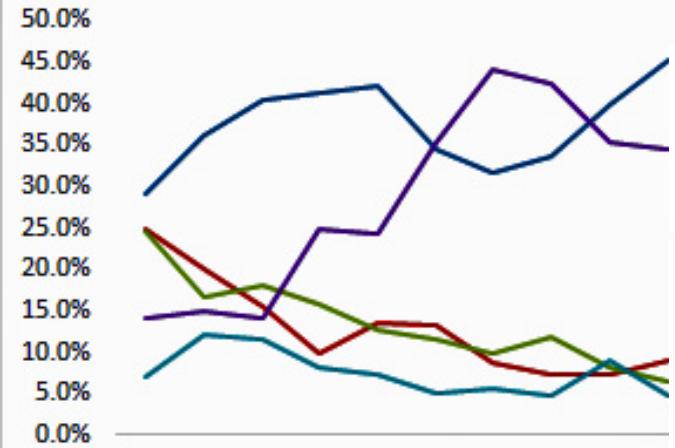


# Practical Statistics for Threat Intelligence

Figure 1: Industry Sectors Percentage of Overall Breaches

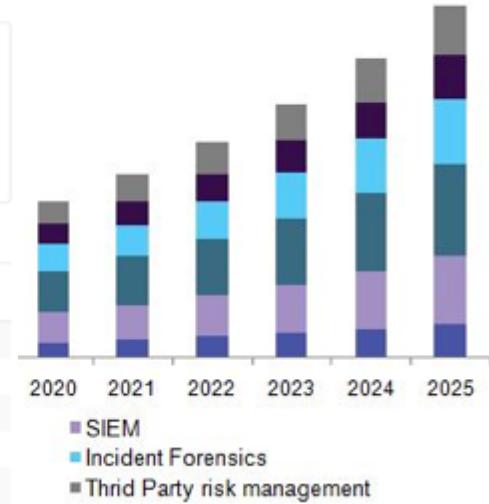
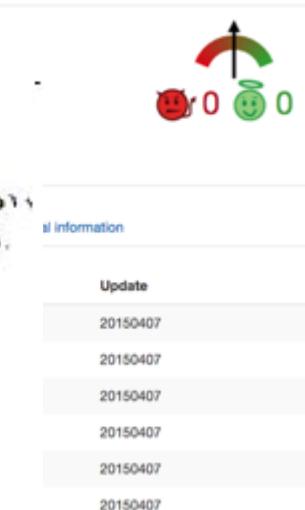


virus total

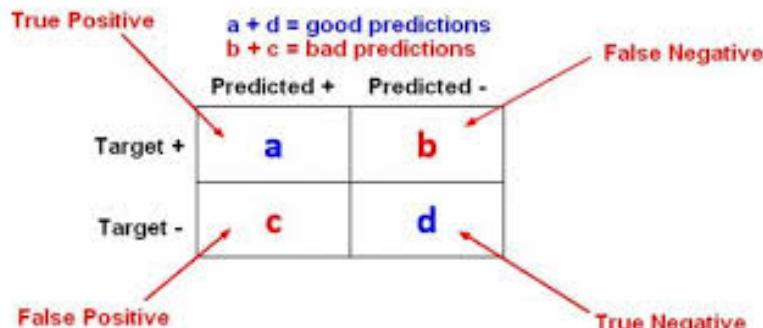
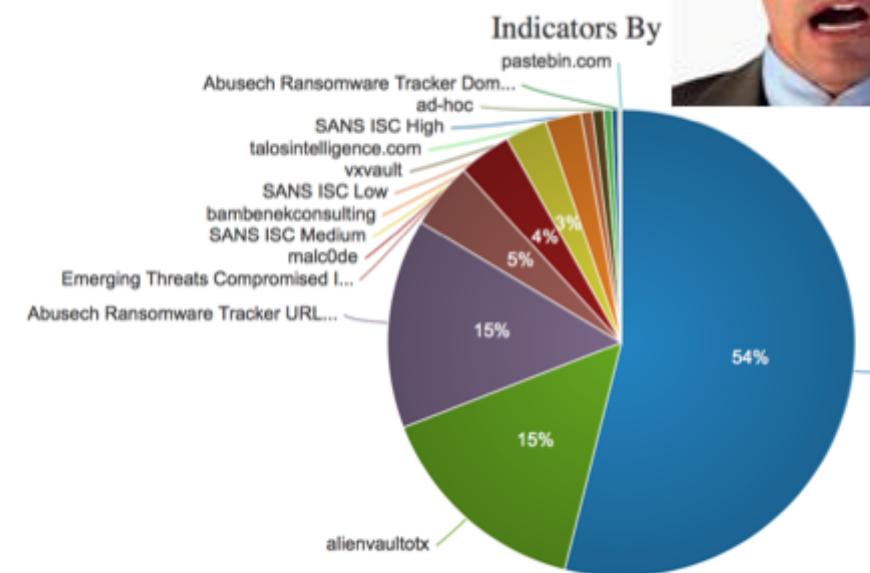
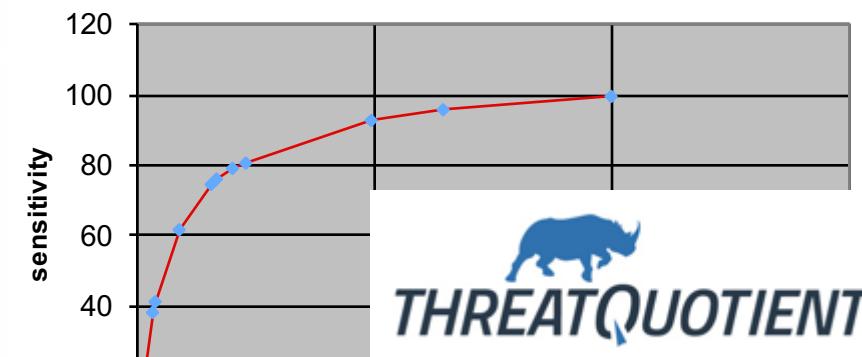
SHA256: b459f10d09a24918704c12d62af9b467e2f7d9726c6ebc632122b88990c1bec1  
File name: 5804993fb083173707dfa1f6911b9d.exe

$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{X_1 X_2} \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$

$= \sqrt{\frac{(n_1 - 1)S_{X_1}^2 + (n_2 - 1)S_{X_2}^2}{n_1 + n_2 - 2}}$



ROC Curve



Nir Yosha – Threat Intel Engineer

# Agenda

- What is Threat Intel?
- Main Challenges with Threat Intel
- Can Statistics Help?
- Confidence Level
- Conditional Probability
- Correlation
- Scoring
- Summary



# Nir Yosha



Israeli Intelligence Corps



**Technion**  
Israel Institute of  
Technology



Threat Intel Platform

Network Security



User Behavior Analytics

# Indicators of Compromise (IOCs)

The collage illustrates the following stages of a cyber attack:

- RECONNAISSANCE:** A screenshot of a Gmail inbox showing an email from PayPal about a changed email address. Below it, a list of domain names is shown: haartezenglish.strangled.net, wallanews.sytes.net, and ynet.sytes.net.
- WEAPONIZATION:** A screenshot of a digital certificate viewer showing an SSL certificate for a domain.
- DELIVERY:** A screenshot of a browser developer tools Network tab showing a request to http://maf.y.2waky.com/ with a status code of 200 OK. The remote address is highlighted with a red box.
- EXPLOITATION:** A screenshot of a terminal or log window containing Base64 strings and the word "MALWARE".
- INSTALLATION:** A screenshot of a terminal or log window containing binary code and the word "MALWARE".
- COMMAND & CONTROL:** A screenshot of a terminal or log window containing binary code and the word "MALWARE".
- ACTION ON OBJECTIVES:** A screenshot of a terminal or log window containing binary code and the word "MALWARE".

RECONNAISSANCE

WEAPONIZATION

DELIVERY

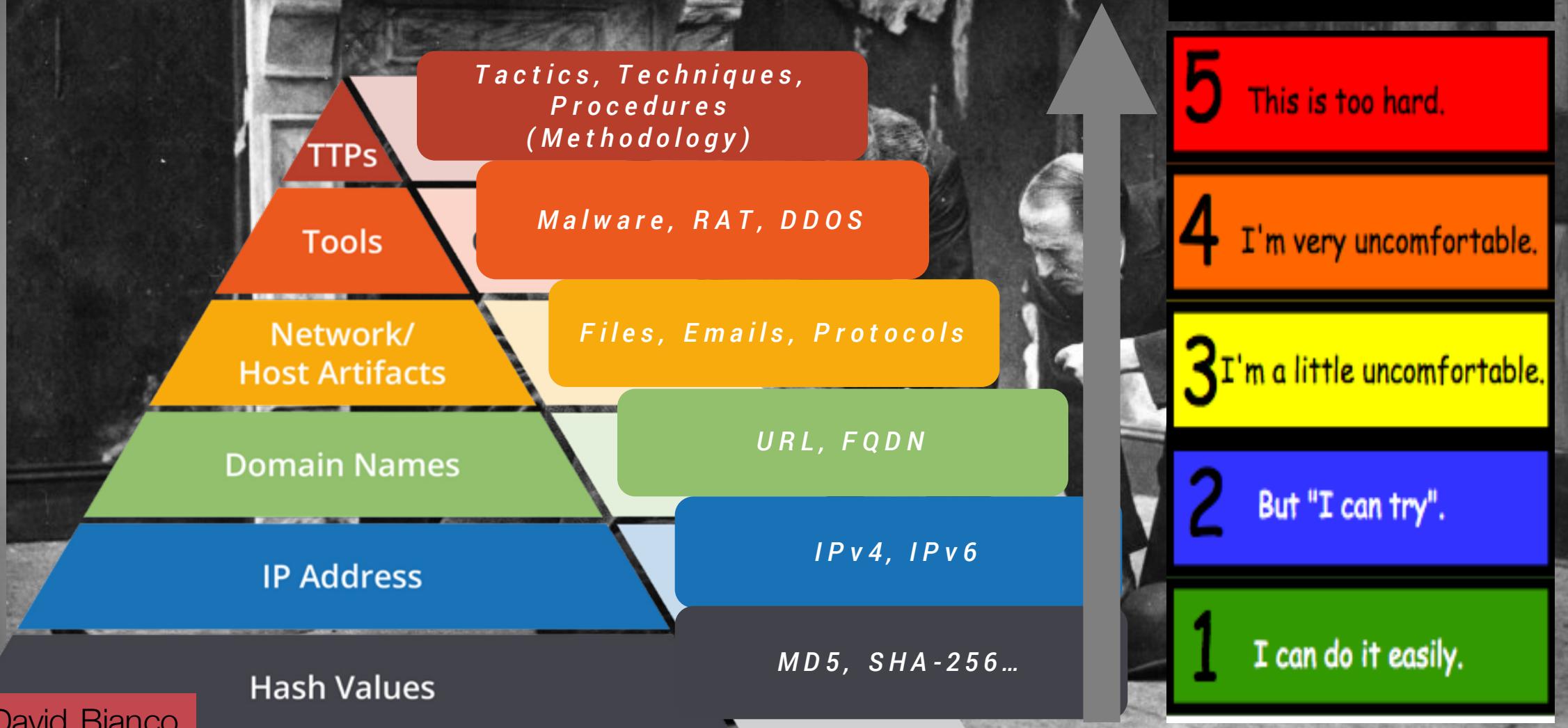
EXPLOITATION

INSTALLATION

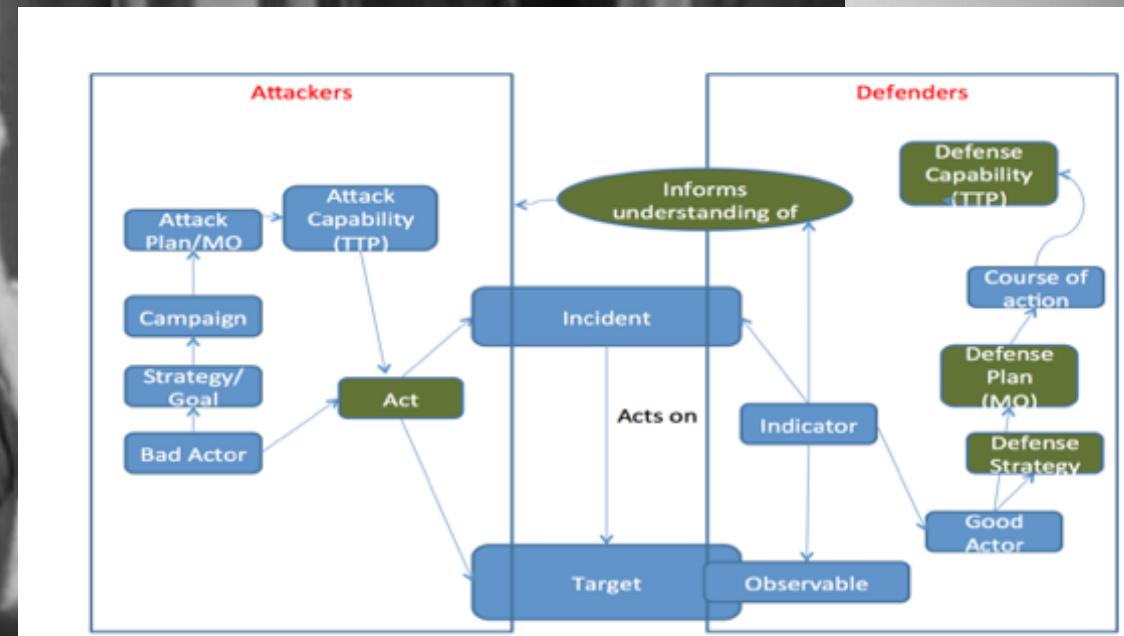
COMMAND & CONTROL

ACTION ON OBJECTIVES

# Pyramid of Pain



# Threat Intelligence



Set of data collected, assessed and applied regarding security threats, threat actors, exploits, malware, vulnerabilities, indicators of compromise.

