

Cybersecurity Data Science (CSDS)

Best Practices in an Emerging Profession

Scott Allen Mongeau

Cybersecurity Data Scientist – SAS Institute

PhD candidate - Nyenrode Business University (Netherlands)

s.mongeau@edp1.nyenrode.nl

scott.mongeau@sas.com

@SARK7 #CSDS2020 #FloCon2020

FloCon 2020

JANUARY 6-9, 2020 | SAVANNAH, GA



PhD academic research / book

- ~June 2020 release

Research on cybersecurity data science (CSDS) as an emerging profession

- I. Literature: What is CSDS and is it a profession?
- II. Interviews: 50 CSDS practitioners
- III. Designs: Approaches to address challenges



Scott Mongeau

Cybersecurity
Data Science:
Best Practices in an
Emerging Profession

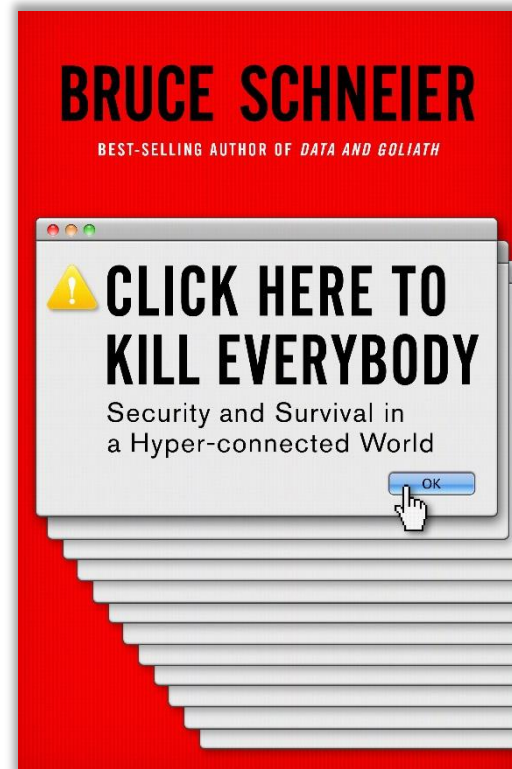
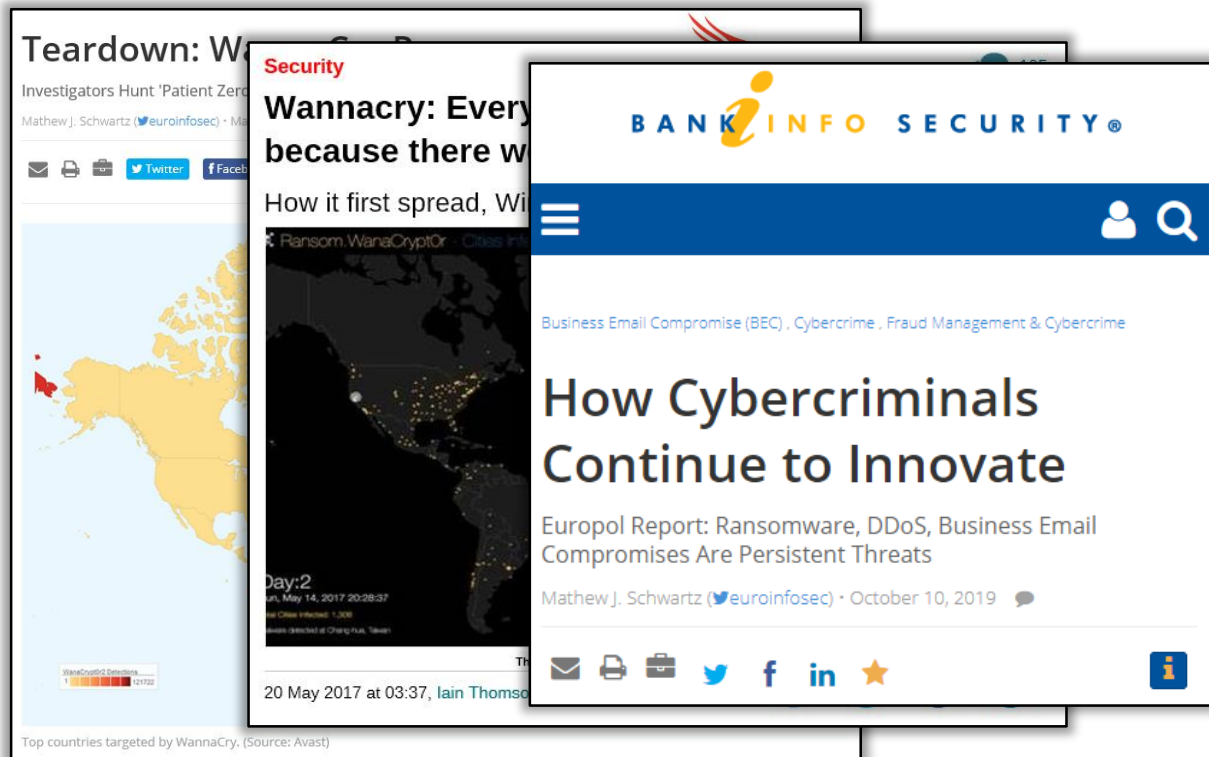
 Springer



I. CSDS Literature

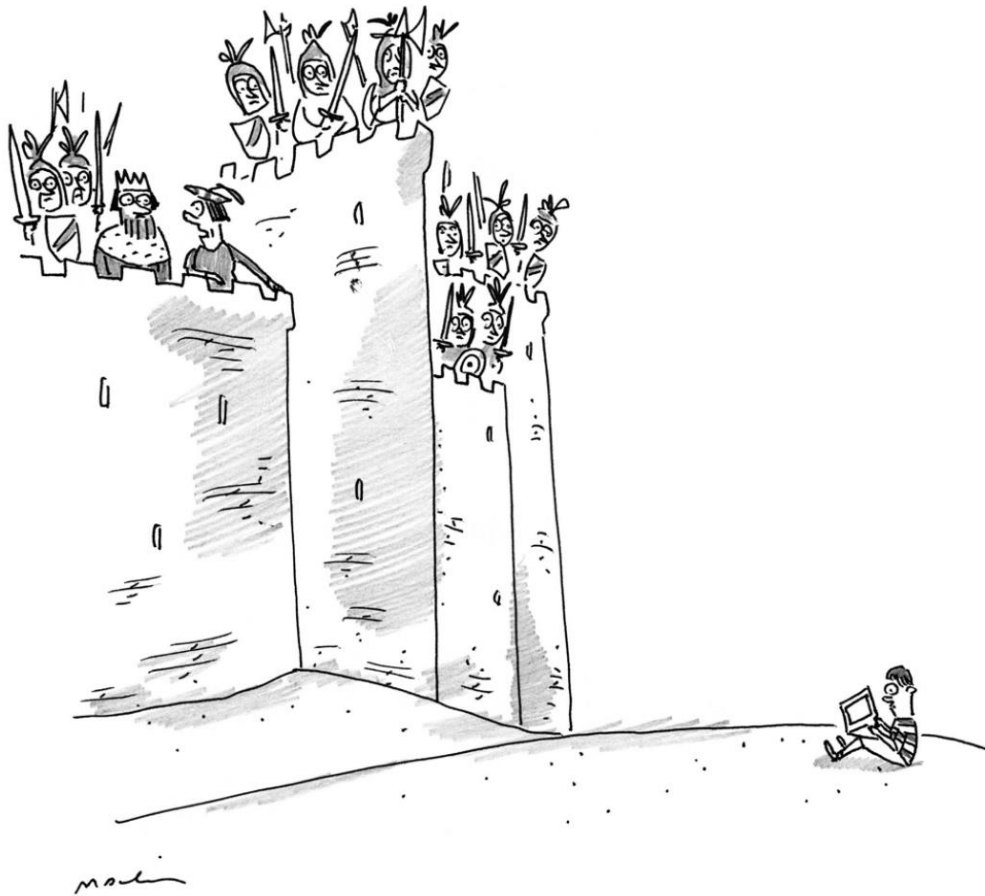
FUD Fear, Uncertainty, Doubt

Expansion of exposure and targets >!< Increasing sophistication, frequency, and speed of attacks

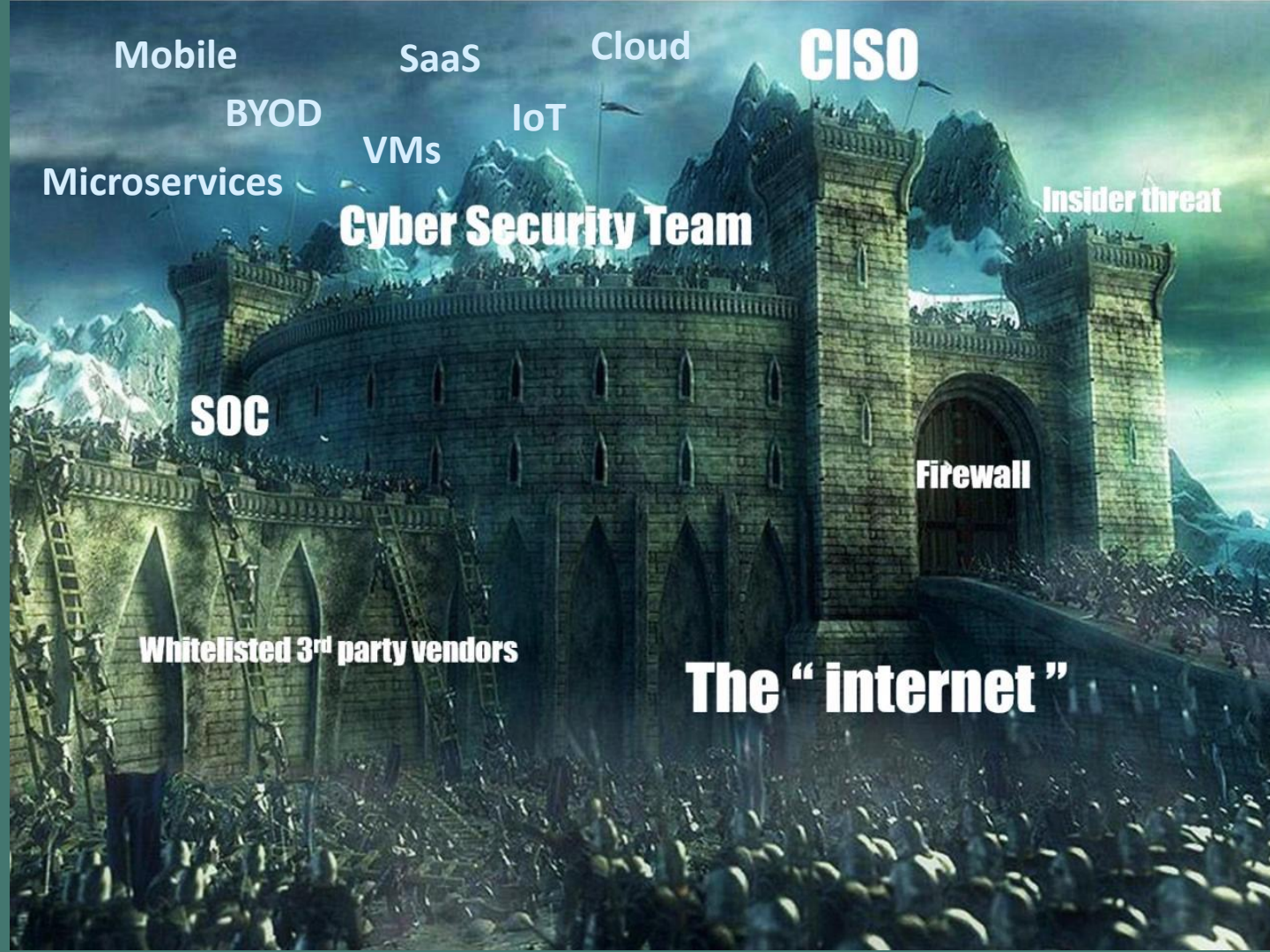


Castle and Moat

How quaint!



