

AI & ML IN CYBERSECURITY

Why Algorithms Are Dangerous

RAFFAEL MARTY

Vice President of Corporate Strategy, Forcepoint

BlackHat USA | August 2018



Copyright © 2018 Forcepoint.

A BRIEF SUMMARY

- ▶ We don't have **artificial intelligence** (yet)
- ▶ Algorithms are getting ‘smarter’, but **experts** are more important
- ▶ Stop throwing algorithms on the wall - they are not spaghetti
- ▶ **Understand** your data and your algorithms
- ▶ Invest in people who **know** security (and have experience)
- ▶ Build systems that capture “**export knowledge**”
- ▶ Think out of the box, history is bad for innovation
- ▶ Focus on advancing **insights**

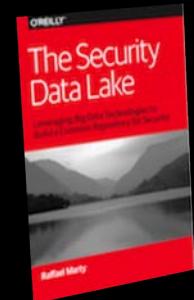
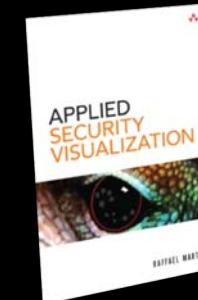
RAFFAEL MARTY

FORCEPOINT

- ▶ Sophos
- ▶ PixlCloud
- ▶ Loggly
- ▶ Splunk
- ▶ ArcSight
- ▶ IBM Research



- ▶ Security Visualization
- ▶ Big Data
- ▶ ML & AI
- ▶ SIEM
- ▶ Corp Strategy
- ▶ Leadership
- ▶ Zen



OUTLINE

01

STATISTICS, MACHINE LEARNING & AI

Defining the concepts

02

THE ALGORITHMIC PROBLEM

Understanding the data and the algorithms

03

AN EXAMPLE

Let's get practical

01

STATISTICS MACHINE LEARNING & 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)”

ML AND AI – WHAT IS IT?

MACHINE LEARNING

Algorithmic 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, dim reduction)

DEEP LEARNING

A “newer” machine learning algorithm

- ▶ Eliminates the feature engineering step
- ▶ Explainability / verifiability issues

DATA MINING

Methods to explore data – automatically

ARTIFICIAL INTELLIGENCE

“Just calling something AI doesn’t make it AI.”

“A program that doesn't simply classify or compute model parameters, but comes up with novel knowledge that a security analyst finds insightful.”

We don't have artificial intelligence (yet)

WHAT “AI” DOES TODAY

KICK A HUMAN'S
ASS AT GO



DESIGN MORE
EFFECTIVE DRUGS



MAKE SIRI
SMARTER



MACHINE LEARNING USES IN SECURITY

SUPERVISED

▶ Malware classification

- ▶ Deep learning on millions of samples - 400k new malware samples a day
- ▶ Has increased true positives and decreased false positives compared to traditional ML

▶ Spam identification

▶ MLSec project on firewall data

- ▶ Analyzing massive amounts of firewall data to predict and score malicious sources (IPs)

UNSUPERVISED

▶ DNS analytics

- ▶ Domain name classification, lookup frequencies, etc.

▶ Threat Intelligence feed curation

- ▶ IOC prioritization, deduplication, ...

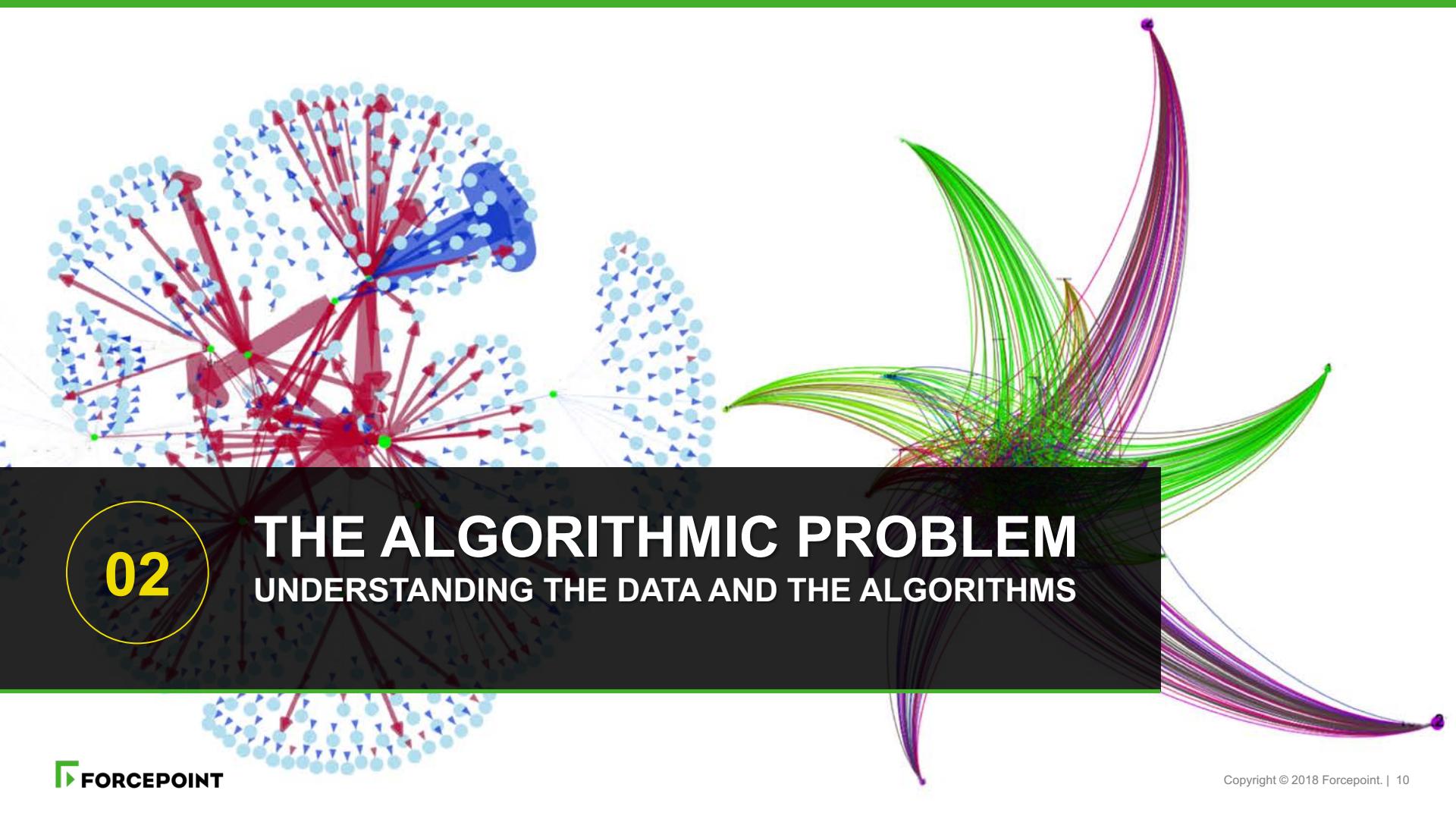
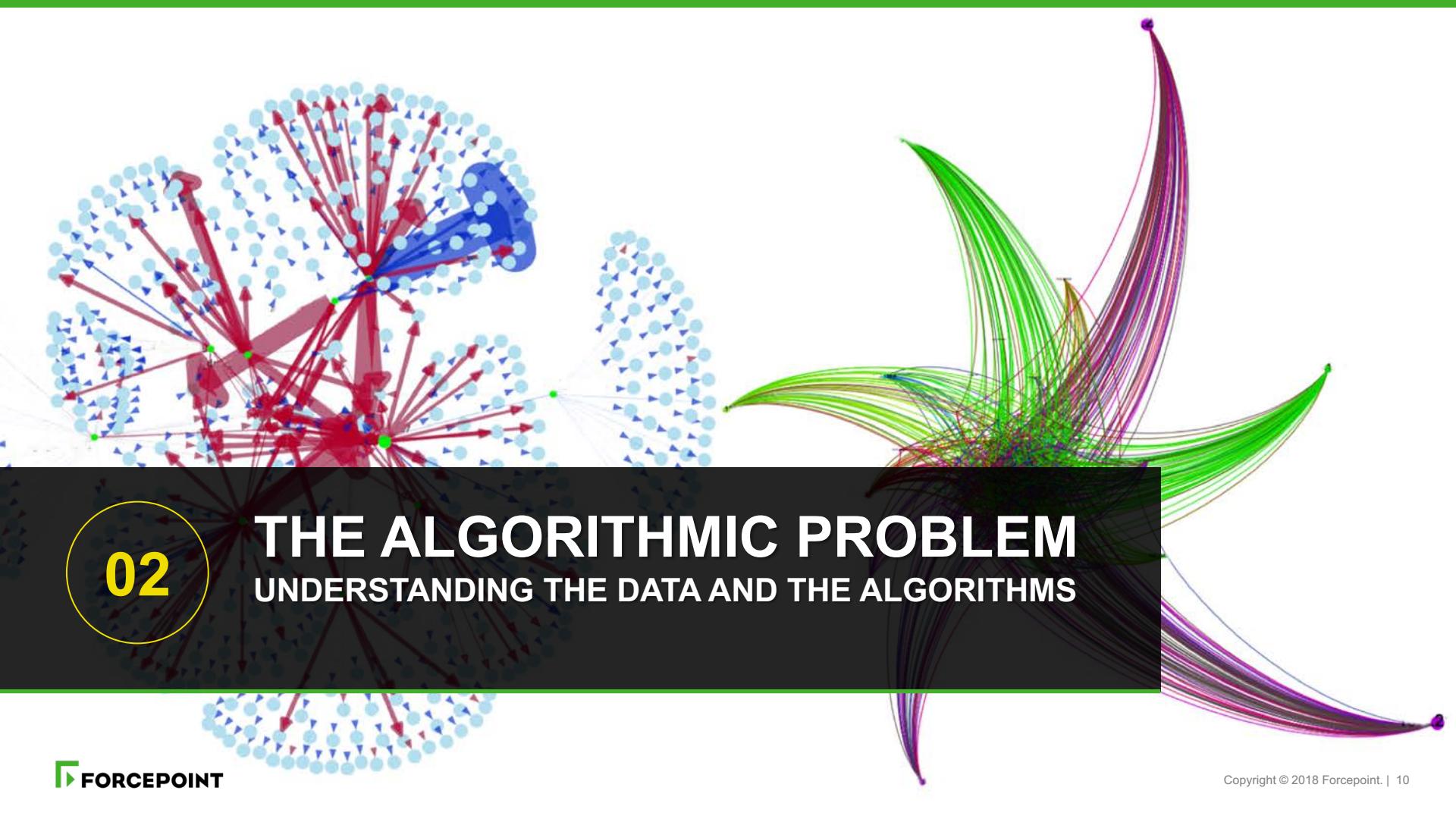
▶ Tier 1 analyst automation