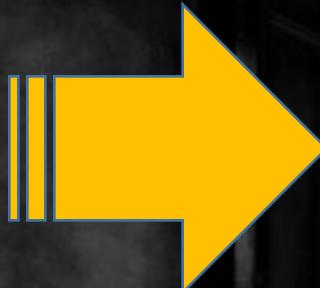
# Threat Info ≠ Threat Intel

Threat Information	Threat Intelligence
Raw, unfiltered feed	Processed information
Not enriched	Normalized & enriched
Not actionable	Actionable

Report name	Domains	Hosts	Samples
Iranian Threat Agent OilRig Delivers Digitally Signed Malware, Impersonates University of Oxford..._ClearSky Cybersecurity.pdf			
(Click name to open as page)			
Iranian Threat Agent OilRig Delivers Digitally Signed Malware, Impersonates University of Oxford..._ClearSky Cybersecurity.pdf	ns2.ayu-update.com digitallySignedMalware.mimic UniversityOfOxford..._ClearSkyCyber security.pdf	151.80.211.156 178.23.94.47 156.64.57.61 germancountyservices.com it-servicen.in www.googlecountyservices.com tecsupport.in www.cuford-careers.com www.windows-dns-resolver.org ns1.windowsupdates.me windows-dns-resolver.org 879e0f78.deckertjen.com update3.i mail.ru ns2.applicationframehost.in main-google-resolver.com te4.dn update-kern.net dockersabin.com ns2.dnresolverservice.tk kernel.ws hellfac.in ns1.windows-dns-resolver.org ns2.windowsupdates.me dnresolverservice.tk impersonationscan.in keswom.us oxford-careers.com sys-updates.com	3620c1179fe07a48fcfa942021d5fd 44 1623e01932d99fb5ab16b7e715 023995ee7952057244029e611 7e 5713ca031967d9771ac79e183e541 8 6af5d721b548cd0c95cde4842a43 c 79e035ed11901bc431a6b24aeebd 2 0bf0e95ac7d3c8943e02c2b5d98 8 0302e729e199ff14394393951e0321 405ba45259e71885794c225ca0c4a62 1 72048732509804024aa309aee200 0 0340c084195992bc289f15141339f10 033995f1c 187f078022371826603730654a0c03 2a 17903cb052536775d1445f73278f145 0 3a9f0ca803116485c4030859447411 8 2520c259802c48ce0fa80a0c8ed5 87 250fa0c49571a0bda442ca2352795

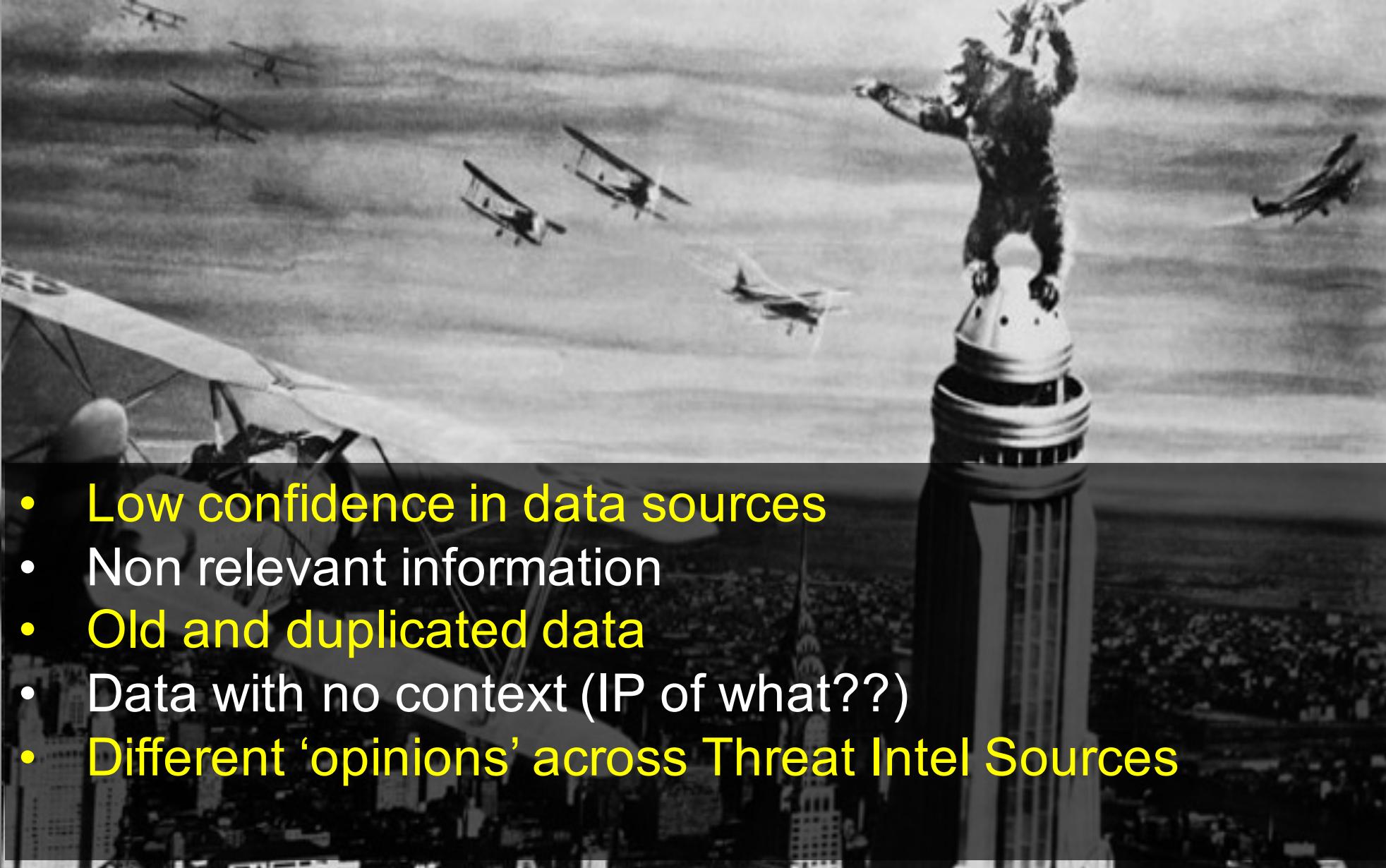
# Multiple sources of data

- OS
- Community
- Commercial
- Enrichment
- Internal



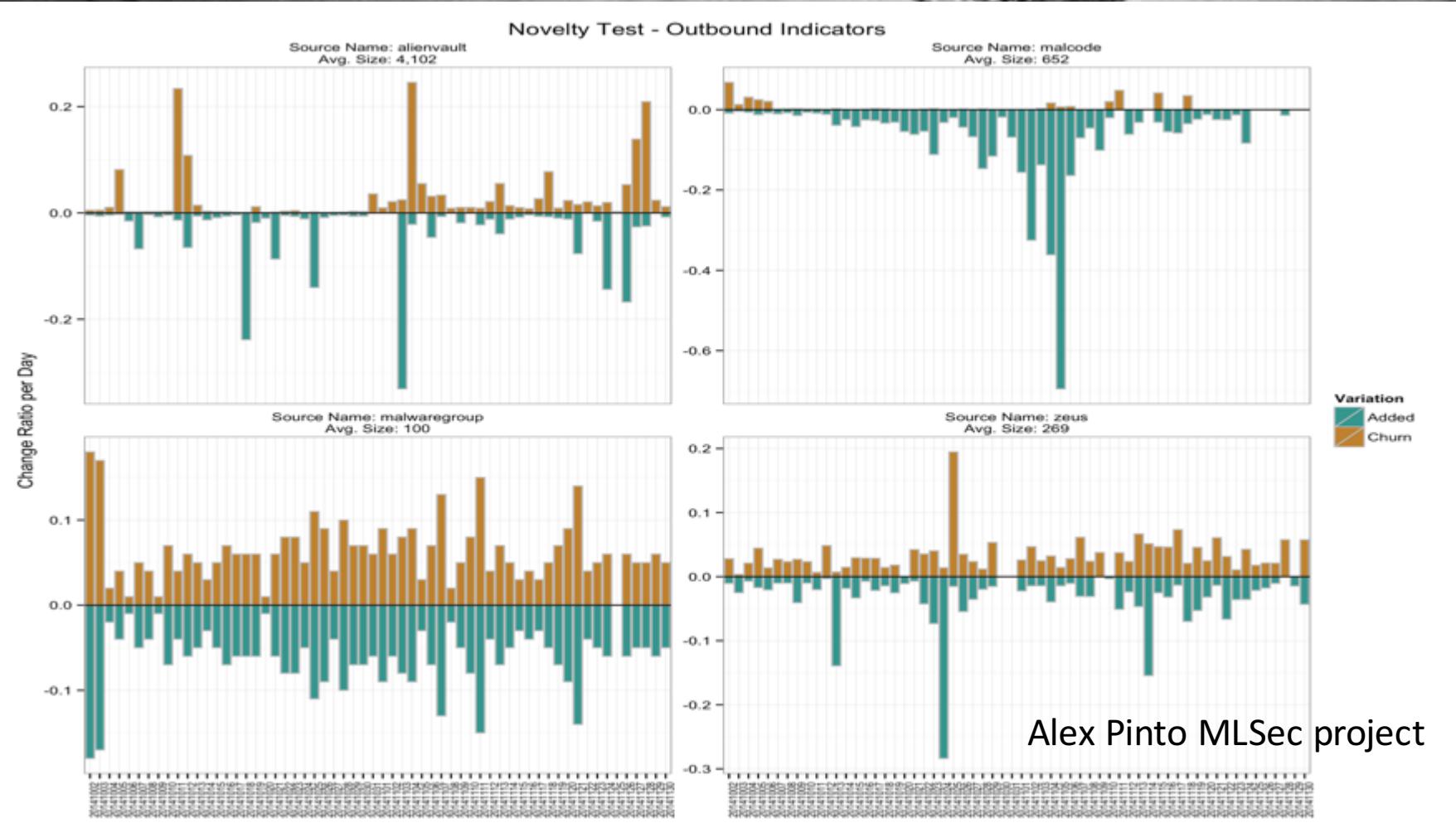
- TTP
- Victims
- Infrastructure
- IOCs (Indicators)
- Adversary (Actor)

