# Applying Data Science to Suricata

## Anomaly Hunting with Suricata & Splunk

splunk>

# Disclaimer

During the course of this presentation, we may make forward looking statements regarding future events or the expected performance of the company. We caution you that such statements reflect our current expectations and estimates based on factors currently known to us and that actual events or results could differ materially. For important factors that may cause actual results to differ from those contained in our forward-looking statements, please review our filings with the SEC. The forward-looking statements made in the this presentation are being made as of the time and date of its live presentation. If reviewed after its live presentation, this presentation may not contain current or accurate information. We do not assume any obligation to update any forward looking statements we may make. In addition, any information about our roadmap outlines our general product direction and is subject to change at any time without notice. It is for informational purposes only and shall not, be incorporated into any contract or other commitment. Splunk undertakes no obligation either to develop the features or functionality described or to include any such feature or functionality in a future release.

splunk> listen to your data

# Speaker's Bio

- Anthony Tellez
  - Splunk Public Sector Federal Team
  - Previously @ NGA
  - Splunkbase App Developer
  - Interests
    - Machine Learning
    - National Security
    - Internet of Things
  - https://github.com/anthonygtellez/
    - https://github.com/anthonygtellez/TA-Suricata
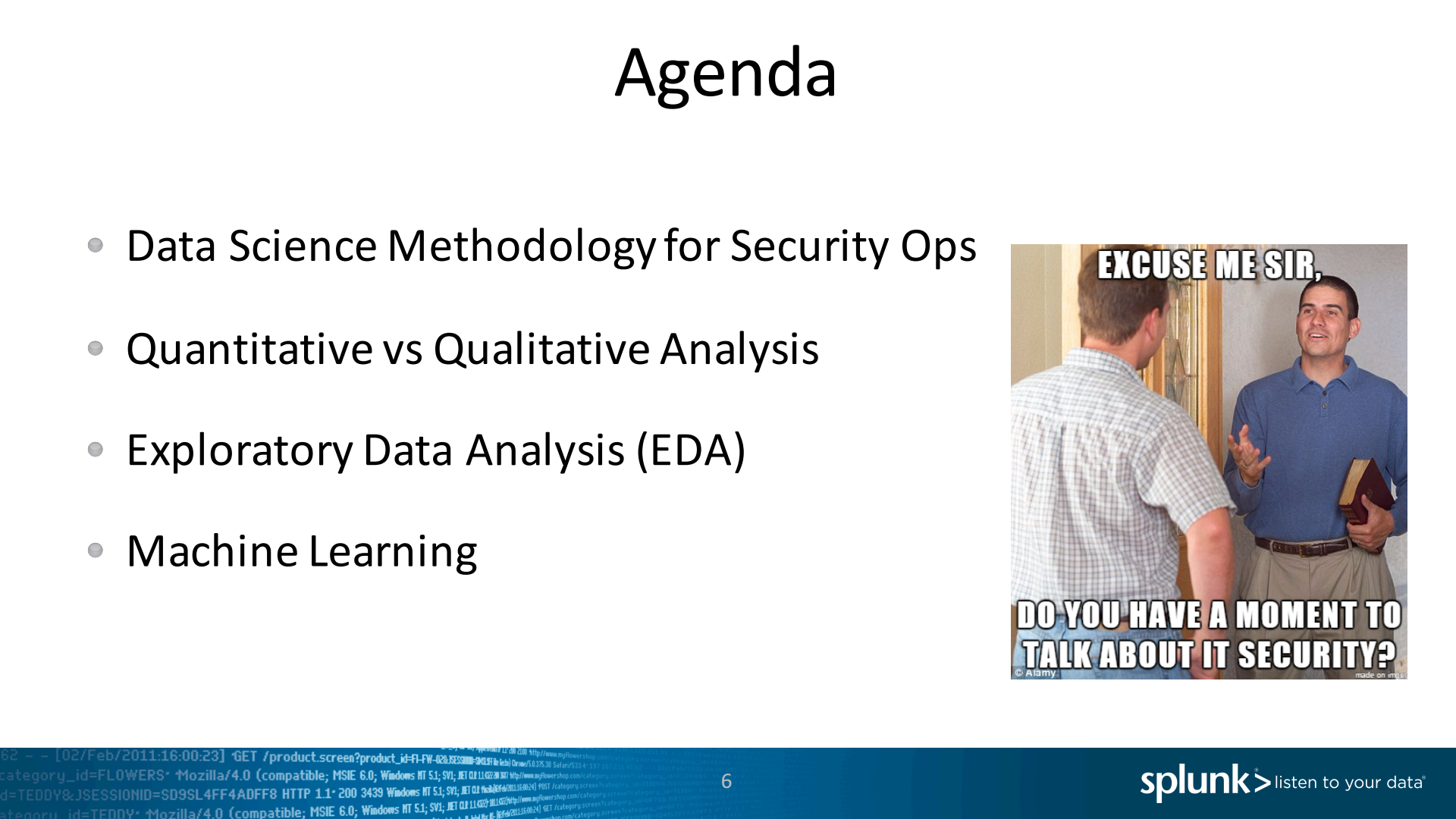    - https://github.com/anthonygtellez/TA-sshd_auth

# What is Data Science?

"Data science is the civil engineering of data. Its acolytes possess a practical knowledge of tools and materials, coupled with a theoretical understanding of what's possible."

*-Mike Driscoll CEO, Metamarket*

# Agenda

- Data Science Methodology for Security Ops

- Quantitative vs Qualitative Analysis

- Exploratory Data Analysis (EDA)

- Machine Learning

62 - - [02/Feb/2011:16:00:23] "GET /product.screen?product_id=FI-FW-02&JSESSIONID=SD3SL9FF4ADFF8 HTTP 1.1" 200 3439 http://www.myflowershop.com/category_id=FLOWERS" Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4322)" 407 http://www.myflowershop.com/category
=TEDDY&JSESSIONID=SD9SL4FF4ADFF8 HTTP 1.1" 200 3439 Windows NT 5.1; SV1; .NET CLR 1.1.4322)http://www.myflowershop.com/category screen?category
category_id=TEDDY" Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4322)" 3[02/Feb/2011:16:00:24] "GET /category.screen?category

splunk® > listen to your data

# Security Data Analysis

- Information Overload
  - IDS alerts, Virus Scans, tools.

- Multidisciplinary approach is needed for next gen problems
  - SIEM alone, ML alone, are not enough without SME.

- Our goal is to empower security analysts to reach the middle using statistical techniques built into many SIEMs.

- **Everyone is capable of becoming a unicorn.**

62 - - [02/Feb/2011:16:00:23] "GET /product.screen?product_id=FI-FW-02&JSESSIONID=SD5SL6FF6ADFF8 HTTP 1.1" 200 2100 http://www.myflowershop...
category_id=FLOWERS" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; .NET CLR 1.1.4322)" 307 http://www.myflowershop.com/category...
d=TEDDY&JSESSIONID=SD9SL4FF4ADFF8 HTTP 1.1" 200 3439 Windows NT 5.1; SV1; .NET CLR 1.1.4322)" 307 /category.screen?category...
category_id=TEDDY" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1;

splunk> listen to your data®

# Correlation != Causation ☹



Worldwide non-commercial space launches correlates with Sociology doctorates awarded (US). Correlation: 78.92% (r=0.78915)

• Correlating some data may be a waste of time if you don't have an understanding of what the data represents.
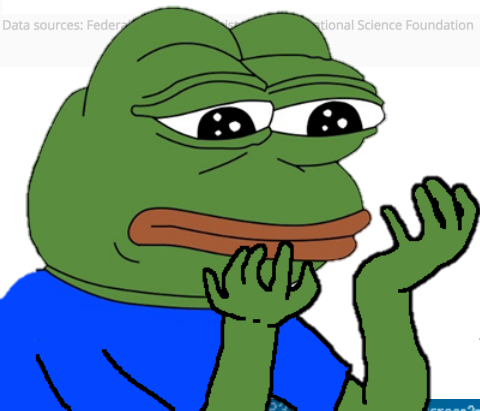


Total revenue generated by arcades correlates with Computer science doctorates awarded in the US. Correlation: 98.51% (r=0.985065)

A good example of why you need a SME

splunk > listen to your data

# 5 Step Data Science Methodology for Security OPS

**Step 1**    Scope relevant machine data to onboard.

**Step 2**    Collect requirements and validate relevant machine data.

**Step 3**    Exploratory Data Analysis. (Searching & Visualizing!)

**Step 4**    Formulate hypothesis working with Domain Experts.

**Step 5**    Test and repeat steps as needed until hypothesis is answered.

splunk > listen to your data®

# Applying Data Science to Security OPS

**Step 1** — Scope relevant machine data to onboard.

Data Sources:
- /var/log/auth.log
- All Network Traffic



THE DATA SUGGESTS.... WAIT, THERE IS NO DATA!

