

Big Lessons in Small Data

SESSION ID: DSP-T09

Wade Baker

RISK Team

Verizon

@wadebaker

Jay Jacobs

RISK Team

Verizon

@jayjacobs



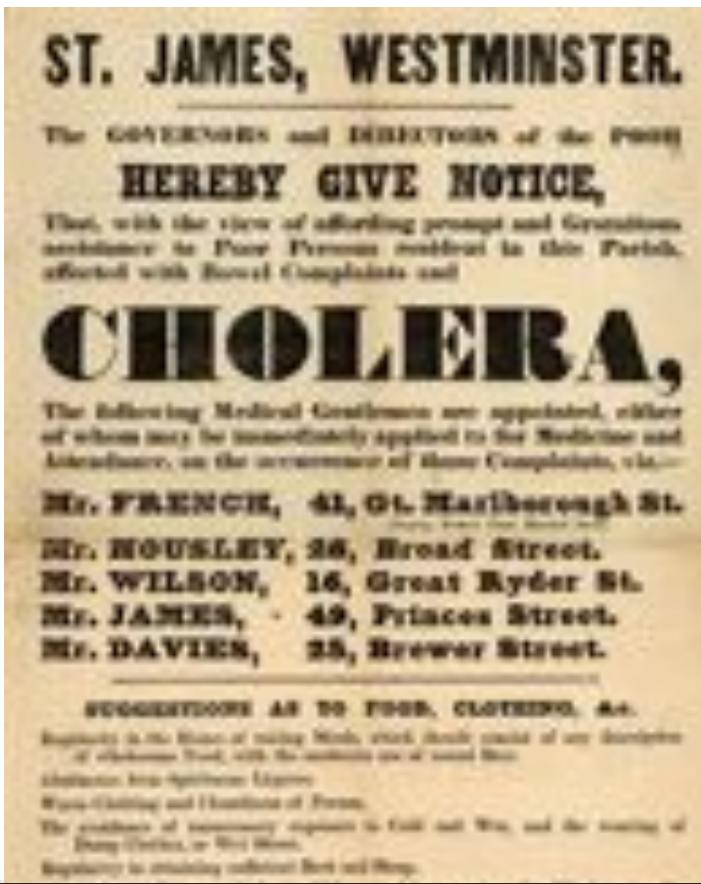
Overview

- ◆ Part 1 : Brief History of Data Analysis
- ◆ Part 2 : Current state of Infosec
- ◆ Part 3 : Putting data analysis into practice

Part 1: Brief History of Data Analysis

Early Data Analysis: Epidemiology

William Farr
1807-1883

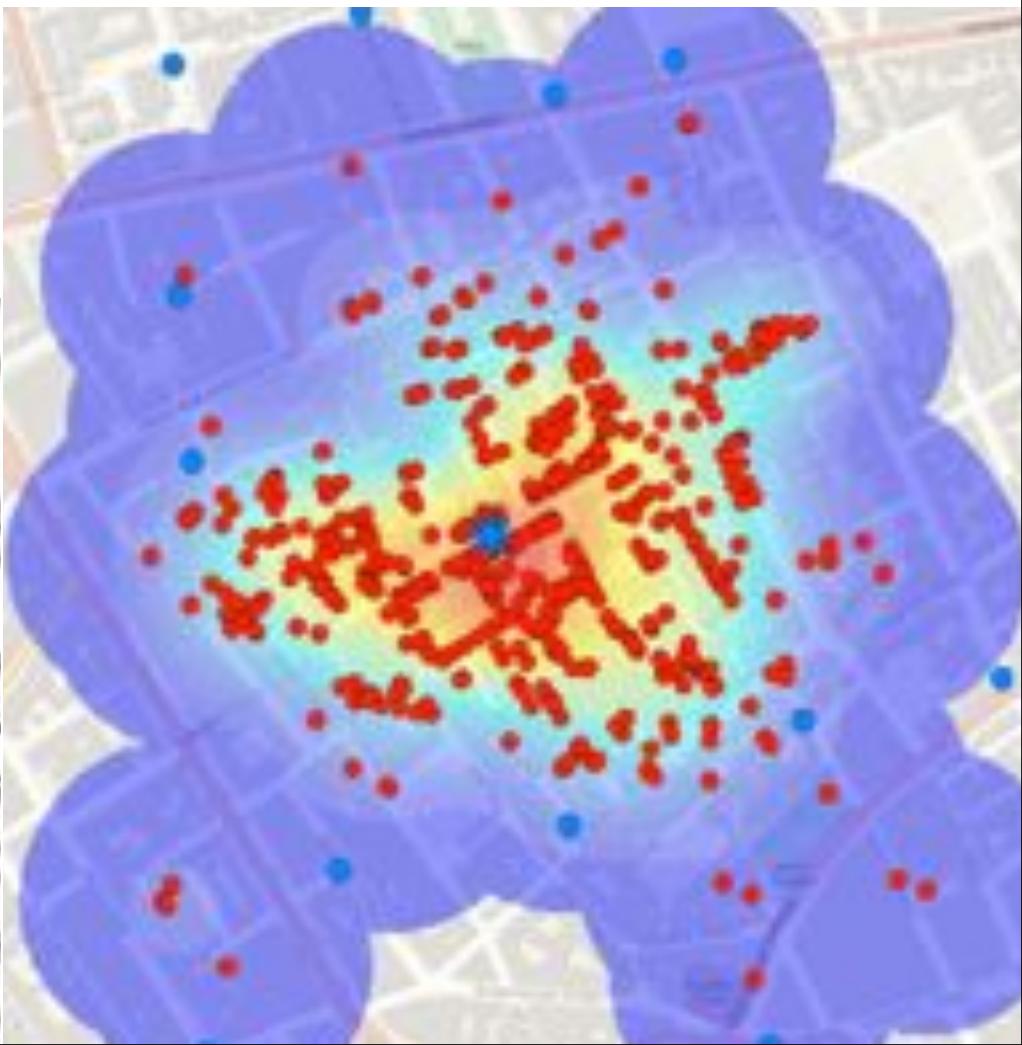


John Snow
1813-1858



John Snow:

Cholera is spread by consuming
water contaminated by a
“special animal poison”



William Farr:

“The elevation of the soil in London has a more constant relation with mortality from cholera than any other known element.”

Table 2 Deaths from cholera in London, registered in 1849 by registration district together with eight possible explanatory variables.

District	Deaths from cholera in 1849 per 10,000 inhabitants	Elevation above high water (feet)	Annual deaths from all causes 1838-1844 per 10,000 inhabitants	Persons per acre	Persons per inhabited house	Average annual value of house (£)	Annual value of house per person (£)	Poor rate percent per pound of house value	Water supply ^a
Harrow	144	-3	232	101	5.8	32	3,788	0.075	1
Rockingham	205	0	277	19	5.8	33	4,238	0.143	1
Bermondsey	161	0	264	66	6.2	18	3,077	0.134	1
St George	164	0	267	181	7.0	32	3,318	0.088	1
Southwark									
St Olave	181	3	281	114	7.9	35	4,559	0.079	1
St Saviour	153	2	292	141	7.1	36	5,291	0.079	1
Westminster	68	2	280	70	6.8	36	4,189	0.076	1
Lambeth	120	3	233	34	6.5	28	4,389	0.039	1
Camberwell	97	4	197	12	5.8	25	4,508	0.072	1
Greenwich	75	8	238	18	6.8	22	3,379	0.038	2
Poplar	71	10	241	15	6.2	44	7,380	0.061	2
Chelsea	46	12	287	62	7.1	29	4,210	0.060	1
Hammersmith	33	12	228	11	6.8	33	9,070	0.067	3

Quote of William Farr as it appears in:

P. Bingham, N.Q. Verlander, M.J. Cheal

John Snow, William Farr and the 1849 outbreak of cholera that affected London: a reworking of the data highlights the importance of the water supply

Bingham, et al. (2004):

“Had logistic regression been available to Farr, its application to his 1852 data set would have changed his conclusion.”

Table 2 Deaths from cholera in London, registered in 1849 by registration district together with eight possible explanatory variables.

