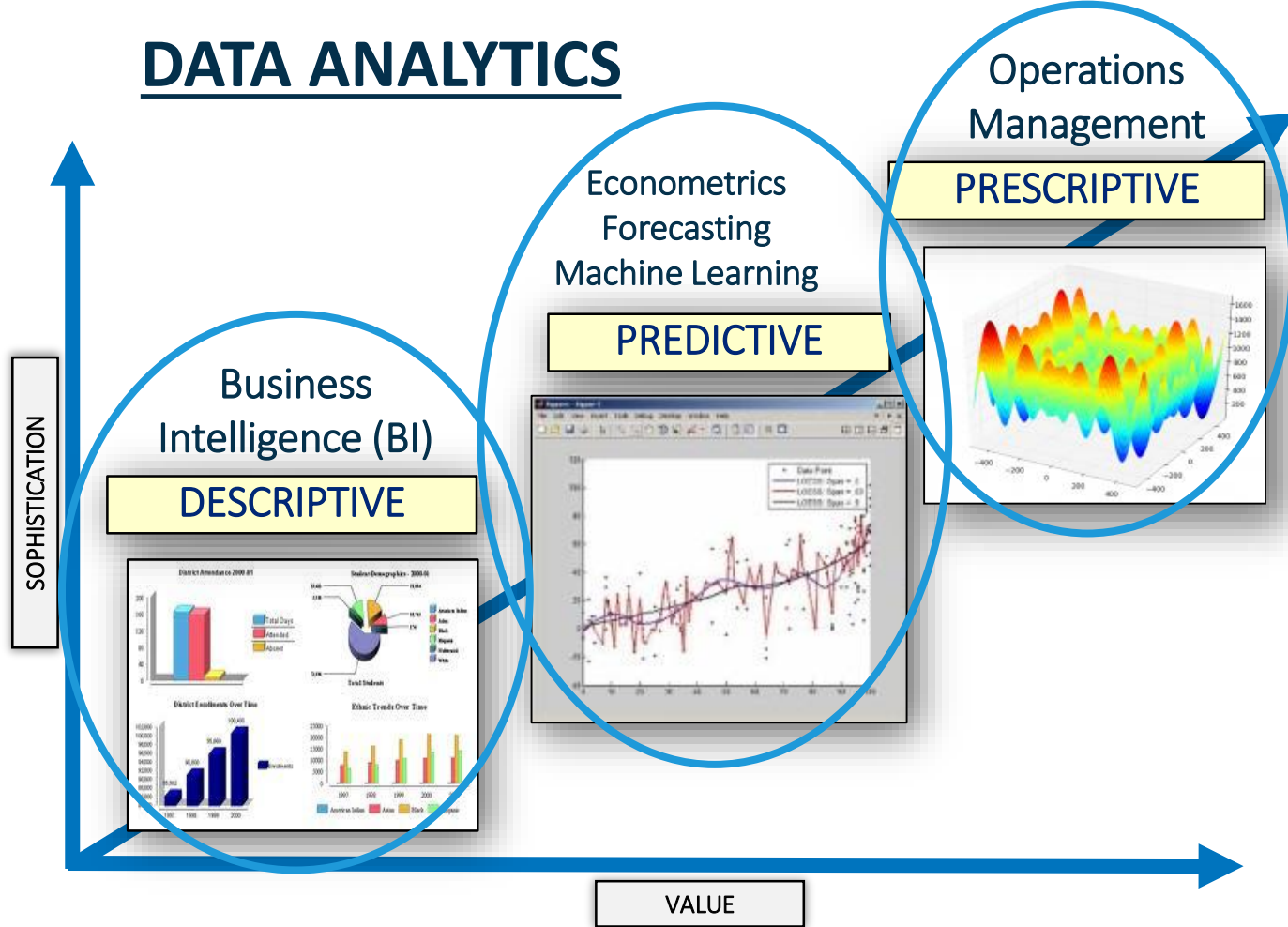
What to do?

PRESCRIPTIVE



VALUE

# DATA ANALYTICS



# DATA ANALYTICS

How can we optimize remediation from focused alerts?

What is the 'normal' behaviour of particular users or devices?

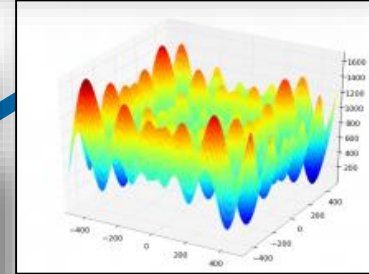
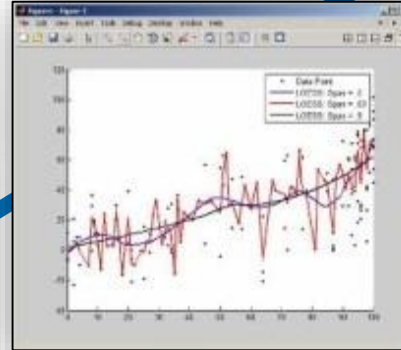
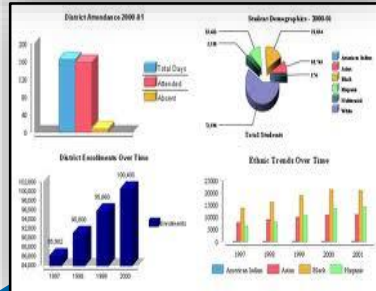
PRESCRIPTIVE

PREDICTIVE

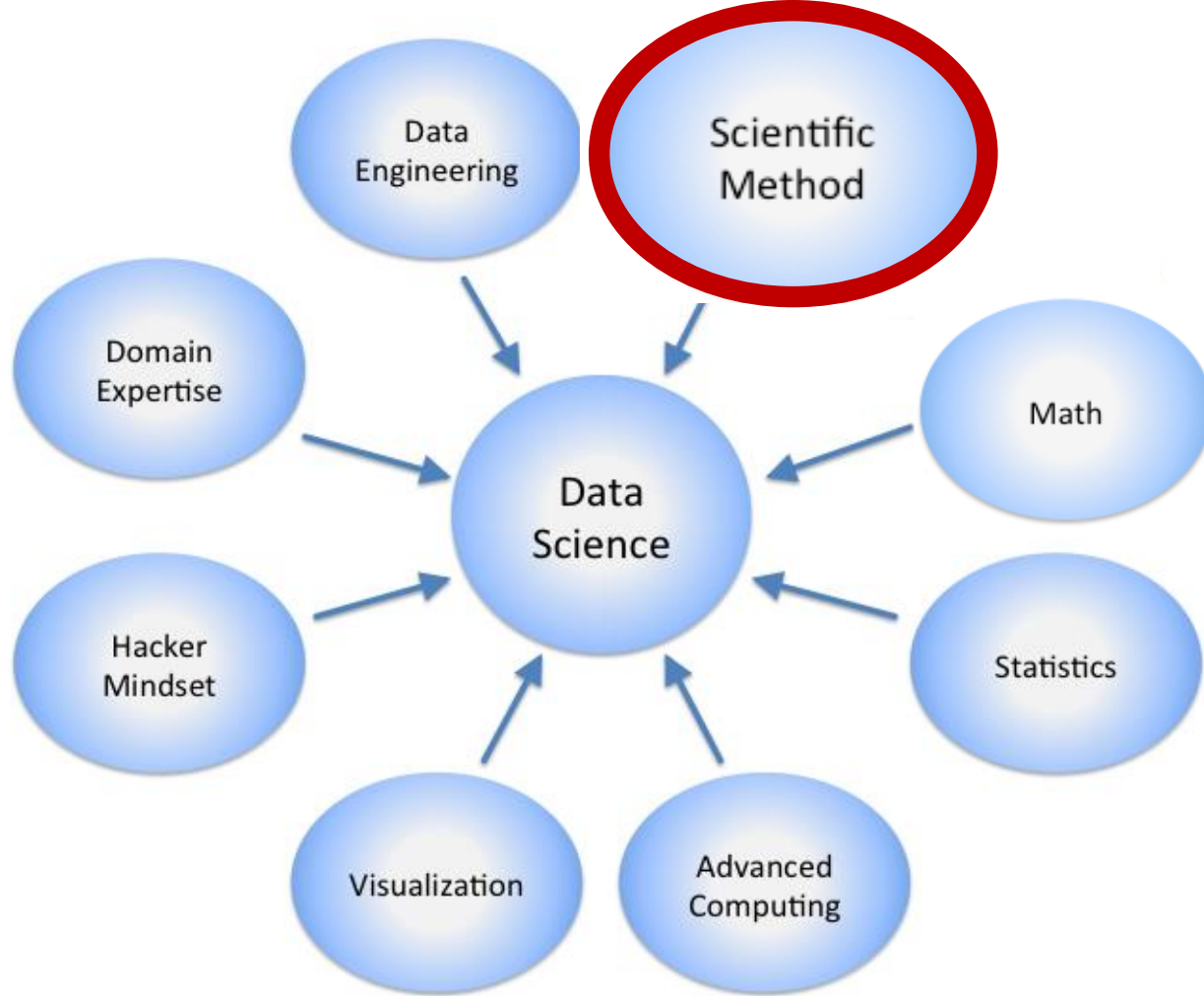
What are trends and patterns in network and device usage?