**Deployment/ On PremTier**

**Search Head, Indexer, Deployment Server**

Digital Ocean

App

SURICATA

Splunk Instance

Suricata

SURICATA

Bro

Switch w/ Port Mirroring

# Security Patterns in Machine Data

| What To Look For | Data Source |
|---|---|
| Abnormally high number of file transfers to USB or CD/DVD | **Operating system** |
| Abnormally high number of files or records downloaded from an internal file store or database containing confidential information | **File server / Database** |
| Abnormally large amount of data emailed to personal webmail accounts or uploaded to external file hosting site | **Email server / web proxy** |
| Unusual physical access attempts (after hours, accessing unauthorized area, etc.) | **Physical badge records / Authentication** |
| Excessive printer activity and employee is on an internal watch list as result of demotion / poor review / impending layoff | **Printer logs / HR systems** |
| User name of terminated employee accessing internal system | **Authentication / HR systems** |
| IT Administrator performing an excessive amount of file deletions on critical servers or password resets on critical applications (rogue IT administrator) | **Operating system /Authentication / Asset DB** |
| Employee not taking any vacation time or logging into critical systems while on vacation (concealing fraud) | **HR systems / Authentications** |
| Long running sessions, bandwidth imbalance between client & server, Bad SSL Configurations | **IPS / IDS / Stream** |
| Known cloud or malware domains, bad SSL Configurations | **Threat Intelligence, Custom Lookups** |
| High Entropy Subdomains | **Web proxy, DNS, Wiredata** |

splunk>

# Applying Data Science to Security OPS

**Step 1**   Scope relevant machine data to onboard.

**Step 2**    Collect requirements and validate relevant machine data.

Example Collection Methods
- Universal Forwarder / Agent on Endpoints
  - /var/log/suricata/eve.json
  - /var/log/auth.log

Example Validation Methods
- Add Ons (TA-Suricata, & TA-sshd_auth) / SIEM Parsers
- Regex to build additional fields
- Common Information Model

```
[suricata]
SHOULD_LINEMERGE = true
TIME_PREFIX=timestamp":
BREAK_ONLY_BEFORE = ^{
KV_MODE = json
FIELDALIAS-suricata_global = proto AS transport src_ip AS src dest_ip AS dest
##Vendor Fields
FIELDALIAS-suricata_vendor_id = alert.signature_id AS vendor_sid alert.gid AS vendor_gid ale
rt.rev AS vendor_rev
EVAL-suricata_signature_id = vendor_gid.":".vendor_sid.":".vendor_rev

##FIELD ALIAS FOR IDS
FIELDALIAS-suricata_ids = alert.action AS action alert.gid AS alert_gid alert.rev AS alert_r
ev alert.severity AS severity_id alert.category AS category alert.signature AS signature hos
t AS dvc

##FIELD ALIAS FOR WEB
FIELDALIAS-suricata_web = http.hostname AS dest http.url AS url http.http_user_agent AS http
_user_agent http.http_content_type AS http_content_type http.cookie AS cookie http.length AS
 bytes http.protocol AS http_protocol http.status AS status http.http_method AS http_method
http.http_refer AS http_referrer

##FIELD ALIAS FOR DNS
FIELDALIAS-suricata_dns = dns.id AS transaction_id dns.rcode AS reply_code dns.rdata AS answ
er dns.rdata AS dest dns.rrname AS query dns.ttl AS ttl dns.tx_id AS tx_id dns.type AS messa
ge_type

##FIELD ALIAS FOR SSL
FIELDALIAS-suricata_ssl = tls.fingerprint AS ssl_publickey tls.issuerdn AS ssl_issuer_common
_name tls.sni AS ssl_server_name_indication tls.subject AS ssl_subject_common_name tls.versi
on AS ssl_version
```
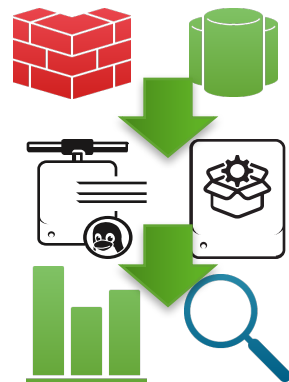
splunk > listen to your data

# Applying Data Science to Security OPS

**Step 1**  Scope relevant machine data to onboard.

**Step 2**  Collect requirements and validate relevant machine data.

**Step 3**  Exploratory Data Analysis. (Searching & Visualizing!)

- Number of connections between src_ip & dest_ip, iplocation
- Torrent activity (dest_port 6881-6889, 6969), connections to Tor Addresses, or Malware domains
- Interesting Fields: http_user_agent, http_method, bytes
- Descriptive Statistics: Producer Consumer Ratio Categories
  Bytes_in/Bytes_Total | Bytes_out /Bytes_total

splunk > listen to your data

# Applying Data Science to Security OPS

**Step 1**  Scope relevant machine data to onboard.

**Step 2**  Collect requirements and validate relevant machine data.

**Step 3**  Exploratory Data Analysis. (Searching & Visualizing!)

**Step 4**  Formulate hypothesis working with Domain Experts.

- Is this real torrent traffic or another application using the same ports?
- Can users install or run TOR Browser onto their desktops in this VLAN?
- Is this SQL injection valid in user_agent field or just bad parsing of data during the onboarding process?

**Can I disprove the activity by adding more data or context?**

# Relevant Data Sources

| Raw Data | Lookups | Context | Value |
|---|---|---|---|
| Firewall Traffic | Username to IP | 10.0.0.12 fails to login to 5 different servers | Determine user responsible |
| Proxy | Username to IP | 10.0.0.12 visits Dropbox and uploads 1TB of data | Determine user responsible |
| Active Directory | User to Group Mapping | SPLUNK\JohnDoe authenticates to 30 different hosts in 30 second period | Determine scope of compromise, domain admin, SQL admin only? |
| DHCP | User to IP, Host to IP | 10.0.0.12, 10.0.0.35 attempt to connect to TOR IP address | Determine user or hosts responsible |
| Email Transport | Baseline Usage | User sends email with large file attachments | Determine normal behavior |
| Exchange / Email | Baseline usage | User sends 40 emails in 60 minute period | Determine normal behavior |
| Packet Capture / Wire Data | Subnet to physical location / priority of asset | 10.0.0.0/27 shows successful SSH connections originating from Russia | Determine where an asset is physically or scope of compromise based on VLAN |

splunk > listen to your data®

# Applying Data Science to Security OPS
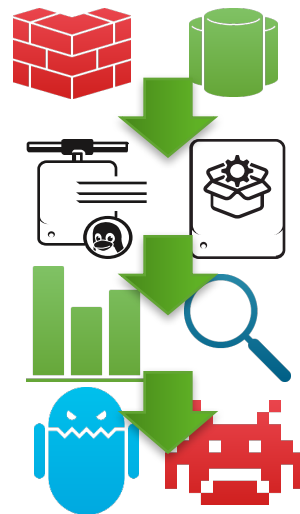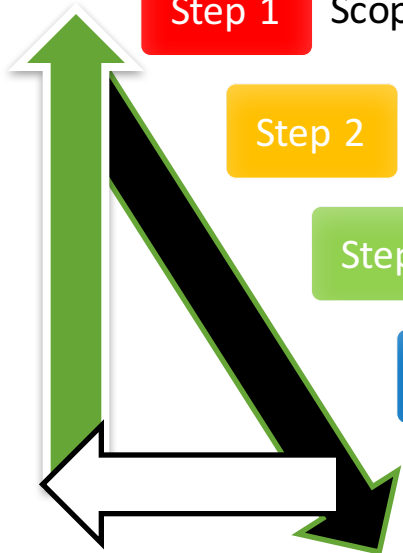
**Step 1**    Scope relevant machine data to onboard.