District	Deaths from cholera in 1849 per 10,000 inhabitants	Elevation above high water (feet)	Annual deaths from all causes 1838-1844 per 10,000 inhabitants	Persons per acre	Persons per inhabited house	Average annual value of house (£)	Annual value of house per person (£)	Poor rate percent per pound of house value	Water supply ^a
Newington	144	-3	232	101	5.8	32	3,788	0.075	1
Rotherhithe	205	0	277	19	5.8	33	4,238	0.143	1
Bermondsey	161	0	264	66	6.2	18	3,077	0.134	1
St George	164	0	267	181	7.0	32	3,318	0.088	1
Southwark									
St Olave	181	3	281	114	7.9	35	4,559	0.079	1
St Saviour	153	2	292	141	7.1	36	5,291	0.079	1
Westminster	68	2	280	70	6.8	36	4,189	0.076	1
Lambeth	120	3	233	34	6.5	28	4,389	0.039	1
Camberwell	97	4	197	12	5.8	25	4,508	0.072	1
Greenwich	75	8	238	18	6.8	22	3,379	0.038	2
Poplar	71	10	241	15	6.2	44	7,380	0.061	2
Chelsea	46	12	287	62	7.1	29	4,210	0.060	1
Hammersmith	33	12	228	11	6.8	33	5,070	0.067	3

Quote of William Farr as it appears in:
P. Bingham, N.Q. Verlander, M.J. Cheal
John Snow, William Farr and the 1849 outbreak of cholera that affected London: a reworking of the data highlights the importance of the water supply

“Modern” statistical inference: Design of experiments

R. A. Fisher
1890-1962



“After two or three months of investigation it will be found possible to understand some of Fisher’s sentences.”

Sir Fred Hoyle, Astronomer

And then the transistor...

John Tukey
1915-2000

and Exploratory Data Analysis



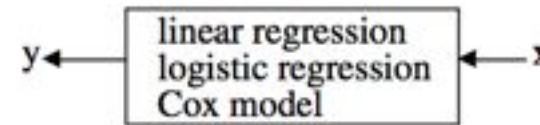
“Numerical quantities focus on expected values, graphical summaries on unexpected values.”

Statistical Modeling: The Two Cultures

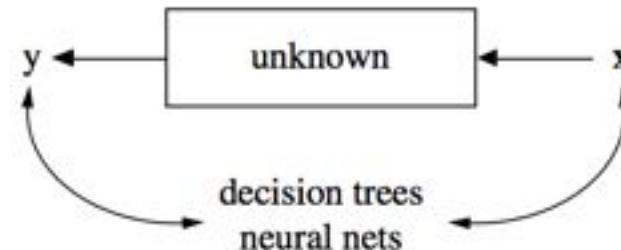
Leo Breiman
1928-2005



Data Modeling Culture



Algorithmic Modeling Culture



Part 2: Current State of Infosec

(Hint: 1800's statistics now with computers!)

Our Goal:

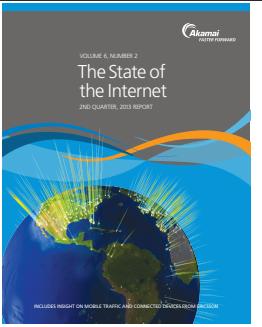
- ◆ How is infosec doing with our data (small or big)?

Our Approach:

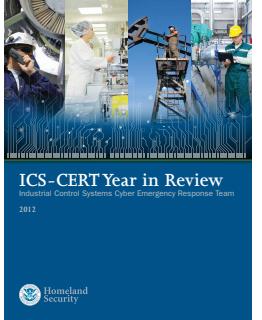
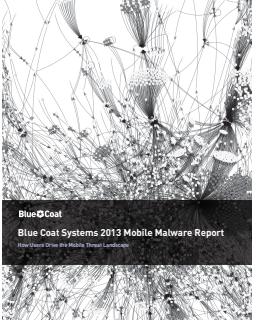
- ◆ Analyze data analysis techniques in (industry) publications

Our Conclusion:

- ◆ We have some “opportunities for improvement”



ARBOR



SOPHOS

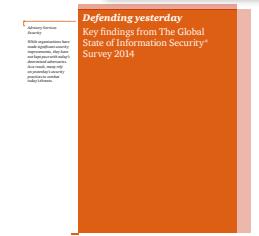
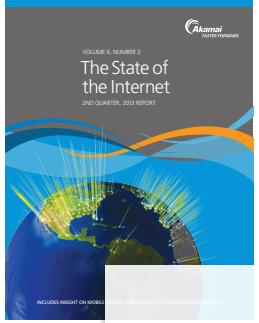
Security Threat Report 2013



Presenter's Company Logo –
replace on master slide

Over 50 Reports, seeking:

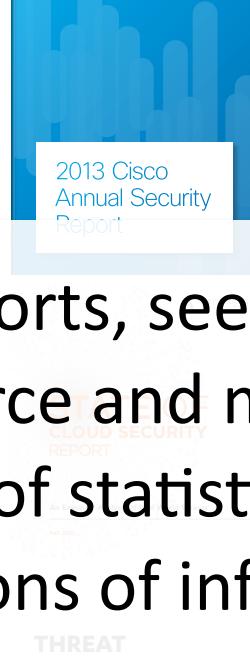
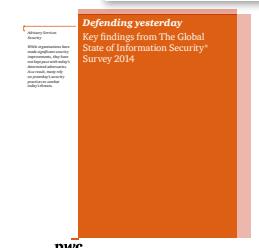
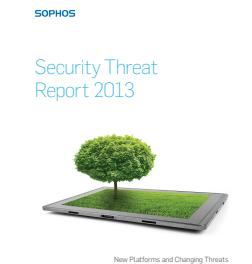
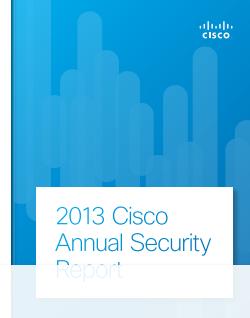
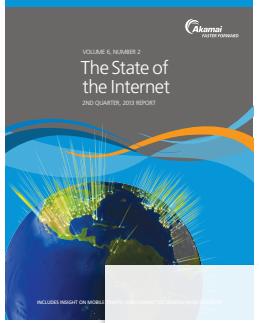
- ◆ Data Source and methodology
- ◆ [Mis]use of statistics
- ◆ Applications of inference



Presenter's Company Logo –
replace on master slide

Over 50 Reports, seeking:

- ◆ Data Source and methodology
- ◆ [Mis]use of statistics
- ◆ Applications of inference



Presenter's Company Logo –
replace on master slide

Results

Note: to RSA reviewers, we are still finalizing the numbers here

- ◆ mm Surveys, mm Device Data, mm Reports
- ◆ mm of 52 defined their methodology
- ◆ mm defined their **sample** of which mm defined their **population**
- ◆ **None** applied a probabilistic sampling method*
- ◆ Heavy use of “count and compare”
- ◆ Conclusion: *We are in the 1800's for data analysis*

* Not a bad thing

Opportunities to Learn...

- ◆ More science! (less FUD)
- ◆ More statistics!
- ◆ More transparency!

What is a Scientist?
‘A scientist is a person who asks questions and tries different ways to answer them.’



Opportunities to Learn 1: More Science!



“My job was to find questions about baseball that have objective answers, that’s all that I do, that’s all that I’ve done.”

-- Bill James, Sabermetrician

Scientific method

1. Research question
2. Design/plan experiment
3. Gather data
4. Make sense of data
5. Communicate results



<http://scienceandstory.files.wordpress.com/2012/04/science.jpg>

Opportunities to Learn 2: More Statistics!

- ◆ Exploratory Data Analysis
- ◆ Sample Error vs. Sample Bias
- ◆ (let's argue about p-values later)



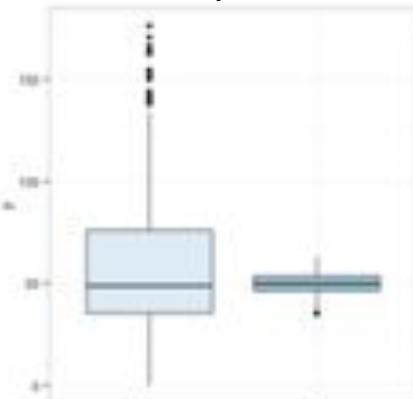
Exploratory Data Analysis (get to know your data)

“Five Number” summary:

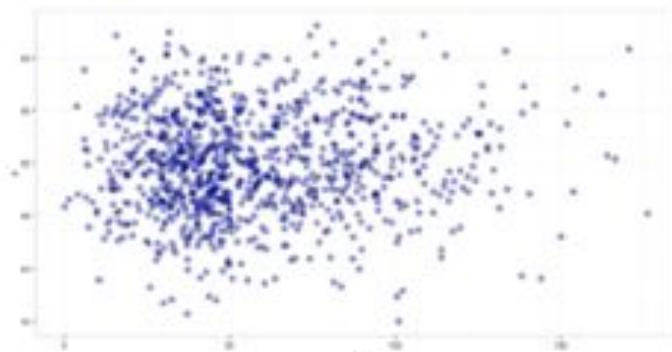
```
> summary(x)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.10	35.42	48.90	57.07	75.97	176.40

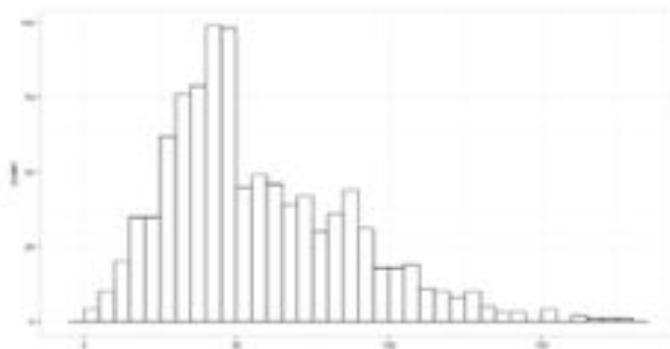
Box plot



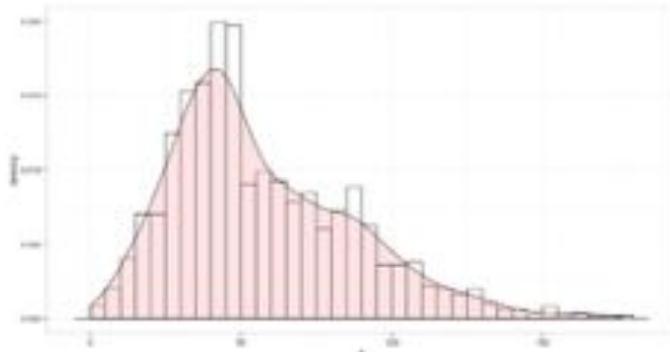
Scatter plot



Histogram



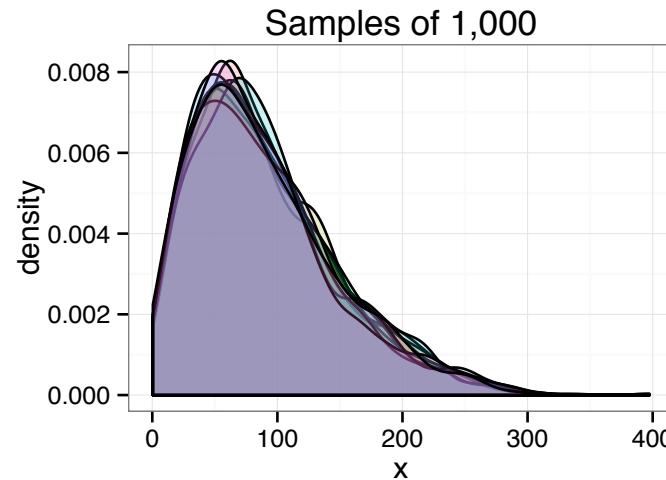
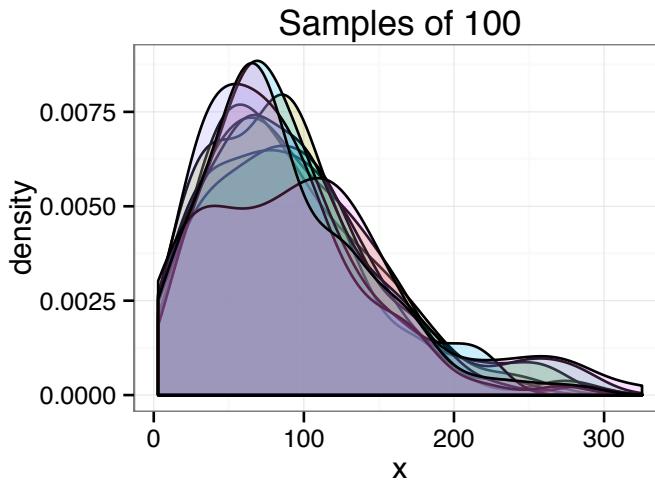
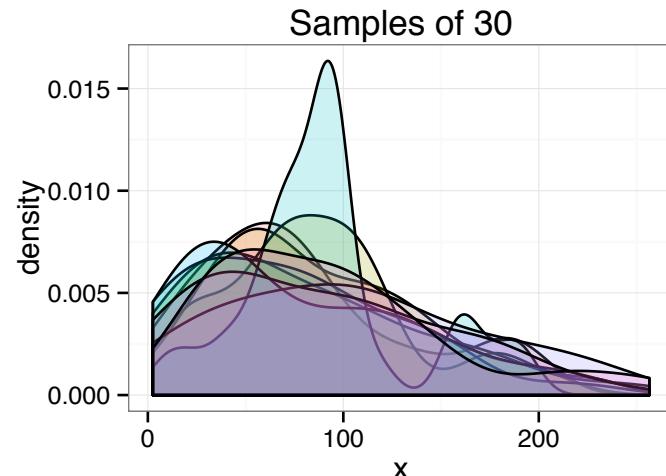
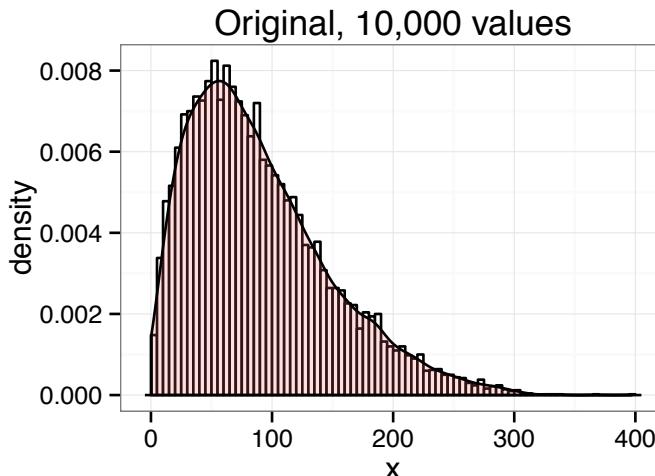