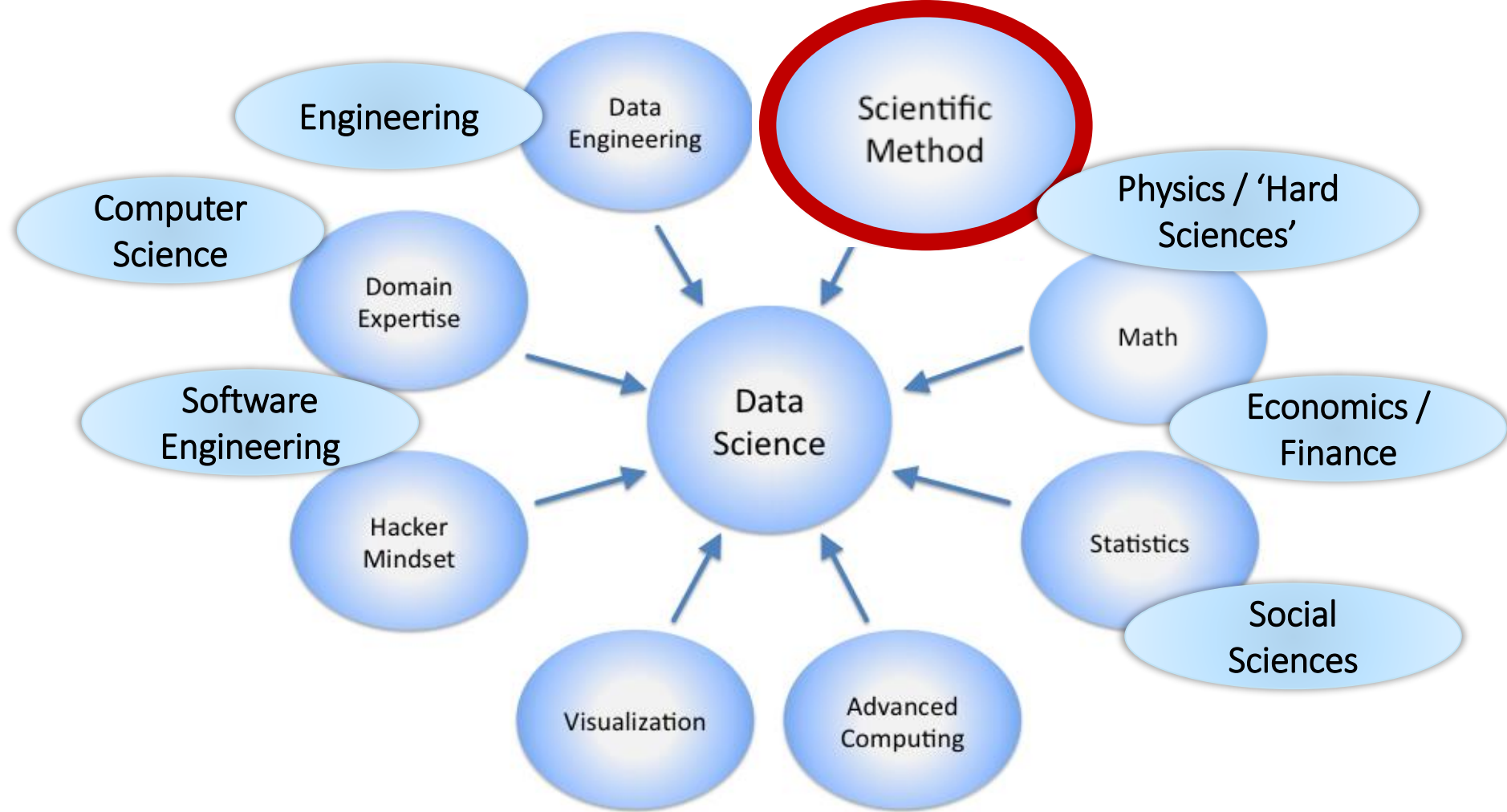
DESCRIPTIVE

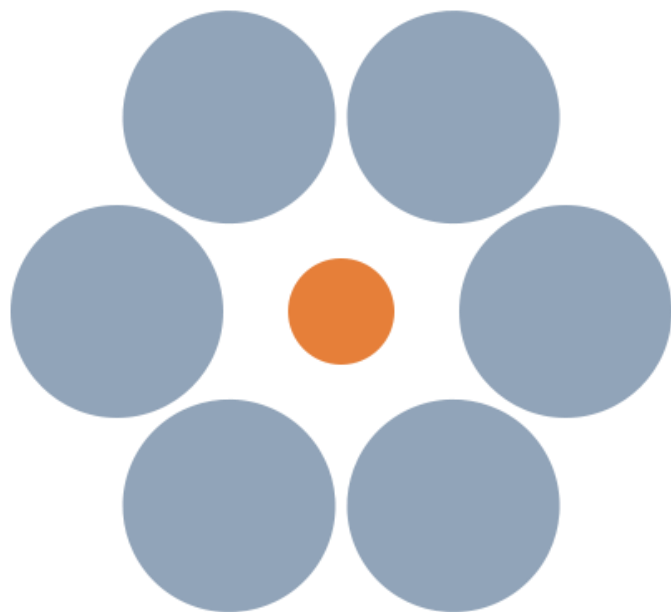
SOPHISTICATION



VALUE



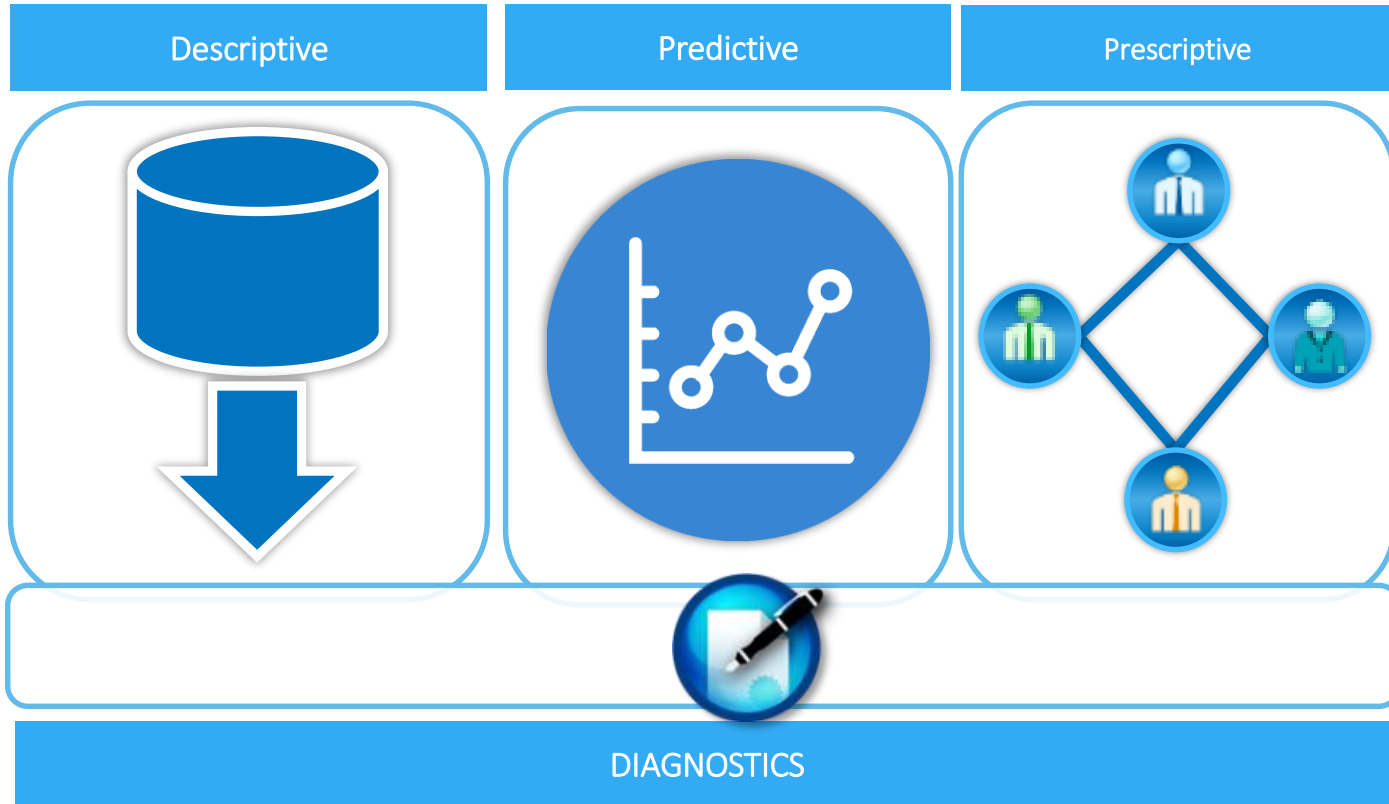




Scientific test...