**Step 2**    Collect requirements and validate relevant machine data.

**Step 3**    Exploratory Data Analysis. (Searching & Visualizing!)

**Step 4**    Formulate hypothesis working with Domain Experts.

**Step 5**    Test and repeat steps as needed until hypothesis is answered.

splunk> listen to your data

# Applying Data Science to Security OPS

**Step 1**  Scope relevant machine data to onboard.

**Step 2**  Collect requirements and validate relevant machine data.

**Step 3**  Exploratory Data Analysis. (Searching & Visualizing!)

**Step 4**  Formulate hypothesis working with Domain Experts.

**Step 5**  Test and repeat steps as needed until hypothesis is answered.

splunk > listen to your data

# **Quan**titative vs **Qual**itative Analysis



- **Quan**titative measure:

-     25 GB of Data uploaded in 60 mins

-     Threshold and periodicity fixed

- **Qual**itative measure:

-     The data uploaded during <span style="color:red">abnormal</span> time periods.

-     Threshold and periodicity is variable

# Quantitative

```
| tstats `summariesonly` sum(Web.bytes_out) as bytes_out
from datamodel=Web where Web.bytes_out>0 by
Web.src,Web.dest
| `drop_dm_object_name("Web")`
| search bytes_out>10485760 | `get_asset(src)`
| search src_priority=high OR src_priority=critical)
```

What does correlation rule this mean??
- Summarize Bytes Out by source, trigger when bytes out exceeds **10485760** and the asset is tagged by the user as **high or critical**.
- Rule fails when asset isn't tagged properly, or bytes is only **10485759**, doesn't take time into context. (Would **10485760** bytes be acceptable over 1 year, 30 days, 1hour?)

splunk> listen to your data

# Qualitative

Enterprise Security 3 - 4+  SA-ExtremeSearch

Create the model in a Context

Count traffic by src in 30m (Takes time into account) √

```
| tstats `summariesonly` dc(All_Traffic.src) as src_count from datamodel=Network_Traffic.All_Traffic by _time span=30m
```

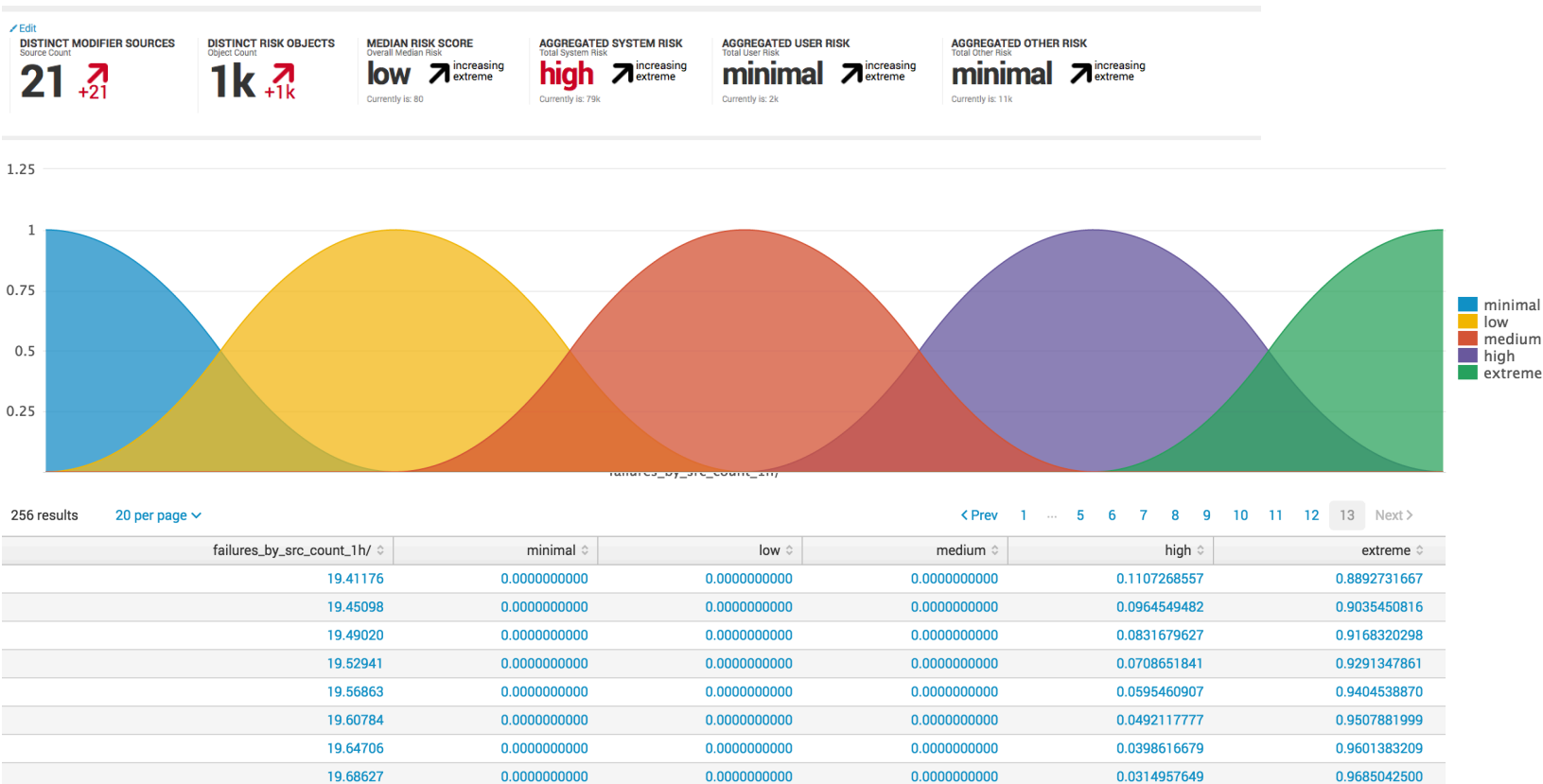Gather stats median, min, max, (descriptive statistics) per src √

```
| stats count, median(total_count) as median, stdev(total_count) as size | search size>0
```

Create a context with current stats per src √

```
| xsupdateddcontext name=count_30m container=network_traffic
terms="minimal,low,medium,high,extreme" type=median_centered width=3 app=SA-
NetworkProtection scope=app
| stats count
```

Time Range -25h to -1h

splunk > listen to your data®

# Visualize Context

# Static to Dynamic Thresholds

- Quantitate v. Qualitative

- Exploratory Data Analysis
  - Descriptive Statistics + Moving Window = Context
  - Visualization
  - Entropy & Correlation

- Machine Learning
  - Supervised v. Unsupervised
  - Security Application of ML
  - Adaptive Thresholding

# EDA - Descriptive Statistics

- In high school statistics you learned about **mean, mode, median, min, max**, & **frequency aka "Descriptive Statistics"**.