Density plot



Sample Error

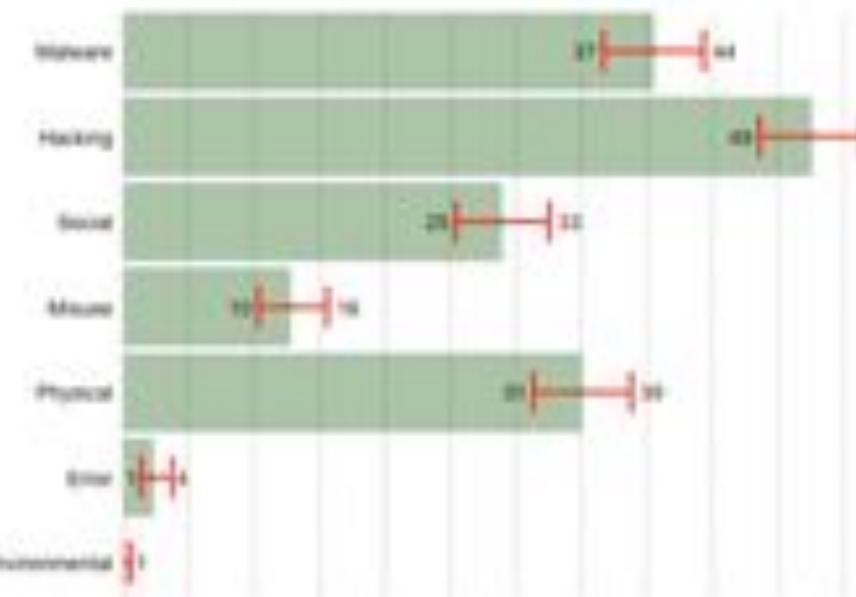
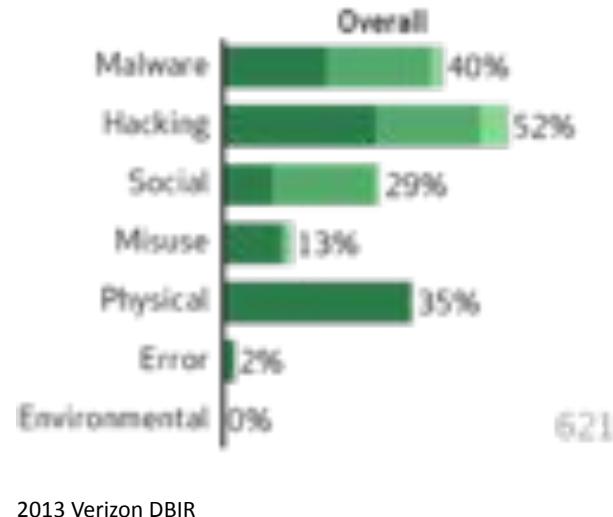
("error" does not mean mistake)

Increasing sample size reduces sample error.



Sample Error

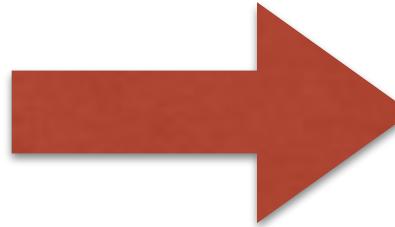
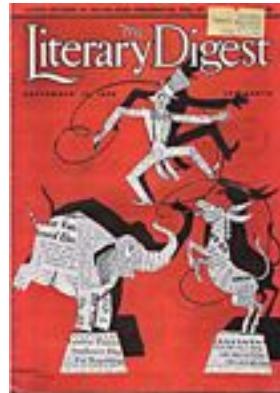
Figure 16: Threat action categories



95% confidence intervals

Sample Bias

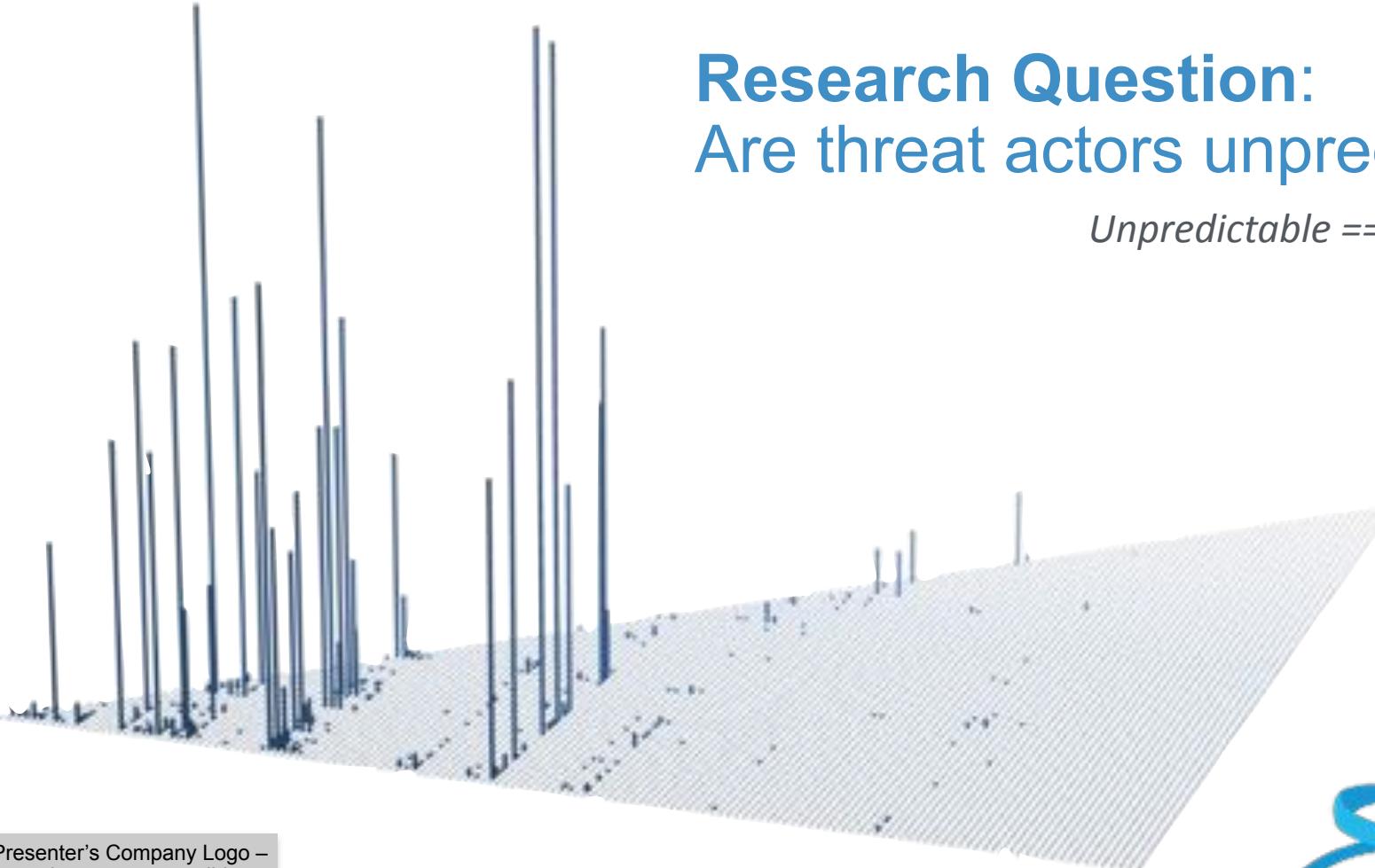