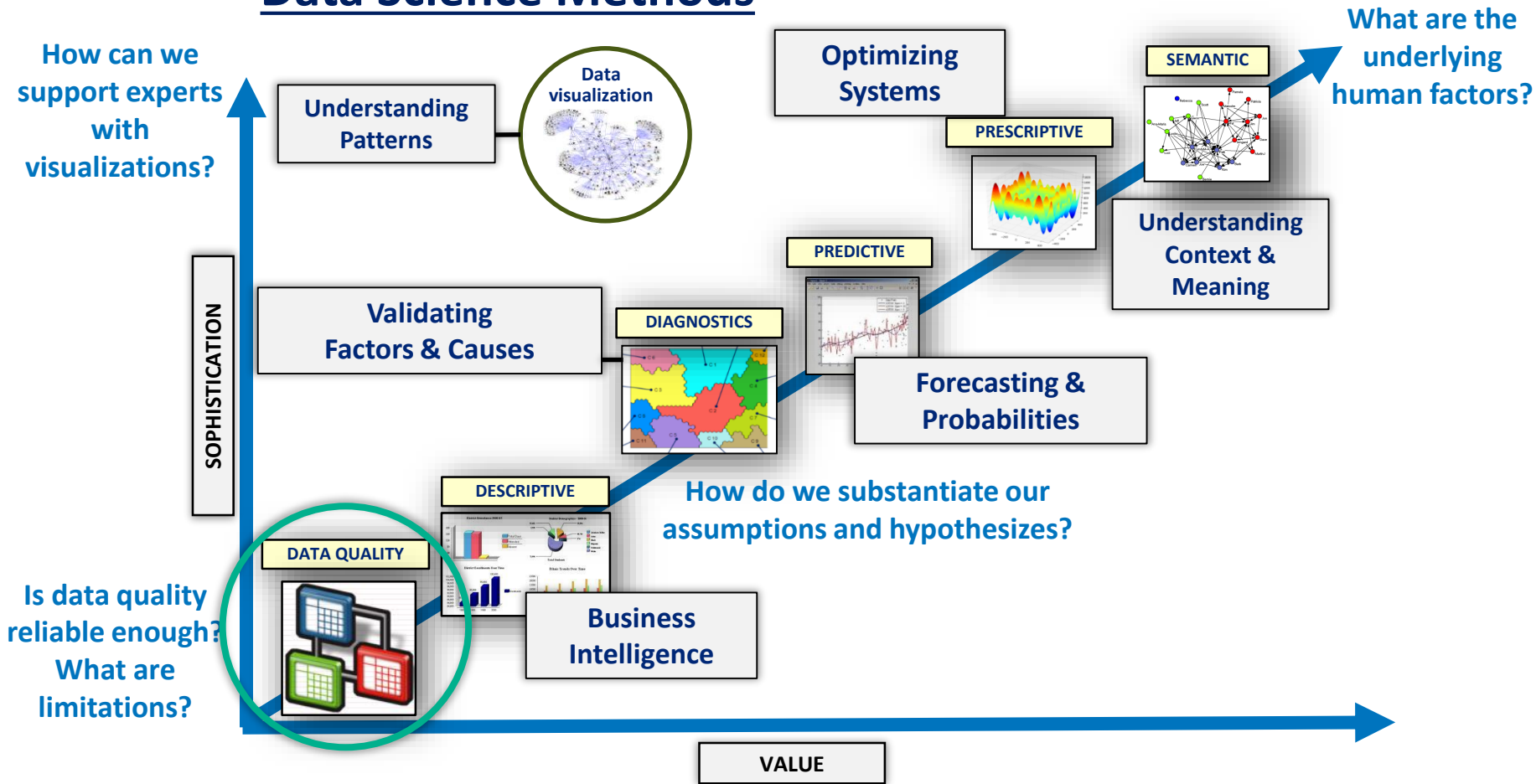


# VALIDATE WITH DIAGNOSTICS



i.e. statistical tests, causal and explanatory models, experiments, validation, model performance, etc.