# Main Challenges with Threat Intel

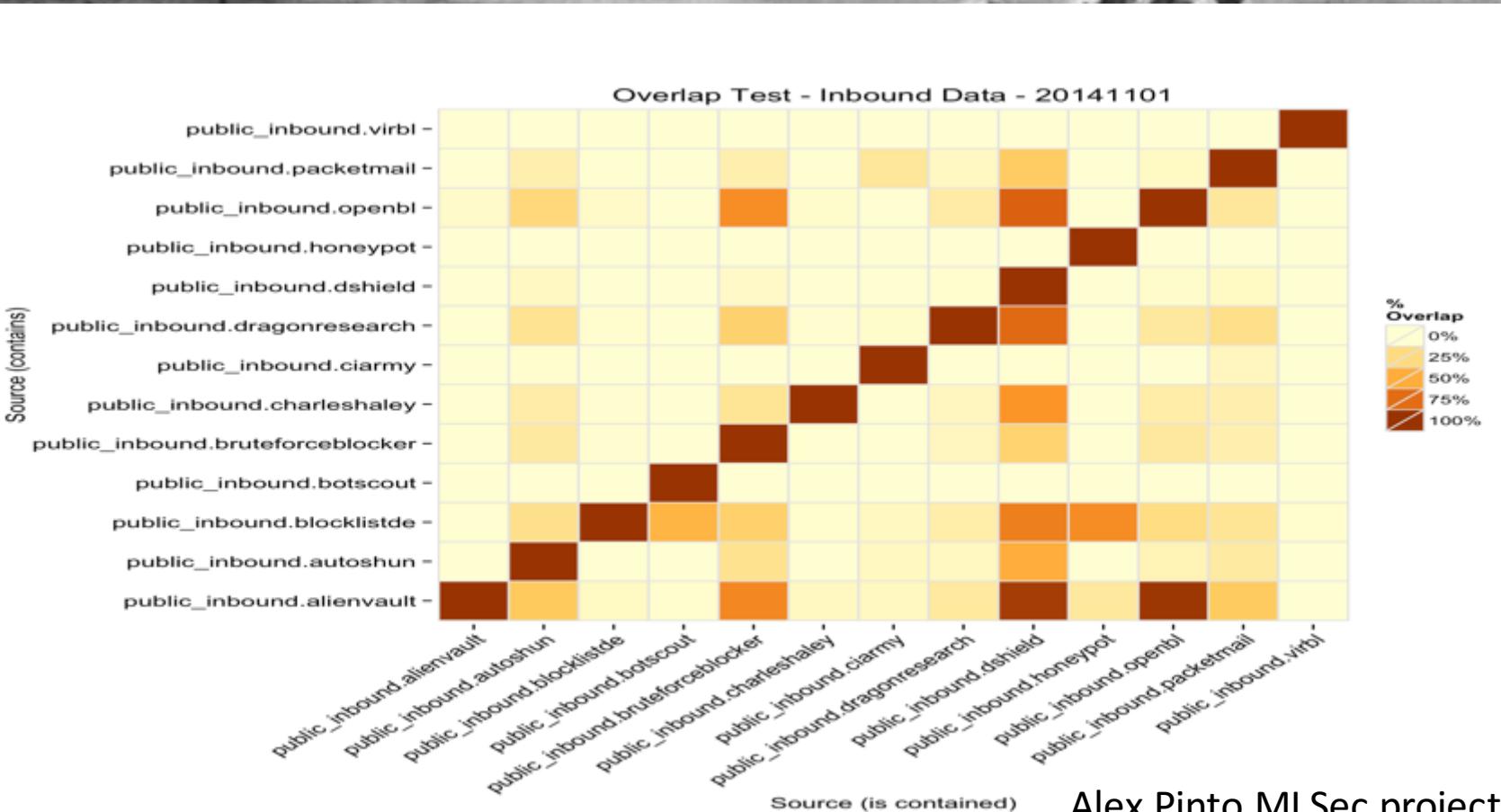


- Low confidence in data sources
- Non relevant information
- Old and duplicated data
- Data with no context (IP of what??)
- Different 'opinions' across Threat Intel Sources

# Novelty test

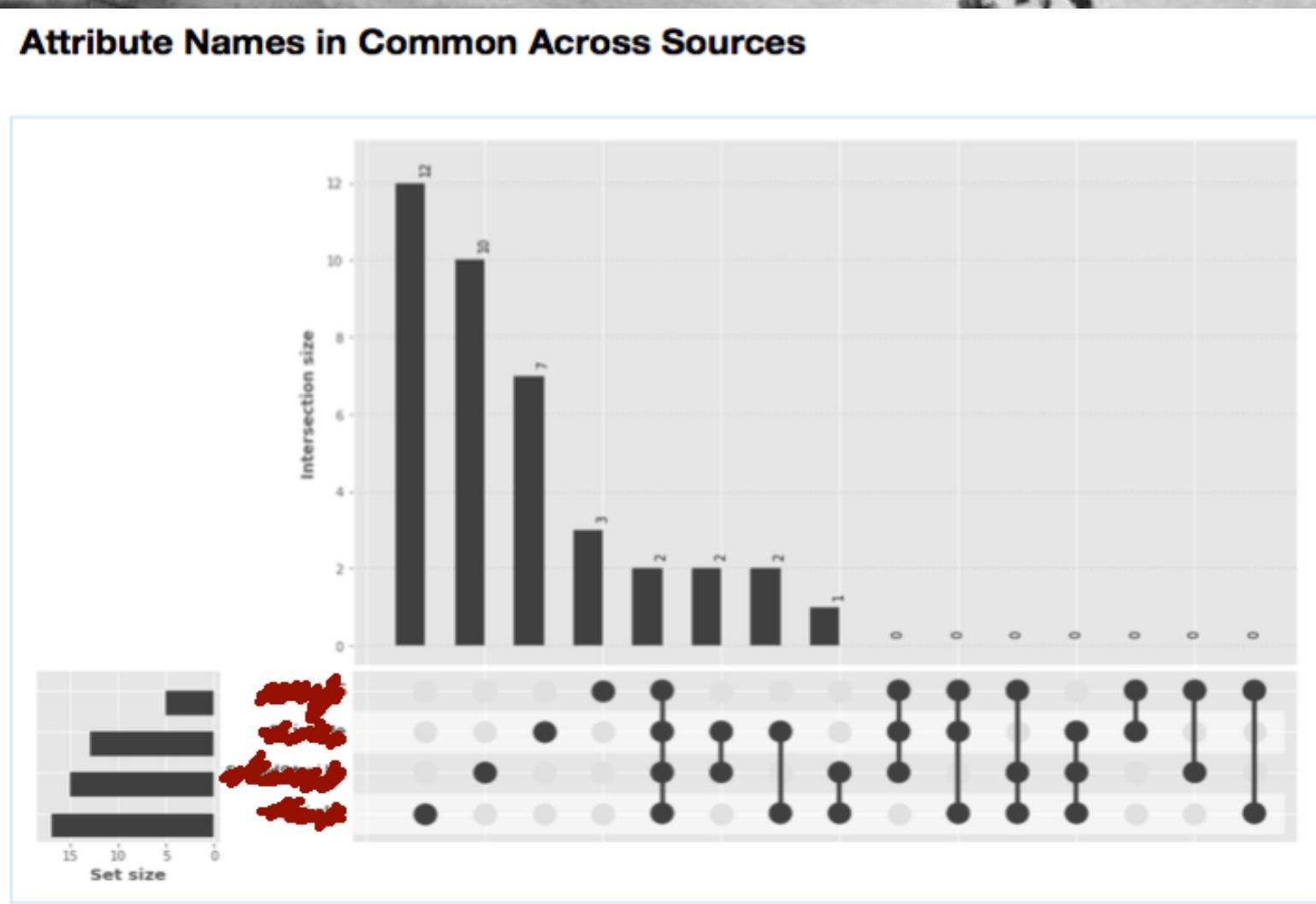


# Overlap test

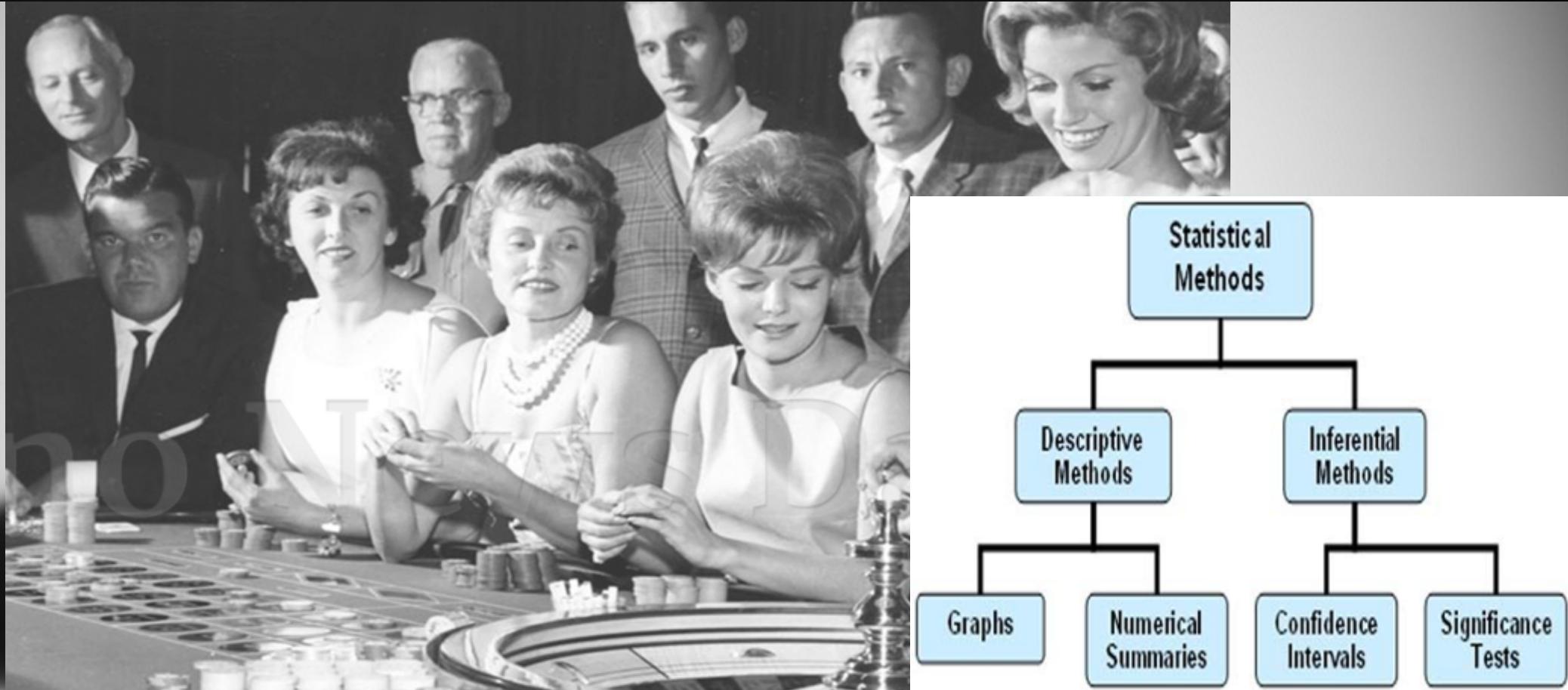


Alex Pinto MLSec project

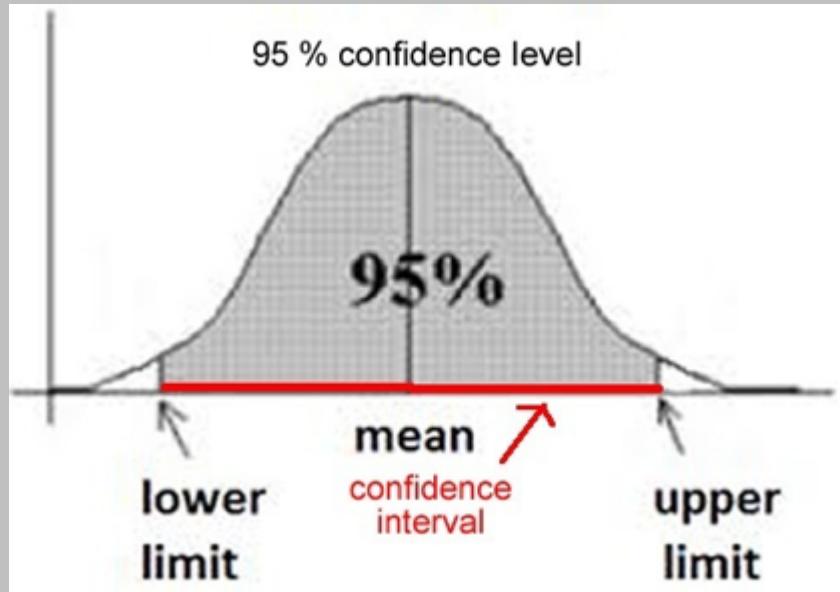
# Vendor specific terminology



# Can Statistics Help Threat Intel?



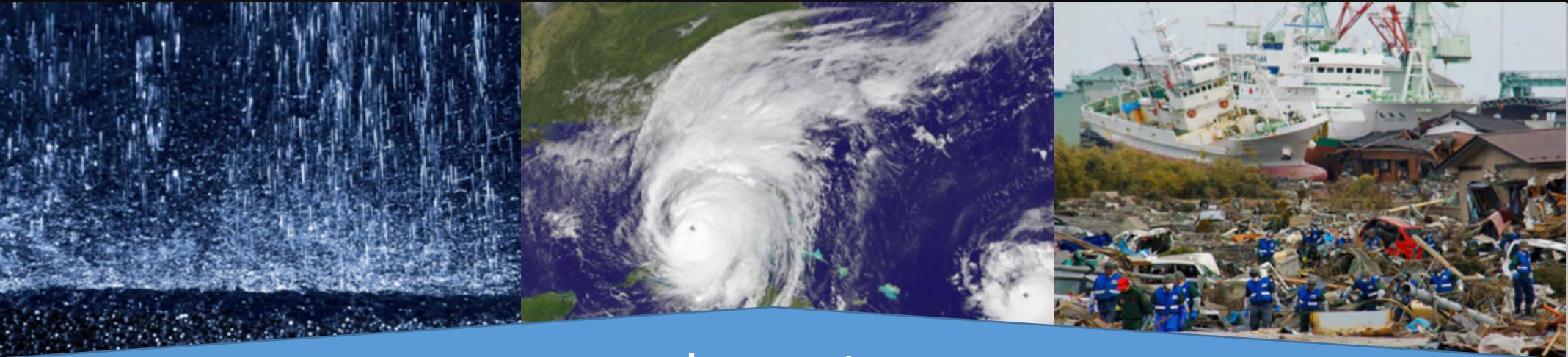
Collection, analysis, interpretation, and presentation of masses of numerical data.