telephone directories +
lists of magazine subscribers +
rosters of clubs and associations
=
2.4 Million people polled and
Landon - 57%
Roosevelt - 43%



<http://www.math.upenn.edu/~deturck/m170/wk4/lecture/case1.html>

Part 3: Putting Data Analysis into Practice

[Examples of big lessons in small data]



Research Question:

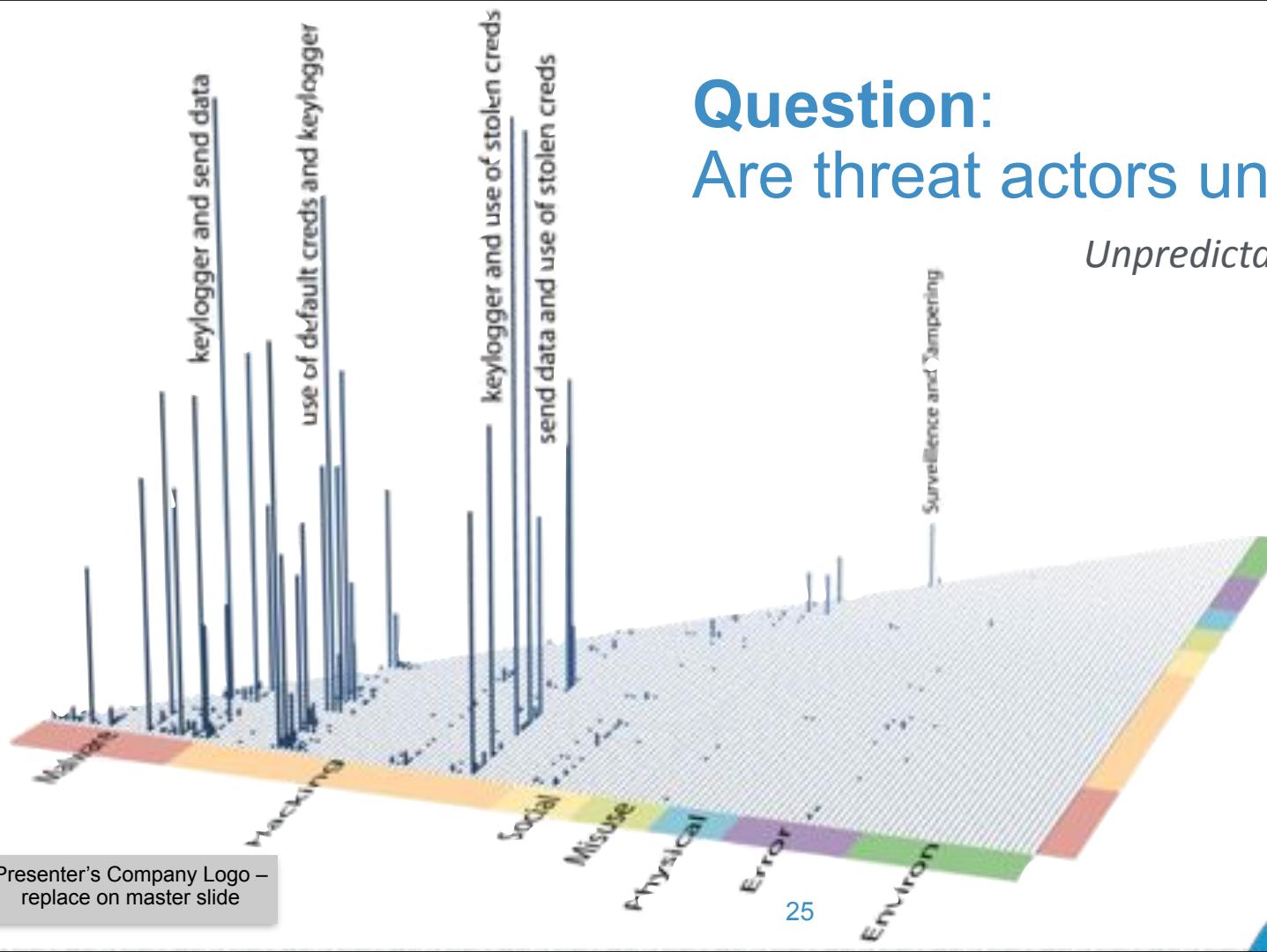
Are threat actors unpredictable?