# Data Science Methods



# 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

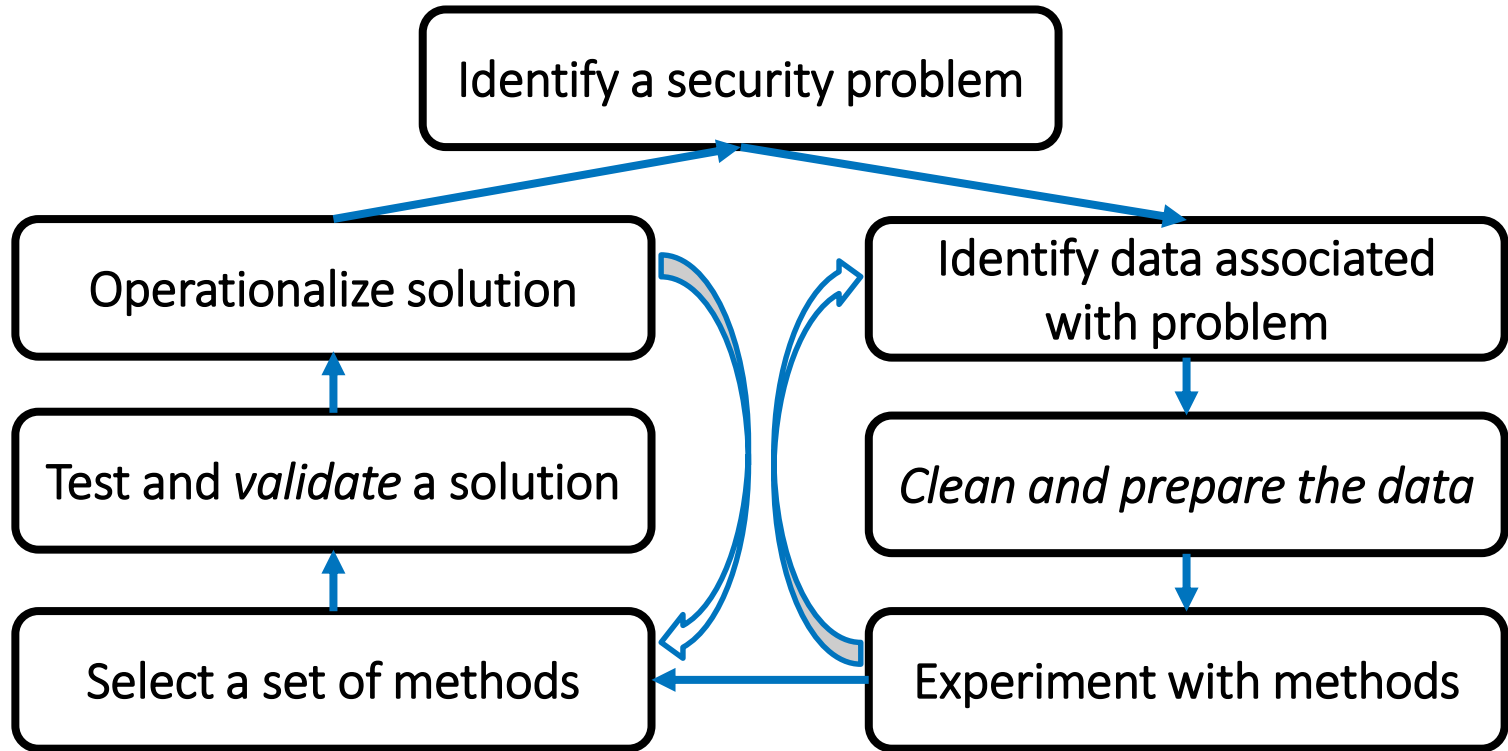


**Data engineering** to structured and integrate distributed big data into '**smart data**'



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

# Cybersecurity Data Science: Typical Sequence



# Data Foundations: Log File Analytics



# The devil is in the data

010001011010001000101100000001101010111100010000110001  
100000110100101000101100100111011100101001000001100100  
00001101101100110011011100110010101111000100001000011  
010101101000100010010001000100101011011101101010011101  
111000010011001100100111100011010111100001001100100001  
111011010010100111000100101000110100110010011001010010  
010111110100100011101000100011010100111011000101100111  
110001001011001000100110100011010101101000010000100001  
0100010110100010000100000001101010101100001010010010  
101010110100010001001111000110101011010100100001000001  
1110010011001000110000100011010101111001010000111001  
011100011010001110000101111100001110011001110101001  
101001011000001000110000100011010101111000111100100010  
010101101011101000101000100010010101011100001000000011  
011101000001100100100011001001010111011010100110101

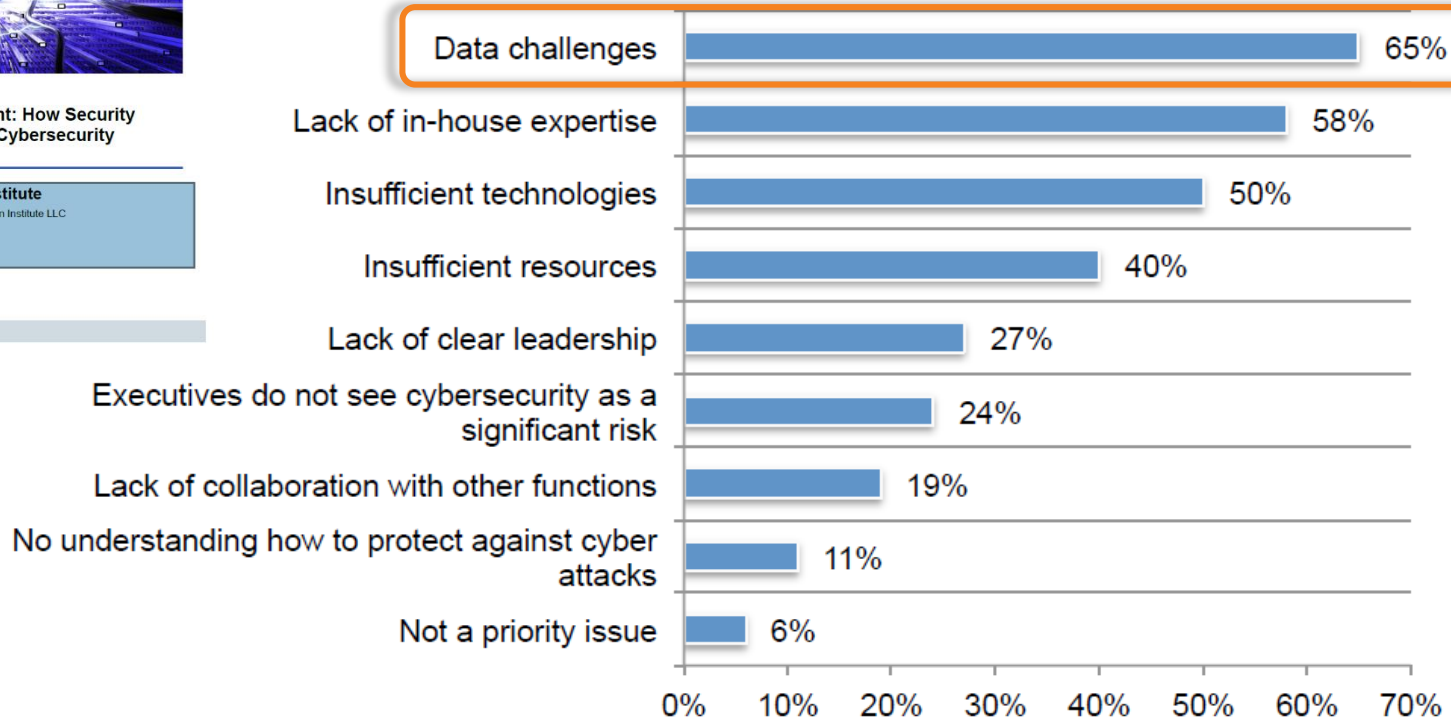


When Seconds Count: How Security Analytics Improves Cybersecurity Defenses

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

Ponemon Institute® Research Report

# Challenges preventing successful use of cybersecurity analytics\*





## 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](https://www.sas.com/en_th/whitepapers/analyze-this-108217.html)

## Log files are the most common source of cybersecurity data...

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

Log files

60%

Network flow

48%

Identity and access management systems

43%

Physical security systems

43%

Endpoint monitoring

40%

Packet capture

39%

SIEM

19%

# Security Log Data Sources

## Operating systems

- Windows & registry events
- Unix logs

## Data access

- File system
- ODBC
- RMDB (i.e. Oracle, SQL Server)

## Authentication and Authorization Reports

- Active Directory Services
- Kerberos
- Proxy logs

## Network logs

- Firewall logs
- Router
- Traffic flows / packet captures

## • Web server

- IIS W3C Logs, Apache
- HTTP Error Logs, URL Scans

## • Resource access and performance

## • Physical security access

## • Malware and endpoint activity

## • Failure and critical error reports

## • Specialized

- Cloud applications (i.e. SaaS)
- *Specialized systems / applications (i.e. ERP, thin client, specialized or home-grown)*

SOURCE: Marty, Raffael. *Applied Security Visualization*. Pearson Education.

# Use Cases: Security Log File Analytics

- User-entity behavioral analytics (UEBA)\*
- Performance insights
- Optimization of resources
- Asset tracking
- Refining focused alerts
- Refining risk indicator metrics
- Enriching SIEM or other repositories
- Data / asset protection
- Preventive insights
- Identify / attribute attacks
- Incident response
- Incident root-cause analysis
- Compliance / risk reporting
- CISO dashboard

\* [Global Behavioral Analytics Market](#) \$3.5 billion projected market value by 2024



# Diving into Cybersecurity Log Data

This exercise illustrates an example of examining cybersecurity log data with a variety of tools.