# Confidence



# Confidence ≠ Impact



Impact



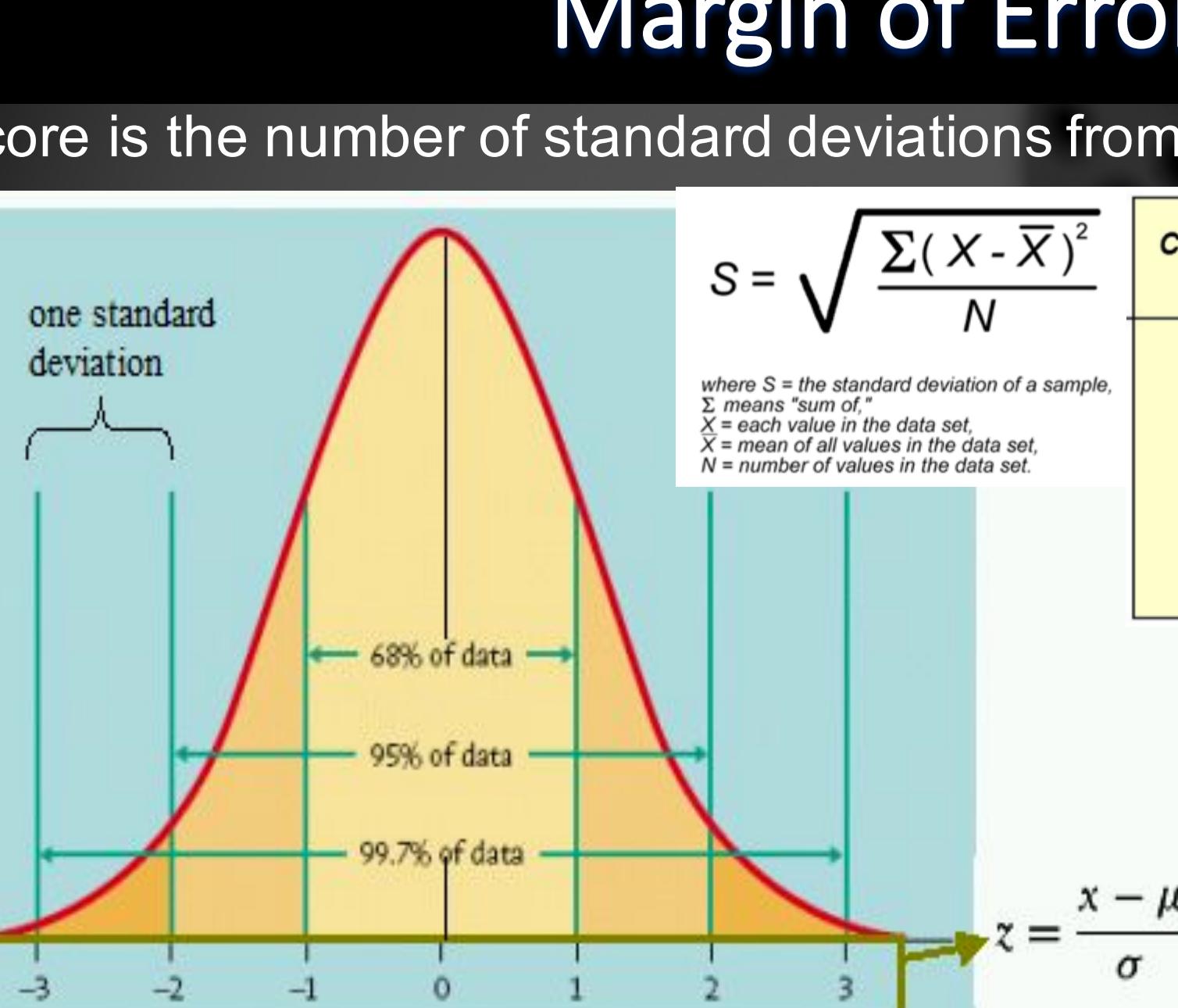
# Margin of Error

Z-score is the number of standard deviations from the mean a data point is.

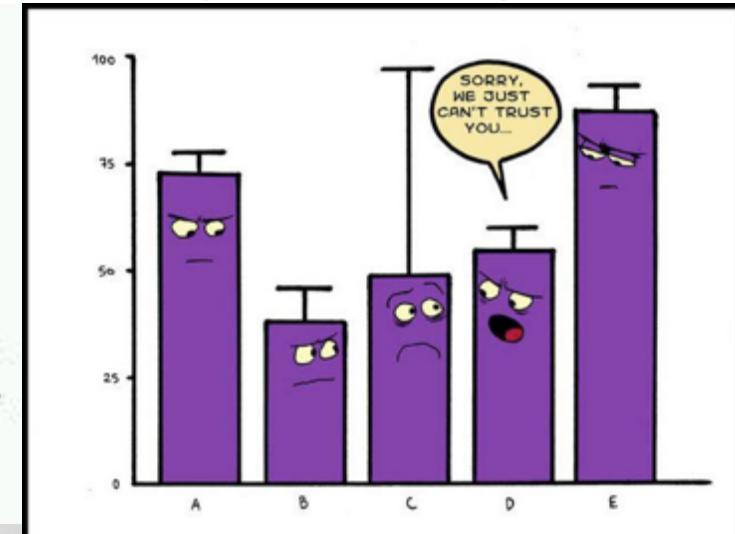
one standard deviation

$$S = \sqrt{\frac{\sum (X - \bar{X})^2}{N}}$$

where  $S$  = the standard deviation of a sample,  
 $\Sigma$  means "sum of,"  
 $X$  = each value in the data set,  
 $\bar{X}$  = mean of all values in the data set,  
 $N$  = number of values in the data set.



Confidence Level	Confidence Coefficient, $1 - \alpha$	$z$ value, $z_{\alpha/2}$
80%	.80	1.28
90%	.90	1.645
95%	.95	1.96
98%	.98	2.33
99%	.99	2.58
99.8%	.998	3.08
99.9%	.999	3.27



## Example 1

A suspicious file is running in VirusTotal  
# of AVs = 20    14/20 positive scans

What is the margin of error with a confidence level of 95%?



SHA256: b459f10d09a24918704c12d62af9b467e2f7d9726c6ebc632122b88990c1bec1

File name: 5804993f8cf083173707dfa1f6911b9d.exe

Detection ratio: 2 / 56

Analysis date: 2015-04-07 05:01:26 UTC (3 days, 11 hours ago)

Analysis File detail Additional information Comments Votes Behavioural information

Antivirus	Result	Update
AVG	Opera Software ASA	20150407
Sophos	Mal/EncPk-AAK	20150407
ALYac	✓	20150407
AVware	✓	20150407
Ad-Aware	✓	20150407
AegisLab	✓	20150407

### Margin of Error Formulas

$$ME = z \sqrt{\left( \frac{p(1-p)}{n} \right)}$$

$$ME = 1.96 \sqrt{\left( \frac{p(1-p)}{n} \right)}$$

for 95% confidence interval



# Size does matter!

## Solution 1

$$P(\text{Malware}) = 14/20 = 0.7$$

Z Score of 95% = 1.96 (normal dis.)

# of AVs = 20

$$ME = 1.96 \sqrt{(0.3 \times 0.7 / 20)} = 0.2 = 20\%$$

50% < P(Malware) < 90%

## What if we increase the sample size?

$$P = 140/200 = 0.7$$

Z Score of 95% = 1.96 (normal dis.)

# of AVs = 200

$$ME = 1.96 \sqrt{(0.3 \times 0.7 / 200)} = 0.063 = 6.3\%$$

63.7% < P(Malware) < 76.3%

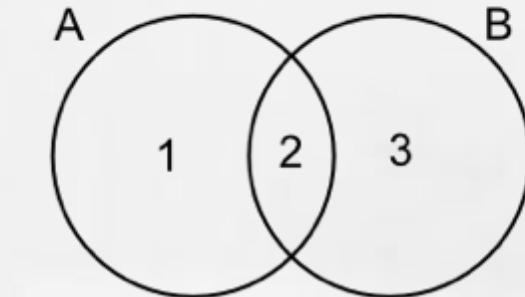
The screenshot shows the VirusTotal analysis page for a file with SHA256: b459f10d09a24918704c12d62af9b467e2f7d9726c6ebc632122b88990c1bec1. The detection ratio is 2/56, with the analysis date being 2015-04-07 05:01:26 UTC (3 days, 11 hours ago). The interface includes tabs for Analysis, File detail, Additional information, Comments, Votes, and Behavioural information. Below these tabs is a table of antivirus results:

Antivirus	Result	Update
AVG	Opera Software ASA	20150407
Sophos	Mal/EncPk-AAK	20150407
ALYac	✓	20150407
AVware	✓	20150407
Ad-Aware	✓	20150407
AegisLab	✓	20150407