- You should make use of these to <u>describe the data</u> you are looking at, <u>explore potential relationships</u> within your data, and <u>ask questions</u> of your data.

- **This iterative process is called "<span style="color:red">Exploratory Data Analysis</span>" it is critical to Machine Learning and Security Analytics.**



DON'T KNOW WHAT A P-VALUE IS

AND AT THIS POINT I'M TOO AFRAID TO ASK

memegenerator.net

# EDA - Descriptive Statistics

- **Compare different duration times of data set for a specific time period.**

- **index**=suricata event_type=flow
  **| stats** count **as** number_events, **min**(duration) **as** min_duration, **max**(duration) **as** max_duration, **avg**(duration) **as** avg_duration, **median**(duration) **as** median_duration, perc95(duration) **as** perc95_duration, **stdev**(duration) **as** stdev_duration

- Are there any long running sessions in the last 60 minutes?

| number_events | min_duration | max_duration | avg_duration | median_duration | perc95_duration | stdev_duration |
|---|---|---|---|---|---|---|
| 3397 | 0 | 3654 | 14.274948 | 0 | 60 | 78.859433 |

splunk> listen to your data

# Applying Descriptive Statistics - PCR

Describing network flows with Producer Consumer  Ratio (PCR)

1. **Create a ratio of bytes_in to bytes_out**

2. **Apply case logic to determine inbound or outbound imbalance between client & server**

- 
```
index=suricata event_type=flow
| eval bytes_total=bytes_in+bytes_out
| eval bytes_ratio= ((bytes_out-bytes_in)/bytes_total)
| eval bytes_pcr_range = case(bytes_ratio > 0.4  "Pure Push", bytes_ratio > 0  "70:30 Export", bytes_ratio == 0
"Balanced Exchange", bytes_ratio >= -0.5  "3:1 Import", bytes_ratio > -1  "Pure Pull"
| stats sparkline(count) AS activity by src_ip src_port dest_ip dest_port bytes_in bytes_out bytes_pcr_range
```
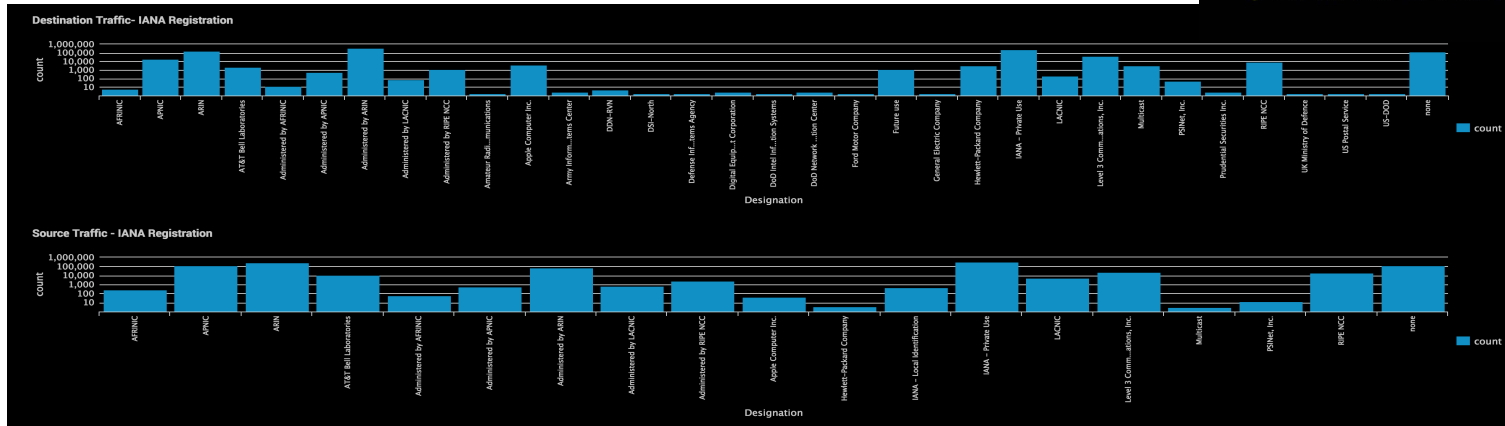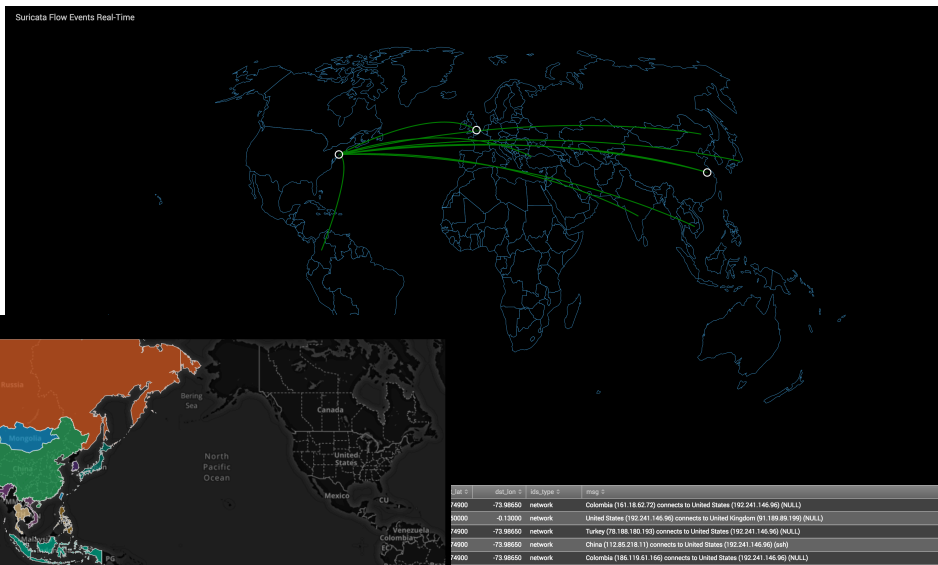
STOP ALL THE DOWNLOADING!

**PCR Ranges:**
**1.0 – Pure Push - FTP**
upload, multicast,
beaconing
**0.4 – 70:30 export -**
Sending Email
**0.0 – Balanced Exchange**
- NTP, ARP probe
**-0.5 – 3:1 import - HTTP**
Browsing
**-1.0 – pure pull - HTTP**
Download

Data Exfiltration, PCR Categories

| src_ip ⌄ | src_port ⌄ | dest_ip ⌄ | dest_port ⌄ | bytes_in ⌄ | bytes_out ⌄ | bytes_pcr_range ⌄ | activity ⌄ |
|---|---|---|---|---|---|---|---|
| 1.196.57.52 | 11595 | 45.79.169.212 | 23 | 54 | 74 | 70:30 Export | |
| 1.34.249.55 | 57909 | 10.10.0.5 | 23 | 54 | 56 | 70:30 Export | |
| 10.0.0.3 | 49488 | 131.253.34.234 | 443 | 5860 | 7253 | 70:30 Export | |
| 10.0.0.3 | 49490 | 65.52.108.231 | 443 | 5904 | 7626 | 70:30 Export | |
| 10.0.0.3 | 49491 | 65.52.108.254 | 443 | 4436 | 3753 | 3:1 Import | |
| 10.0.0.3 | 49492 | 65.52.108.213 | 443 | 5283 | 5724 | 70:30 Export | |
| 10.0.0.3 | 49493 | 131.253.34.230 | 443 | 4436 | 3753 | 3:1 Import | |
| 10.0.0.3 | 49495 | 131.253.34.230 | 443 | 4436 | 3753 | 3:1 Import | |
| 10.0.0.3 | 49782 | 75.75.75.75 | 53 | 210 | 82 | 3:1 Import | |
| 10.0.0.3 | 50185 | 75.75.75.75 | 53 | 255 | 82 | Pure Pull | |

« prev  1  2  3  4  5  6  7  8  9  10  next »

splunk > listen to your data

# Visualization & Creating Context (EDA)

- **Visualization** is a powerful EDA tool
  - Not everything can be described as bits, bytes, plaintext or pie charts.