Unpredictable == Random

Question:

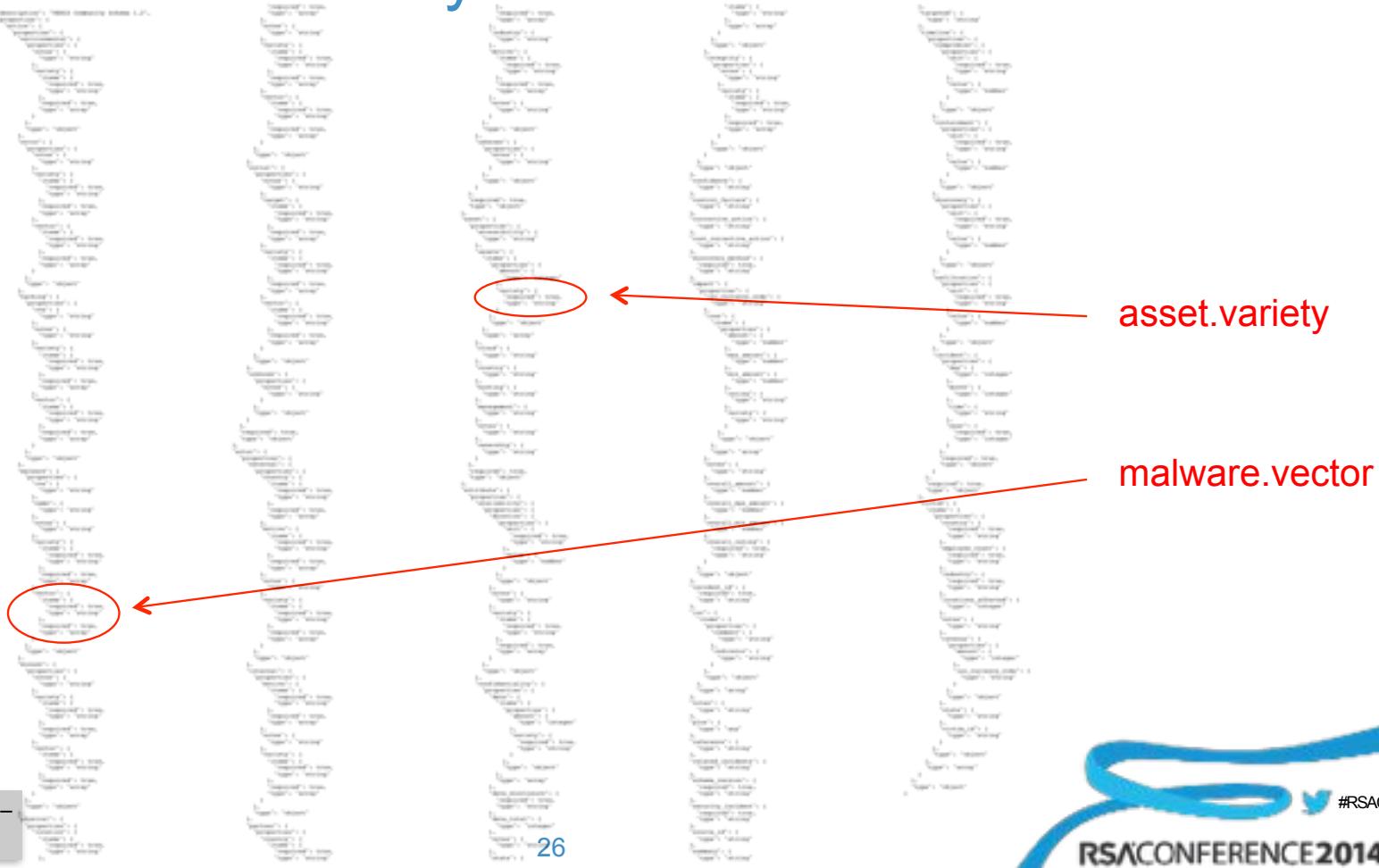
Are threat actors unpredictable?

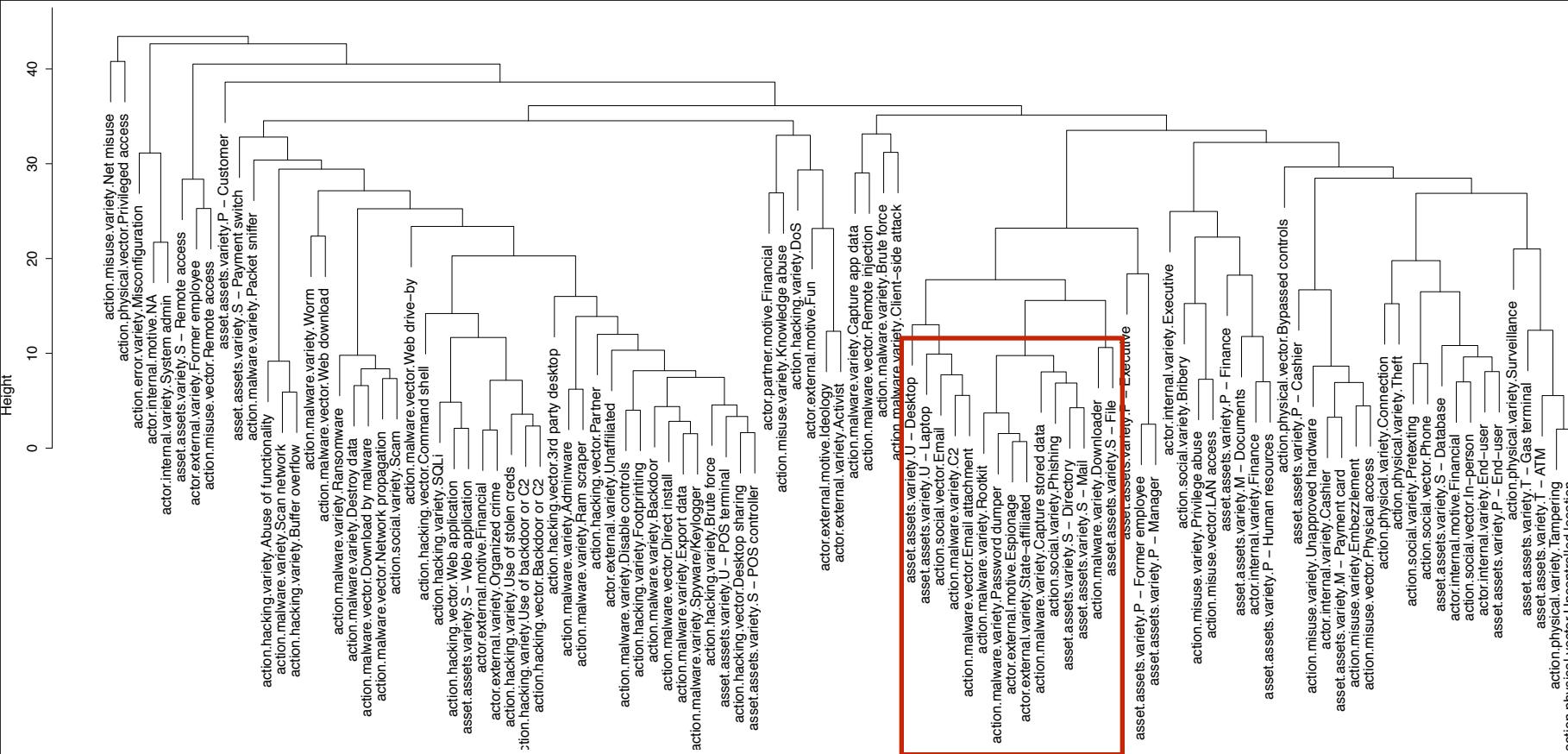
Unpredictable == Random



Presenter's Company Logo –
replace on master slide

The “DNA” of a Security Incident

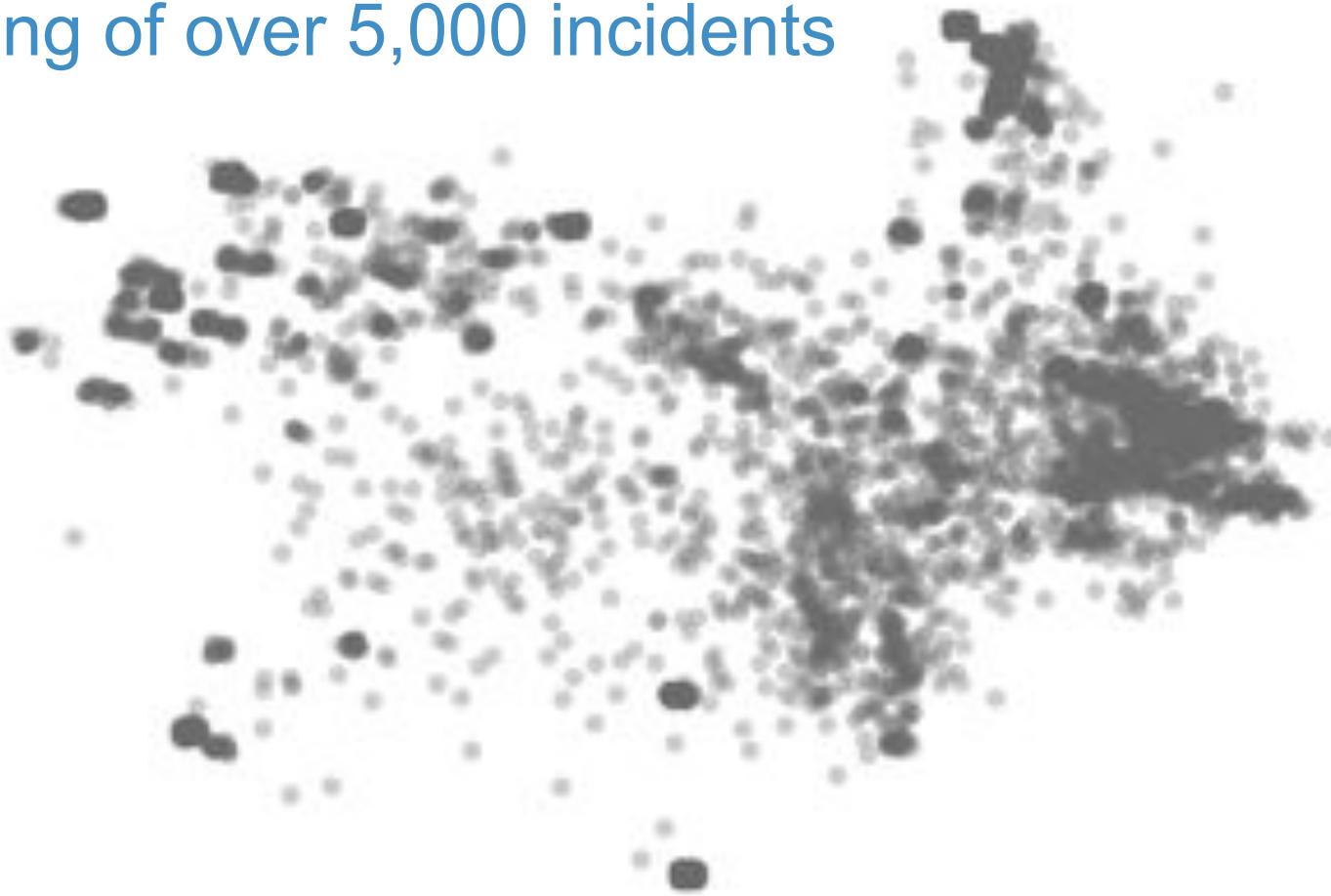




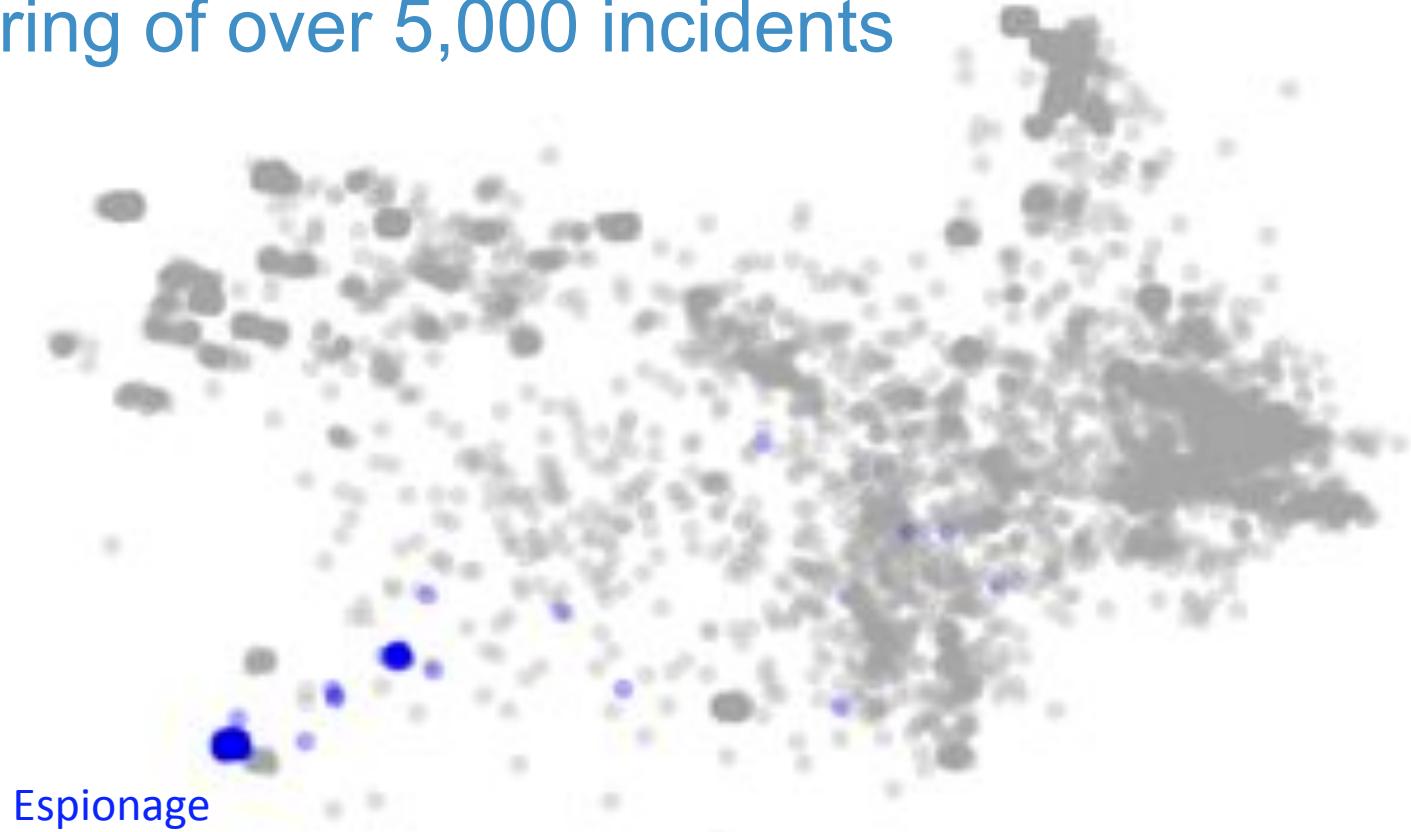
Key variables in “Espionage” pattern