# Diving Right In!

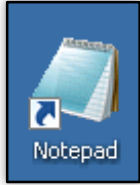
## EXERCISE: Assessing Raw Log Data

```
authlog - Notepad
File Edit Format View Help

Mar 16 08:12:04 app-1 login[4659]: pam_unix(login:session): session opened for user user3 by LOGIN(uid=0)Mar 16 08:12:09 app-1 sudo: user3 : TTY=ttty1 ; PW
0]: removed group 'user4' owned by 'user4' Mar 16 08:12:38 app-1 groupadd[4702]: new group: name=user1, GID=1001Mar 16 08:12:38 app-1 useradd[4703]: new user:
name=sshd, UID=104, GID=65534, home=/var/run/ssh, shell=/usr/sbin/nologinMar 16 08:25:22 app-1 usermod[4846]: change user 'sshd' passwordMar 16 08:25:22 app-1
med for user root by user3(uid=0)Mar 16 09:17:01 app-1 CRON[5085]: pam_unix(cron:session): session opened for user root by (uid=0)Mar 16 09:17:01 app-1 CRON[54
(cron:session): session opened for user root by (uid=0)Mar 16 10:17:01 app-1 CRON[5184]: pam_unix(cron:session): session closed for user rootMar 16 10:39:54 aj
session closed for user rootMar 16 15:17:01 app-1 CRON[5435]: pam_unix(cron:session): session opened for user root by (uid=0)Mar 16 15:17:01 app-1 CRON[5435]:
t by user3(uid=0)Mar 16 17:12:41 app-1 sudo: pam_unix(sudo:session): session closed for user rootMar 16 17:12:41 app-1 su[5535]: Successful su for root by root
(cron:session): session closed for user rootMar 16 17:32:58 app-1 su[4679]: pam_unix(su:session): session closed for user rootMar 16 09:43:06 app-1 sudo: pam_unix(sudo:session): se
D=/bin/suMar 18 09:43:06 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 09:43:06 app-1 sudo: pam_unix(sudo:session): se
r rootMar 18 09:49:52 app-1 su[4673]: Successful su for root by rootMar 18 09:49:52 app-1 su[4673]: + ttty1 root:rootMar 18 09:49:52 app-1 su[4673]: pam_unix(s
sued for user root by user3(uid=0)Mar 18 09:54:25 app-1 sshd[4614]: Server listening on :: port 22,Mar 18 09:54:26 app-1 sshd[4614]: error: Bind to port 22 on
r 18 10:00:06 app-1 passwd[4763]: pam_unix(passwd:chautok): password changed for user1Mar 18 10:00:10 app-1 sshd[4764]: Accepted password for user1 from 76.
TTYpts/0 ; PwD=/home/user1 ; USER=root ; COMMAND=/bin/su -Mar 18 10:02:09 app-1 sudo: user1 : TTY=pts/0 ; PwD=/home/user1 ; USER=root ; COMMAND=/bin/su -Mar
ql, UID=106, GID=115, home=/var/lib/mysql, shell=/bin/falseMar 18 10:18:26 app-1 chage[6967]: changed password expiry for mysqlMar 18 10:18:26 app-1 chfn[6968]
for root by rootMar 18 10:35:25 app-1 su[8909]: + pts/1 root:rootMar 18 10:35:25 app-1 su[8909]: pam_unix(su:session): session opened for user root by user3(ui
for user rootMar 18 10:36:14 app-1 su[8921]: pam_unix(su:session): session closed for user rootMar 18 10:38:07 app-1 sudo: user1 : TTY=pts/0 ; PwD=/opt/soft
8 10:41:09 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 10:41:09 app-1 sudo: pam_unix(sudo:session): session closed f
44:09 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 10:44:09 app-1 sudo: pam_unix(sudo:session): session closed for u
Mar 18 10:55:31 app-1 sudo: pam_unix(sudo:session): session closed for user rootMar 18 10:56:12 app-1 sudo: user1 : TTY=pts/0 ; PwD=/opt/software/web/app ; l
do:session): session closed for user rootMar 18 10:59:49 app-1 sudo: user1 : TTY=pts/0 ; PwD=/opt/software/web/app ; USER=root ; COMMAND=/usr/bin/vi /opt/softwa
sion): session closed for user rootMar 18 11:01:49 app-1 sudo: user1 : TTY=pts/0 ; PwD=/opt/software/web/app ; USER=root ; COMMAND=/usr/bin/vi /opt/software
udo:session): session closed for user rootMar 18 11:04:52 app-1 sudo: user1 : TTY=pts/0 ; PwD=/opt/software/web/app ; USER=root ; COMMAND=/usr/bin/vi /home/
ened for user root by (uid=0)<78>Mar 18 11:09:01 /usr/sbin/cron[9398]: (root) CMD ( [ -x /usr/lib/phps/maxlifetime ] && [ -d /var/lib/phps ] && find /var/lib/
rootMar 18 11:20:19 app-1 su[9504]: pam_authenticate: Authentication failureMar 18 11:20:19 app-1 su[9504]: FAILED su for root by user1Mar 18 11:20:19 app-1 su
su[9507]: pam_unix(su:session): session closed for user rootMar 18 11:21:02 app-1 sudo: user1 : TTY=pts/0 ; PwD=/var/log ; USER=root ; COMMAND=/bin/su - user
com/$DOMAIN/g /etc/apache2/sites-available/wwwMar 18 11:23:26 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 11:23:26 a
root ; COMMAND=/etc/init.d/apache2 restartMar 18 11:23:56 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 11:23:56 app-1
125:00 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 11:25:08 app-1 sudo: pam_unix(sudo:session): session closed for
root by user1(uid=0)Mar 18 11:27:11 app-1 sudo: pam_unix(sudo:session): session closed for user rootMar 18 11:27:19 app-1 sudo: user1 : TTY=pts/0 ; PwD=/etc/
session): session closed for user rootMar 18 11:28:28 app-1 sudo: user1 : TTY=pts/0 ; PwD=/etc/apache2 ; USER=root ; COMMAND=/bin/cp domain.org.key domain.o
18 11:29:12 app-1 sudo: user1 : TTY=pts/0 ; PwD=/etc/apache2 ; USER=root ; COMMAND=/usr/bin/openssl x509 -req -days 365 -in domain.org.csr -signkey domain.o
ession): session closed for user rootMar 18 11:32:02 app-1 sudo: user1 : TTY=pts/0 ; PwD=/etc/apache2 ; USER=root ; COMMAND=/etc/init.d/apache2 restartMar 11
/apache2 ; USER=root ; COMMAND=/usr/bin/vi /opt/software/base/vmscripts/app/base_setup.sshMar 18 11:33:51 app-1 sudo: pam_unix(sudo:session): session opened fo
2/domain.org.crt /etc/apache2/domain.org.csr /etc/apache2/domain.org.key .Mar 18 11:35:10 app-1 sudo: pam_unix(sudo:session): session opened for user root by u
TTYpts/0 ; PwD=/opt/bin/g ; USER=root ; COMMAND=/usr/bin/openssl x509 -req -days 365 -in domain.org.csr -signkey domain.org.csr .Mar 18 11:35:10 app-1 sudo: pam_unix(
me=uid=0 euid=0 tty=ssh ruser= rhost=192.171.132.129.212.dsl.pltn13.pacbell.net. user=user2Mar 18 11:38:59 app-1 sshd[10158]: Failed password for user2 from :
2 from 71.132.129.212 port 40961 ssh2Mar 18 11:40:56 app-1 sshd[10202]: pam_unix(sshd:session): session opened for user user2 by (uid=0)Mar 18 11:41:38 app-1 :
port 41296 ssh2Mar 18 11:51:31 app-1 sshd[10296]: pam_unix(sshd:session): session opened for user user2 by (uid=0)Mar 18 11:59:43 app-1 sshd[10333]: Accepted f
rhost=d192-24-91-113.tty.wideopenwest.com
59413 Apr 19 13:11:14 app-1 sshd[31673]: Failed password for invalid user amanda from 24.192.113.91 port 58114 ssh2
59414 Apr 19 13:11:12 app-1 sshd[31675]: Invalid user rpm from 24.192.113.91
59415 Apr 19 13:11:12 app-1 sshd[31675]: pam_unix(sshd:auth): check pass; user unknown
59416 Apr 19 13:11:12 app-1 sshd[31675]: pam_unix(sshd:auth): authentication failure; logname= uid=0 euid=0 tty=ssh ruser=
rhost=d192-24-91-113.tty.wideopenwest.com
59417 Apr 19 13:11:14 app-1 sshd[31675]: Failed password for invalid user rpm from 24.192.113.91 port 58237 ssh2
59418 Apr 19 13:11:14 app-1 sshd[31677]: Invalid user operator from 24.192.113.91
59419 Apr 19 13:11:14 app-1 sshd[31677]: pam_unix(sshd:auth): check pass; user unknown
59420 Apr 19 13:11:14 app-1 sshd[31677]: pam_unix(sshd:auth): authentication failure; logname= uid=0 euid=0 tty=ssh ruser=
rhost=d192-24-91-113.tty.wideopenwest.com
59421 Apr 19 13:11:16 app-1 sshd[31677]: Failed password for invalid user operator from 24.192.113.91 port 58352 ssh2
59422 Apr 19 13:11:17 app-1 sshd[31679]: Invalid user sgi from 24.192.113.91
```