- **Correlation** to add context to your data during the EDA process or test hypothesis.



PERCENTAGE OF PIE CHART RESEMBLING PAC MAN
PERCENTAGE OF PIE CHART NOT RESEMBLING PAC MAN

# Geographical EDA - Visualization

- Visualization useful for exploring multi-dimensional relationships.

- Tells a story about the data you can't describe in text or tables.

- "Where are connections 'originating', and how often am I seeing this activity?"



Suricata Flow Events Real-Time



Attempted SSH Access by Country

SSH Attempts - Numerical Outliers

Number of Connections: 13
Number of Connections: 15
Number of Connections: 3
Number of Connections: 231
Number of Connections: 47616
Number of Connections: 7
Number of Connections: 95
Number of Connections: 1098
Number of Connections: 891
Number of Connections: 22

Number of Connections: 47616

- I don't remember hiring any remote employees in China.

splunk > listen to your data

# Cool. So how do I operationalize this?

**DNS Water Torture**
- **Botnet sends queries with 16 letters randomly prepended to the victim's domain.**
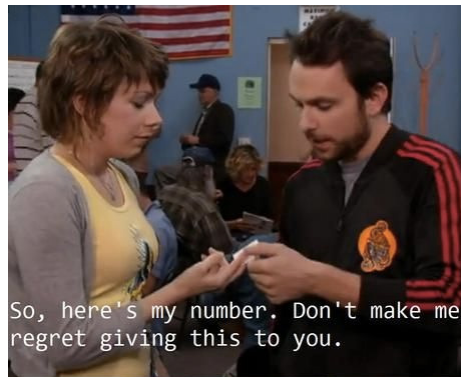  ‣ xyuiasdfcosic.www.halpme.com
  ‣ alkdfejenjasd.www.halpme.com

**C&C Beaconing Activity (Dynamic DNS)**
- **Advanced malware uses a Domain Generation Algorithm (Random Subdomain)**
  ‣ d0290d00xasdf.no-ip[.]org

**Data Exfiltration**
- **DNS Tunneling (Query)**
  ‣ **dnscat.912701a98e9bde415c4ad70007beaf54d2**
  ‣ **dnscat.925401a98ebe0cf540b20d001a4b5e726494b001bb4c192bb68fe73c000bf7c1c0e**

- **Two Techniques to detect this activity in Suricata:**
  - **Shannon Entropy of DNS Query, HTTP destination**
  - **Character Length of DNS Query, HTTP destination**



So, here's my number. Don't make me regret giving this to you.

Wow, so that's your number? Huh, I was so close.

splunk> listen to your data

# Shannon Entropy for EDA & Hunting

- **What is Shannon Entropy?**
  - "… a measure of uncertainty in a random variable"
- **How does it help us find malware and anomalous activity?**

$$H = -\sum p(x) \log p(x)$$

  - The more random a string is, the higher its calculation of randomness.
    ‣ *aaaaa.com (Score 1.8)*
    ‣ *Google.com (Score 2.6)*
    ‣ Ic49f66b73141b5c1.com (Score 4.1)
  - Domains and subdomains with high entropy are good indicators of malicious behavior.
  - **We can filter to domains or subdomains with a score above 3 or 4.**
- **Cons:**
  - **False positives**
    ‣ CDNs like Amazon, Akamai, and others use pseudorandom generated subdomains
    ‣ Requires to you to keep a blacklist or whitelist of domains to reduce noise when hunting (but, relatively easy to do in Splunk)
  - **Malware evolves**
    ‣ Locky & others using shorter subdomains or domains to reduce randomness, reducing entropy score

splunk > listen to your data®

- Python Lookups - Entropy Analysis of DNS / HTTP

- # Full Query for Suricata HTTP

```
index=suricata host=suricata event_type=http
| lookup ut_parse_extended_lookup url AS dest
| lookup ut_shannon_lookup word AS ut_subdomain OUTPUT ut_shannon AS ut_shannon_subdomain
| lookup ut_shannon_lookup word AS dest OUTPUT ut_shannon AS ut_shannon_dest | search ut_shannon_dest > 4
OR ut_shannon_subdomain > 4
| table ut_subdomain ut_shannon_subdomain dest ut_shannon_dest
| dedup dest ut_subdomain
```

- # Results of Suricata HTTP Entropy Scoring

**Subdomain & Domain Entropy Scoring**

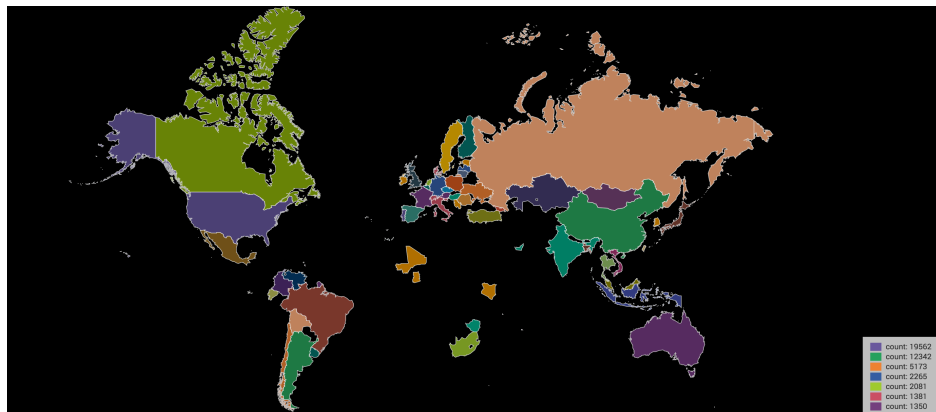| ut_subdomain ⌄ | ut_shannon_subdomain ⌄ | dest ⌄ | ut_shannon_dest ⌄ |
|---|---|---|---|
| ic.49f66b73.141b5c.1.msxbassets.loris | 4.1086680695965025 | ic.49f66b73.141b5c.1.msxbassets.loris.llnwd.net | 4.288082736032309 |
| ic.49f66b73.13d264.1.msxbassets.loris | 4.1831244885738945 | ic.49f66b73.13d264.1.msxbassets.loris.llnwd.net | 4.3041441722488552 |
| ic.49f66b73.020b6e.1.msxbassets.loris | 4.162722123650557 | ic.49f66b73.020b6e.1.msxbassets.loris.llnwd.net | 4.314574491305427 |
| ic.49f66b73.0cdf21.1.xboxone.loris | 4.19438848899739 | ic.49f66b73.0cdf21.1.xboxone.loris.llnwd.net | 4.279519187707896 |
| ic.49f66b73.0fd207.1.xboxone.loris | 4.194388448997389 | ic.49f66b73.0fd207.1.xboxone.loris.llnwd.net | 4.279519187707896 |
| srv-2016-07-31-21.pixel | 3.7950885863977324 | srv-2016-07-31-21.pixel.parsely.com | 4.229003731107054 |
| d1ai9qtk9p41kl | 3.378783493486176 | d1ai9qtk9p41kl.cloudfront.net | 4.142295219190902 |
| srv-2016-07-31-21.config | 3.8868421881310122 | srv-2016-07-31-21.config.parsely.com | 4.350209029099896 |
| d2b3uqm49lqeua | 3.521640636343319 | d2b3uqm49lqeua.cloudfront.net | 4.142295219190901 |
| async-lb-2129785755.us-east-1.elb | 4.028946391954607 | async-lb-2129785755.us-east-1.elb.amazonaws.com | 4.270237192601036 |