"Bad news, Your Majesty—it's a cyberattack."



Mobile

SaaS

Cloud

CISO

BYOD

IoT

VMs

Microservices

Cyber Security Team

Insider threat

SOC

Firewall

Whitelisted 3rd party vendors

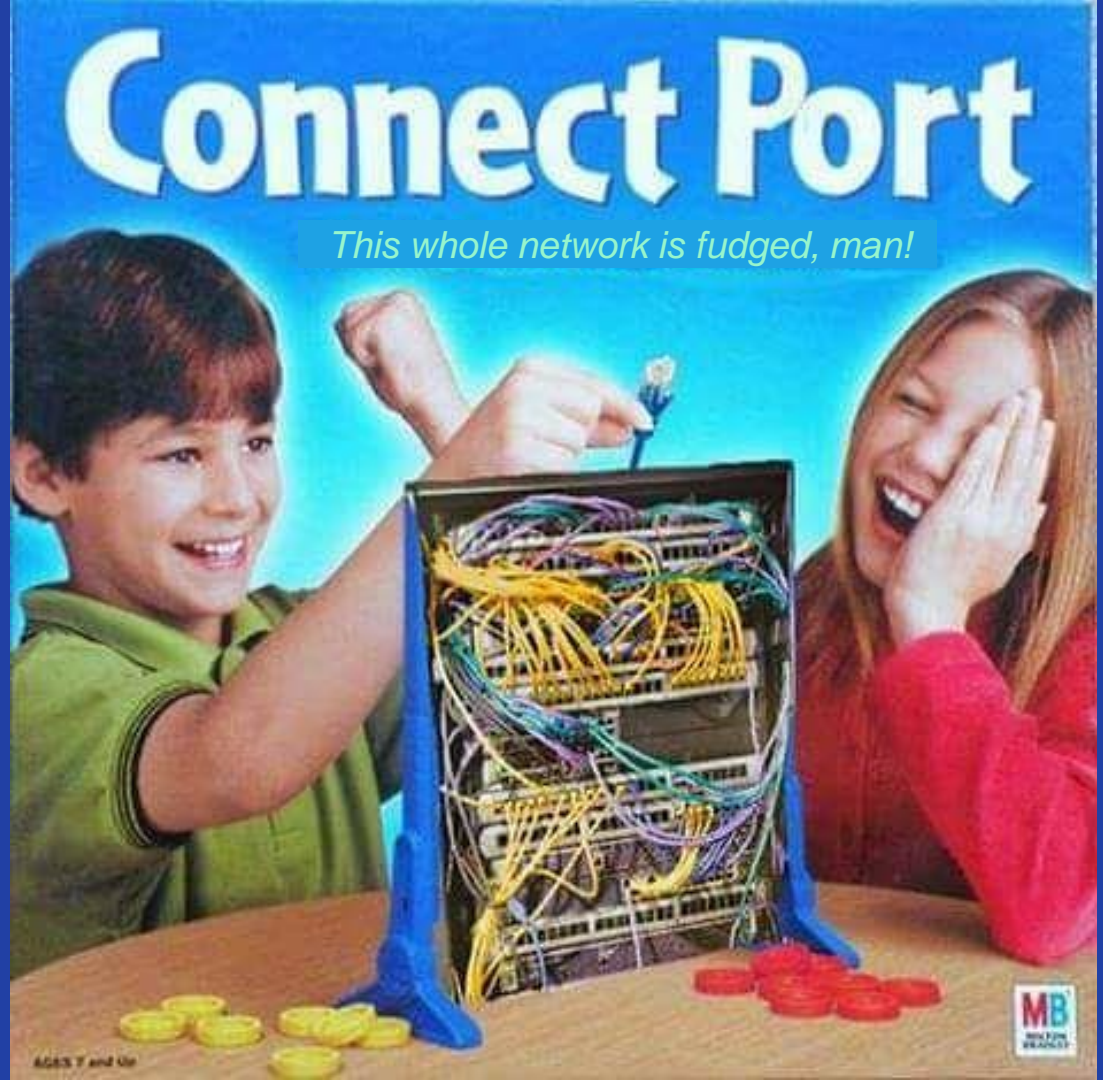
The “internet”

Cybersecurity Challenges

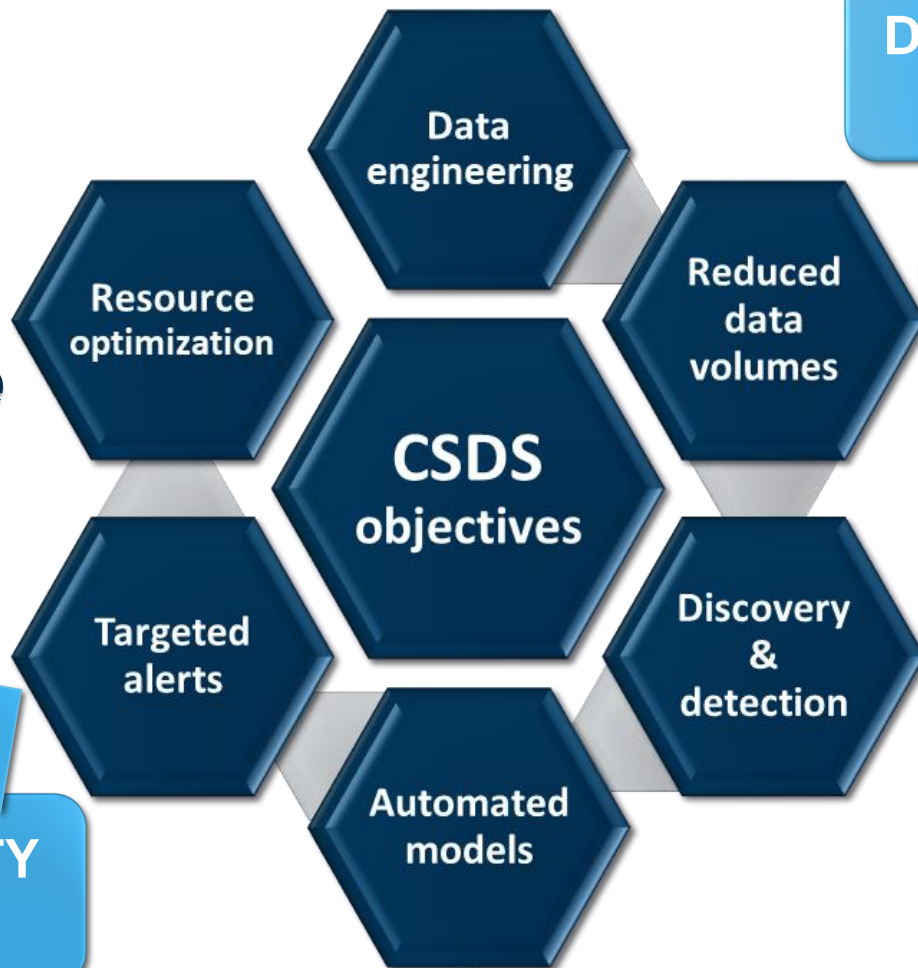


Data Science

New hope amidst
complexity and
confusion...



CSDS
***Cyber
Security
Data
Science***



**DATA SCIENCE
METHODS**

**CYBERSECURITY
GOALS**

CSDS: Existing Professionals + Demonstrated Efficacy

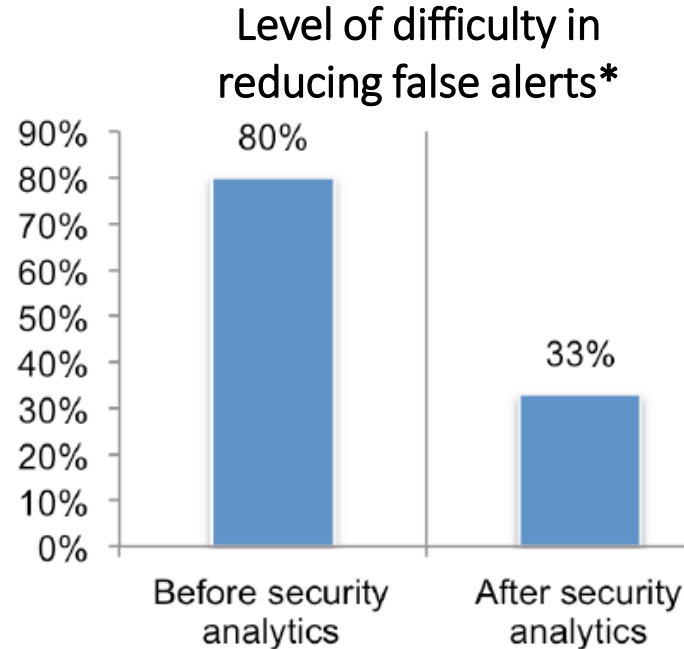


When Seconds Count: How Security Analytics Improves Cybersecurity Defenses

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

Ponemon Institute® Research Report

https://www.sas.com/en_us/whitepapers/ponemon-how-security-analytics-improves-cybersecurity-defenses-108679.html



EXAMPLE CSDS PRACTICAL APPLICATIONS

- Spam filtering
- Phishing email detection
- Malware & virus detection
- Network monitoring
- Endpoint protection

* Survey of 621 global IT security practitioners

'Professional Maturity' Comparison

#	CRITERIA	CYBER	DS	CSDS
1	Broad interest	●	●	●
2	People employed	●	◐	◐
3	Informal training	●	●	◐
4	Informal groups	●	●	◐
5	Professional literature	●	●	◐
6	Research literature	◐	◐	◐
7	Formal training	●	◐	◐
8	Formal prof. groups	●	◐	○
9	Professional certificates	◐	◐	○
10	Standards bodies	●	◐	○
11	Academic discipline	◐	◐	○

CYBER =
Growing challenges +
rapid paradigm shift

DATA SCIENCE =
Poorly defined standards
“whatever you want it to be!”