# Demonstration / Exercise

## EXERCISE: Assessing Raw Log Data

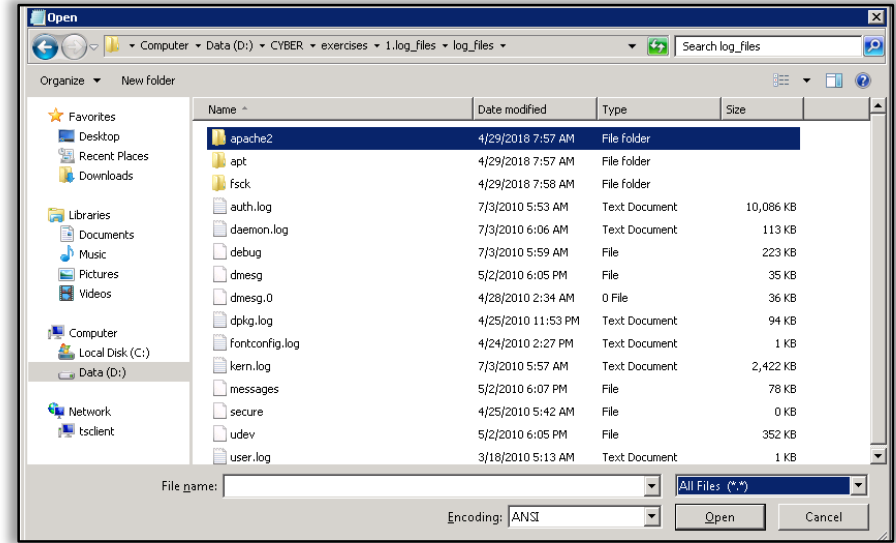


1. Open 'Notepad'

2. File => Open

3. Select 'All Files (\*.\*)' (lower right)

3. D: => @CYBER => 1.FRAME =>  
B.Log\_Files => log\_files => *auth.log*



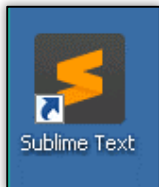
```

Mar 16 08:12:04 app-1 login[4659]: pam_unix(login:session): session opened for user user3 by LOGIN(uid=0)Mar 16 08:12:09 app-1 sudo: user3 : TTY=ttty1 ;
0]: removed group `user4' owned by `user4' Mar 16 08:12:38 app-1 groupadd[4702]: new group: name=user1, GID=1001Mar 16 08:12:38 app-1 useradd[4703]: new use
name=sshd, UID=104, GID=65534, home=/var/run/sshd, shell=/usr/sbin/nologinMar 16 08:25:22 app-1 usermod[4846]: change user `sshd' passwordMar 16 08:25:22 ap
ned for user root by user3(uid=0)Mar 16 09:17:01 app-1 CRON[5085]: pam_unix(cron:session): session opened for user root by (uid=0)Mar 16 09:17:01 app-1 CRO
(cron:session): session opened for user root by (uid=0)Mar 16 10:17:01 app-1 CRON[5184]: pam_unix(cron:session): session closed for user rootMar 16 10:39:54
session closed for user rootMar 16 15:17:01 app-1 CRON[5435]: pam_unix(cron:session): session opened for user root by (uid=0)Mar 16 15:17:01 app-1 CRON[5435]
t by user3(uid=0)Mar 16 17:12:41 app-1 sudo: pam_unix(sudo:session): session closed for user rootMar 16 17:12:41 app-1 su[5535]: Successful su for root by
(cron:session): session closed for user rootMar 16 17:32:58 app-1 su[4679]: pam_unix(su:session): session closed for user rootMar 18 09:41:44 app-1 sshd[46
D=/bin/suMar 18 09:43:06 app-1 sudo: pam_unix(sudo:session): session opened for user root by user3(uid=0)Mar 18 09:43:06 app-1 sudo: pam_unix(sudo:session)
r rootMar 18 09:49:52 app-1 su[4673]: Successful su for root by rootMar 18 09:49:52 app-1 su[4673]: + tty1 root:rootMar 18 09:49:52 app-1 su[4673]: pam_uni
ened for user root by user3(uid=0)Mar 18 09:54:25 app-1 sshd[4614]: Server listening on :: port 22.Mar 18 09:54:26 app-1 sshd[4614]: error: Bind to port 22
r 18 10:00:06 app-1 passwd[4763]: pam_unix(passwd:chauthtok): password changed for user1Mar 18 10:00:10 app-1 sshd[4764]: Accepted password for user1 from
TTY=pts/0 ; PWD=/home/user1 ; USER=root ; COMMAND=/bin/su -Mar 18 10:02:09 app-1 sudo: user1 : TTY=pts/0 ; PWD=/home/user1 ; USER=root ; COMMAND=/bin/su
ql, UID=106, GID=115, home=/var/lib/mysql, shell=/bin/falseMar 18 10:18:26 app-1 chage[6967]: changed password expiry for mysqlMar 18 10:18:26 app-1 chfn[69
for root by rootMar 18 10:35:25 app-1 su[8909]: + pts/1 root:rootMar 18 10:35:25 app-1 su[8909]: pam_unix(su:session): session opened for user root by user1
for user rootMar 18 10:36:14 app-1 su[8921]: pam_unix(su:session): session closed for user rootMar 18 10:38:07 app-1 sudo: user1 : TTY=pts/0 ; PWD=/opt/sc
8 10:41:09 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 10:41:09 app-1 sudo: pam_unix(sudo:session): session close
44:09 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 10:44:09 app-1 sudo: pam_unix(sudo:session): session closed for
Mar 18 10:55:31 app-1 sudo: pam_unix(sudo:session): session closed for user rootMar 18 10:56:12 app-1 sudo: user1 : TTY=pts/0 ; PWD=/opt/software/web/app
do:session): session closed for user rootMar 18 10:59:49 app-1 sudo: user1 : TTY=pts/0 ; PWD=/opt/software/web/app ; USER=root ; COMMAND=/usr/bin/vi /opt,
sion): session closed for user rootMar 18 11:01:49 app-1 sudo: user1 : TTY=pts/0 ; PWD=/opt/software/web/app ; USER=root ; COMMAND=/usr/bin/vi /opt/softwa
udo:session): session closed for user rootMar 18 11:04:52 app-1 sudo: user1 : TTY=pts/0 ; PWD=/opt/software/web/app ; USER=root ; COMMAND=/usr/bin/vi /hor
ened for user root by (uid=0)<78>Mar 18 11:09:01 /USR/SBIN/CRON[9398]: (root) CMD ( [ -x /usr/lib/php5/maxlifetime ] && [ -d /var/lib/php5 ] && find /var/
rootMar 18 11:20:19 app-1 su[9504]: pam_authenticate: Authentication failureMar 18 11:20:19 app-1 su[9504]: FAILED su for root by user1Mar 18 11:20:19 app-
su[9507]: pam_unix(su:session): session closed for user rootMar 18 11:21:02 app-1 sudo: user1 : TTY=pts/0 ; PWD=/var/log ; USER=root ; COMMAND=/bin/su -
com/$DOMAIN/g /etc/apache2/sites-available/wwwMar 18 11:23:26 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 11:23:
root ; COMMAND=/etc/init.d/apache2 restartMar 18 11:23:56 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 11:23:56 ap
:25:08 app-1 sudo: pam_unix(sudo:session): session opened for user root by user1(uid=0)Mar 18 11:25:08 app-1 sudo: pam_unix(sudo:session): session closed fo
root by user1(uid=0)Mar 18 11:27:11 app-1 sudo: pam_unix(sudo:session): session closed for user rootMar 18 11:27:19 app-1 sudo: user1 : TTY=pts/0 ; PWD=/e
session): session closed for user rootMar 18 11:28:28 app-1 sudo: user1 : TTY=pts/0 ; PWD=/etc/apache2 ; USER=root ; COMMAND=/bin/cp domain.org.key domain
18 11:29:12 app-1 sudo: user1 : TTY=pts/0 ; PWD=/etc/apache2 ; USER=root ; COMMAND=/usr/bin/openssl x509 -req -days 365 -in domain.org.csr -signkey domain
session): session closed for user rootMar 18 11:32:02 app-1 sudo: user1 : TTY=pts/0 ; PWD=/etc/apache2 ; USER=root ; COMMAND=/etc/init.d/apache2 restartMar
/apache2 ; USER=root ; COMMAND=/usr/bin/vi /opt/software/base/vmscripts/app/base_setup.shMar 18 11:33:51 app-1 sudo: pam_unix(sudo:session): session opened
2/domain.org.crt /etc/apache2/domain.org.csr /etc/apache2/domain.org.key .Mar 18 11:35:10 app-1 sudo: pam_unix(sudo:session): session opened for user root
TTY=pts/0 ; PWD=/opt/software/web/config ; USER=root ; COMMAND=/usr/bin/vi /opt/software/base/vmscripts/app/apache_setup.shMar 18 11:35:50 app-1 sudo: pam_u

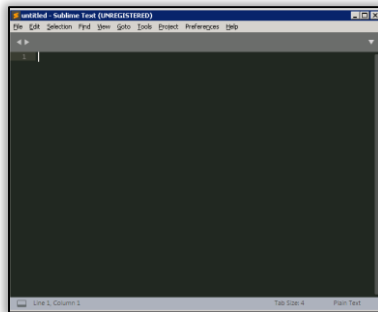
```