« prev 1 2 3 4 5 6 7 8 9 10 next »

# Correlation – Finding Mirai

- **Technique**
  - Default credentials hard-coded in the Scanner.C module give us a **behavioral signature** to look for.
  - Telnet/SSH attempts using invalid users (tech, mother, ubnt, 666666, 888888) are unique to Mirai, & other botnets (post source code leak).
  - **Correlate** list of IPs with Suricata to find other activity from these IoT nodes attempting to breach my network.



```
add_auth_entry("\x14\x14\x14\x14\x14\x14", "\x14\x14\x14\x14\x14\x14", 1);
add_auth_entry("\x1A\x1A\x1A\x1A\x1A\x1A", "\x1A\x1A\x1A\x1A\x1A\x1A", 1);       // 666666    666666
add_auth_entry("\x57\x40\x4C\x56", "\x57\x40\x4C\x56", 1);                       // 888888    888888
add_auth_entry("\x50\x4D\x4D\x56", "\x49\x4E\x54\x13\x10\x11\x16", 1);           // ubnt      ubnt
add_auth_entry("\x50\x4D\x4D\x56", "\x78\x56\x47\x17\x10\x13", 1);               // root      Zte521
add_auth_entry("\x50\x4D\x4D\x56", "\x4A\x4B\x11\x17\x13\x1A", 1);               // root      hi3518
add_auth_entry("\x50\x4D\x4D\x56", "\x48\x54\x40\x58\x46", 1);                   // root      jvbzd
add_auth_entry("\x50\x4D\x4D\x56", "\x43\x4C\x49\x4D", 4);                       // root      anko
add_auth_entry("\x50\x4D\x4D\x56", "\x58\x4E\x5A\x5A\x0C", 1);                   // root      zlxx.
add_auth_entry("\x50\x4D\x4D\x56", "\x15\x57\x48\x6F\x49\x4D\x12\x54\x4B\x58\x5A\x54", 1); // root   7ujMko0vizxv
add_auth_entry("\x50\x4D\x4D\x56", "\x15\x57\x48\x6F\x49\x4D\x12\x43\x46\x4F\x4B\x4C", 1); // root   7ujMko0admin
add_auth_entry("\x50\x4D\x4D\x56", "\x51\x5B\x51\x56\x47\x4F", 1);               // root      system
add_auth_entry("\x50\x4D\x4D\x56", "\x4B\x49\x55\x40", 1);                       // root      ikwb
add_auth_entry("\x50\x4D\x4D\x56", "\x46\x50\x47\x43\x4F\x40\x4D\x5A", 1);       // root      dreambox
add_auth_entry("\x50\x4D\x4D\x56", "\x57\x51\x47\x50", 1);                       // root      user
add_auth_entry("\x50\x4D\x4D\x56", "\x50\x47\x43\x4E\x56\x47\x49", 1);           // root      realtek
add_auth_entry("\x50\x4D\x4D\x56", "\x12\x12\x12\x12\x12\x12\x12\x12", 1);       // root      00000000
add_auth_entry("\x43\x46\x4F\x4B\x4C", "\x13\x13\x13\x13\x13\x13\x13", 1);       // admin     1111111
add_auth_entry("\x43\x46\x4F\x4B\x4C", "\x13\x10\x11\x16\x17", 1);               // admin     1234
add_auth_entry("\x43\x46\x4F\x4B\x4C", "\x17\x16\x11\x10\x13", 1);               // admin     12345
add_auth_entry("\x43\x46\x4F\x4B\x4C", "\x13\x10\x11\x14\x14", 1);               // admin     54321
add_auth_entry("\x43\x46\x4F\x4B\x4C", "\x15\x57\x48\x6F\x49\x4D\x12\x43\x46\x4F\x4B\x4C", 1); // admin  7ujMko0admin
add_auth_entry("\x43\x46\x4F\x4B\x4C", "\x16\x11\x10\x13", 1);                   // admin     1234
add_auth_entry("\x43\x46\x4F\x4B\x4C", "\x52\x43\x51\x51", 1);                   // admin     pass
add_auth_entry("\x43\x46\x4F\x4B\x4C", "\x4F\x47\x4B\x4C\x51\x4F", 1);           // admin     meinsm
add_auth_entry("\x56\x47\x41\x4A", "\x56\x47\x41\x4A", 1);                       // tech      tech
add_auth_entry("\x4F\x4D\x56\x4A\x47\x50", "\x44\x57\x41\x49\x47\x50", 1);       // mother    fucker
```

- Mirai Scanner.C module

splunk>listen to your data

# Finding Mirai – Behavioral & Contextual

- Create CSV Lookup – Invalid users in /var/log/auth.log using information found in Source Code to create a Behavioral Signature

- index=os operation="invalid user"
  | stats count by user src_ip
  | fields user src_ip
  | outputlookup all_invalid_logins.csv

- Filter our CSV to invalid users unique to Mirai

- |inputlookup all_invalid_logins.csv where user="ubnt" OR user="mother" OR user="666666" OR user="888888" OR user="supervisor" OR user="tech"

● **Mirai Scanner.C Adapted to Scan for ARM?**

| | |
|---|---|
| 11/3/16 2:26:26.000 AM | Nov 3 02:26:26 digitalocean sshd[25418]: Invalid user raspberry from 91.200.12.153 |
| | host = ml-bootcamp  source = /var/log/auth.log  sourcetype = ssh:invalid:user |
| 11/3/16 2:26:18.000 AM | Nov 3 02:26:18 digitalocean sshd[25410]: Invalid user raspberry from 91.200.12.153 |
| | host = ml-bootcamp  source = /var/log/auth.log  sourcetype = ssh:invalid:user |
| 11/3/16 2:26:14.000 AM | Nov 3 02:26:14 digitalocean sshd[25407]: Invalid user pi from 91.200.12.153 |
| | host = ml-bootcamp  source = /var/log/auth.log  sourcetype = ssh:invalid:user |
| 11/3/16 2:26:06.000 AM | Nov 3 02:26:06 digitalocean sshd[25399]: Invalid user pi from 91.200.12.153 |
| | host = ml-bootcamp  source = /var/log/auth.log  sourcetype = ssh:invalid:user |
| 11/3/16 2:26:03.000 AM | Nov 3 02:26:03 digitalocean sshd[25396]: Invalid user ubnt from 91.200.12.153 |
| | host = ml-bootcamp  source = /var/log/auth.log  sourcetype = ssh:invalid:user |
| 11/3/16 2:25:55.000 AM | Nov 3 02:25:55 digitalocean sshd[25389]: Invalid user ubnt from 91.200.12.153 |
| | host = ml-bootcamp  source = /var/log/auth.log  sourcetype = ssh:invalid:user |
| 11/3/16 2:25:52.000 AM | Nov 3 02:25:52 digitalocean sshd[25386]: Invalid user admin from 91.200.12.153 |
| | host = ml-bootcamp  source = /var/log/auth.log  sourcetype = ssh:invalid:user |
| 11/3/16 2:25:44.000 AM | Nov 3 02:25:44 digitalocean sshd[25379]: Invalid user admin from 91.200.12.153 |
| | host = ml-bootcamp  source = /var/log/auth.log  sourcetype = ssh:invalid:user |

| src_ip | user |
|---|---|
| 195.22.126.193 | mother |
| 184.173.118.110 | supervisor |
| 193.201.225.113 | supervisor |
| 193.201.225.82 | supervisor |
| 195.22.126.193 | supervisor |
| 201.144.228.137 | supervisor |
| 212.33.200.73 | supervisor |
| 220.178.13.70 | supervisor |
| 95.46.140.178 | supervisor |
| 114.108.150.118 | tech |
| 116.12.146.226 | tech |
| 119.254.162.204 | tech |
| 195.22.126.193 | tech |

from 58.108.219.65

om 123.31.31.95

Is anyone on the internet, the person they say they are?!