$$ME = z \sqrt{\left( \frac{p(1-p)}{n} \right)}$$

$$ME = 1.96 \sqrt{\left( \frac{p(1-p)}{n} \right)}$$

for 95% confidence interval

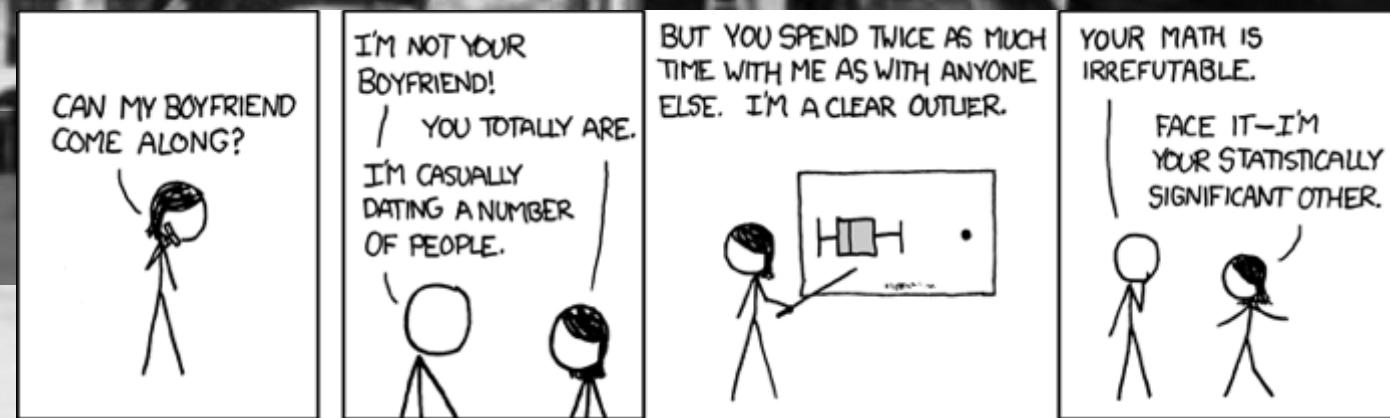


4

$P(A|B)$  is A given B

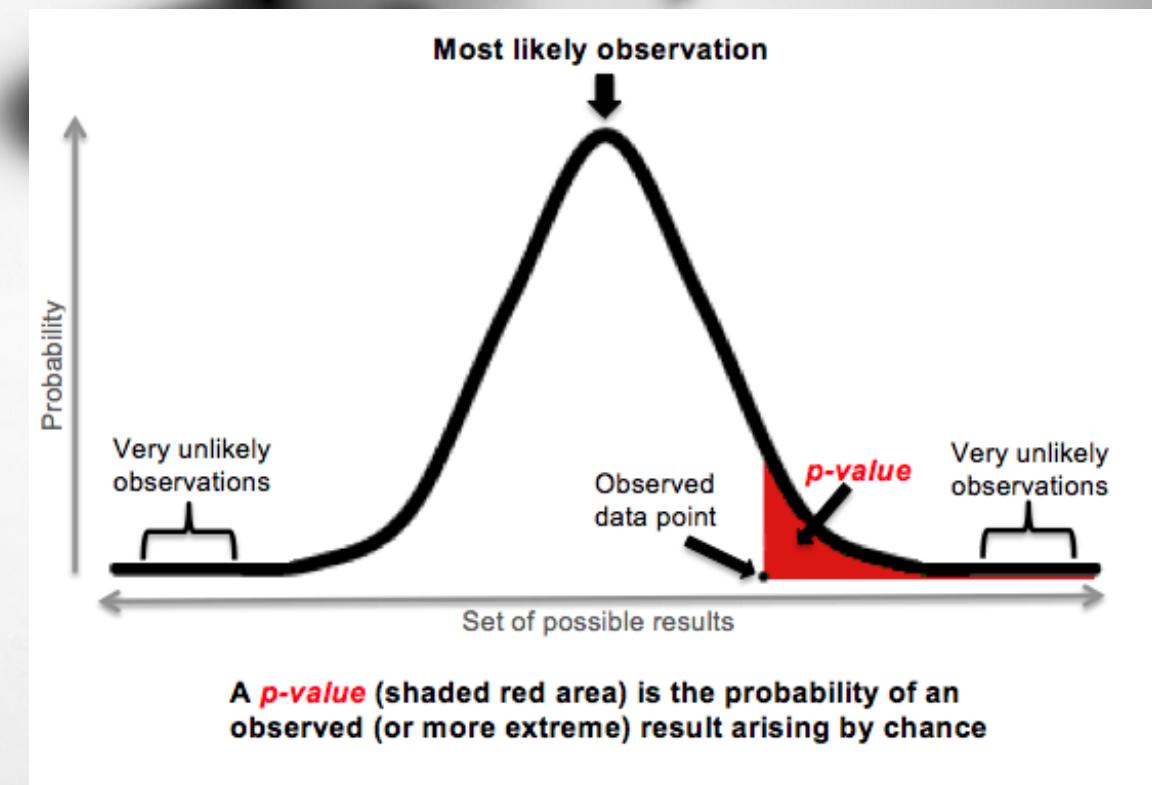
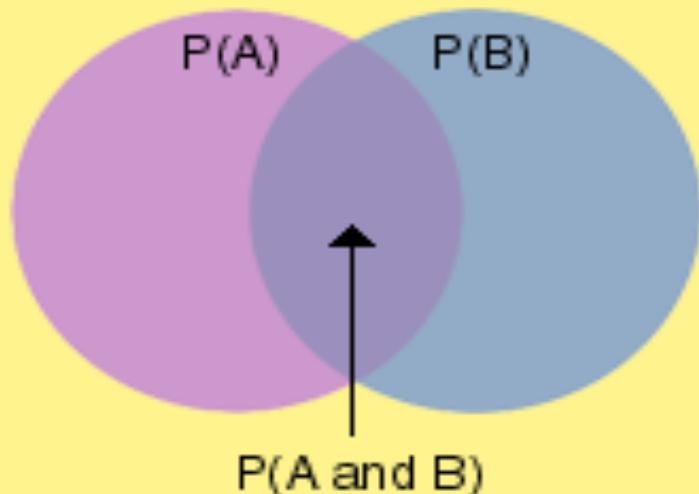
$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{2}{2+1} = \frac{2}{3}$$

# Probability



# Conditional Probability

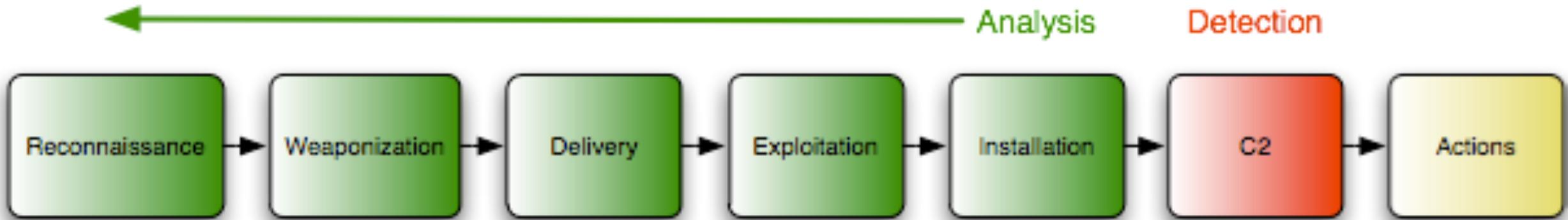
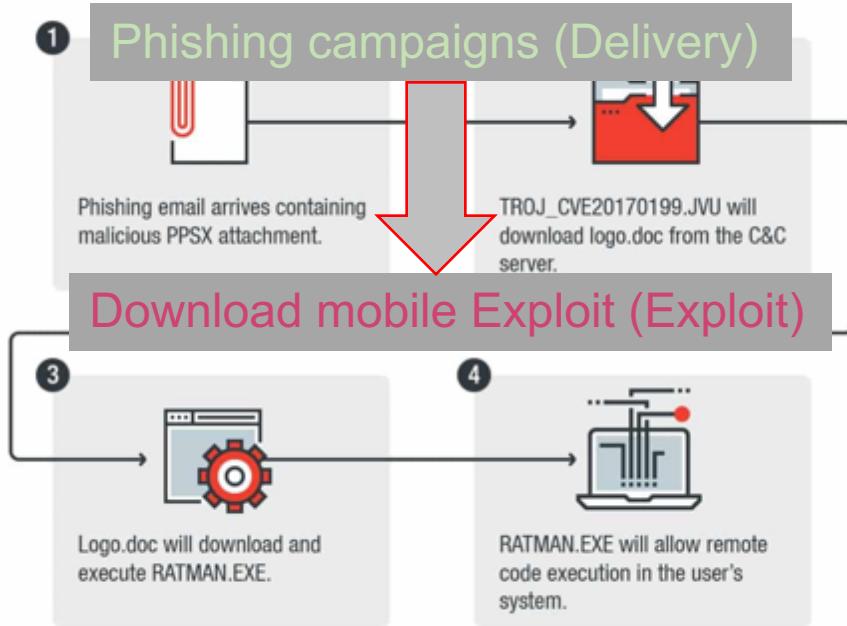
$$P(A|B) = \frac{P(A \text{ and } B)}{P(B)}$$



## Example 2

Out of 100 incidents under investigation 40 are a result of **phishing campaigns (Delivery)** and 30 are **mobile Exploit (Exploit)**.

What is the probability of a mobile malware giving that its delivery method is phishing campaign?



# Improve accuracy with time

## Solution 2

$$P(A \text{ and } B) = P(A) \times P(B) \Rightarrow$$

$$0.3 \times 0.4 = 0.12 = 12\%$$

$$P(A|B) = P(A \text{ AND } B)/P(B) \Rightarrow$$

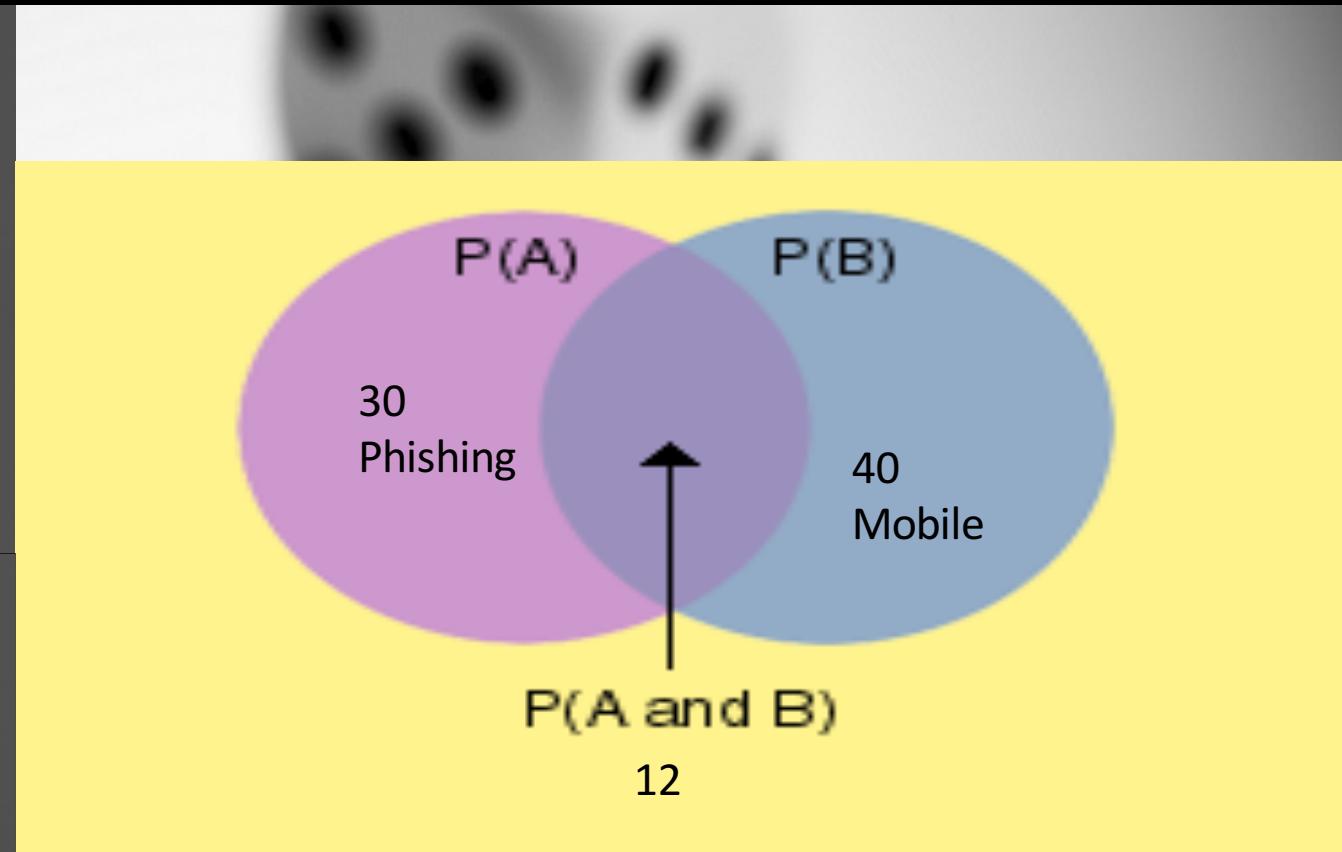
$$12/40 = 30\%$$

What if a sample data includes the following?

$$P(A \text{ AND } B) = 20\% \text{ (known)}$$

$$P(A|B) = P(A \text{ AND } B)/P(B) \Rightarrow$$

$$20\%/40\% = 50\%$$

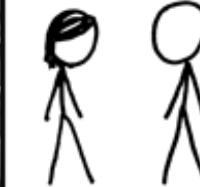


$$P(A|B) = \frac{P(A \text{ and } B)}{P(B)}$$



# Correlation

I USED TO THINK  
CORRELATION IMPLIED  
CAUSATION.

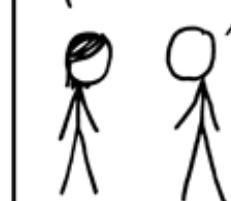


THEN I TOOK A  
STATISTICS CLASS.  
NOW I DON'T.



SOUNDS LIKE THE  
CLASS HELPED.

WELL, MAYBE.