- ▶ Reducing workload from 600M raw events to 100 incidents*

▶ User and Entity Behavior Analytics (UEBA)

- ▶ Uses mostly regular statistics and rule-based approaches

* See Respond Software Inc.



02

THE ALGORITHMIC PROBLEM

UNDERSTANDING THE DATA AND THE ALGORITHMS

WARNING

ALGORITHMS ARE DANGEROUS

FAMOUS AI (ALGORITHM) FAILURES



December 2009 | Hewlett-Packard investigates instances of so-called "racist camera software" which had trouble recognizing dark-skinned people

March 2015 | A Carnegie Mellon University study determines that some personalized ads from sites such as Google and Facebook are gender-biased

July 2015 | A Google algorithm mistakenly captions photos of black people as "Gorillas"

March 2016 | Microsoft shuts down AI chatbot Tay after it starts using racist language

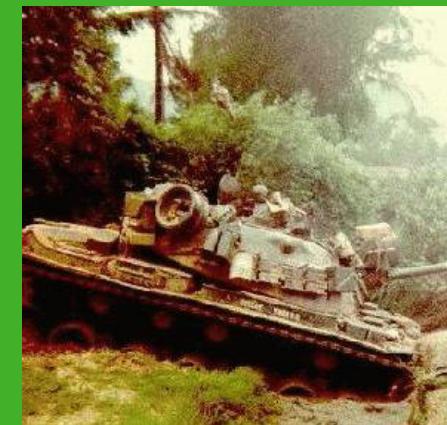
May 2016 | ProPublica investigation finds that a computer program used to track future criminals in the US is racially biased

September 2016 | Machine-learning algorithms used to judge an international beauty contest displays bias against dark-skinned contestants



PENTAGON - AI FAIL

<http://neil.fraser.name/writing/tank/>



WHAT MAKES ALGORITHMS DANGEROUS?

ALGORITHMS MAKE ASSUMPTIONS ABOUT THE DATA

- ▶ Assume ‘clean’ data (src/dst confusion, user feedback, etc.)
- ▶ Often assume a certain type of data and its distribution
- ▶ Generally don’t deal with outliers
- ▶ Machine learning assumes enough, representative data
- ▶ Need contextual features (e.g., not just IP addresses)
- ▶ Assume all input features are ‘normalized’ the same way

ALGORITHMS ARE TOO EASY TO USE THESE DAYS (TENSORFLOW, TORCH, ML ON AWS, ETC.)

- ▶ The process is more important than the algorithm (e.g., feature engineering, supervision, drop outs, parameter choices, etc.)

ALGORITHMS DO NOT TAKE DOMAIN KNOWLEDGE INTO ACCOUNT

- ▶ Defining meaningful and representative distance functions, for example
- ▶ e.g., each L4 protocol exhibits different behavior. Train it separately.
- ▶ e.g., interpretation is often unvalidated - beware of overfitting and biased models.
- ▶ Ports look like integers, they are not, same is true for IPs, processIDs, HTTP return codes, etc.

WHAT MAKES ALGORITHMS DANGEROUS?

SAMPLE BIAS

What Should We Remember?

We favor simple-looking options and complete information over complex, ambiguous options

To avoid mistakes, we aim to preserve autonomy and group status, and avoid irreversible decisions

To get things done, we tend to complete things we've invested time & energy in

To stay focused, we favor the immediate, relatable thing in front of us

Need To Act Fast

To act, we must be confident we can make an impact and feel what we do is important

We store memories differently based on how they were experienced

We reduce events and lists to their key elements

We discard specifics to form generalities

We edit and reinforce some memories after the fact

Too Much Information

We notice when something has changed

We are drawn to details that confirm our own existing beliefs

We notice flaws in others more easily than we notice flaws in ourselves

KNOWLEDGE BIAS

Less-is-better effect • Conjunction fallacy • Law of triviality • Rhyme as reason effect • Belief effect • Ambiguity bias

Status quo bias • Social comparison bias • Decoy effect • Reactance • Reverse psychology • System justification

Backfire effect • Endowment effect • Processing difficulty effect • Pseudocertainty effect • Dispositionality effect • Zero-risk bias • Unit bias • IKEA effect • Loss aversion • Generation effect • Intentional communication • Escalation of commitment • Irrational escalation • Sunk cost fallacy

Identifiable victim effect • Appeal to novelty • Hyperbolic discounting

Risk-perturbation effect • Effort justification bias • Fundamental attribution error • Defensive attribution bias • Just-world belief • Self-handicapping effect • Optimism bias • Confirmation bias • Self-reinforcing effect • Overconfidence bias • Self-licensing effect • Self-justification effect • Self-justification bias

Optimism bias • Confirmation bias • Self-licensing effect • Self-justification effect • Self-justification bias

We project our current mindset and assumptions onto the past and future

CONFIRMATION BIAS

We tend to find stories and patterns even when looking at sparse data

AVAILABILITY BIAS

We fill in characteristics from stereotypes, generalities, and prior histories

We imagine things and people we're familiar with or fond of as better

We simplify probabilities and numbers to make them easier to think about

We think we know what other people are thinking

HISTORY BIAS

Our prior homogeneity bias • Our prior heterogeneity bias • Confirmation bias • Confirmation effect • Confirmation bias • Confirmation effect • Confirmation bias • Confirmation bias

Outcome bias • Bias towards recent news • Confirmation bias • Confirmation bias • Confirmation bias • Confirmation bias • Confirmation bias

Confirmation bias • Confirmation bias

Confirmation bias • Confirmation bias

Confirmation bias • Confirmation bias

CORRELATION BIAS

COGNITIVE BIASES

How biased is your data set? How do you know?

- ▶ Only a single customer's data
- ▶ Learning from an 'infected' data set
- ▶ Collection errors
- ▶ Missing data (e.g., due to misconfiguration)
- ▶ What's the context the data operates in?
- ▶ FTP although generally considered old and insecure, isn't always problematic
- ▶ Don't trust your IDS (e.g. "UDP bomb")

English ▾ he is a nurse. she is a doctor. Edit Hungarian ▾ ō ápolónő. ō egy orvos.