CSDS =
At risk problem child?

The Blessing and Curse of Data Science

PROS

- Commercial interest
 - Range of methods
- Freedom to experiment
 - Delivers efficiencies
- Big data engineering
 - Insightful questions
- Power of machine learning



CONS

- Hype & noise
- Befuddling array of approaches
- Lack of standards
- Myth of automation
- Big data ipso facto is not solution
- Wait, what is the question?
- “Throwing the statistical baby out with grampa’s bathwater?”



II. CSDS Interviews

CSDS Practitioner Interviews

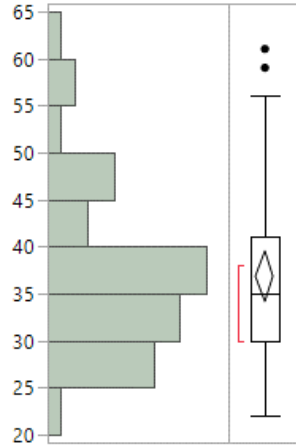
30 minutes per interviewee

- ENTRY: How did you become involved in domain?
- What are perceived central CHALLENGES?
- What are key BEST PRACTICES?

Demographic Profile (n=50)

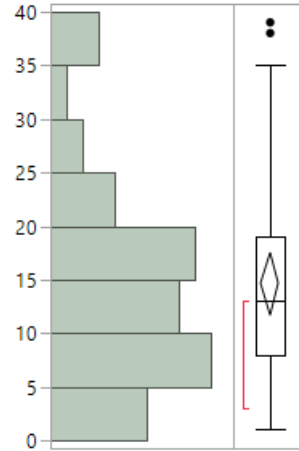
LinkedIn => 350 candidates => 50 participants

Age*



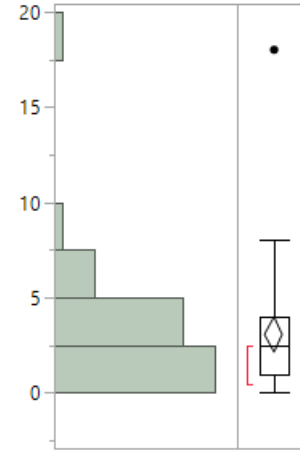
Mean	36.8
StdDev	9.1

Yrs Employed*



Mean	14.2
StdDev	9.5

Yrs CSDS*

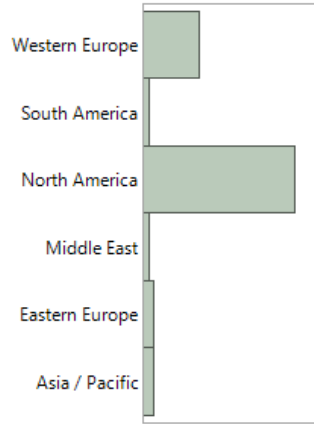


Mean	2.9
StdDev	1.9

** Estimates inferred from LinkedIn profile data*

Demographic Profile (n=50)

Current Region



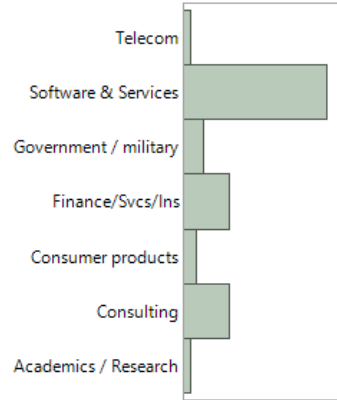
Current Region ¹	n	%
North America	35	70%
Western Europe	10	20%
Eastern Europe	2	4%
Middle East	2	4%
South America	1	2%

22% (n=11) relocated from native region

18% (n=9) relocated to US specifically

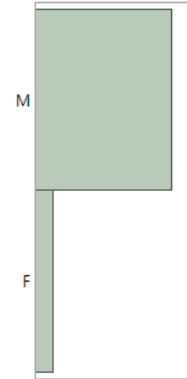
10% (n=5) relocated specifically from Asia/Pacific to US

Current Industry



Industry	n	%
Software and services	28	56%
Consulting	7	14%
Finance/financial services/insurance	7	14%
Government / military	3	6%
Consumer products	2	4%
Academics / research	2	4%
Telecom	1	2%

Gender



Gender	n	%
Male	43	86%
Female	7	14%

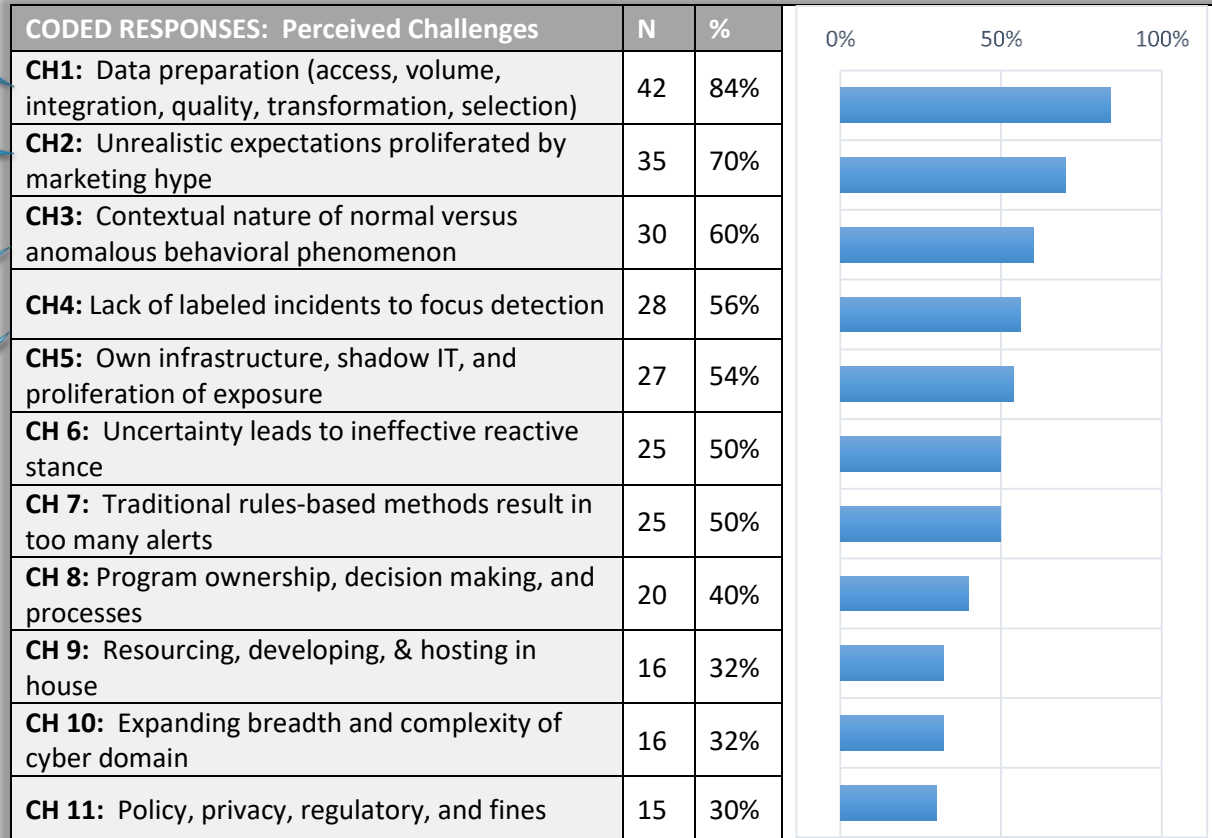
CSDS 'CHALLENGES': 11

DATA PREPARATION!
84%

Marketing hype 70%

Establishing context
60%

Labeled incidents
(evidence) 56%



CSDS 'BEST PRACTICES': 26

DATA PREPARATION!
84%

Cross-domain
collaboration 76%

Scientific rigor 68%

RESPONSES: Advocated best practices	Family	N	%	0%	50%	100%
BP1: Structured data preparation, discovery, engineering process	Proc	42	84%			
BP2: Building process focused cross-functional team	Org	38	76%			
BP3: Cross-training team in data science, cyber, engineering	Org	37	74%			
BP4: Scientific method as a process	Proc	34	68%			
BP5: Instill core cyber domain knowledge	Org	33	66%			
BP6: Vulnerability, anomaly & decision automation to operational capacity	Tech	33	66%			
BP7: Data normalization, frameworks & ontologies	Tech	32	64%			
BP8: Model validation and transparency	Proc	31	62%			
BP9: Data-driven paradigm shift away from rules & signatures	Org	29	58%			
BP10: Track and label incidents and exploits	Proc	28	56%			
BP11: Cyclical unsupervised and supervised machine learning	Proc	25	50%			
BP12: Address AI hype and unrealistic expectations directly	Org	23	46%			
BP13: Understand own infrastructure & environment	Org	23	46%			

RESPONSES: Advocated best practices	Family	N	%	0%	50%	100%
BP14: Cloud and container-based tools and data storage	Tech	22	44%			
BP15: Distinct exploration and detection architectures	Tech	22	44%			
BP16: Participate in data sharing consortiums and initiatives	Tech	21	42%			
BP17: Deriving probabilistic and risk models	Org	20	40%			
BP18: Upper management buy in and support	Org	16	32%			
BP19: Human-in-the-loop reinforcement	Proc	14	28%			
BP20: Survey academic methods and techniques	Org	13	26%			
BP21: Cyber risk as general enterprise risk & reward	Org	12	24%			
BP22: Segment risk programmatically and outsource components	Org	9	18%			
BP23: Adding machine learning to SIEM	Tech	5	10%			
BP24: Preventative threat intelligence	Org	4	8%			
BP25: Hosting and pushing detection to endpoints	Tech	4	8%			
BP26: Honeypots to track and observe adversaries	Tech	2	4%			

KEY CSDS GAPS: Factor-to-Factor Fitting

CH F1 Expansive complexity
CH F2 Tracking & context
CH F3 Data management
CH F4 Expectations versus limitations
CH F5 Unclear ownership
CH F6 Data policies