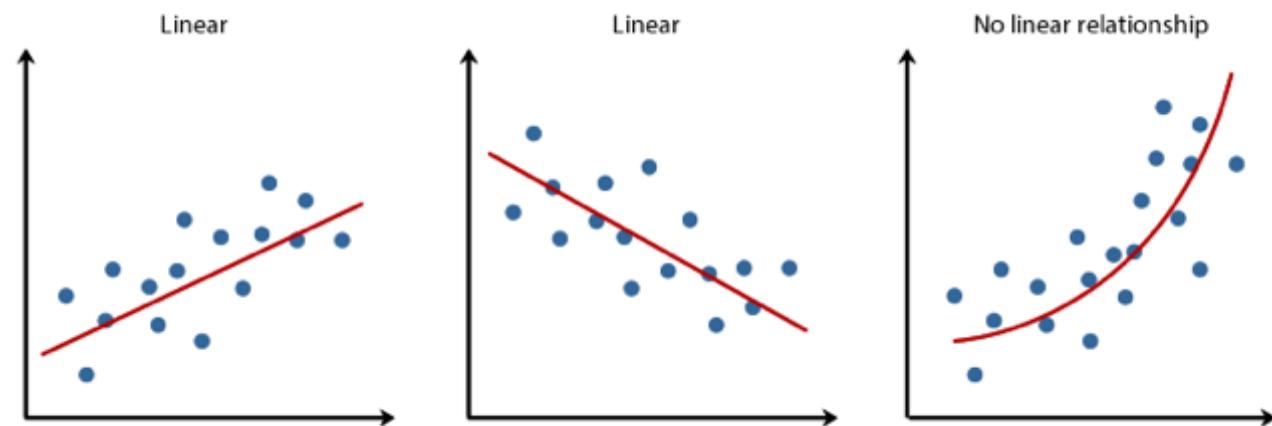


# Regression and Forecasting

Regression Statistics	
Multiple R	0.85
R Square	0.73
Adjusted R Square	0.68
Standard Error	7.83
Observations	15.00

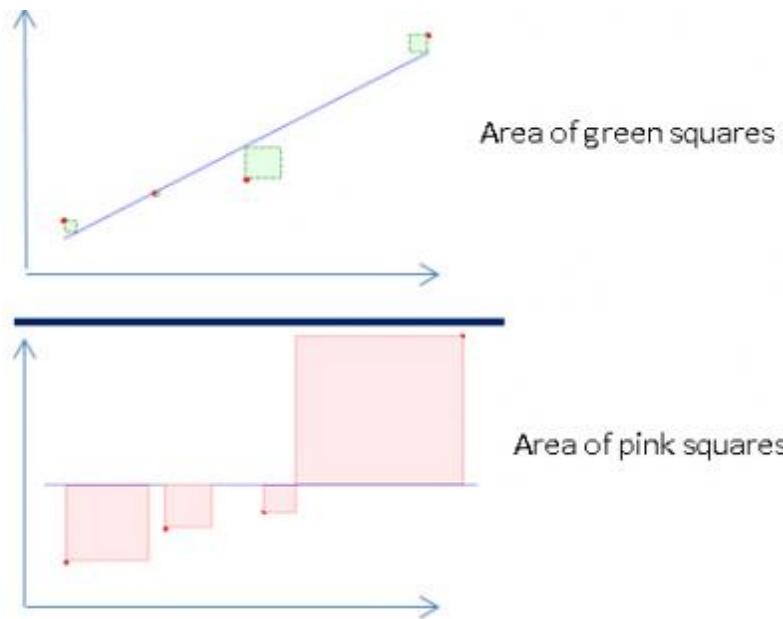
ANOVA					
	df	SS	MS	F	Significance F
Regression	2.00	1984.27	992.13	16.19	0.00
Residual	12.00	735.33	61.28		
Total	14.00	2719.60			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%
Intercept	419.95	37.40	11.23	0.00	338.46	501.45
Price	-42.80	8.30	-5.15	0.00	-60.89	-24.70
Cprice	4.39	9.74	0.45	0.66	-16.82	25.60



Copyright 2014. Laerd Statistics.

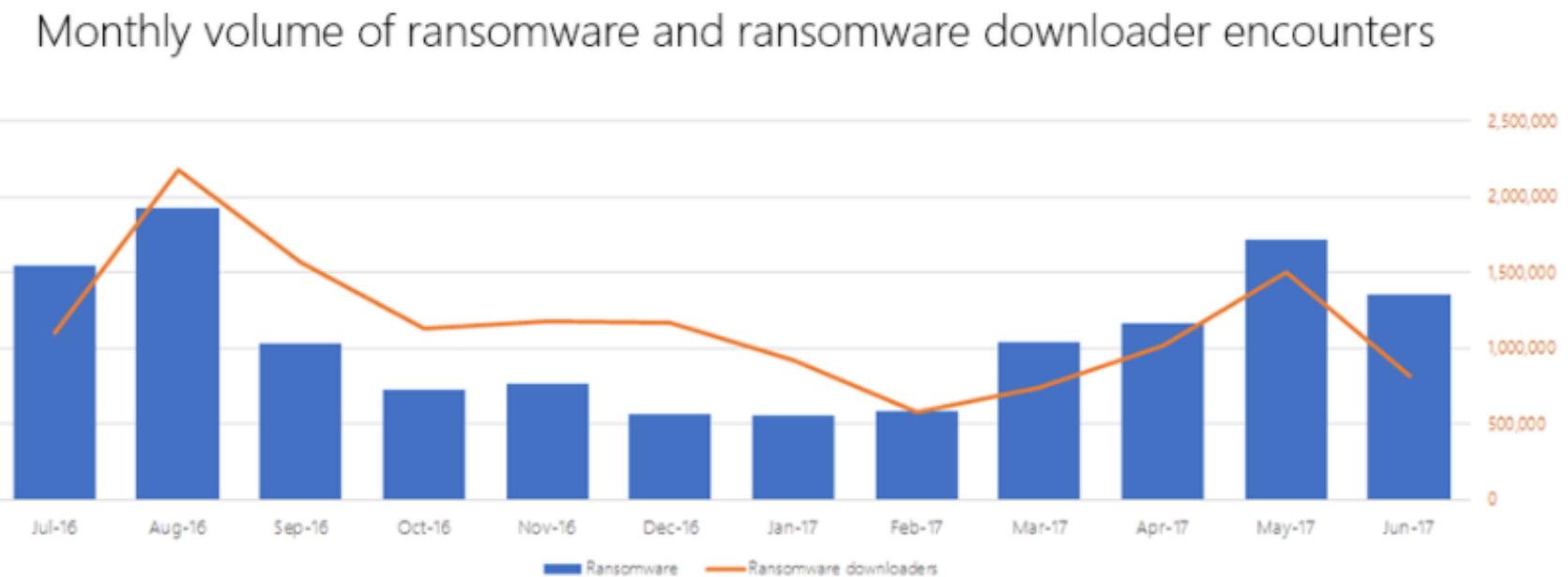
$$R^2 = 1 - \frac{\text{Area of green squares}}{\text{Total Area}}$$



## Example 3

Predict the number of incidents in July based on the past 6 month. Use moving average, smoothing, seasoning and regression.

Month	Incidents
Jan	1,325
Feb	1,353
Mar	1,305
Apr	1,275
May	1,210
Jun	1,195
Jul	?



## Solution 3

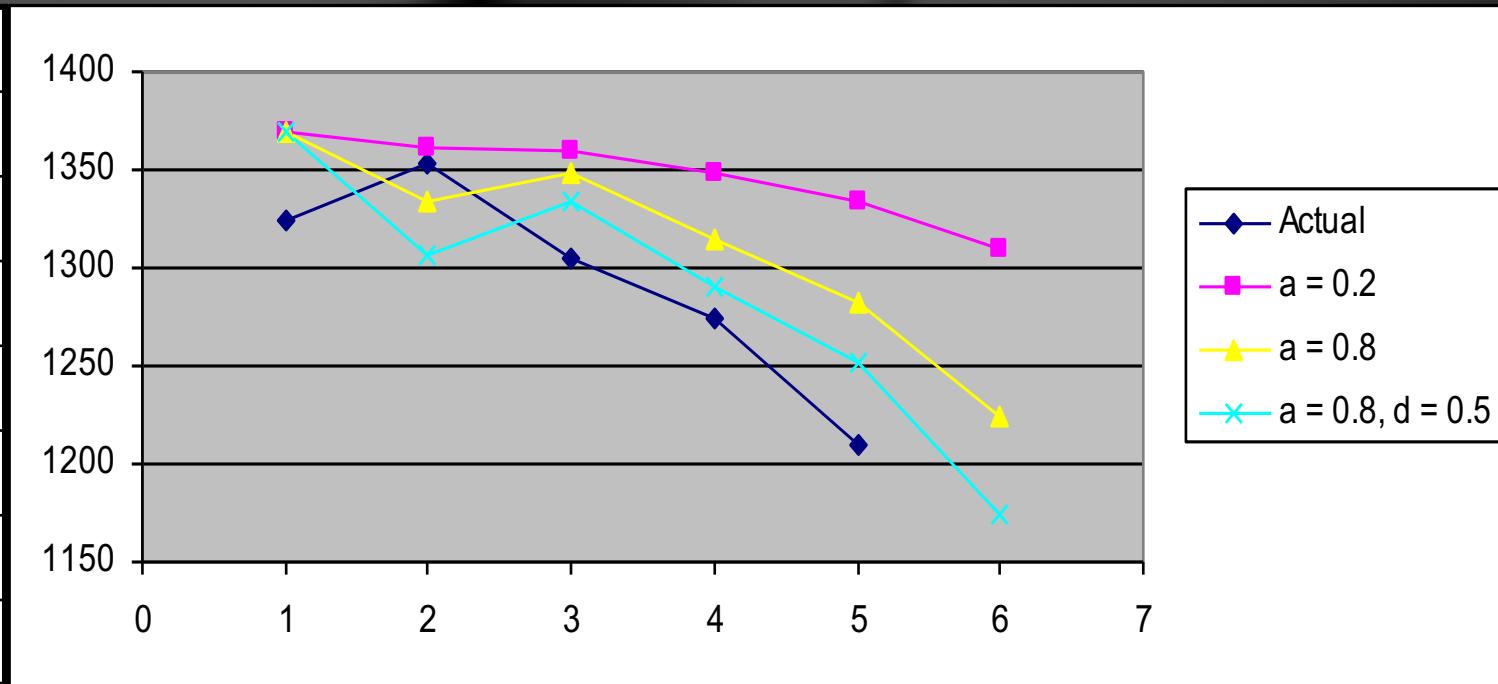
include trend analysis in exponential smoothing.

$$T_t = T_{t-1} + \delta(F_t - FIT_{t-1})$$

FIT: Forecast including trend  
 $\delta$ : Trend smoothing constant

$$FIT_t = F_t + T_t \quad F_t = FIT_{t-1} + \alpha(A_{t-1} - FIT_{t-1})$$

Month	Incidents
Jan	1,325
Feb	1,353
Mar	1,305
Apr	1,275
May	1,210
Jun	1,195
Jul	?



Regular exponential smoothing will always lag behind the trend.

$F$  is the forecast without trend and  $T$  is the trend component

# Regression

Influence Score based on:

Indicator Type	Indicator Source	Attributes	Adversary Relationship
----------------	------------------	------------	------------------------

Regression Statistics	
Multiple R	0.645712437
R Square	0.416944552
Adjusted R Square	0.412105917
Standard Error	1.120229993
Observations	487

ANOVA					
	df	SS	MS	F	Significance F
Regression	4	432.5435862	108.135897	86.16988086	3.48716E-55
Residual	482	604.8691448	1.25491524		
Total	486	1037.412731			

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	6.32248956	0.280393065	22.548666	2.10687E-77	5.771545821	6.8734333	5.77154582	6.8734333
attrb_total	-0.583131574	0.036721225	-15.87996	5.59042E-46	-0.655285031	-0.5109781	-0.655285	-0.5109781
related_sources	3.140066752	0.250619197	12.5292347	2.25903E-31	2.647625623	3.63250788	2.64762562	3.63250788
related_adv	-0.304786318	0.053510731	-5.695798	2.13846E-08	-0.409929441	-0.1996432	-0.4099294	-0.1996432
related_ind	0.005577054	0.00177405	3.1436849	0.001771244	0.002091227	0.00906288	0.00209123	0.00906288