Hungarian ▾ ō ápolónő. ō egy orvos. English ▾ she's a nurse. he is a doctor.

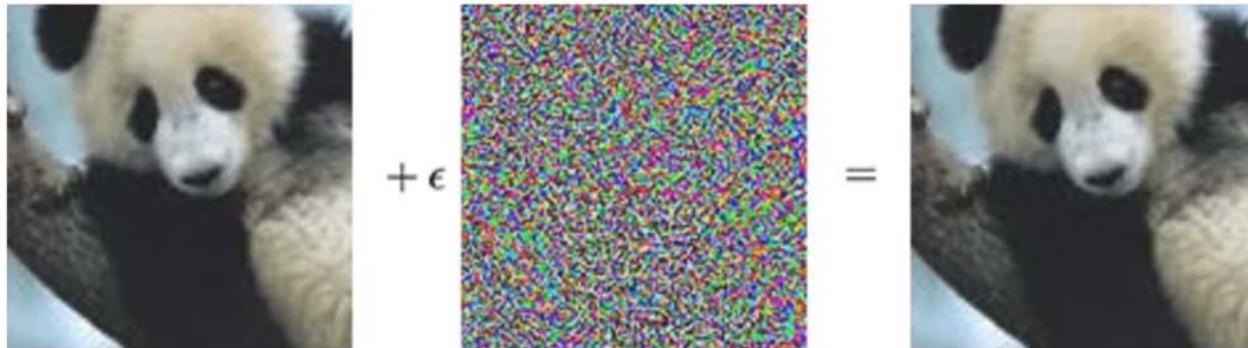
THE DANGERS WITH DEEP LEARNING – WHEN NOT TO USE IT



- ▶ Not enough or no quality **labelled data**
- ▶ Data **cleanliness** issues (timestamps, normalization across fields, etc.)
- ▶ Bad understanding of the data to engineer meaningful **features** (e.g., byte stream for binaries)
- ▶ Data is prone to **adversarial** input
- ▶ No well trained **domain experts** and data scientists to oversee the implementation
- ▶ A need to understand what ML actually learned (**explainability**)
- ▶ **Verifiability** of output
- ▶ **Interpretation** of output

ADVERSARIAL MACHINE LEARNING

An example of an attack on deep learning



"panda"

57.7% confidence

Above: Image Credit: Ian Goodfellow

"gibbon"

99.3% confidence

03

EXAMPLE LET'S GET PRACTICAL

FINDING ANOMALIES / ATTACKS IN NETWORK TRAFFIC

- ▶ Given: Network communications (i.e., netflow)

► Task: Find anomalies / attacks

```
07:40:37.300221 IP 172.20.5.103.5353 > 224.0.0.291.5353: 0 PTR (QM)? _googlecast._tcp.local. (40)
07:40:37.382675 IP 172.20.7.238.5353 > 224.0.0.251.5353: 0 PTR (QM)? _googlecast._tcp.local. (40)
07:40:37.387937 IP 172.20.5.103.5353 > 224.0.0.251.5353: 0 PTR (QM)? _googlecast._tcp.local. (40)
07:40:37.437678 IP 172.20.7.238.5353 > 224.0.0.251.5353: 0 PTR (QM)? _googlecast._tcp.local. (40)
07:40:37.741591 IP6 fe80::878:ec91:c717:8aab > ff02::2: ICMP6, router solicitation, length 16
07:40:37.741825 IP 172.20.8.228.59338 > 224.0.0.1.8612: UDP, length 16
07:40:37.906586 IP 172.20.12.5.49248 > 224.0.0.253.3544: UDP, length 40
07:40:38.100075 IP 172.20.6.254.5353 > 224.0.0.251.5353: 0 PTR (QM)? _googlecast._tcp.local. (40)
07:40:38.169207 IP 172.20.6.254.5353 > 224.0.0.251.5353: 0 PTR (QM)? _googlecast._tcp.local. (40)
07:40:38.663287 IP 172.20.3.85.49993 > 224.0.0.253.3544: UDP, length 40
07:40:38.829006 IP 172.20.8.221.5353 > 224.0.0.251.5353: 0*- [0q] 1/0/0 TXT "model=N51AP" (85)
07:40:38.859024 IP 172.20.5.192.5353 > 224.0.0.251.5353: 0*- [0q] 4/0/4 (Cache flush) PTR micky-laptop.local., (Cache flush) PTR micky-laptop.local., (Cache flush) A 172.20.5.192, (Cache flush) AAAA fe80::38f7:1ff:fe00:19:f1c8 (374)
07:40:38.940469 IP 172.20.15.116.5353 > 224.0.0.251.5353: 0 [6q] PTR (QU)? _ipp._tcp.local. PTR (QU)? _scanner._tcp.local. PTR (QU)? _pdl-datastream._tcp.local. PTR (QU)? _printer._tcp.local. PTR (QU)? _uscan._tcp.local. PTR (QU)? _ptp._tcp.local. (109)
07:40:39.086181 IP 172.20.9.92.5353 > 224.0.0.251.5353: 0*- [0q] 1/0/1 TXT "model=N42AP" (103)
07:40:39.892283 IP 172.20.6.254 > 224.0.0.22: igmp v3 report, 1 group record(s)
07:40:39.947767 IP 172.20.15.116.5353 > 224.0.0.251.5353: 0 [6q] PTR (QU)? _ipp._tcp.local. PTR (QU)? _scanner._tcp.local. PTR (QU)? _pdl-datastream._tcp.local. PTR (QU)? _printer._tcp.local. PTR (QU)? _uscan._tcp.local. PTR (QU)? _ptp._tcp.local. (109)
07:40:40.199461 IP 172.20.15.236.5353 > 224.0.0.251.5353: 0 PTR (QM)? _googlecast._tcp.local. (40)
07:40:40.274673 IP 172.20.3.157.61059 > 239.255.255.250:1900: UDP, length 97
07:40:40.817470 IP 172.20.9.196.64236 > 224.0.0.1.8612: UDP, length 1
07:40:40.849139 IP 172.20.2.8.5353 > 224.0.0.251.5353: 0*- [0q] 12/0/5 (Cache flush) TXT "", PTR _apple-mobdev2._tcp.local., PTR f8:27:93:4b:78:16@fe80::fa27:93ff:fe4b:7816._apple-mobdev2._tcp.local., PTR f8:27:93:4b:78:16@fe80::fa27:93ff:fe4b:7816._apple-mobdev2._tcp.local., PTR f8:27:93:4b:78:16@fe80::fa27:93ff:fe4b:7816._apple-mobdev2._tcp.local., (Cache flush) SRV iPhone-6.local.:62078 0 0, TXT "model=N51AP", (Cache flush) PTR iPhone-6.local., (Cache flush) PTR iPhone-6.local., (Cache flush) AAAA fe80::1035:f6f3:82d6:23ad, (Cache flush) A 172.20.2.8 (563)
07:40:41.223219 IP 172.20.9.195.53599 > 224.0.0.252.5355: UDP, length 29
```

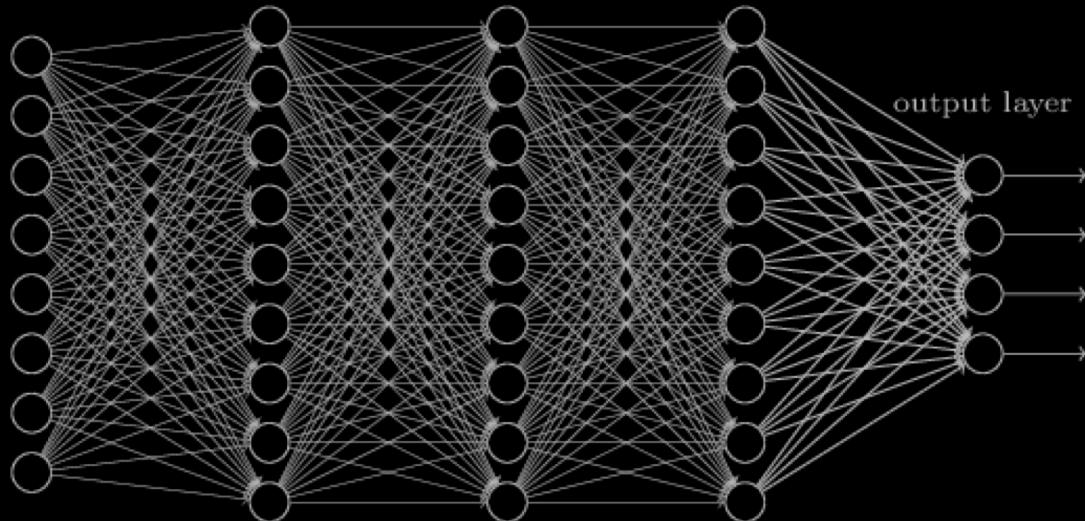
DEEP LEARNING – THE SOLUTION TO EVERYTHING