Challenge
Factor Score

	FACTOR1	FACTOR2
1	-1.10951	-1.2847
2	-0.65954	0.82659
3	-1.14351	0.85817
4	0.27474	0.98433
5	0.185896	1.06243
6	-0.98246	-1.3272
7	-1.19556	-1.3651
8	-1.08428	0.62937
9	0.19231	-1.19096
10	-0.19805	0.990378
11	0.771806	-1.22723
12	-0.93501	0.76347
13	1.374426	-1.3837
14	0.740622	0.65038
15	-0.95034	0.96529
16	0.889892	0.78447
17	-0.03689	0.8046
18	-0.9646	-0.8116
19	0.97118	-1.163
20	-1.17033	0.54904
21	1.328284	-1.243
22	0.092641	0.91744
23	-0.13444	-1.0019
24	0.402174	-1.1042
25	-0.37696	-1.208
26	-0.26951	-1.2847
27	0.827517	0.75920
28	1.460472	-1.3033
29	-1.16343	0.927441
30	-0.16308	0.875596
31	0.558327	0.780959
32	0.024778	-1.0038
33	-0.61325	0.827283
34	-0.15817	0.49019
35	1.399657	0.53047
36	0.175996	-1.0535
37	0.624724	-1.3192
38	-0.64063	-1.146
39	0.978056	0.58732
40	-0.88673	-1.0306
41	-0.7452	-1.2999
42	1.333037	0.78515
43	1.246992	0.70482
44	-1.02385	0.84158
45	1.333037	0.78515
46	-0.95034	0.96529
47	-1.02385	0.84158
48	1.277203	0.705515
49	1.333037	0.78515
50	-1.02385	0.841588

I. Data Management

II. Scientific Processes

III. Cross-Domain Collaboration

Estimated Factor
(Respondent)

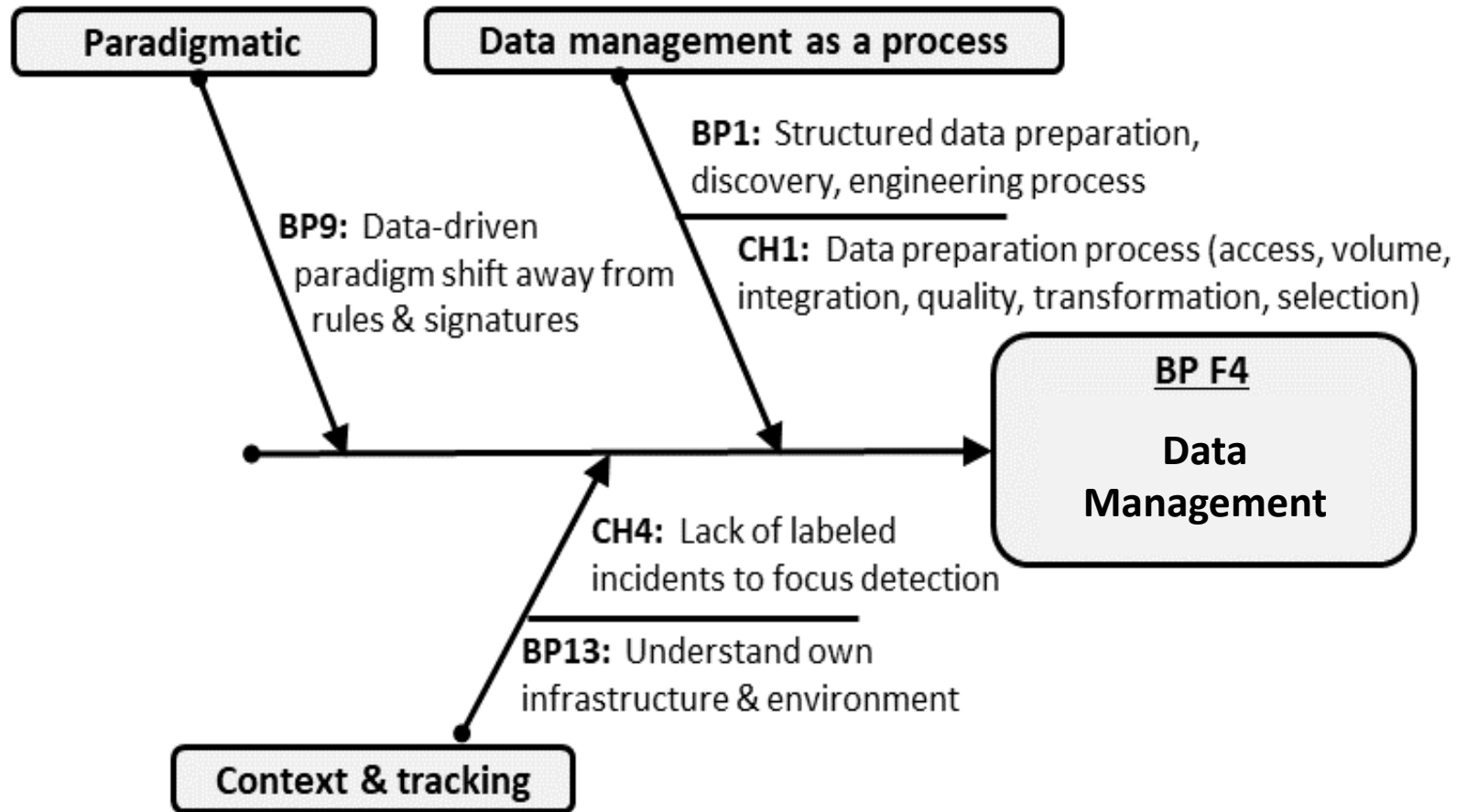
FACTOR1	FACTOR2	FACTOR3
706186	-0.61424	-0.65698
658827	1.416078	0.069793
582112	-0.54697	0.54427
142545	1.240915	1.300297
0.78656	1.188467	1.298155
616202	1.437153	-0.4533
598131	1.492765	-0.64768
-1.55677	-0.47443	-0.98028
-1.4047	-0.66512	0.249976
-1.34841	-0.72229	0.914774
-1.47342	-0.0147	-0.6551
-1.4738	-0.52214	-0.26985
0.505144	0.860759	-0.83992
0.499477	1.085276	-0.92164
816931	-0.60279	0.385758
582528	1.037561	-0.21421
624868	1.37377	-0.07003
0.85632	-1.06151	1.808688
490955	0.94173	1.072211
1.66482	1.632577	-0.97099
427606	0.860188	0.634095
694619	1.223261	0.527547
1.29747	-0.80965	1.240221
1.43948	1.145655	1.109696
1.47246	-0.72714	-0.24082
706186	-0.61424	-0.65698
1.53194	-0.95627	-0.35647
0.6132	-1.24625	-0.84922
0.798861	-0.54718	0.191376
0.89624	1.42463	-0.72156
0.558327	0.780959	0.319014
0.856401	0.81968	0.818443
0.545722	-0.6712	-0.60713
-0.31032	-0.7599	-1.45588
1.75781	1.000571	-1.16323
758866	1.129298	0.965849
629087	-1.03894	-0.90378
1.54722	-0.67359	1.207946
729118	1.031358	-0.54349
867878	-0.78845	0.711205
626722	1.376429	-0.4549
627234	-1.15099	-0.65862
519178	0.956019	-0.64932
738291	-0.57459	-0.272
627234	-1.15099	-0.65862
816931	-0.60279	0.385758
738291	-0.57459	-0.272
0.406074	-1.13126	-0.3584
0.558327	0.780959	0.319014
627234	-1.15099	-0.65862
738291	-0.57459	-0.272

BP F1 Scientific process
BP F2 Cross-domain collaboration
BP F3 Risk management focus
BP F4 Data-driven / data management
BP F5 Focused tools
BP F6 Structured discovery process