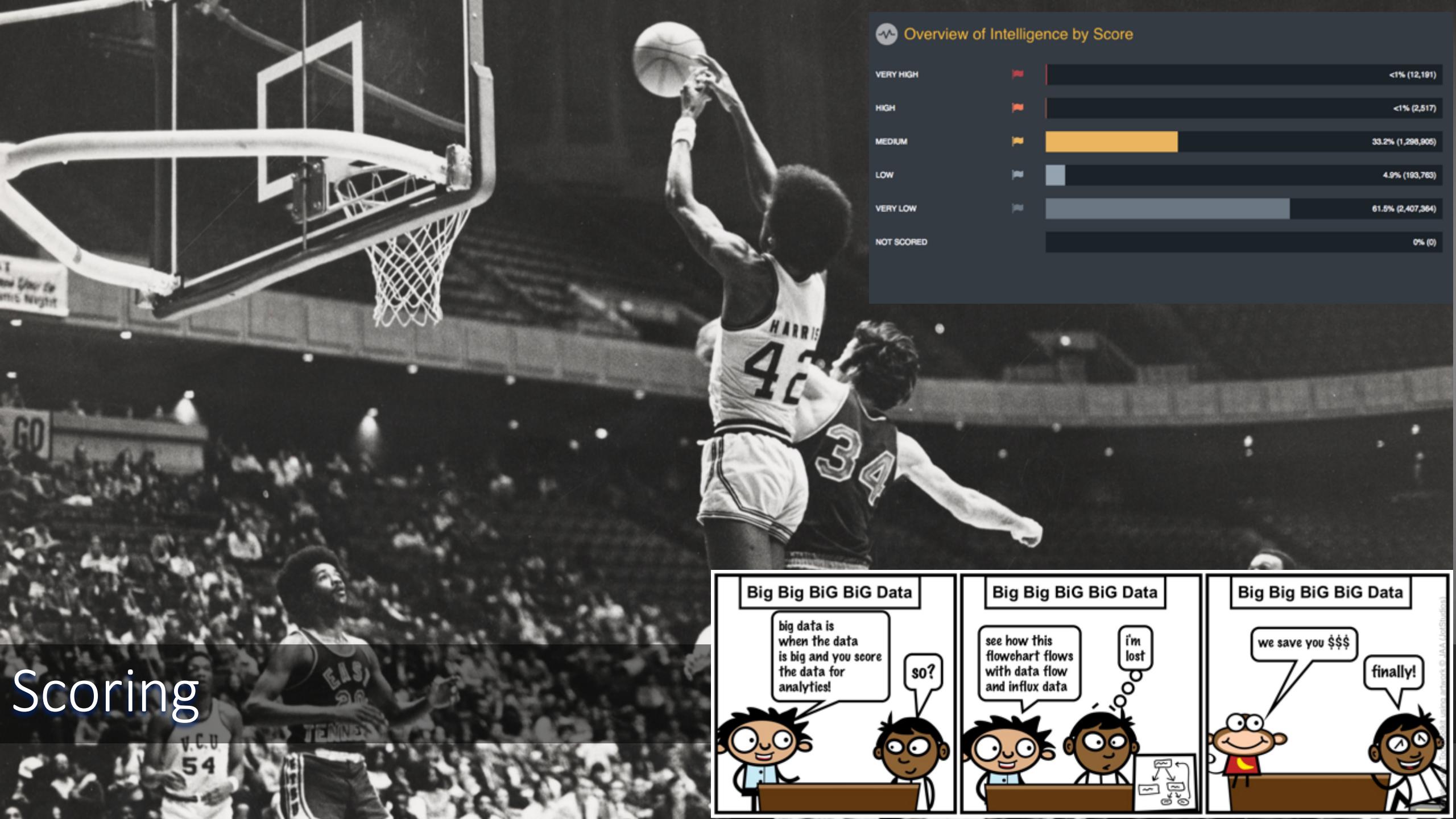
Eg. Estimated score for indicator with 6 attributes, 2 sources and 20 related indicators ->  $Y = 6.32 - 0.58X6 + 3.14X2 + 0.005X20 \Rightarrow$

Observations.

Confidence intervals for the regression line.

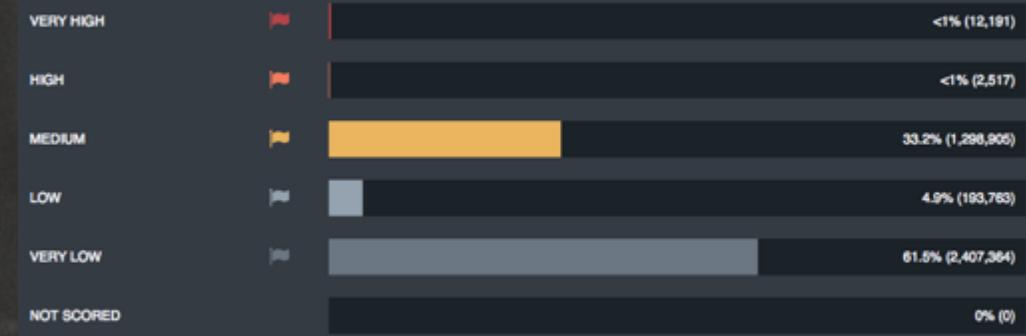
Fitted regression line.

Prediction intervals.



# Scoring

## Overview of Intelligence by Score



### Big Big BiG BiG Data

big data is  
when the data  
is big and you score  
the data for  
analytics!

so?

### Big Big BiG BiG Data

see how this  
flowchart flows  
with data flow  
and influx data

i'm  
lost

### Big Big BiG BiG Data

we save you \$\$\$

finally!

# Scoring Categories

Threat agent factors				Vulnerability factors				
Skill level	Motive	Opportunity	Size		Ease of discovery	Ease of exploit	Awareness	Intrusion detection
5	2	7	1		3	6	9	2
Overall likelihood=4.375 (MEDIUM)								

Technical Impact				Business Impact				
Loss of confidentiality	Loss of integrity	Loss of availability	Loss of accountability		Financial damage	Reputation damage	Non-compliance	Privacy violation
9	7	5	8		1	2	1	5
Overall technical impact=7.25 (HIGH)					Overall business impact=2.25 (LOW)			

# Scoring IOCs

The screenshot shows a web-based interface for managing Indicators of Compromise (IOCs). At the top, there are status indicators: "Being Watched" (green), "Score: 10 - Very High" (red), and "Status: Active". A red box highlights a tooltip for "Attributes". The tooltip lists the following attributes with their values and sources:

Attribute	Value	Source	Last Updated
Attack Phase	Delivery	abuse.ch SSLBL (Extended)	02/07/17 04:59pm
Malware Family	Locky	abuse.ch SSLBL (Extended)	02/07/17 04:58pm
Port	443	abuse.ch SSLBL (Extended)	01/10/17 05:46pm
Malware Family	Quakbot C&C	abuse.ch SSLBL (Extended)	01/10/17 05:46pm

Below the attributes, there is a section for "Sources" which lists "abuse.ch SSLBL (Extended)" with a timestamp of "01/10/17 05:46pm".

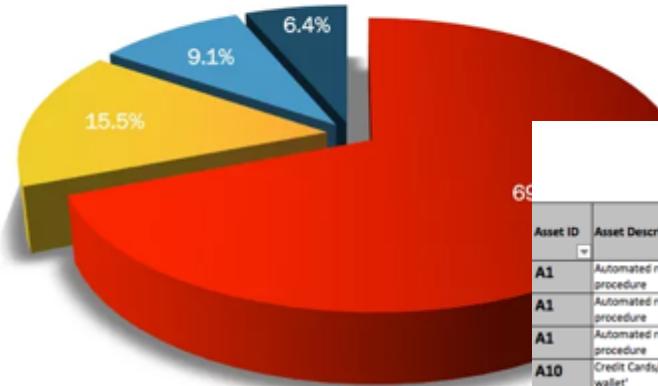
# IOC Classification

Classifiers	Example	Why
Indicator Type	IP address, URL, Domain, File	Support different security devices
Attributes	Language, Country, Malware family	Used by most vendors
Adversary	Industry, Government, PII, PHI, PCI	Relevant adversaries to your vertical / endpoints
Incidents or Events	Age, Owner, ASN, Day of Week, User ID	Connect the dots across the kill chain
Relevance	CVE, OS, User, Brand	Relevant to your environment

# Threat Intel Final Scoring

Motivations Behind Attacks

June 2016



Cyber Crime

Identification & Assessment of Risks  
(Asset-Threat-Vulnerability)

Asset ID	Asset Description	Asset Value	Vulnerability Description	A/V Value	Threats		Threat Value	A/V/T ID	Risk Value (Final)
					T	V			
A1	Automated reservation, check-in and boarding procedure	4	V38. Lack of common or harmonised legislation in EU Member States	4	T32. Profiling	4	A1.V38.T32	9	
A1	Automated reservation, check-in and boarding procedure	4	V43. Lack of respect to the data conservation principle	3	T12. Non-compliance with data protection legislation	4	A1.V41.T12	8	
A1	Automated reservation, check-in and boarding procedure	4	V42. Lack of respect to the rights of the data subject (such as the right for rectification, blocking or deletion of data).	4	T12. Non-compliance with data protection legislation	4	A1.V42.T12	9	
A10	Credit Cards/Debit cards/Payment cards/e-wallet*	4	V21. Inappropriate / inadequate identity management	4	T13. Denial of service attack / Flood / Buffer overflow	3	A10.V21.T1	9	
A10	Credit Cards/Debit cards/Payment cards/e-wallet*	4	V21. Inappropriate / inadequate identity management	4	T2. Spoofing of credentials / bypass authentication	5	A10.V21.T2	10	
A10	Credit Cards/Debit cards/Payment cards/e-wallet*	4	V21. Inappropriate / inadequate identity management	4	T4. Traffic analysis / scan / probe	3	A10.V21.T4	8	
A10	Credit Cards/Debit cards/Payment cards/e-wallet*	4	V39. Insufficient protection of wireless networks and communication (weak or no encryption etc.)	4	T24. Worms, viruses & malicious code	3	A10.V39.T24	9	
A10	Credit Cards/Debit cards/Payment cards/e-wallet*	4	V39. Insufficient protection of wireless networks and communication (weak or no encryption etc.)	4	T29. MANET/Adhoc network routing attack	2	A10.V39.T29	7	
A10	Credit Cards/Debit cards/Payment cards/e-wallet*	4	V39. Insufficient protection of wireless networks and communication (weak or no encryption etc.)	4	T2. Spoofing of credentials / bypass authentication	5	A10.V39.T2	10	
A11	Other RFID cards	3	V21. Inappropriate / inadequate identity management	4	T13. Denial of service attack / Flood / Buffer overflow	3	A11.V21.T1	8	
A11	Other RFID cards	3	V39. Insufficient protection of wireless networks and communication (weak or no encryption etc.)	3	T20. Fake / rogue RFID readers / scanning of RFID reader and / or tag	3	A11.V39.T20	7	
A11	Other RFID cards	3	V39. Insufficient protection of wireless networks and communication (weak or no encryption etc.)	3	T24. Worms, viruses & malicious code	3	A11.V39.T24	7	
A11	Other RFID cards	3	V39. Insufficient protection of wireless networks and communication (weak or no encryption etc.)	3	T29. MANET/Adhoc network routing attack	2	A11.V39.T29	5	
A11	Other RFID cards	3	V39. Insufficient protection of wireless networks and communication (weak or no encryption etc.)	3	T2. Spoofing of credentials / bypass authentication	5	A11.V39.T2	8	
A12	Scanners & detectors	3	V1. Inappropriate design of procedures - includes: lack of accountability, high complexity of procedures, assigning extensive responsibilities to end-users (in critical parts of the procedures) etc.	2	T6. Social engineering attack	4	A12.V1.T6	6	

## Equation group victims

Legend: Finance (blue), Government (red), Research institution (green), University (orange), Energy / Infrastructure (purple), Military (pink), Aerospace (grey)

### High infection rate

Iran, Russia Federation, Pakistan, Afghanistan, India, China, Syria, Mali

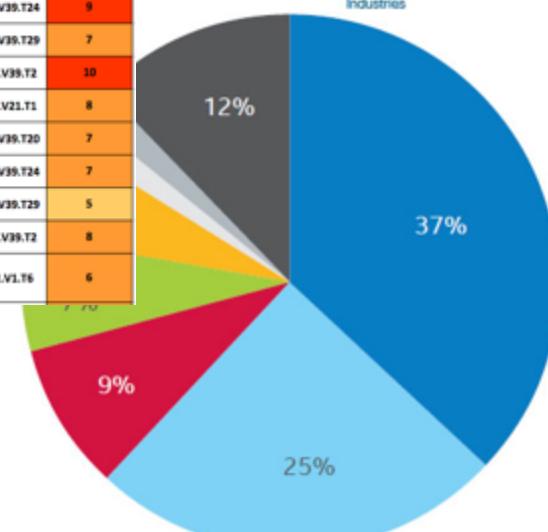
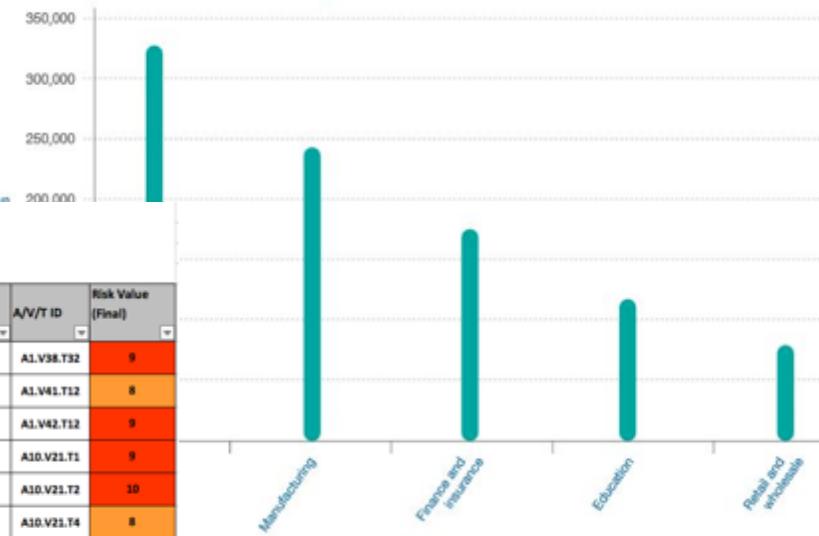
### Medium-level infection rate