# Demonstration / Exercise

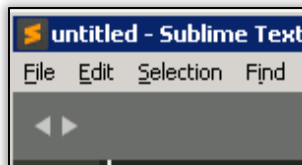
## EXERCISE: Assessing Semi-Structured Log Data



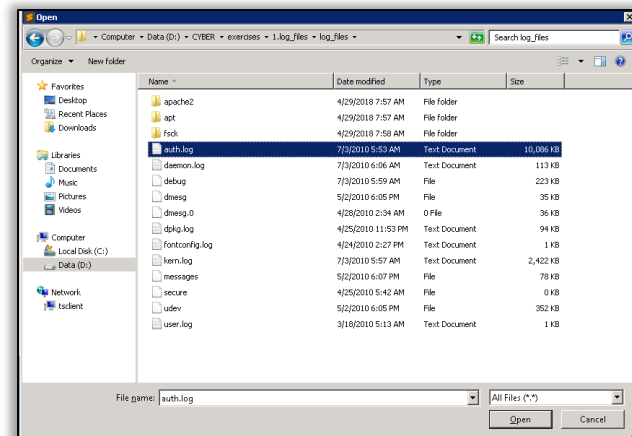
1. Open 'Sublime Text'



2. File => Open File



3. D: => @CYBER => 1.FRAME => B.Log\_Files  
=> log\_files => *auth.log*



```
auth.log
1 Mar 16 08:12:04 app-1 login[4659]: pam_unix(login:session): session opened for user
  user3 by LOGIN(uid=0)
2 Mar 16 08:12:09 app-1 sudo:      user3 : TTY=tty1 ; PWD=/home/user3 ; USER=root ;
  COMMAND=/bin/su
3 Mar 16 08:12:09 app-1 sudo: pam_unix(sudo:session): session opened for user root by
  user3(uid=0)
4 Mar 16 08:12:09 app-1 sudo: pam_unix(sudo:session): session closed for user root
5 Mar 16 08:12:09 app-1 su[4679]: Successful su for root by root
6 Mar 16 08:12:09 app-1 su[4679]: + tty1 root:root
7 Mar 16 08:12:09 app-1 su[4679]: pam_unix(su:session): session opened for user root by
  user3(uid=0)
8 Mar 16 08:12:13 app-1 groupadd[4691]: new group: name=user4, GID=1001
9 Mar 16 08:12:13 app-1 useradd[4692]: new user: name=user4, UID=1001, GID=1001, home=/home
  /user4, shell=/bin/bash
10 Mar 16 08:12:17 app-1 passwd[4695]: pam_unix(passwd:chauthtok): password changed for
  user4
11 Mar 16 08:12:22 app-1 chfn[4696]: changed user `user4' information
12 Mar 16 08:12:31 app-1 userdel[4700]: delete user `user4'
13 Mar 16 08:12:31 app-1 userdel[4700]: removed group `user4' owned by `user4'
14 Mar 16 08:12:38 app-1 groupadd[4702]: new group: name=user1, GID=1001
15 Mar 16 08:12:38 app-1 useradd[4703]: new user: name=user1, UID=1001, GID=1001, home=/home
  /user1, shell=/bin/bash
16 Mar 16 08:12:44 app-1 passwd[4706]: pam_unix(passwd:chauthtok): password changed for
  user1
17 Mar 16 08:12:46 app-1 chfn[4707]: changed user `user1' information
18 Mar 16 08:12:49 app-1 chfn[4708]: changed user `user1' information
19 Mar 16 08:12:55 app-1 groupadd[4710]: new group: name=user2, GID=1002
20 Mar 16 08:12:55 app-1 useradd[4711]: new user: name=user2, UID=1002, GID=1002, home=/home
  /user2, shell=/bin/bash
21 Mar 16 08:13:00 app-1 passwd[4714]: pam_unix(passwd:chauthtok): password changed for
  user2
```

# Processing Raw Cybersecurity Data – Example

## EXERCISE: Assessing Semi-Structured Log Data

### 1. What are we looking at – what is this?

- What system? [http://honeynet.org/challenges/2010\\_5\\_log\\_mysteries](http://honeynet.org/challenges/2010_5_log_mysteries)
- What types of logs / format? <https://help.ubuntu.com/community/LinuxLogFiles>
- How ‘verbose’? Inherent structure (or lack thereof)?
- How might we build context (structure more granularly)?
  - Schema availability? Data dictionary? Guidance from experts?
  - For example, <https://help.ubuntu.com/community/LinuxLogFiles>
  - Codes and indicators - descriptive versus opaque?