FX

splunk > listen to your data®

# Finding Mirai – Behavioral & Contextual

- **Let's correlate the invalid user IPs using Mirai Creds with our Suricata eve.json logs to see if there are any matches on our network!**

- index=suricata
  [ | inputlookup all_invalid_logins.csv  where user="ubnt" OR user="mother" OR user="666666" OR user="888888" OR user="supervisor" OR user="tech" | table src_ip | dedup src_ip ]

- **Using MaxMind we can "geo-locate" the IoT devices trying to gain access:**

- | iplocation src_ip | stats count by Country | geom geo_countries featureIdField=Country

| | |
|---|---|
| 0 - 600 | |
| 600 - 1,200 | |
| 1,200 - 1,800 | |
| 1,800 - 2,400 | |
| 2,400 - 3,000 | |
| 3,000 - 3,600 | |
| 3,600 - 4,200 | |
| 4,200 - 4,800 | |
| 4,800 - 5,400 | |

11/2/16
2:38:04.001 PM

{ [-]
   app_proto: ssh
   dest_ip:
   dest_port: 22
   event_type: flow
   flow: { [-]
      age: 13
      bytes_toclient: 3717
      bytes_toserver: 2943
      end: 2016-11-02T14:37:03.365353+0000
      pkts_toclient: 22
      pkts_toserver: 22
      reason: timeout
      start: 2016-11-02T14:36:50.930487+0000
      state: closed
   }
   flow_id: 1028078808
   proto: TCP
   src_ip: 91.224.160.184
   src_port: 41381
   tcp: { [+]
   }
   timestamp: 2016-11-02T14:38:04.001945+0000
}
Show as raw text

host = ml-bootcamp    source = /var/log/suricata/eve.json    sourcetype = suricata

splunk > listen to your data

# Machine Learning

## Supervised

▸ Classification (Nearest Neighbors, Support Vector Machines, Naïve Bayes, Decision Tree)

 – Group "like" things together based on selected features.

▸ Regression (Linear & Logistic)

 – Infer a relationship between two variables (x) & result (y).

## Unsupervised

▸ Clustering (K Means)

 – Partition events with multiple numeric fields into clusters

▸ Decomposition(PCA, SVD)

 – Dimension Reduction, explains the maximum variance of the higher dimension

# Machine Learning – Security Application

- A toolset for asking research questions which we want to operationalize.

- **Problem:** BotNet DDoS attacks are problematic for all size of organizations. They take a time, money and manpower to resolve. The IP addresses are dynamic making simple whitelist/blacklist mitigation not feasible.

- **Hypothesis:** *"Are there patterns in botnet network activity that can be leveraged to identify the specific botnet and mitigate the threat posed by that botnet?"*

# Machine Learning – Security Application

- 50k random Suricata flow events, dest_port=22

  - Features: packet_ratio, packets_in, packets_out, packets_total

  - Labels: isMirai = 1 or 0

  - Kmeans Cluster

    - K=5



| City | Country | Designation | Region | _time | isMirai | lat | lon | packet_pcr_range | packet_ratio | packets_in | packets_out | packets_total | src_ip | src_port | tcp_flag_hex_to_client | tcp_flag_hex_to_server |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Thailand | APNIC | | 2016-10-13 | 0 | 13.75000 | 100.46670 | 3:1 Import | -0.111111 | 5 | 4 | 9 | 1.10.214.109 | p_20599 | 1e | 1a |
| | China | APNIC | | 2016-08-23 | 1 | 35.00000 | 105.00000 | 3:1 Import | -0.037736 | 55 | 51 | 106 | 1.119.7.35 | p_9224 | 19 | 1b |
| | China | APNIC | | 2016-08-23 | 1 | 35.00000 | 105.00000 | 3:1 Import | -0.070000 | 107 | 93 | 200 | 1.119.7.35 | p_9224 | 1b | 1b |
| Seoul | Republic of Korea | APNIC | Seoul | 2016-09-19 | 0 | 37.59850 | 126.97830 | Balanced Exchange | 0 | 18 | 18 | 36 | 1.209.148.34 | p_44084 | 1b | 1b |
| Seoul | Republic of Korea | APNIC | Seoul | 2016-09-19 | 0 | 37.59850 | 126.97830 | 3:1 Import | -0.058824 | 18 | 16 | 34 | 1.209.148.34 | p_45435 | 1b | 1b |

# Machine Learning – Security Application

- **Cluster 4** is an outlier

- Characteristics of **Cluster 4**



| packet_pcr_range | packet_ratio | packets_in | packets_out | packets_total | isMirai | count |
|---|---|---|---|---|---|---|
| 3:1 Import | -0.008696 | 116 | 114 | 230 | 1 | 1 |
| 3:1 Import | -0.014286 | 142 | 138 | 280 | 1 | 1 |
| 3:1 Import | -0.014749 | 172 | 167 | 339 | 1 | 1 |
| 3:1 Import | -0.017143 | 178 | 172 | 350 | 0 | 1 |
| 3:1 Import | -0.026316 | 156 | 148 | 304 | 0 | 1 |
| 3:1 Import | -0.039301 | 238 | 220 | 458 | 0 | 1 |
| 3:1 Import | -0.047319 | 166 | 151 | 317 | 0 | 1 |
| 3:1 Import | -0.047619 | 154 | 140 | 294 | 0 | 1 |
| 3:1 Import | -0.04797 | 142 | 129 | 271 | 1 | 1 |
| 3:1 Import | -0.054545 | 116 | 104 | 220 | 1 | 1 |
| 3:1 Import | -0.058824 | 108 | 96 | 204 | 0 | 1 |
| 3:1 Import | -0.064615 | 173 | 152 | 325 | 1 | 1 |
| 3:1 Import | -0.07 | 107 | 93 | 200 | 1 | 1 |
| 3:1 Import | -0.081761 | 172 | 146 | 318 | 1 | 1 |
| 3:1 Import | -0.09465 | 133 | 110 | 243 | 1 | 1 |
| 3:1 Import | -0.115789 | 106 | 84 | 190 | 0 | 1 |
| 3:1 Import | -0.122995 | 105 | 82 | 187 | 0 | 1 |
| 3:1 Import | -0.129032 | 105 | 81 | 186 | 0 | 1 |
| 3:1 Import | -0.130435 | 104 | 80 | 184 | 0 | 1 |
| 3:1 Import | -0.132275 | 107 | 82 | 189 | 0 | 1 |

splunk > listen to your data

# Machine Learning – Security Application

- We now have a model which describes different Botnet populations.
- Let's use this model to predict if the connection is Mirai based on Cluster Distance
- Cluster Distance
    - Describes the distance from a centroid
- Prediction Algorithms:
    - Linear Regression
    - Decision Tree
    - Random Forest

Actual vs. Predicted Scatter Chart [↗]

# Machine Learning- Predict Mirai

- Linear Regression

- Results
  - High precision at predicting 0
  - 13.% False positives
  - 41.9% False Negatives
  - 58.1% Correct at getting Mirai Traffic correct

- Summary
  - IPS sensor allowed all of these connections (Not Blocked), while we missed 41.9% of these attacks.
  - We now have a model which we can further refine to identify malicious SSH traffic to investigate.
  - Adds a new layer to our security stack

| Precision ⬈ | Recall ⬈ | Accuracy ⬈ | F1 ⬈ |
|---|---|---|---|
| 0.98 | 0.86 | 0.86 | 0.92 |