(*actor.external.state-affiliated*, *malware.variety.rootkit*,
malware.variety.pwd dumper, *social.variety.phishing*,
asset.variety.S – directory, etc)

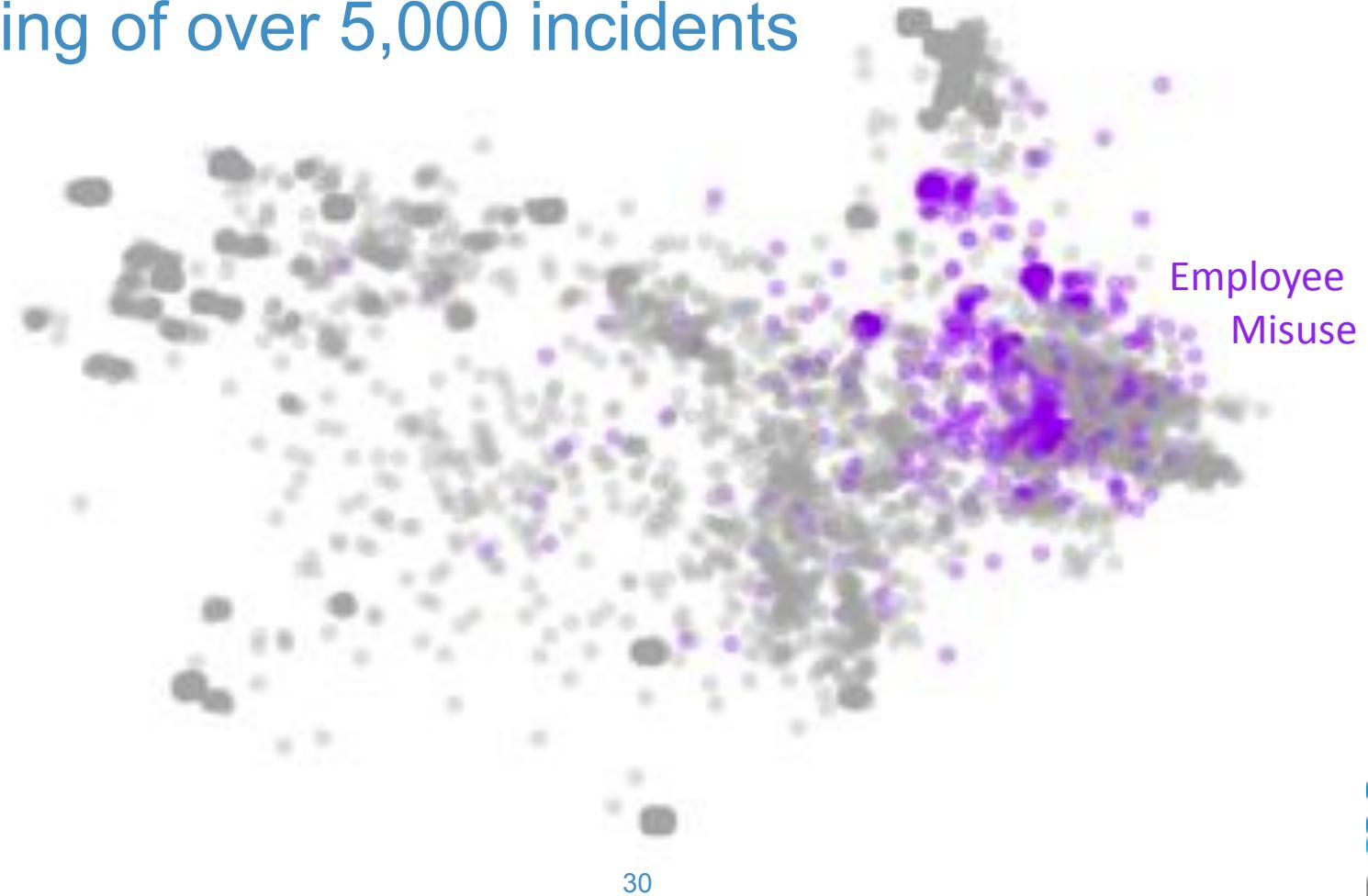
Clustering of over 5,000 incidents



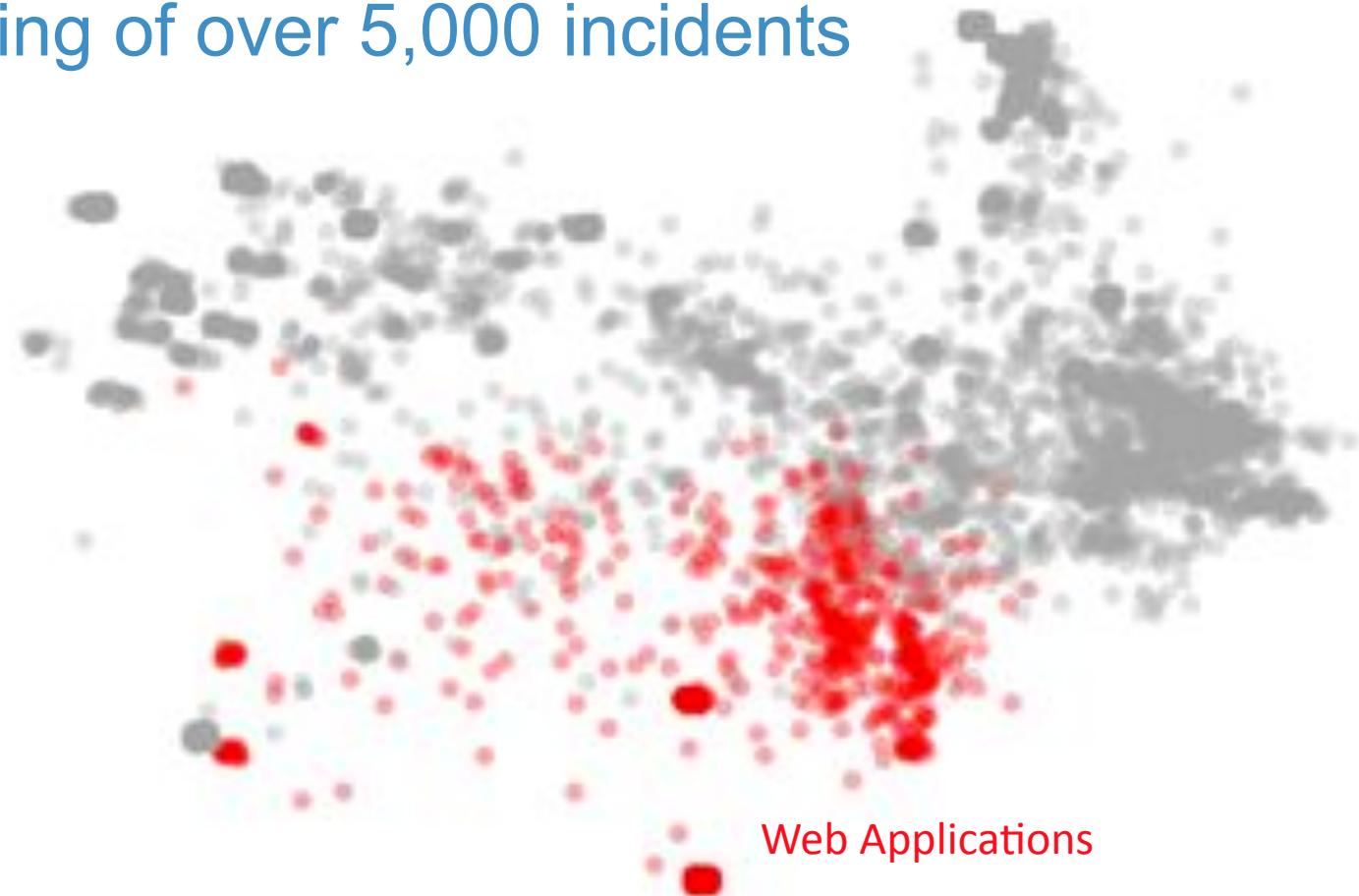
Clustering of over 5,000 incidents



Clustering of over 5,000 incidents

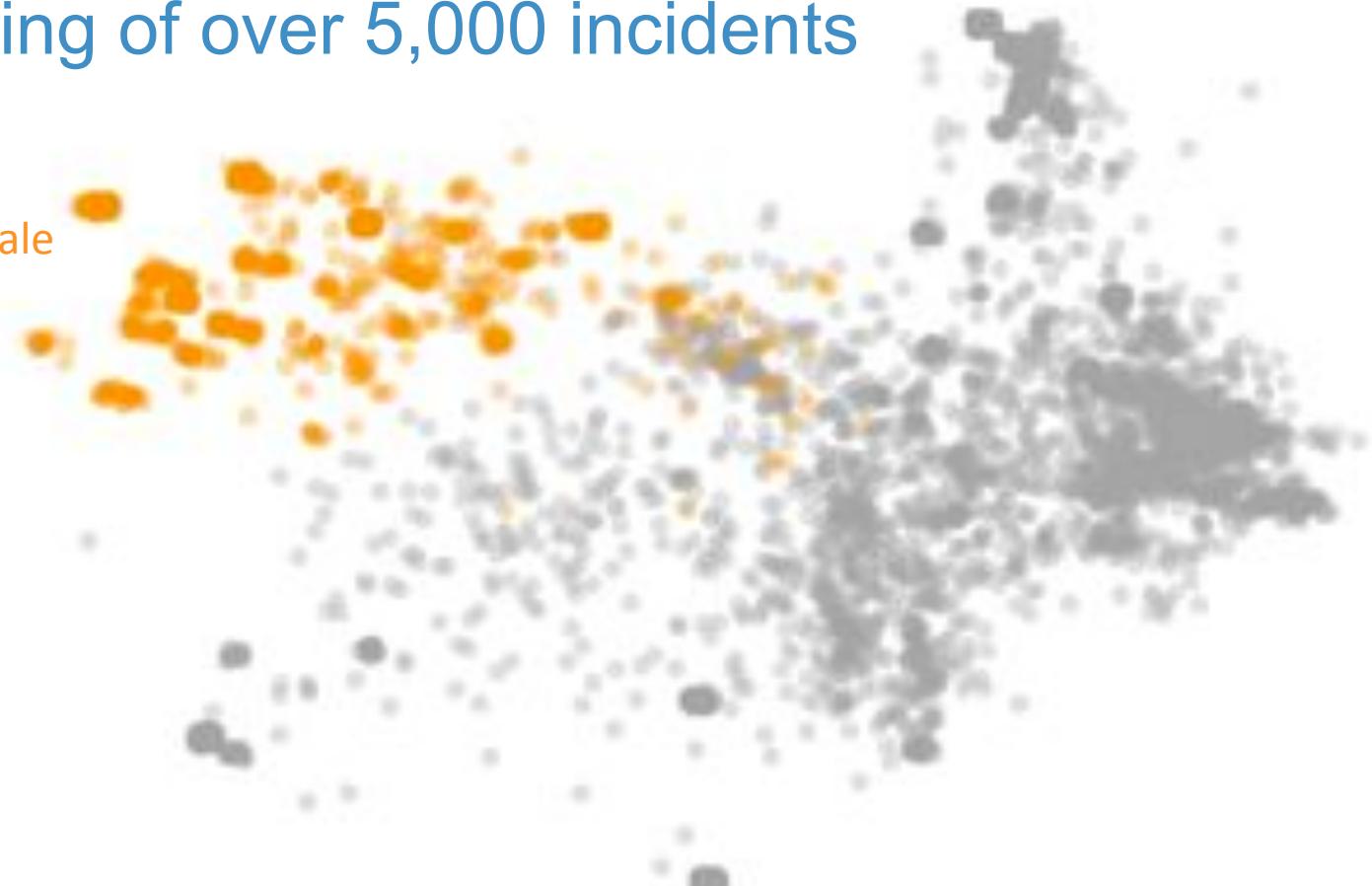


Clustering of over 5,000 incidents