DEEP LEARNING PROMISES A FEW THINGS:

- ▶ ‘Auto’ feature extraction
- ▶ High accuracy of detections

AND WE SATISFY SOME REQUIREMENTS:

- ▶ Lots of data available
- ▶ **BUT:** A single record does not indicate good/bad
- ▶ **BUT:** Not enough ‘information’ within flows – need **context**
- ▶ **BUT:** No **labels** available



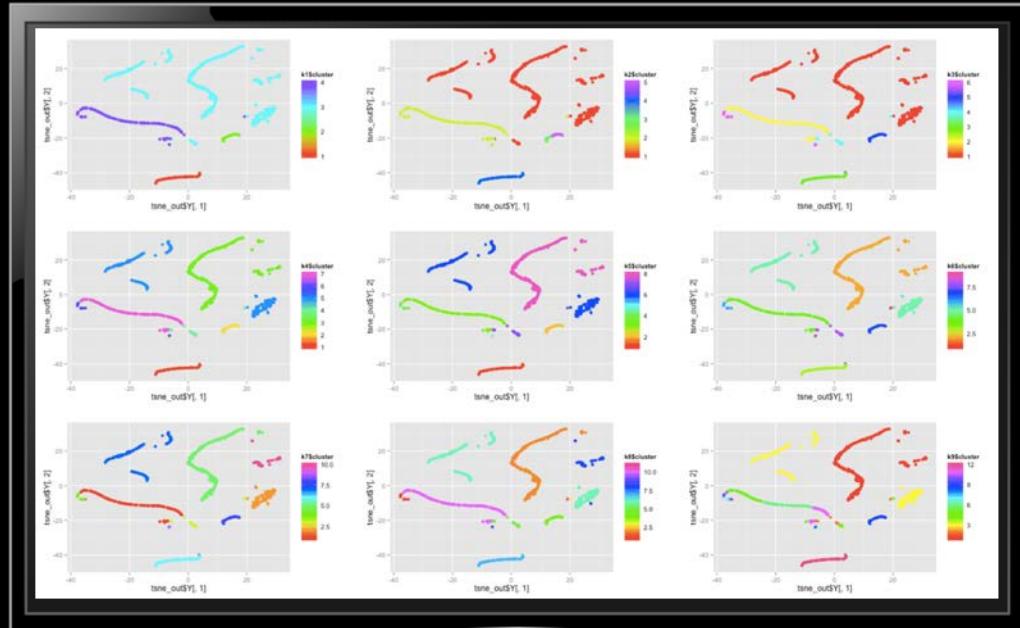
MOST SECURITY PROBLEMS CAN'T
BE SOLVED WITH DEEP LEARNING
or supervised methods in general

UNSUPERVISED TO THE RESCUE?

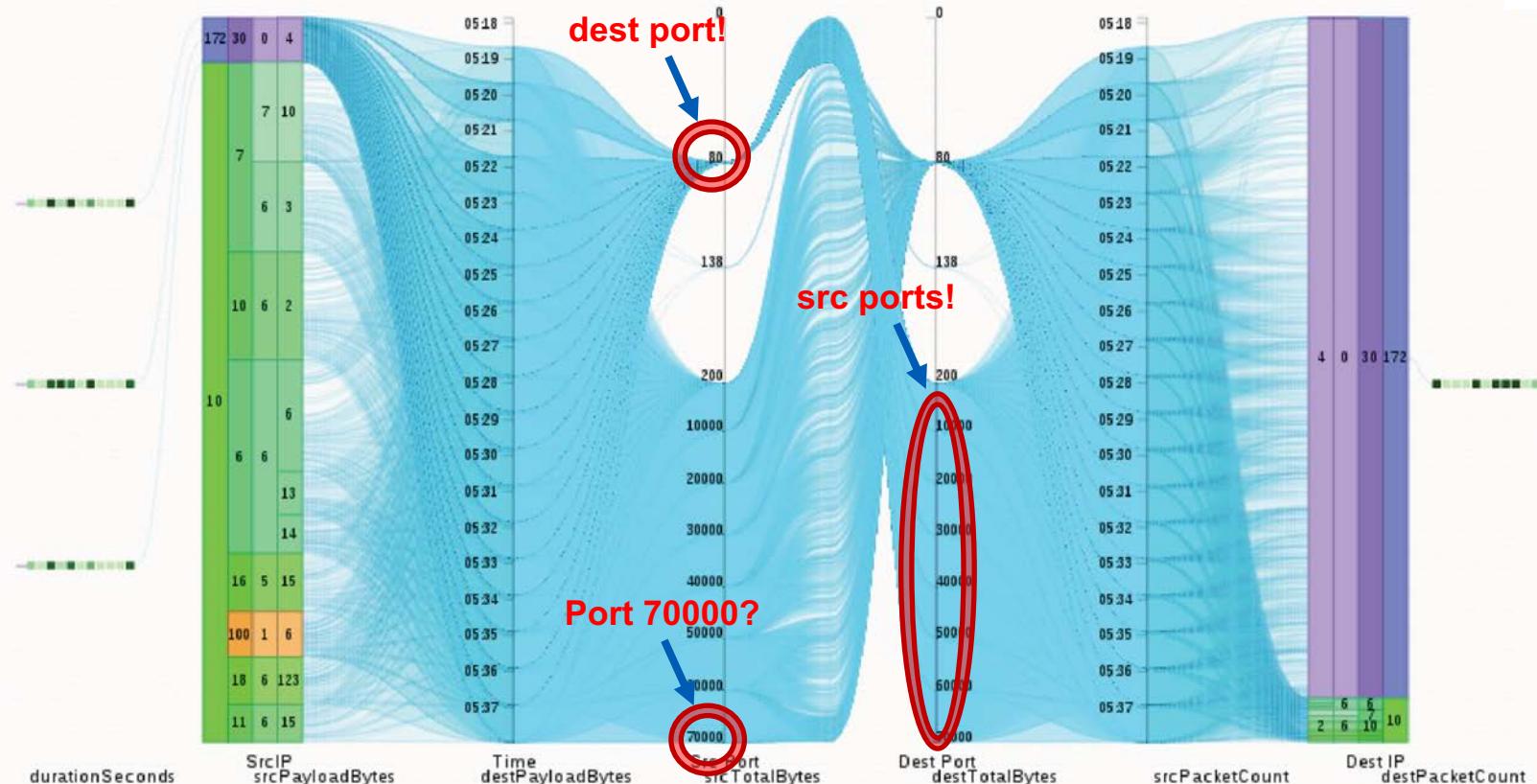
Can we exploit the inherent structure within the data to find anomalies and attacks?

Clustering traffic to find outliers

1. Clean the data
2. Engineer distance functions
3. Figure out the right algorithm
4. Apply the correct algorithmic parameters
5. Data interpretation

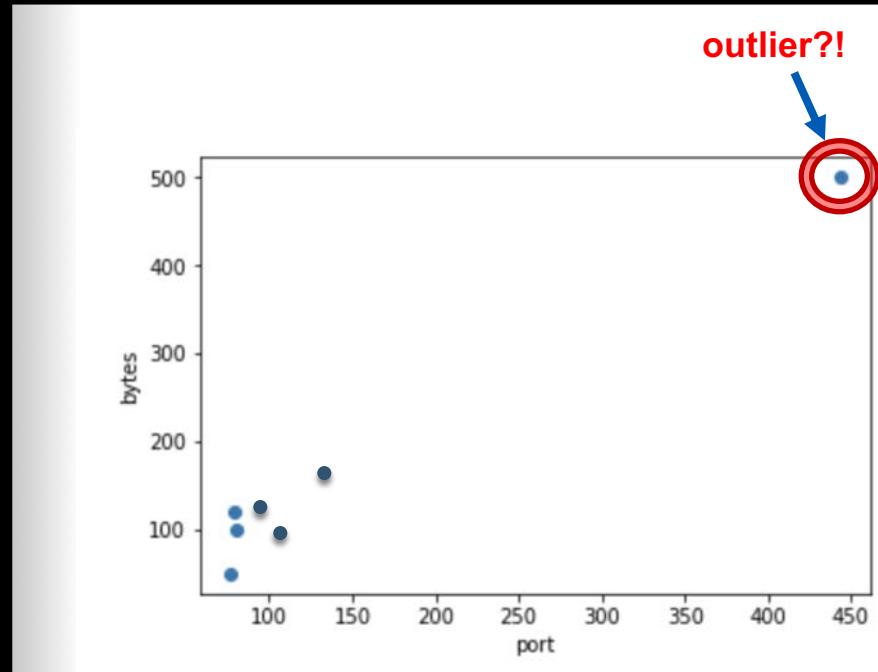


1. UNDERSTAND AND CLEAN THE DATA



2. ENGINEERING DISTANCE FUNCTIONS

- ▶ Distance functions define the similarity of data objects
- ▶ Need domain-specific similarity functions
 - ▶ URLs (simple levenshtein distance versus domain based?)
 - ▶ Ports (and IPs, ASNs) are NOT integers
 - ▶ Treat user names as categorical, not as strings