Open in Search    Show SPL

**Classification Results (Confusion Matrix)** ⬈

| Predicted actual ⇕ | Predicted 0 ⇕ | Predicted 1 ⇕ |
|---|---|---|
| 0 | 21318 (86.7%) | 3271 (13.3%) |
| 1 | 117 (41.9%) | 162 (58.1%) |

Open in Search    Show SPL

gifak-net

splunk > listen to your data

# Machine Learning- Predict Mirai

- Random Forest

- Results
  - High precision at predicting 0
  - Small false positive (8/25,000)
  - 10.6% False Negatives
  - 89.4% Correct at getting Mirai Traffic correct

- Summary
  - IPS sensor allowed all of these connections (Not Blocked), while we missed 10.6%
  - We now have a model which we can further refine to identify malicious SSH traffic to investigate.
  - Adds a new layer to our security stack

**Algorithm**
RandomForestClassifier ▾

**Field to predict**
isMirai ▾

**Fields to use for predicting**
✕ cluster_distance   ✕ packet_pcr_range   ✕ packet_ratio
✕ packets_total

**Split for training / test: 50 / 50**

**N Estimators**
(optional)

**Max Depth**
(optional)

**Max Features**
(optional)

**Min Samples Split**
(optional)

**Max Leaf Nodes**
(optional)

**Save the model as**
predict_mirai_botnet_kmeans_model

[ Fit Model  ⏱ ]  [ Open in Search ]  [ Show SPL ]

### Classification Results (Confusion Matrix) ⬈

| Predicted actual ⇕ | Predicted 0 ⇕ | Predicted 1 ⇕ |
|---|---|---|
| 0 | 24659 (100%) | 8 (0%) |
| 1 | 28 (10.6%) | 235 (89.4%) |

[ Open in Search ]  [ Show SPL ]

### Prediction Results ⬈

| isMirai ⇕ | predicted(isMirai) ⇕ | cluster_distance ⇕ | packet_pcr_range ⇕ | packet_ratio ⇕ | packets_total |
|---|---|---|---|---|---|
| 0 | 0 | 8.00774122578 | 3:1 Import | -0.111111 | 9 |
| 1 | 0 | 2133.4800996 | 3:1 Import | -0.037736 | 106 |
| 1 | 0 | 13991.7377947 | 3:1 Import | -0.07 | 200 |
| 0 | 0 | 210.153851029 | 70:30 Export | 0.083333 | 24 |
| 0 | 1 | 170.967448746 | 70:30 Export | 0.111111 | 45 |
| 0 | 0 | 2.79967846195 | 3:1 Import | -0.028571 | 35 |
| 0 | 0 | 5.71179214833 | Balanced Exchange | 0.0 | 36 |
| 0 | 0 | 210.153851029 | 70:30 Export | 0.083333 | 24 |
| 0 | 0 | 5.71179214833 | Balanced Exchange | 0.0 | 36 |
| 0 | 0 | 84.9325049989 | 3:1 Import | -0.147541 | 61 |

« prev  1  2  3  4  5  6  7  8  9  10  next »

# Machine Learning – Next Steps

- Model is quite accurate because there **\*may\*** be an indicator of compromise it has found.

- How to validate:
  - Assume Null Hypothesis
    - Add more data
    - Validate Variance & Entropy
    - Work with peers to cross validate model
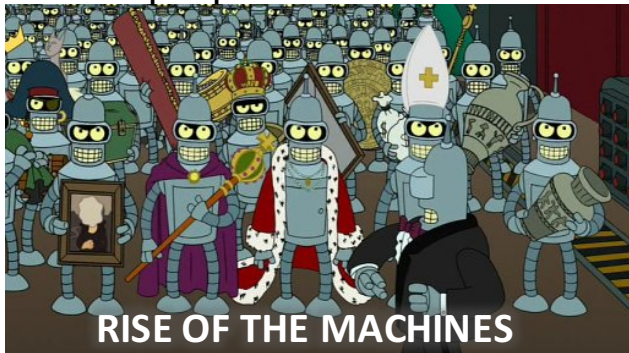


THECUTECORNER.COM

splunk> listen to your data

# Machine Learning - Adaptive Thresholding

- Make use of eval to create bytes_total & bytes_ratio for Producer Consumer Ratio (PCR) for KPI Base Search & NetFLOW

```
index=suricata event_type=flow
| eval bytes_total=bytes_in+bytes_out
| eval bytes_ratio= ((bytes_out-bytes_in)/bytes_total)
```

- Thresholding score compares the current traffic against a rolling hourly average and standard deviation from mean for last 30 days of data.

- Bytes Ratio Thresholds based on PCR Static Ratios
  - **1.0 – pure push - FTP upload, multicast, beaconing**
  - **0.4 – 70:30 export - Sending Email**
  - **0.0 – Balanced Exchange - NTP, ARP probe**
  - **-0.5 – 3:1 import - HTTP Browsing**
  - **-1.0 – pure pull - HTTP Download**



**RISE OF THE MACHINES**

splunk > listen to your data

# Machine Learning - Adaptive Thresholding

- Visualization of the same PCR Suricata Flow

- Health score based on 5 KPIs. The current traffic (bytes_in, bytes_out, bytes_total, packets_in, & packets_out) compared to a rolling hourly average, and standard deviation from mean.

- Attempting to define "What is normal and when is something deviating from the norm I've seen for 30 days?"

- Bytes Ratio based on PCR Ratio for thresholding.

splunk> listen to your data

# Recap

- √ 5 Step Data Science Methodology for Security

- √Descriptive Statistics

- √Quantitative vs Qualitative Analysis

- √Exploratory Data Analysis (EDA)

- √Machine Learning

splunk> listen to your data

Thank You

splunk>

# Glossary

- Descriptive Statistics
  - Min, Max, Median, Average(Mean), Standard Deviation, Mode
  - Z-Scores

- Exploratory Data Analysis
  - Searching the data and looking for relationships
  - Leveraging knowledge ( lookups , reference tables )

- Entropy
  - Measurement of how mixed up something is
    ‣ e.g. non-numerical field such as query compared against wordlist

- P-Values
  - "The p-value is defined as the probability of obtaining a result equal to or "more extreme" than what was actually observed, when the null hypothesis is true."

splunk > listen to your data

# Explore Splunk Analytics

- Anomalies
  - Analyzes numeric fields for their ability to predict another discrete field.

- Anomalousvalue
  - Computes an "unexpectedness" score for an event.

- Anomalydetection
  - Finds and summarizes irregular, or uncommon, search results.

- Cluster
  - Computes a probability for each event and detects unusually small probabilities.

- Kmeans
  - Groups similar events together.

- Outlier
  - Removes outlying numerical values.

- Rare
  - Displays the least common values of a field.

splunk> listen to your data®

# References & Resources

- Spurious Correlations http://www.tylervigen.com/spurious-correlations

- PCR – A New Flow Metric http://qosient.com/argus/presentations/Argus.FloCon.2014.PCR.Presentation.pdf

- Data Driven Security http://datadrivensecurity.info/

- Splunk Syntax Highlighting http://blog.metasyn.pw/splunk-syntax-highlighting/

- Doing Data Science http://shop.oreilly.com/product/0636920028529.do

- Hunting the Known Unknowns (with DNS) https://conf.splunk.com/speakers/2015.html#search=Kovar&

- Lookups, and other goodies https://github.com/anthonygtellez/conf2016_extras

- IDS Evasion w TTL - http://insecure.org/stf/secnet_ids/secnet_ids.html

- Applying Machine Learning to Network Security Monitoring http://www.mlsecproject.org/#conference-presentations

- Scikit-Learn http://scikit-learn.org/

- Machine Learning Toolkit https://splunkbase.splunk.com/app/2890/

- URL Toolbox https://splunkbase.splunk.com/app/2734/

splunk > listen to your data