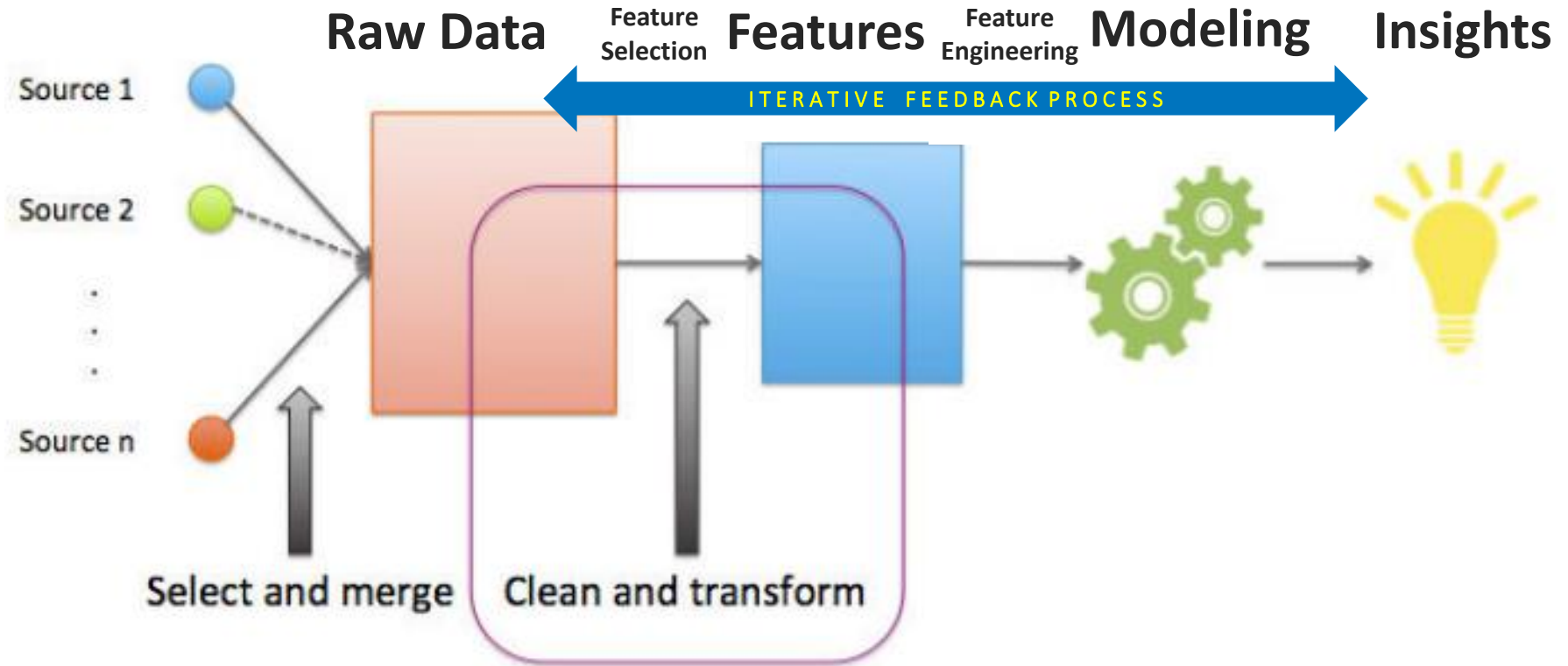


III. CSDS Designs



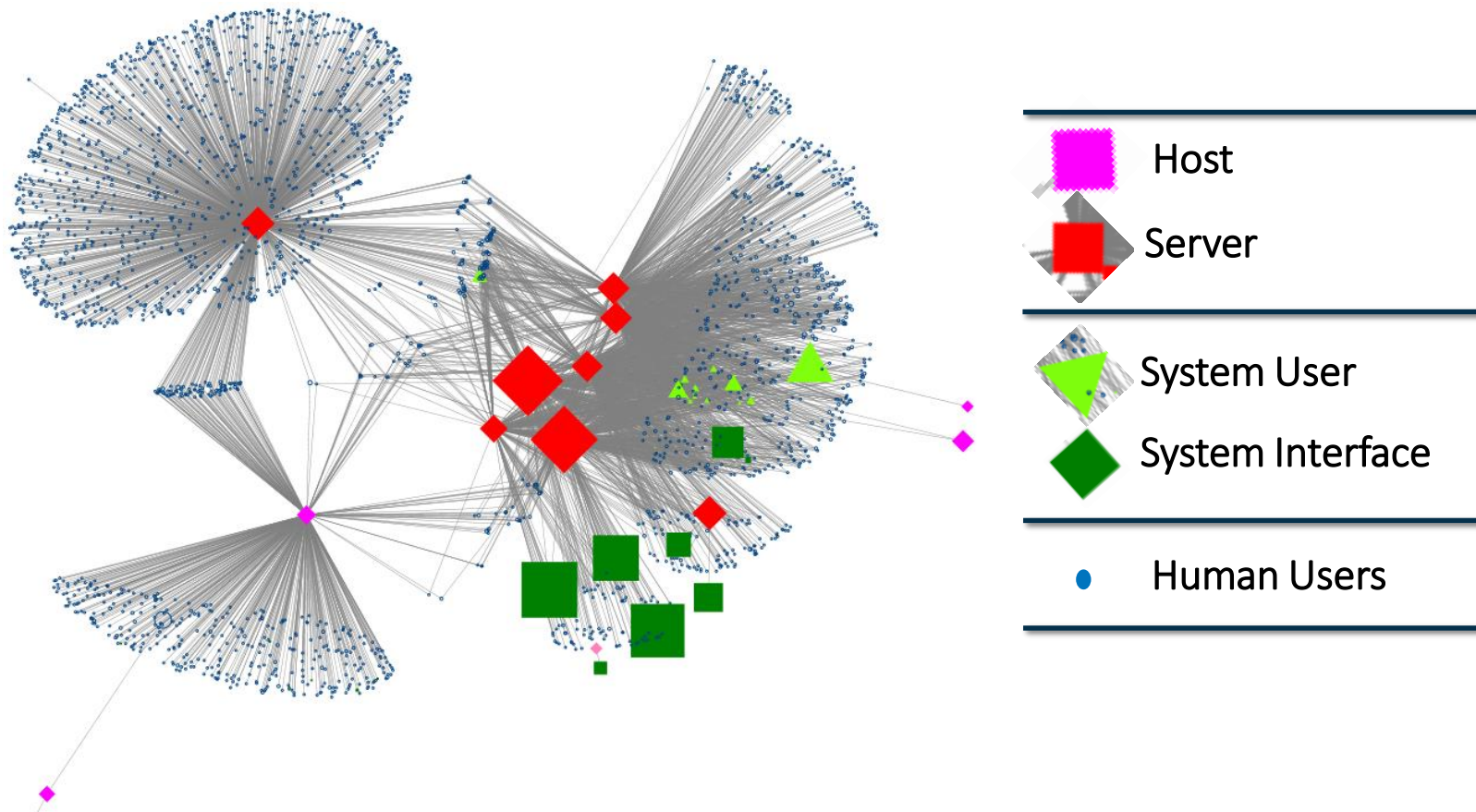


Data Management: EDA Process + Feature Engineering

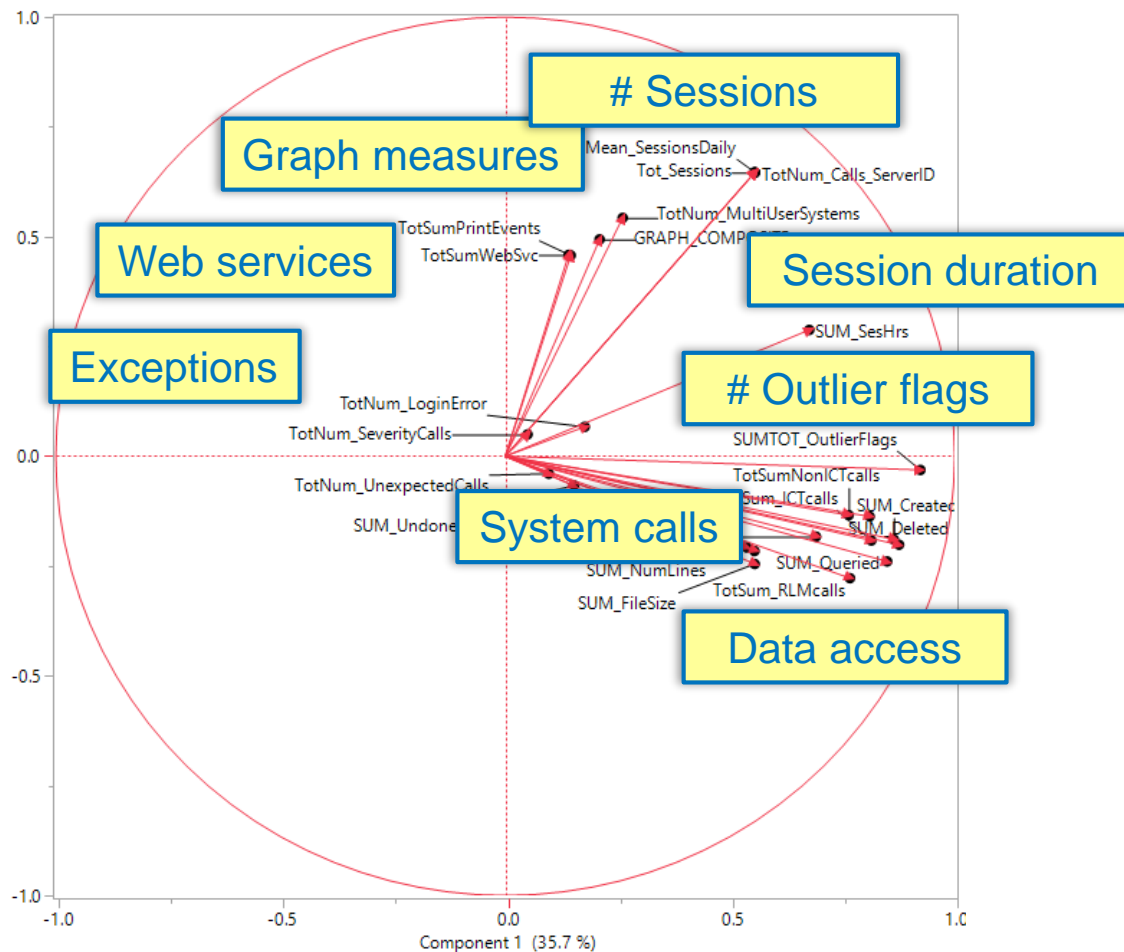
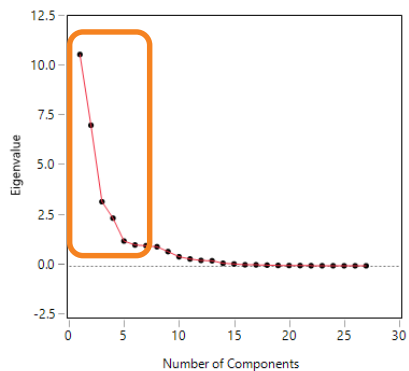


SOURCE: Alice Zheng, Amanda Casari. 2016. [Feature Engineering for Machine Learning Models](#). O'Reilly Media.

Featurization: Example - Graph Analytics



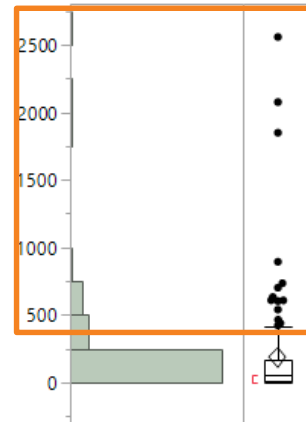
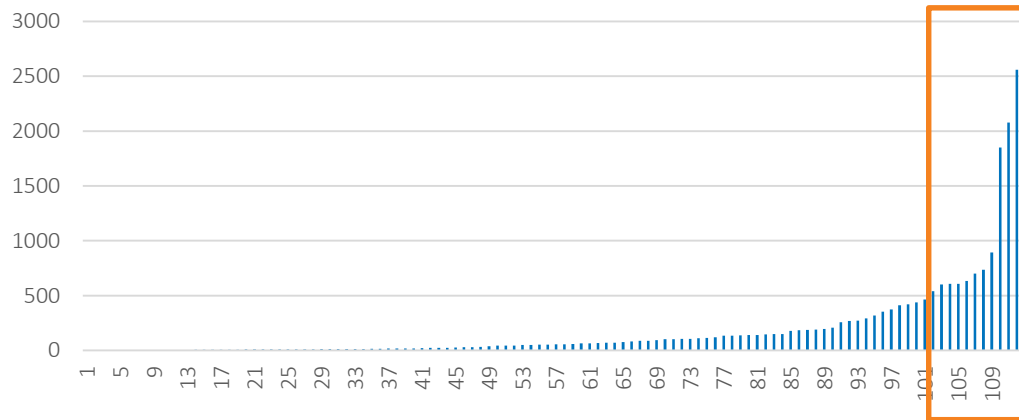
Feature Reduction: Example - Principal Component Analysis (PCA)



Exploratory Data Analysis (EDA): Example – Probabilistic Analysis

Exception Events

Exception messages per user (ranked)