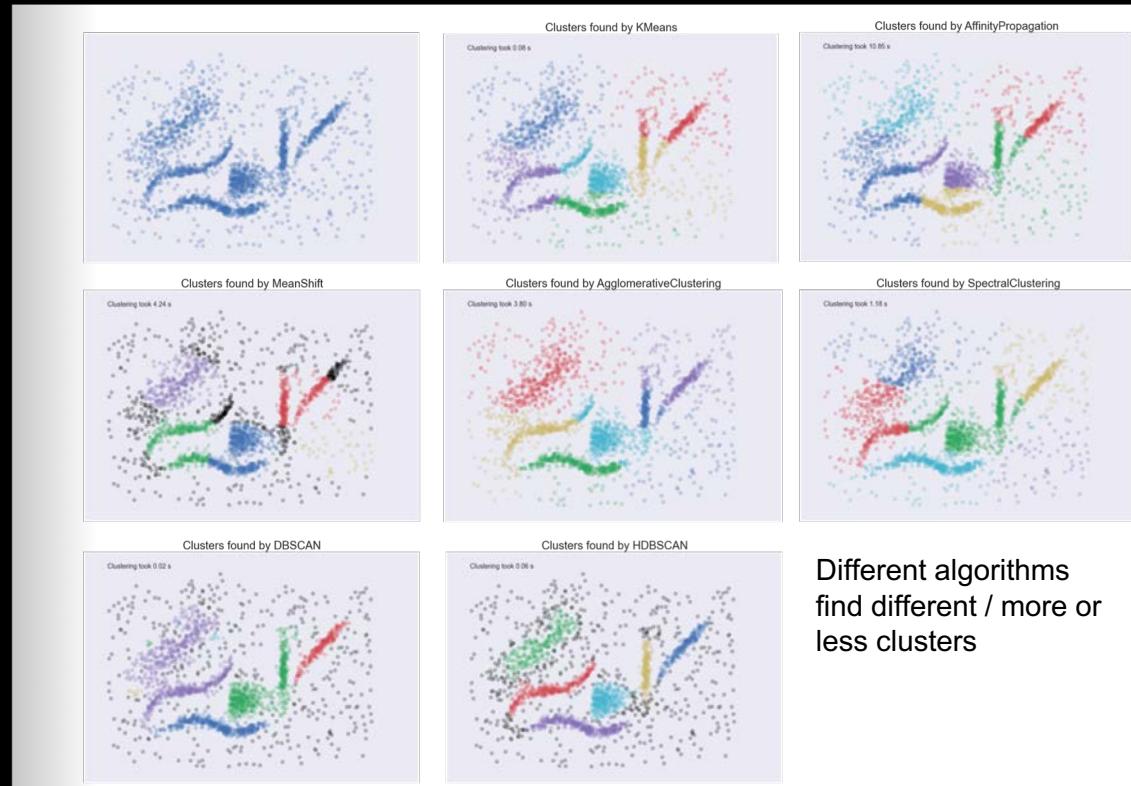
3. CHOOSING THE RIGHT UNSUPERVISED ALGORITHM

CLUSTERING ALGORITHMS

- ▶ K-means
- ▶ Affinity Propagation (AP)
- ▶ DBScan
- ▶ t-SNE

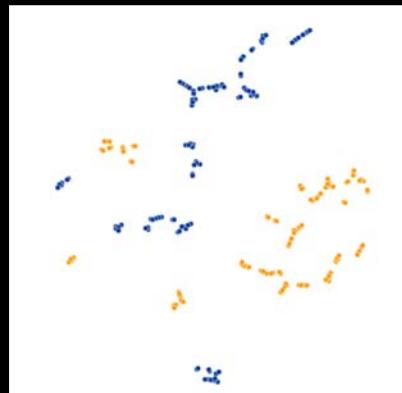
CRITERIA TO CHOOSE AN ALGORITHM

- ▶ Dimensionality of data
- ▶ “Shape” of data
- ▶ Intrinsic algorithm workings
- ▶ Algorithm convergence or speed



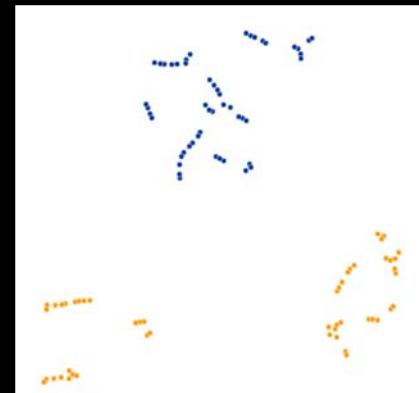
4. CHOOSING THE CORRECT ALGORITHM PARAMETERS

- ▶ The dangers of not understanding algorithmic parameters
- ▶ t-SNE clustering of network traffic from two types of machines



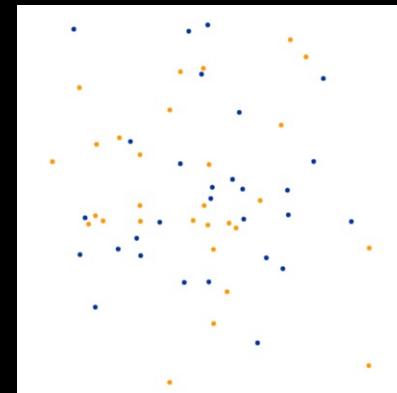
perplexity = 3
epsilon = 3

No clear separation



perplexity = 3
epsilon = 19

3 clusters instead of 2

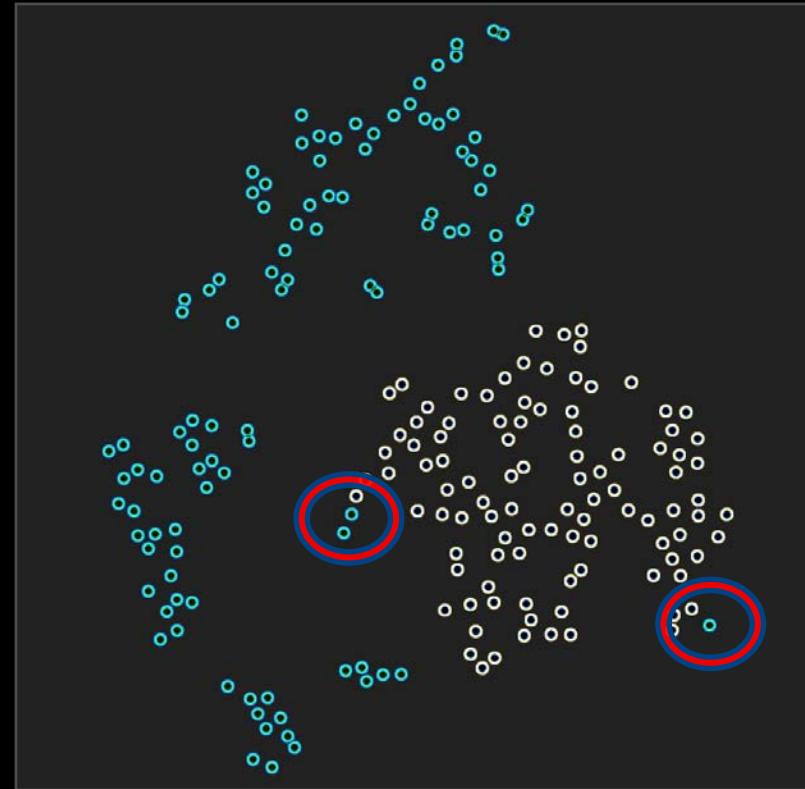


perplexity = 93
epsilon = 19

What a mess

4. CHOOSING THE CORRECT ALGORITHM PARAMETERS

- ▶ And this is when it gets dangerous
- ▶ Access decisions / enforcements based on cluster membership