# Cyber Forensics: Log File Analysis – HoneyNet Challenge

## Virtual Server Log File Analysis (Virtual Ubuntu Linux Server)

1. **Was the system compromised** and when? How do you know that for sure?
2. If the system was compromised, **what was the method used**?
3. Can you locate how many attackers failed? If some succeeded, how many were they?
4. How many stopped attacking after the first success?
5. What happened after the brute force attack?
6. Locate the authentication logs. Was a brute force attack performed? if yes, how many?
7. What is the timeline of significant events? How certain are you of the timing?
8. Anything else that looks suspicious in the logs? Any misconfigurations? Other issues?
9. Was an automatic tool used to perform the attack? if yes, which one?
10. What can you say about the attacker's goals and methods?

[http://honeynet.org/challenges/2010\\_5\\_log\\_mysteries](http://honeynet.org/challenges/2010_5_log_mysteries)

# Processing Raw Cybersecurity Data – Example

## EXERCISE: Assessing Semi-Structured Log Data

### 2. How can we reduce and structure?

- Parsing and extracting – many options
  - Commercial tools (e.g. Splunk, Excel) and free tools (e.g. MS Log Parser)
  - Scripting/programmatic (e.g. Perl, R, Python, SAS DS2, UNIX GREP, PowerShell)

### 3. Where do we store and transform subsequently?

- Flat files
- Database
- SIEM or other specialized repository
- ‘Data Lake’ (note: also can dump the raw logfile there too)
- Analytics platform



# Demonstration / Exercise

## EXERCISE: Assessing Semi-Structured Log Data

1. Open 'Log Parser 2.2'
2. Examples of structured queries to try:

### ALL FAILED AUTHENTICATION EVENTS:

```
logparser -i:TEXTLINE -o:csv "SELECT * INTO  
'D:\@CYBER\1.FRAME\B.Log_Files\results\log_extract_failed.txt' FROM  
'D:\@CYBER\1.FRAME\B.Log_Files\log_files\auth.log' WHERE Text LIKE '%Failed password%'"
```

### ALL FAILED ROOT ACCESS AUTHENTICATION EVENTS – STRUCTURED:

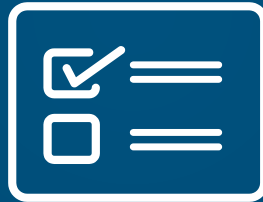
```
logparser -i:TEXTLINE -o:csv "SELECT EXTRACT_TOKEN (Text, 8, ' ') AS UserName,  
EXTRACT_TOKEN (Text, 10, ' ') AS IPSource, EXTRACT_TOKEN (Text, 12, ' ') AS Port,  
EXTRACT_TOKEN (Text, 10, ' ') AS Protocol INTO  
'D:\@CYBER\1.FRAME\B.Log_Files\results\log_extract_failed_root.txt' FROM  
'D:\@CYBER\1.FRAME\B.Log_Files\log_files\auth.log' WHERE Text LIKE '%Failed password for root%'"
```

RESULTS    *D:\@CYBER\1.FRAME\B.Log\_Files\results\*

More fun with LogParser!

<https://mlichtenberg.wordpress.com/2011/02/03/log-parser-rocks-more-than-50-examples/>

# Exercise Review



# Conclusion



# Section Review



# Cybersecurity Data Science as a Process

## Data Engineering



## Advanced Analytics

Diagnostics & patterns    Establishing baselines    Predictive modelling    Anomaly detection    Behavioral insights



## Triage / Validate



## Remediate



Data Manager



Data Scientist



Cyber Investigator



Infosec Response

# Cybersecurity Analytics Maturity

## Anomaly Detection

- Big data management
- Flags, rules, and alerts

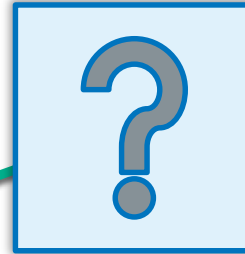
CHASING  
PHANTOM  
PATTERNS



## Data-aware Investigations



## Predictive Detection



## Risk Awareness / Resource Optimization



# Cybersecurity Analytics Maturity

## Anomaly Detection

- Big data management
  - Flags, rules, and alerts
- 
- Structured data
  - Counts of key measures
  - Foundation for comparisons across entities and over time



## Data-aware Investigations



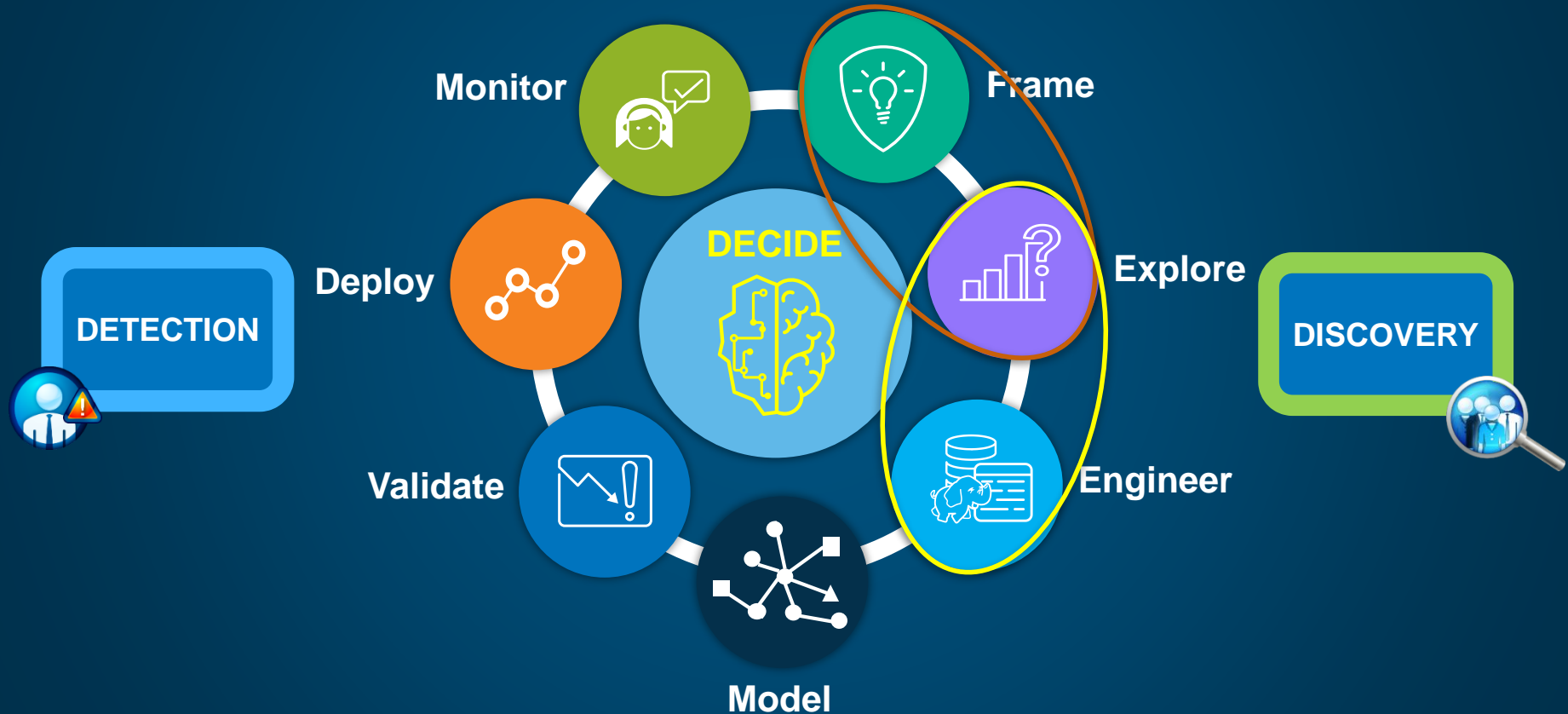
## Predictive Detection



## Risk Awareness / Resource Optimization



# Cybersecurity Data Science (CSDS) Lifecycle



# Cybersecurity Data Science (CSDS)

TOPIC
1. FRAME
2. DATA
3. DISCOVER
4. DETECT
5. DEPLOY