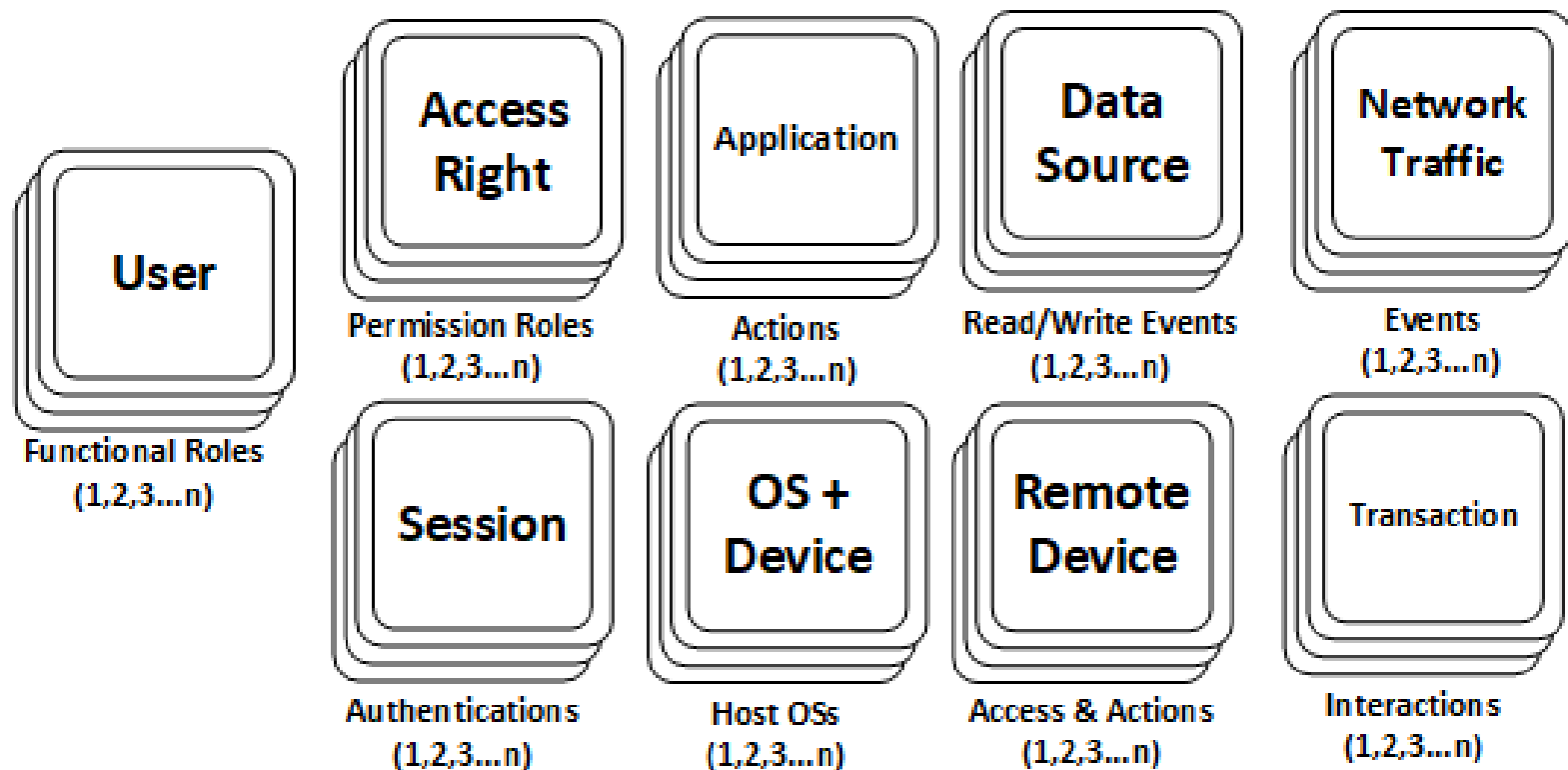
Quantiles

100.0%	maximum	2559
99.5%		2559
97.5%		1889.725
90.0%		517.5
75.0%	quartile	172.75
50.0%	median	55.5
25.0%	quartile	9.75
10.0%		3.3
2.5%		1.825
0.5%		1
0.0%	minimum	1

Summary Statistics

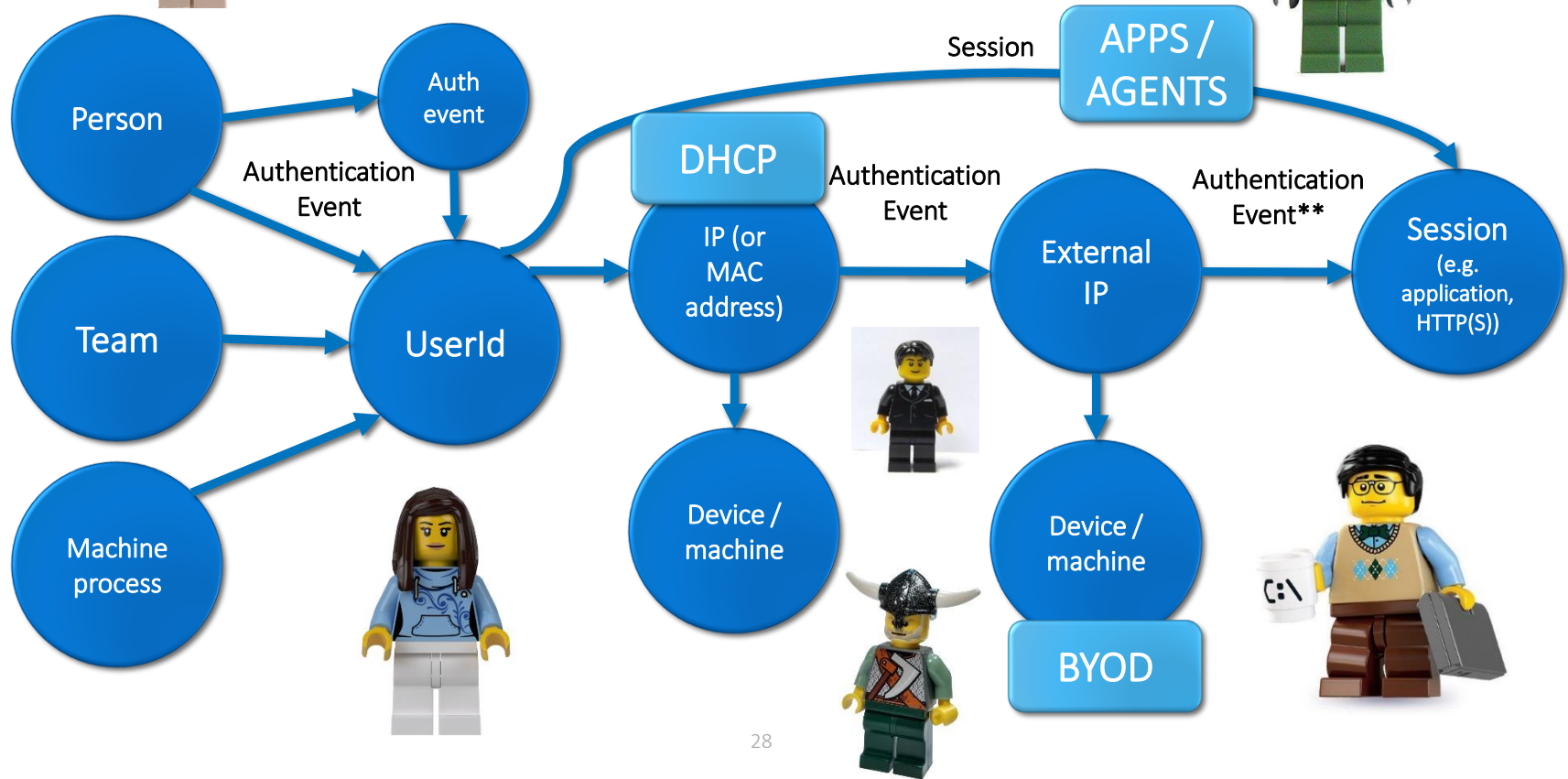
Mean	184.01786
Std Dev	380.96684
Std Err Mean	35.997982
Upper 95% Mean	255.35026
Lower 95% Mean	112.68545
N	112

Entity Resolution

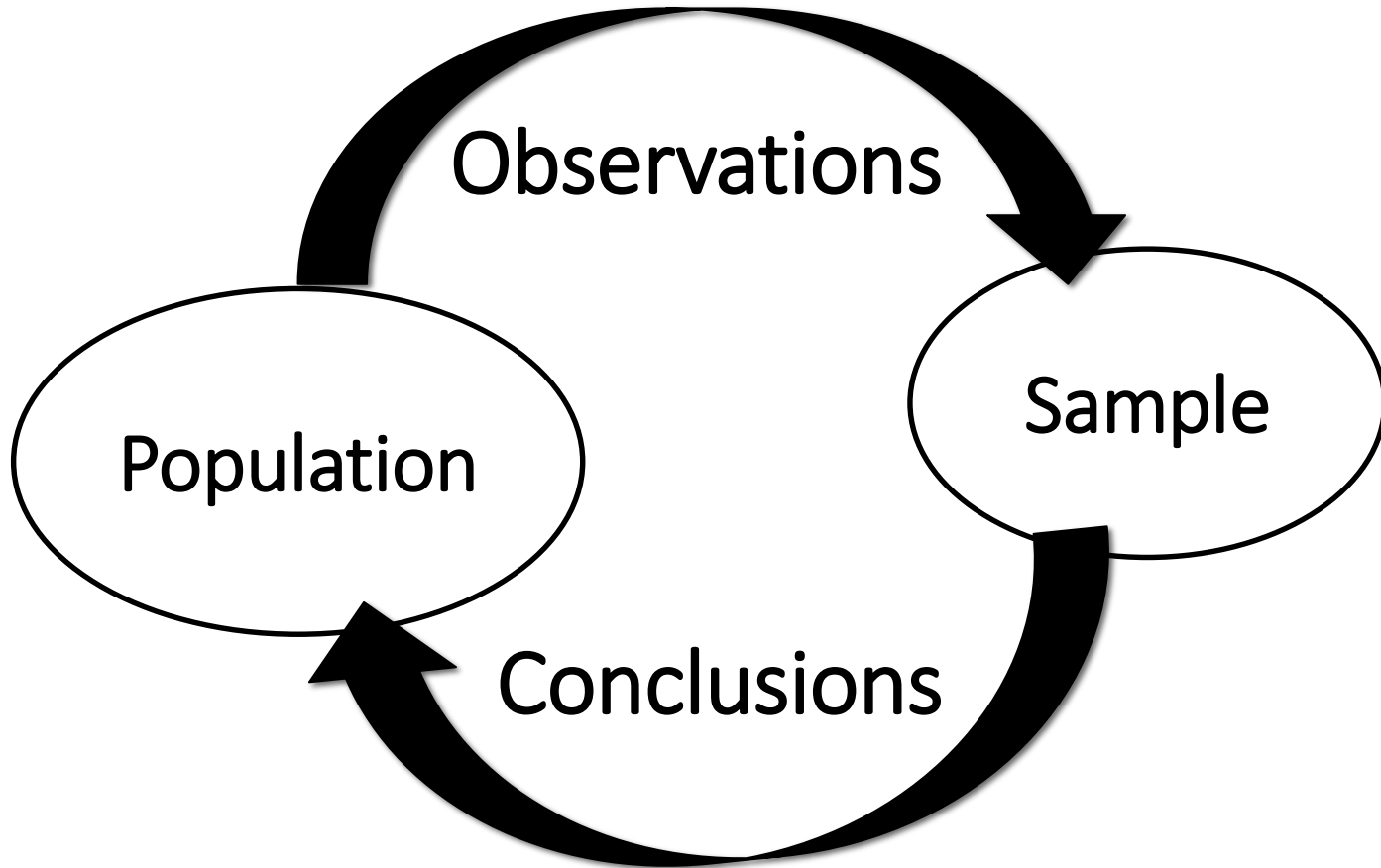


What is a User, anyway?