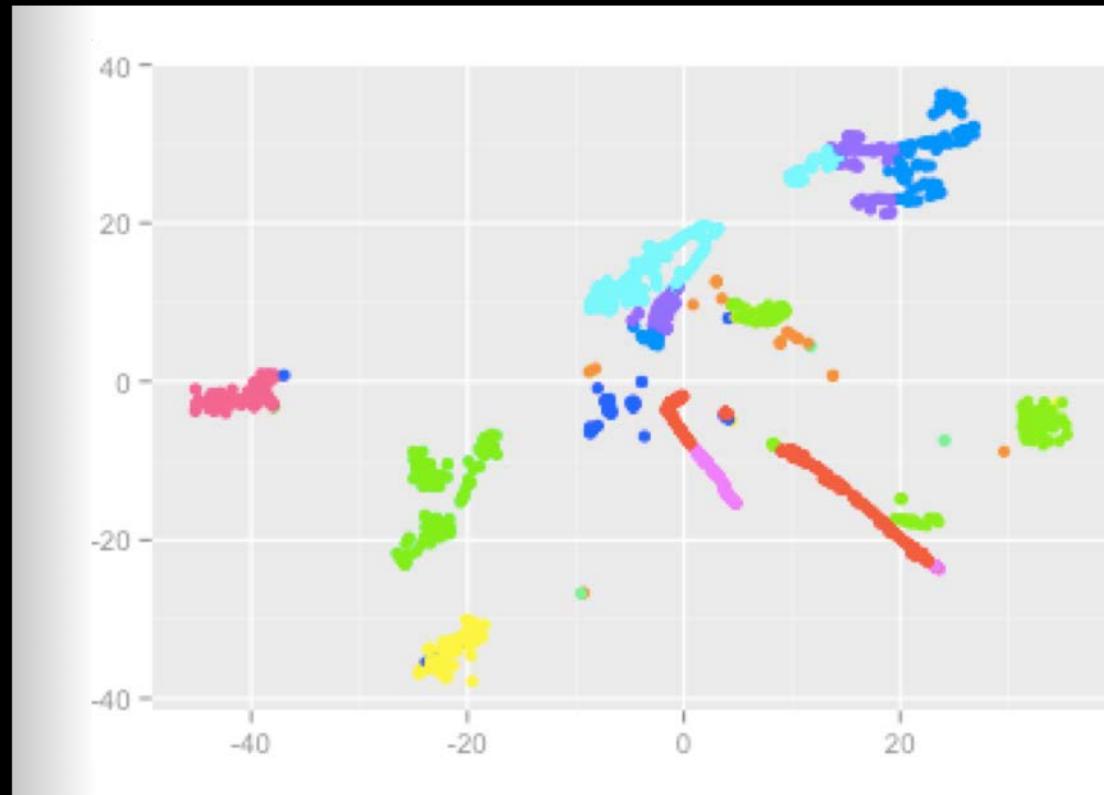


5. INTERPRETING THE DATA

We analyze network traffic. The graph shows an abstract space (X and Y axes have no specific meaning). Each dot represents a device on the network. Colors represent machine-identified clusters.

Interpretation questions:

- ▶ What are these clusters?
- ▶ What are good clusters?
- ▶ What's anomalous?
- ▶ Where are compromised machines / attackers?



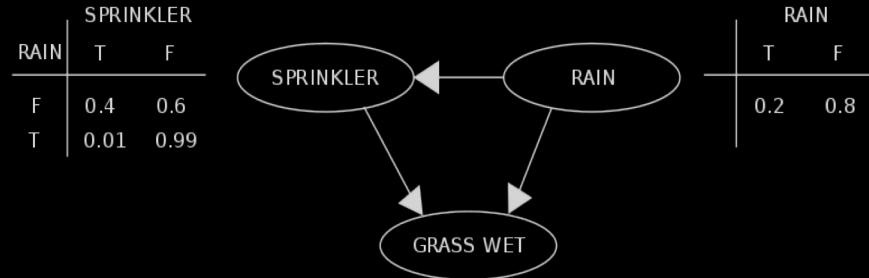
A DIFFERENT APPROACH - PROBABILISTIC INFERENCE

Rather than running algorithms that model the shape of data, we need to take **expert knowledge / domain expertise** into account

Introducing Belief Networks

- ▶ Models that represent the **state** of the ‘world’
- ▶ Helps us make predictions and **reason** about the world
- ▶ A graph rather than huge joint distribution tables across all states
- ▶ Using **Bayes** theorem to calculate ‘belief’
- ▶ Could use ML to learn graph structure (nodes and edges), but it’ll get too unwieldy and non-interpretable!

RAIN	SPRINKLER	
	T	F
F	0.4	0.6
T	0.01	0.99



SPRINKLER	RAIN	GRASS WET	
		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

“What is the probability that it is raining, given the grass is wet?”: 35.77%

BAYESIAN BELIEF NETWORK 1ST STEP – BUILD THE GRAPH

1. What's our objective?

2. What behaviors can we observe?

- ▶ What are observable factors that reduce uncertainty of the central inference (of device compromised)
- ▶ Observations should not be locally dependent – they should be true across all customers / environments
- ▶ Do we have that data?
- ▶ Do we need context for it?



Open port 53

Is using port 23?