Lebanon, Yemen, United Arab Emirates, Algeria, Kenya, United Kingdom, Libya, Mexico, Qatar, Egypt



Legend: Morocco, Malaysia, Kazakhstan, Iraq, Brazil, Uganda, Switzerland, Singapore, Philippines, Peru, France, Ecuador, Belgium, Bahrain

Top industries attacked



# Score Evaluation

TO ESTIMATE CONFIDENCE INTERVALS  
FOR SENSITIVITY, SPECIFICITY AND TWO-  
LEVEL LIKELIHOOD RATIOS:

	Found Relevant	Found Not Relevant
Above Cut-off Score	200	200
Below Cut-off Score	20	450

Enter the required confidence interval: 90

RESULT:

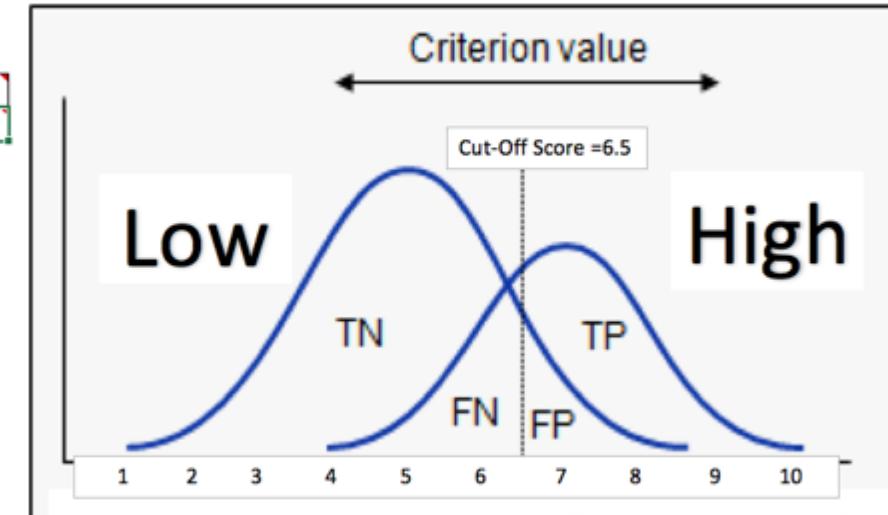
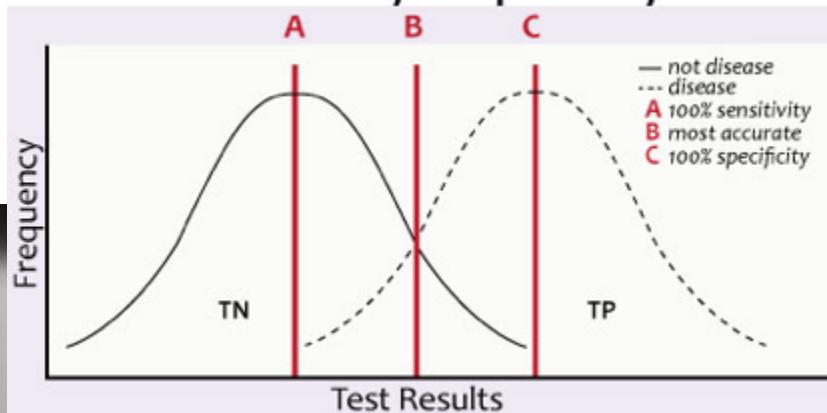
Sensitivity: 0.9091 CI: 0.872 to 0.9362

Specificity: 0.6923 CI: 0.6618 to 0.7212

Positive likelihood ratio: 2.955 CI: 2.666 to 3.275

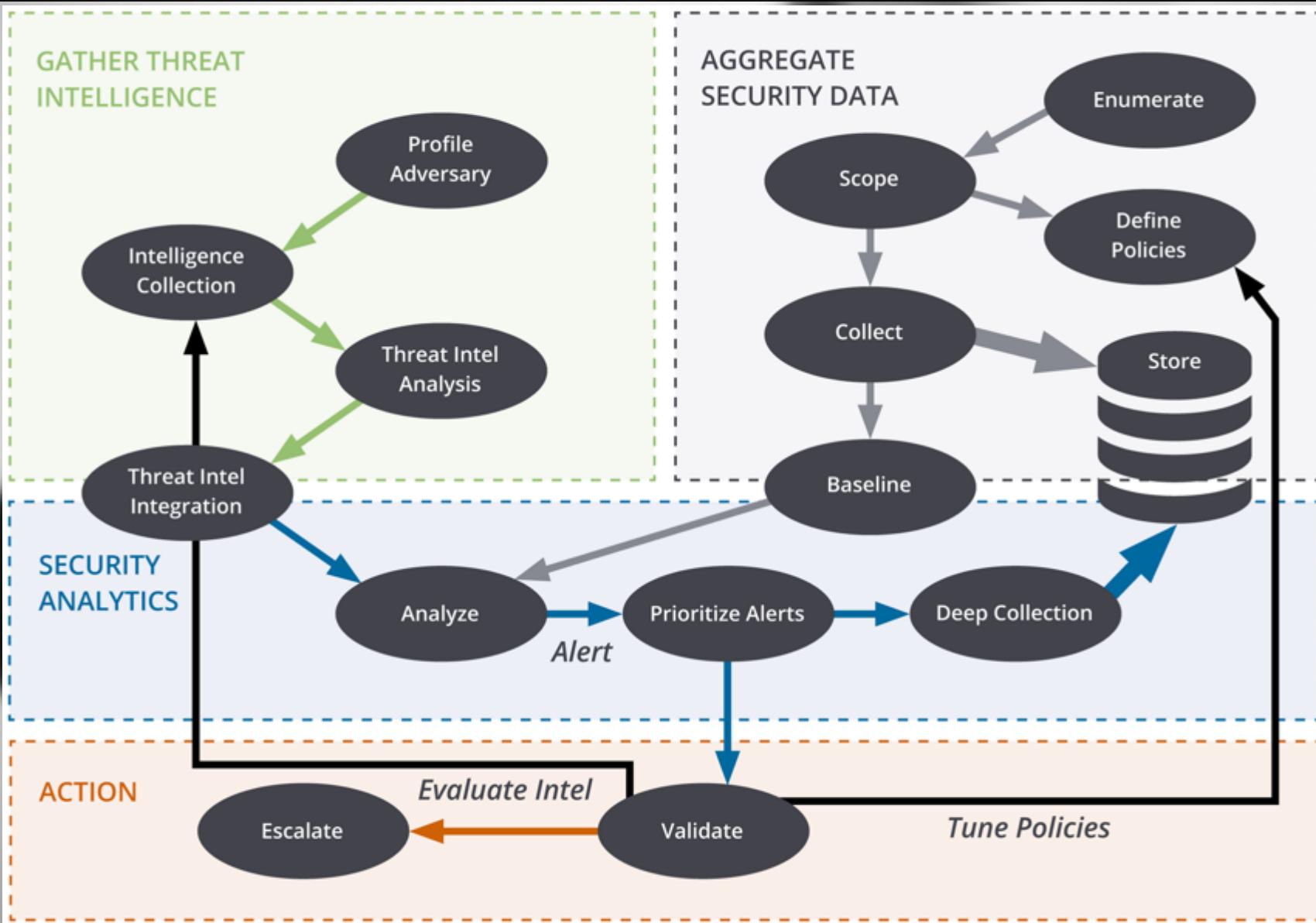
Negative likelihood ratio: 0.131 CI: 0.092 to 0.187

Diagnostic odds ratio: 22.500 CI: 14.928 to 33.914



Threat Score (1 to 10)

# Putting it all together



# Where should we start?

- Aggregate, Normalize and Deduplicate your threat intelligence data (TIP)
- Consult with a data scientist on your current dataset
- Use your intuition and experience when building a model (Scoring)
- Use false positive as feedback to your model (False positives)
- Keep your statistical models aligned with your organizational goals



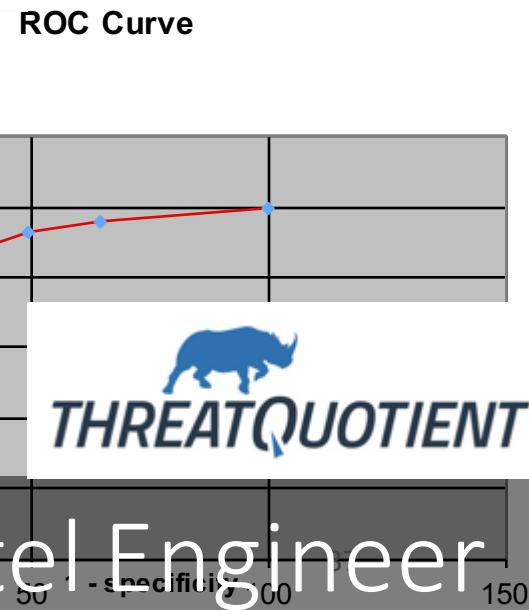
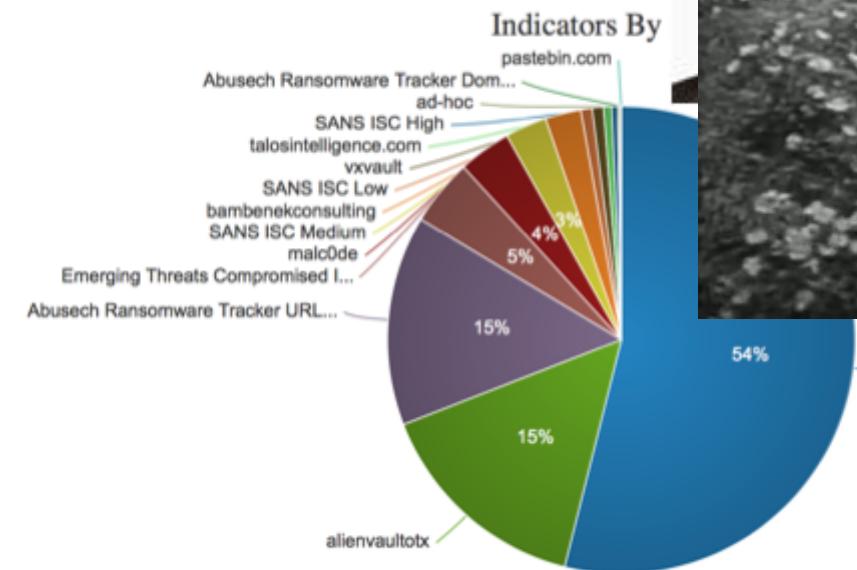
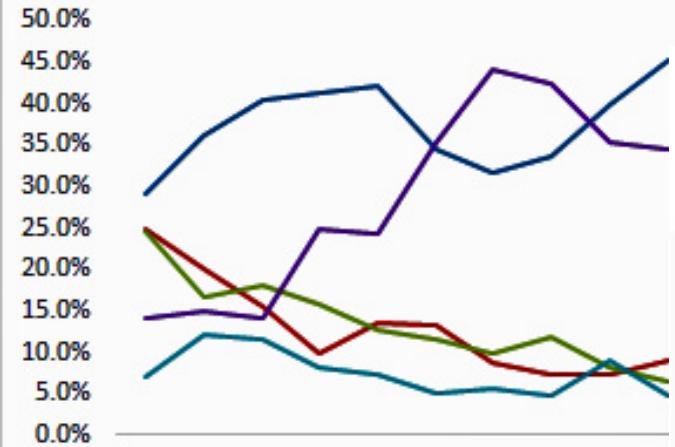
# Summary

- The higher on the Pyramid of Pain, the longer indicators stay relevant.
- Indicators can help not only detect, but also prevent future attacks.
- Current threat intelligence feeds has lots of noise.
- Statistics can help reduce the noise by:
  - Move from subjective to objective decisions
  - Correlate internal and external attributes
  - Reduce the number of false positives
  - Identify threat related trends



# Practical Statistics for Threat Intelligence

Figure 1: Industry Sectors Percentage of Overall Breaches



Nir Yosha – Threat Intel Engineer