What is an IP address, anyway?

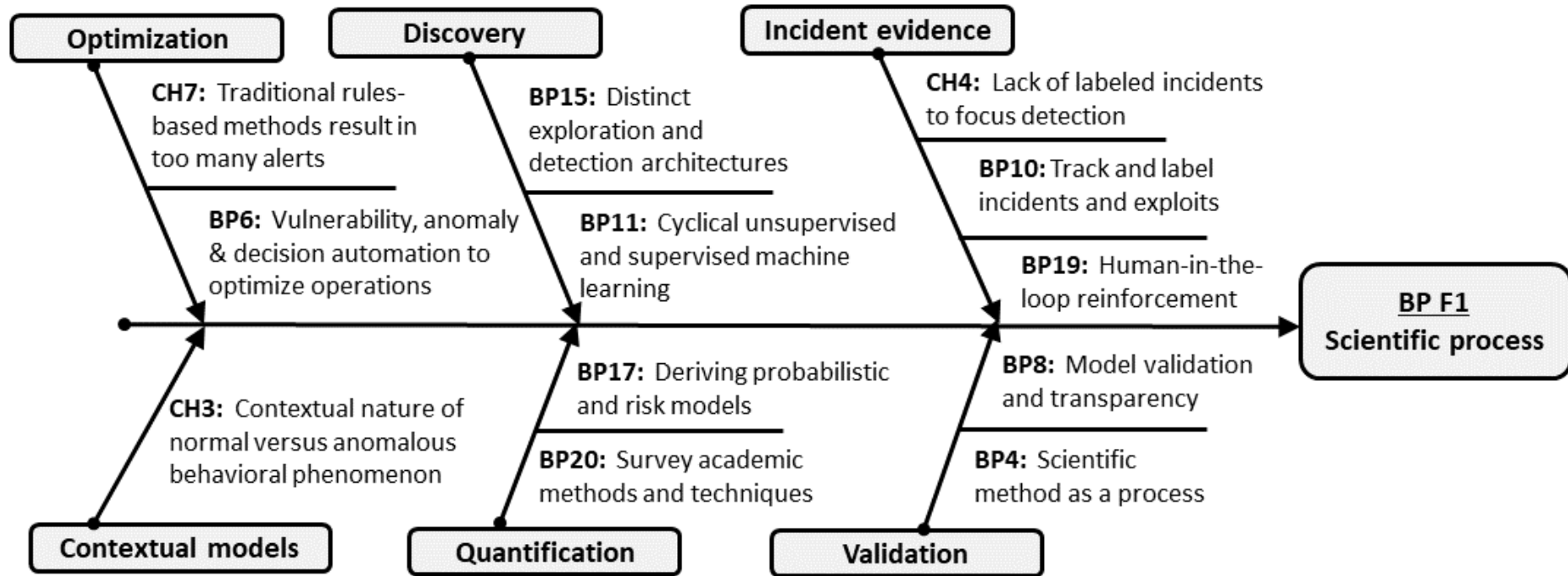


Inferential Statistics





Root Cause Analysis: Fishbone / Ishikawa Diagram



** Resulting from factor analysis and factor-to-factor fitting*

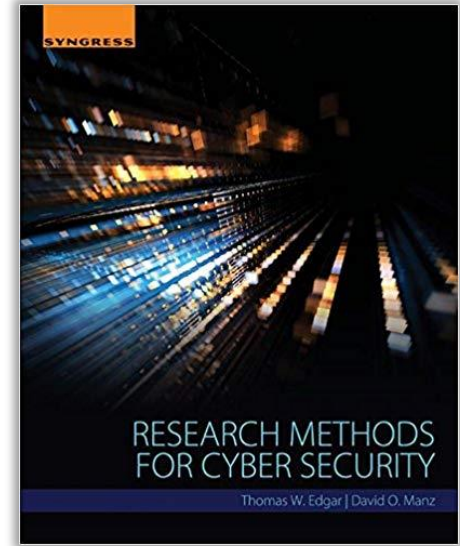
CSDS: What type of science is it?

Controlled experiments
versus
Pattern extrapolation



Research Methods for Cybersecurity

- *Experimental*
 - i.e. hypothetical-deductive and quasi-experimental
- *Applied*
 - i.e. applied experiments and observational studies
- *Mathematical*
 - i.e. theoretical and simulation-based
- *Observational*
 - i.e. exploratory, descriptive, machine learning-based

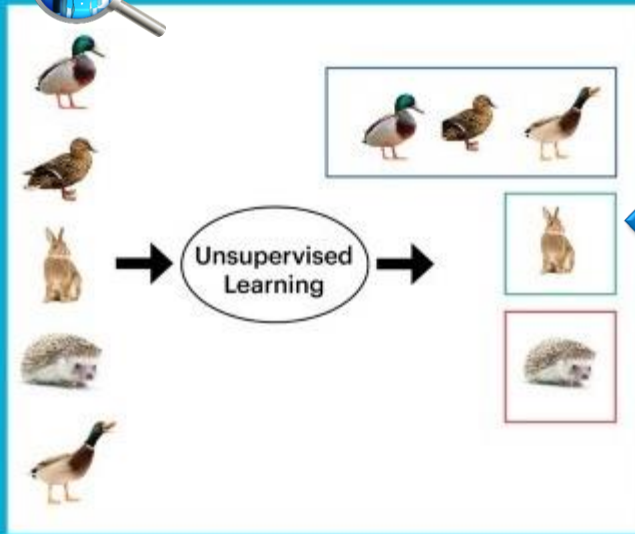


Manz, D. and Edgar, T. (2017)
Research Methods for Cyber Security

Discovery ⇔ Detection

Exploration and
Insights

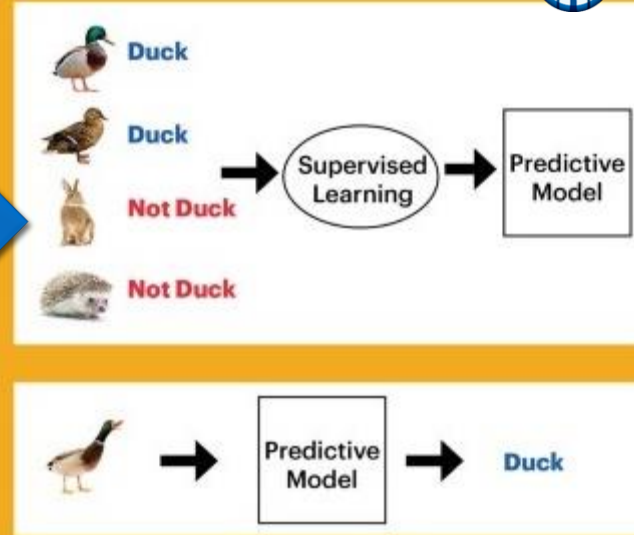
Unsupervised Learning
(Clustering Algorithm)



SEGMENTATION

Supervised Learning
(Classification Algorithm)

Pattern
Detection



CATEGORIZATION

Labels: What constitutes 'evidence'?

EXAMPLES OF SECURITY EVIDENCE

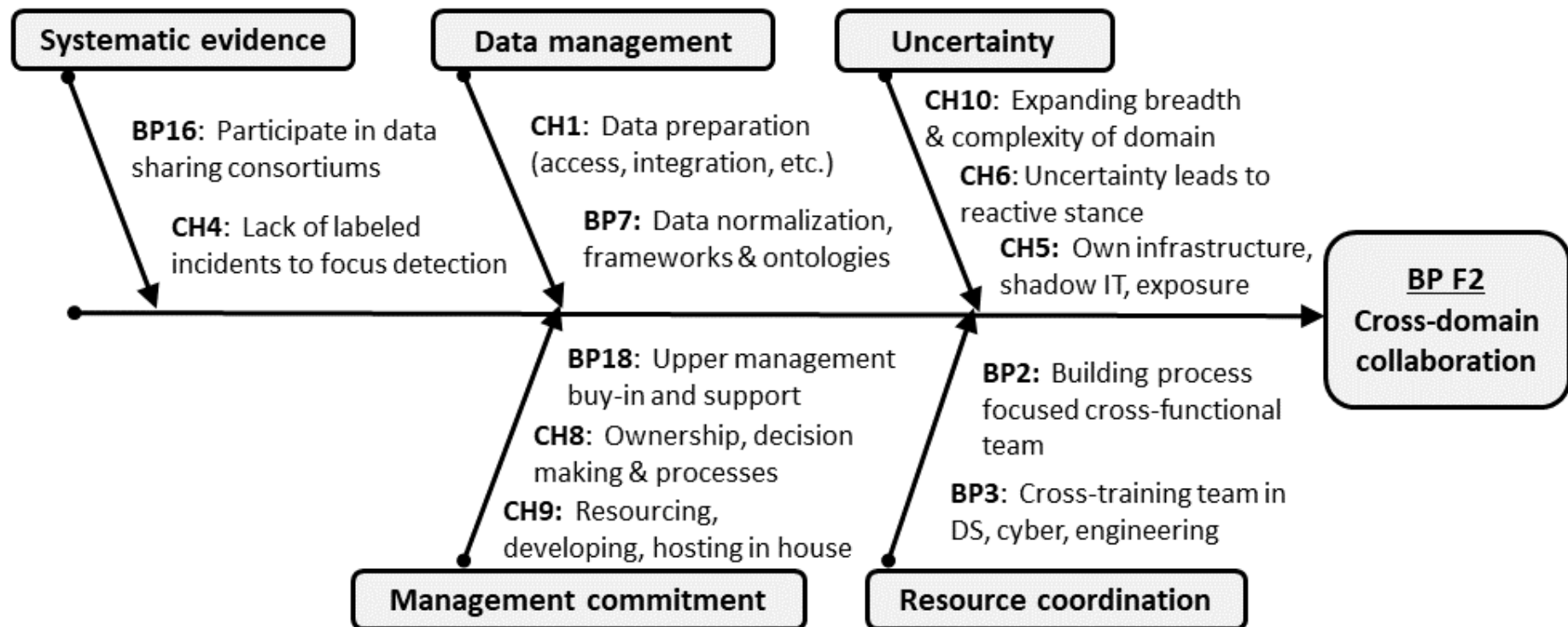
Synthesized Collected	- Field evidence - Probing & testing - 3 rd party sourced	- Rules & signatures - Research & threat intelligence
	- Red Teaming - Simulations - Laboratory	- Expert opinion - Thought experiments
	Inductive	Deductive