Protocol mismatch

New protocol seen

Not seen for a week

Shows up with new OS

Has known vulnerabilities

Mistake in IP classification

Connecting to suspicious IP

Machine got update to new OS

Connecting from suspicious IP

Device is in maintenance mode

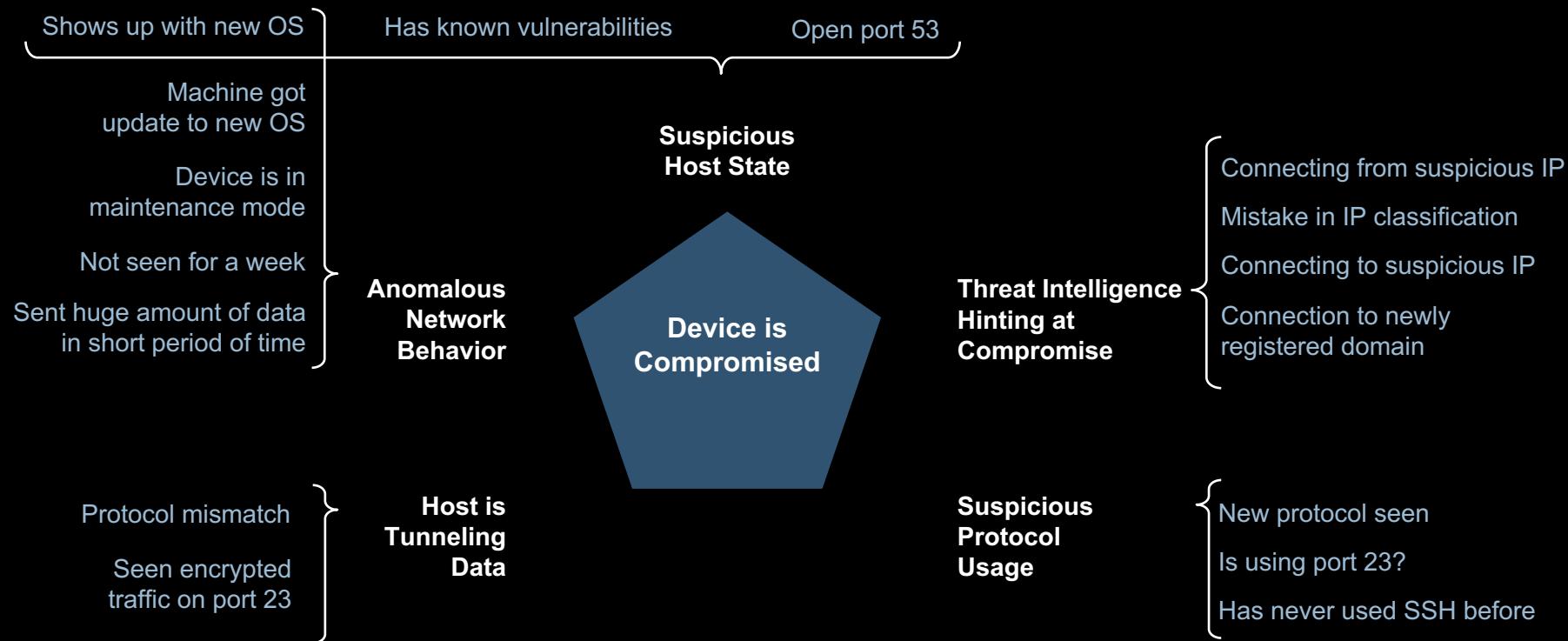
Seen encrypted traffic on port 23

Connection to newly registered domain

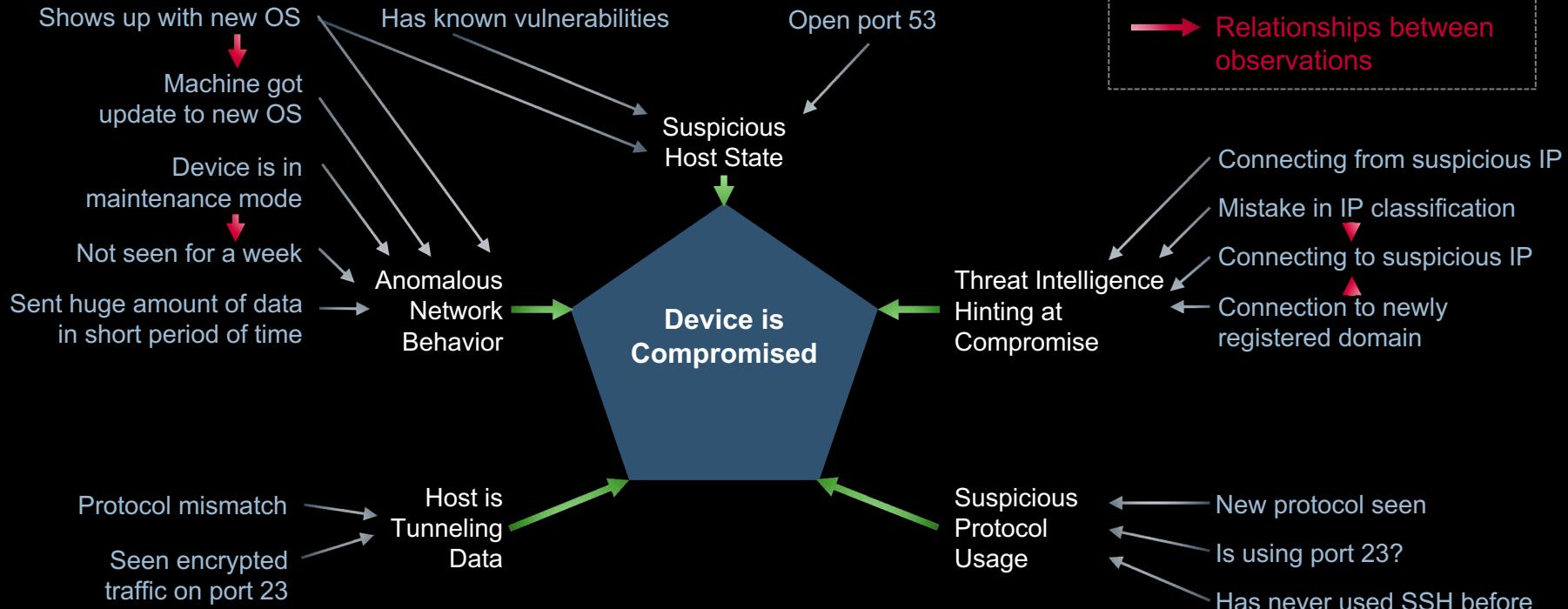
Sent huge amount of data in short period of time

BAYESIAN BELIEF NETWORK 2ND STEP – GROUP NODES

Complexity of this network is too high. We cannot computer all the conditional probabilities.
Therefore we need to introduce “grouping nodes”.



BAYESIAN BELIEF NETWORK 3RD STEP – INTRODUCE DEPENDENCIES



BAYESIAN BELIEF NETWORK 4TH STEP – ESTIMATE PROBABILITIES

NODE PROBABILITIES

- ▶ $P(\text{Protocol mismatch}) = 0.01 \text{ OR "very low"}$
- ▶ $P(\text{Seen encrypted traffic on port 23}) = 0.01 \text{ OR "very low"}$
- ▶ $P(\text{Host is Tunnelling Data}) = 0.01 \text{ OR "very low"}$

}

Expert Knowledge

CONDITIONAL PROBABILITIES

- ▶ Our belief network teaches us: “Tunnelling is not independent of seeing port 23 traffic”
- ▶ $P(\text{Tunnelling} | \text{Enc. Port 23 Traffic}) = (P(\text{Enc. Port 23} | \text{Tunnelling}) * P(\text{Tunnelling})) / P(\text{Enc. Port 23})$

JOINT PROBABILITIES

- ▶ Multiple factors lead to Tunnelling, not just one
- ▶ $P(\text{Tunnelling} | \text{Enc. Port 23 AND Proto mismatch}) = (P(\text{Enc. Port 23 AND Proto mismatch} | \text{Tunnelling}) * P(\text{Tunnelling})) / P(\text{Enc. Port 23 AND Proto mismatch})$

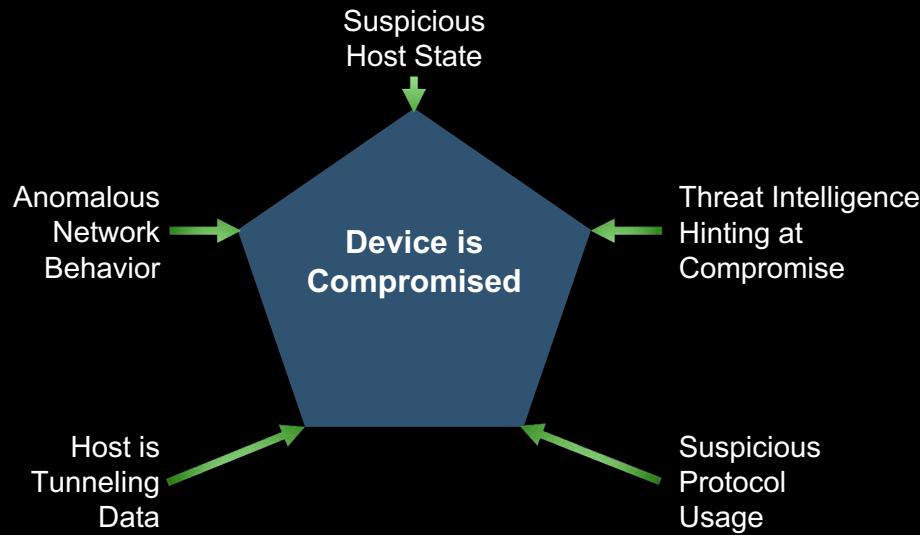
}

More precise than in previous formula

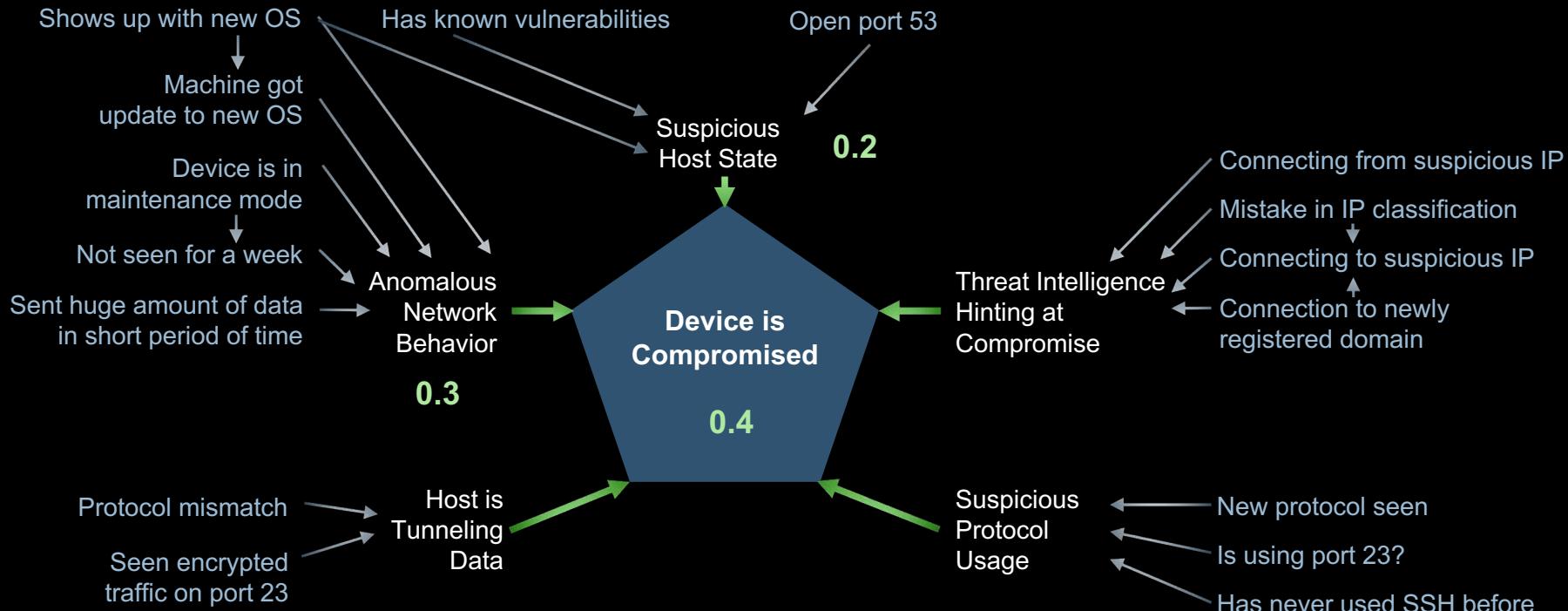


BAYESIAN BELIEF NETWORK 5TH STEP – GOAL COMPUTATION

The probability that we have a compromised device is the joint and conditional probability over all the ‘group nodes’

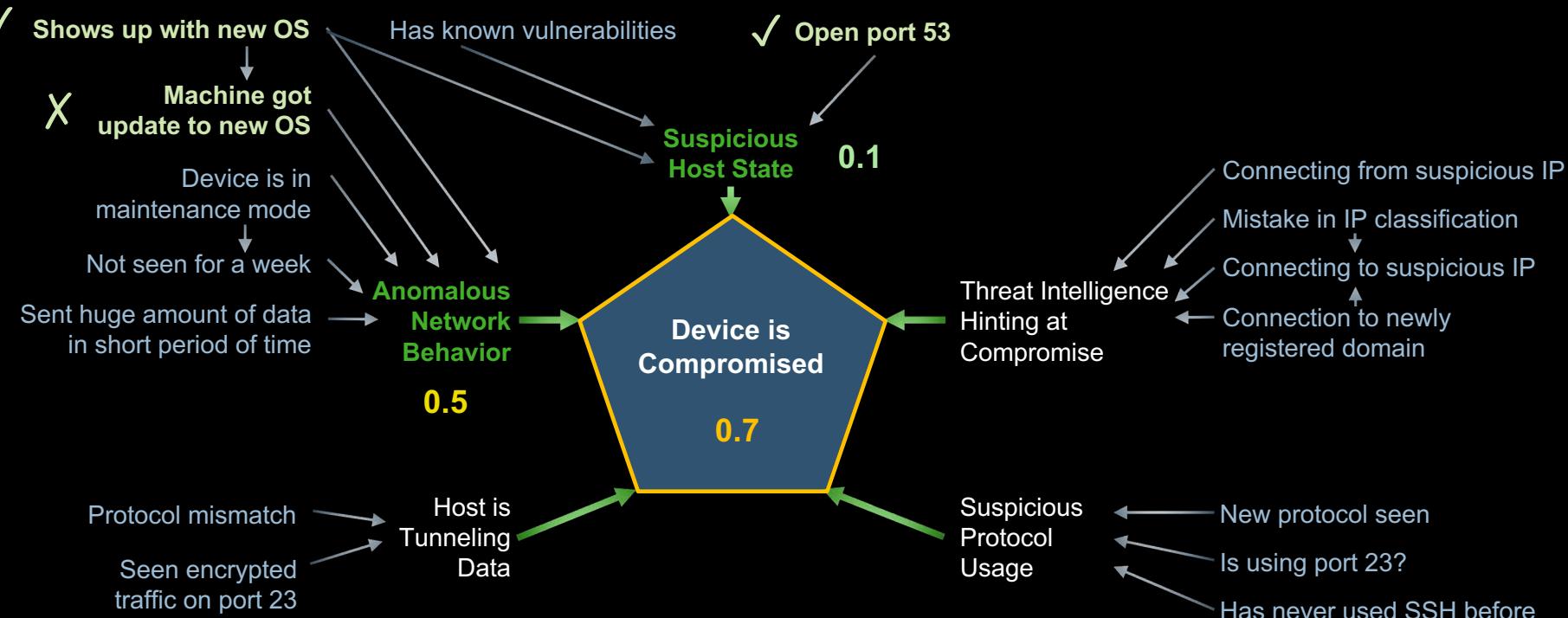


BAYESIAN BELIEF NETWORK 6TH STEP – OBSERVE ACTIVITIES



BAYESIAN BELIEF NETWORK 6TH STEP – OBSERVE ACTIVITIES

1. Update the ‘**observation nodes**’ in the network with observation (what we find in the logs)
2. **Re-compute** probabilities on the connected nodes



BAYESIAN BELIEF NETWORK 7TH STEP – EXPERT INPUT

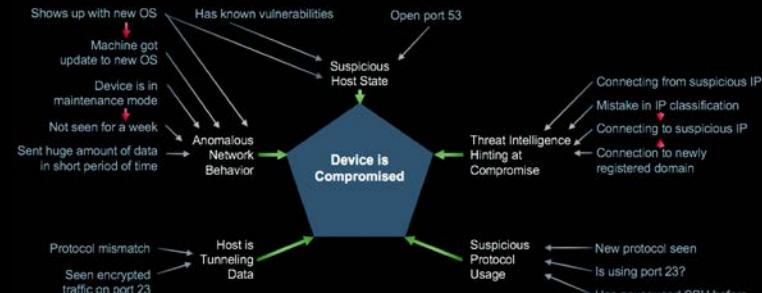
Strengthen the network by introducing **expert knowledge**

- ▶ Pose any combinations of ‘observations’ and ‘group’ nodes as questions to experts
- ▶ Asking *meaningful* questions is an art and requires expert knowledge
- ▶ You will find that it matters how you named your nodes to define good questions

Question	Expert Answer
What's the probability that device is compromised and I have highly suspicious network behavior and nothing on threat intelligence	0.3
Probability that host is in suspicious state, given that port 53 is open, brand new OS	0.1
How likely is it that we see a connection to a newly registered domain and we see port 23 traffic?	0.01
Etc.	

Note how this is not a full joint probability over only a subset of the group nodes.

We can have questions across observational nodes of different groups as well



BELIEF NETWORKS – SOME OBSERVATIONS

Biggest benefit of belief networks is that the learned knowledge can be **verified** and **extracted**!

- ▶ Iterative process of adding more nodes, grouping, adding expert input, etc.
- ▶ Graph allows for answering many questions – e.g., sensitivity analysis
- ▶ Do not determine the probabilities on the observation nodes with historic data. It is only accurate for scenarios that were included in data – how do you know your data covered all scenarios?
- ▶ Each problem requires the definition of a graphs based on expert input
- ▶ A generic “Network Traffic” graph is hard to build and train
 - ▶ Not every FTP is bad
 - ▶ Poor network practice -> e.g., using unencrypted protocols like FTP

IN SUMMARY

RECOMMENDATIONS

- ▶ Start with defining your **use-cases**, not choosing an algorithm
- ▶ ML is **barely ever** the solution to your problem
- ▶ Use **ensembles** of algorithms
- ▶ Teach the algos to **ask for input** – if it's unsure, have it ask an expert rather than making a decision on its own
- ▶ Make sure models keep up with **change** and forget old facts that are not relevant anymore
- ▶ Do you need **white lists / black lists** for your algos to not go haywire?
- ▶ **Verify** your models - use visualization to help with that
- ▶ **Share** your insights with your peers – security is not your competitive advantage
- ▶ GDPR – transparency on what data is collected and used for decisions

 *"The data subject shall have the right not to be subject to a decision based solely on automated processing, including profiling, which produces legal effects concerning him or her or similarly significantly affects him or her."* 

BLACK HAT SOUNDBITES



“Algorithms are getting ‘smarter’,
but experts are more important”

“Understand your data, your algorithms,
and your data science process”

“History is not a predictor
– but knowledge is”



QUESTIONSONS?

[@raffaelmarty](http://slideshare.net/zrlram)