1. Field evidence (e.g. observed incidents)
2. Sourcing own data from field testing (e.g. local experiments)
3. Honeypots
4. IDSs (Intrusion Detection Systems)
5. Simulation findings
6. Laboratory testing (e.g. malware in a staged environment)
7. Stepwise discovery (iterative interventions)
8. Pen testing (attempts to penetrate the network)
9. Red teaming (staged attacks to achieve particular goals)
10. Incidents (records associated with confirmed incidents)
11. Reinforcement learning (self-improving ML to achieve a goal)
12. Research examples (datasets recording attacks from research)
13. Expert review (opinion and guidance from experts)
14. Intelligence feed (indications from a 3rd party service)
15. Thought experiments (e.g. boundary conditions, counterfactuals)

CSDS as a Process: Discovery and Detection





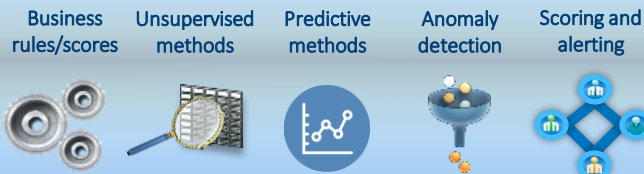


CSDS: High-Level Functional Process

Data management



Advanced Analytics



Triage



Investigation



ALERT ANALYTICS PROCESS



Data Manager



Data Scientist



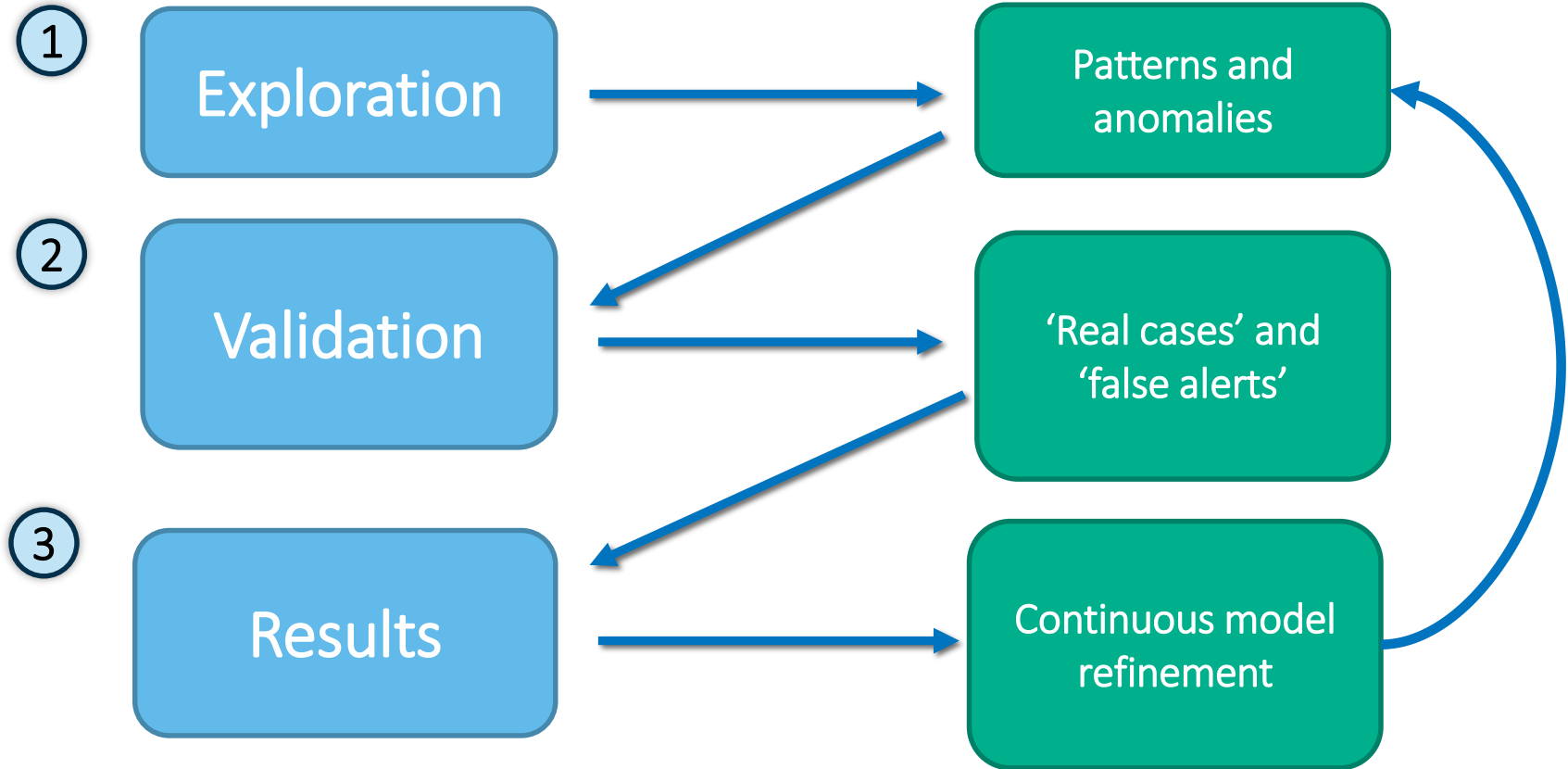
Investigator



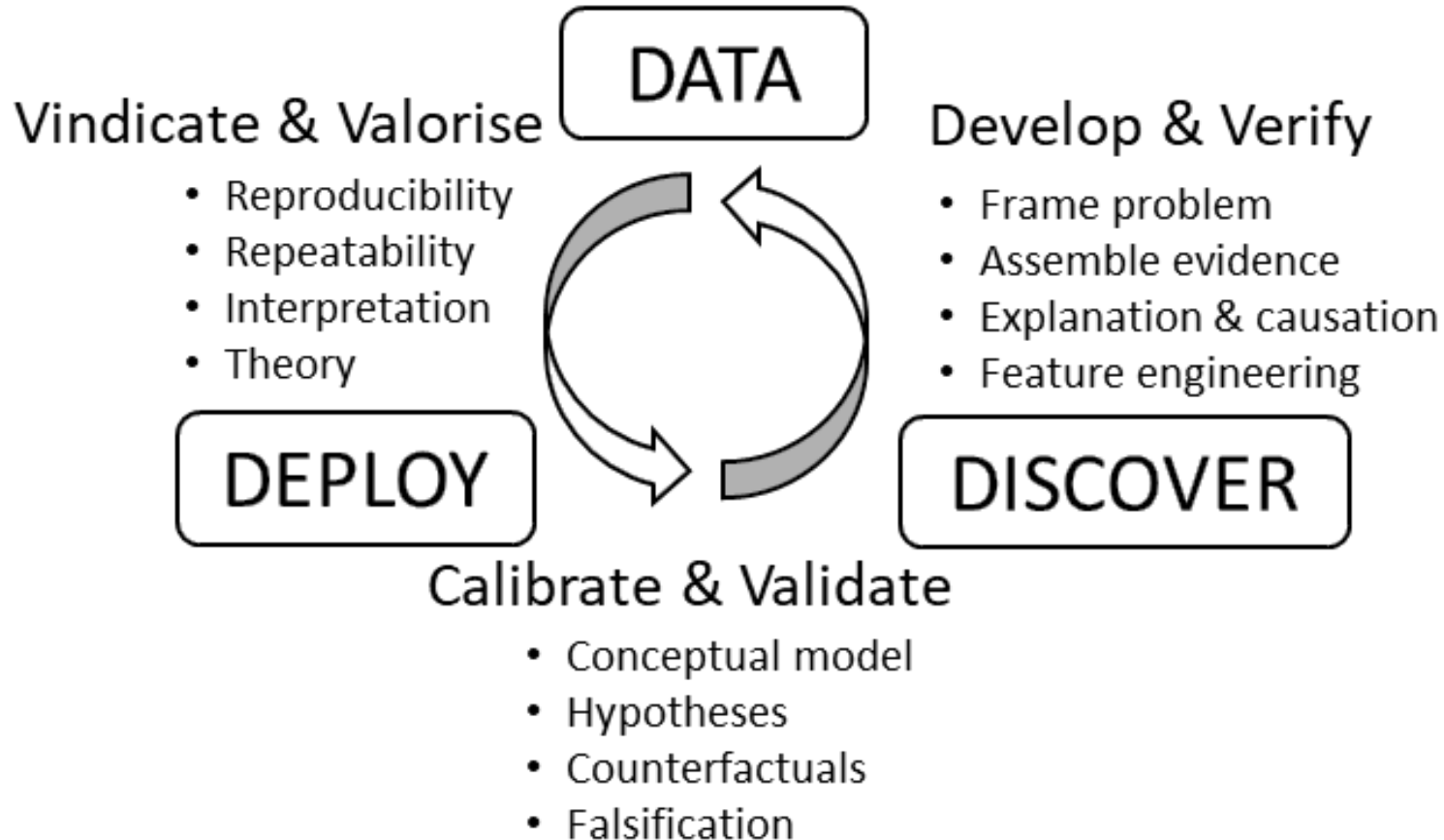
Case
Remediation

RECURSIVE FEEDBACK

Continuous Detection Improvement Process



CSDS Model Development Process





Conclusions

Cybersecurity ✓

Data ✓

Science ?

Not so much...
but, ASPIRATIONAL!





CSDS: A Work in Progress

- **Process of Professionalization**

- Named professionals
- Set of methods and techniques
- Standards, best practices

Training programs

Certifications

Academic degree programs

Focused research journals

Formal sub-specialization



Specialist Surgeon Researcher Diagnostician Primary Care Emergency Care



APPENDIX

References

- Aggarwal, C. (2013). "Outlier Analysis." Springer. <http://www.springer.com/la/book/9781461463955>
- 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://www.sas.com/en_us/whitepapers/hbr-defending-your-networks-108030.html
- Longitude Research. (2014). "Cyberrisk in banking." Available at https://www.sas.com/content/dam/SAS/bp_de/doc/studie/ff-st-longitude-research-cyberrisk-in-banking-2316865.pdf
- 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
- SAS Institute. (2019). "Data Management for Artificial Intelligence." Available at www.sas.com/en_us/whitepapers/data-management-artificial-intelligence-109860.html
- 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
- 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