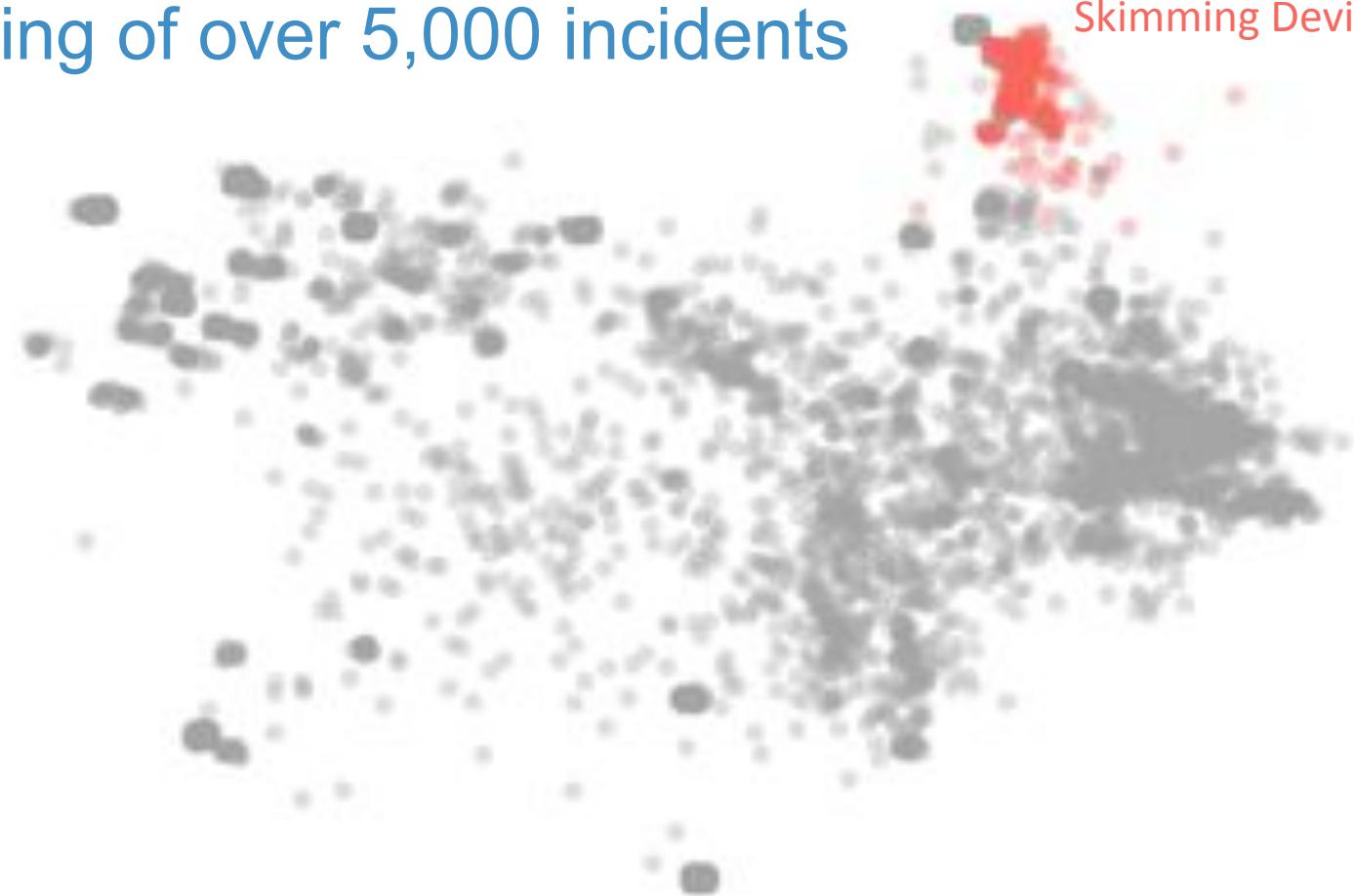


Clustering of over 5,000 incidents

Point of Sale

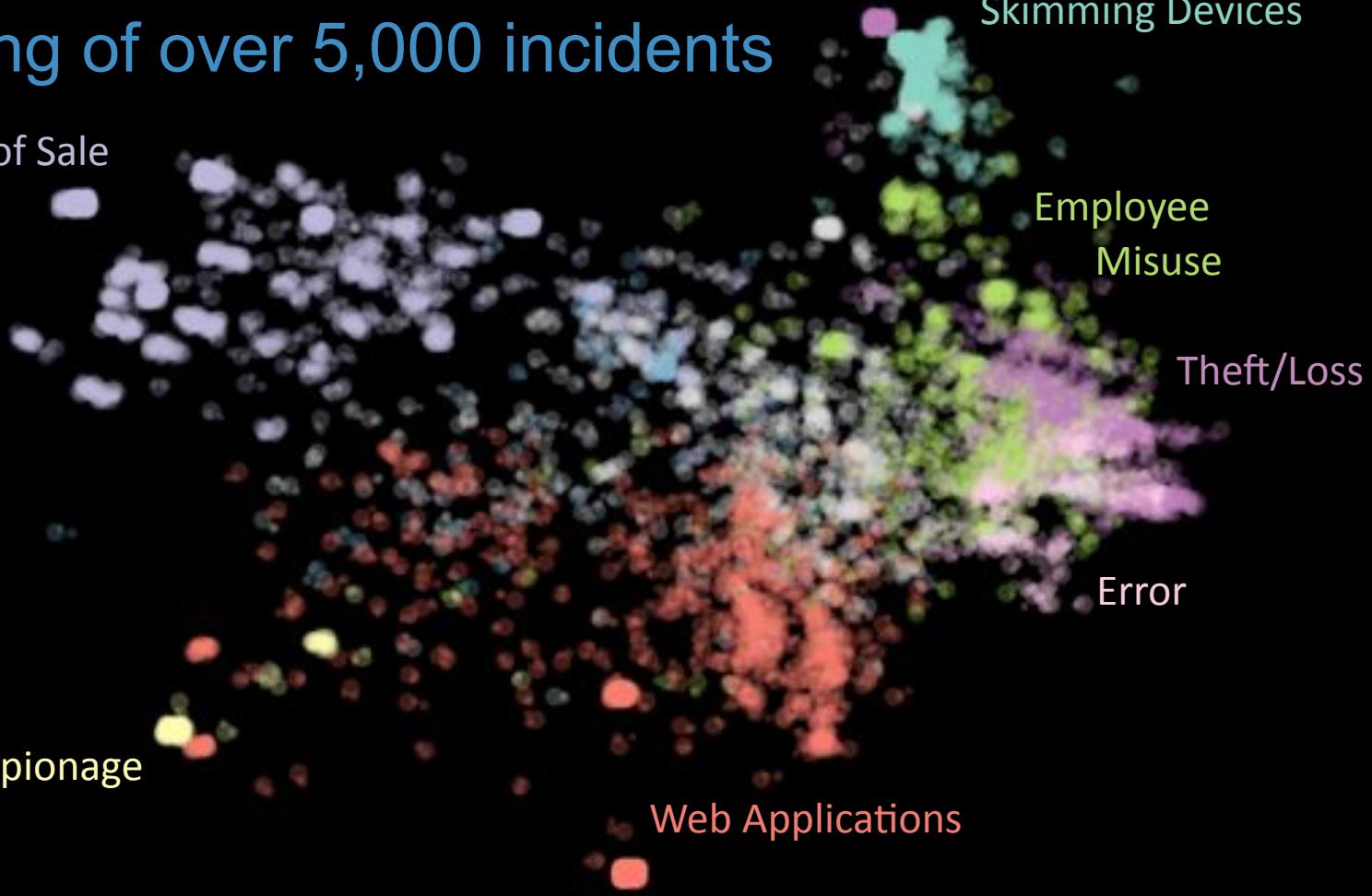


Clustering of over 5,000 incidents



Clustering of over 5,000 incidents

Point of Sale



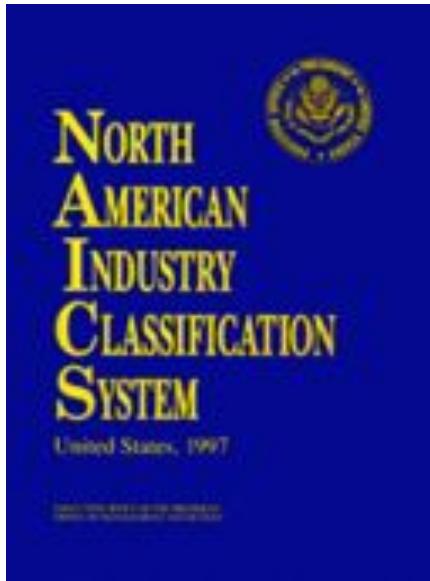
Research Question:

Is there such thing as a “Top 10” list of controls?

Note to RSA reviewers,
We are still working through the data on this
and we'll fill this in.

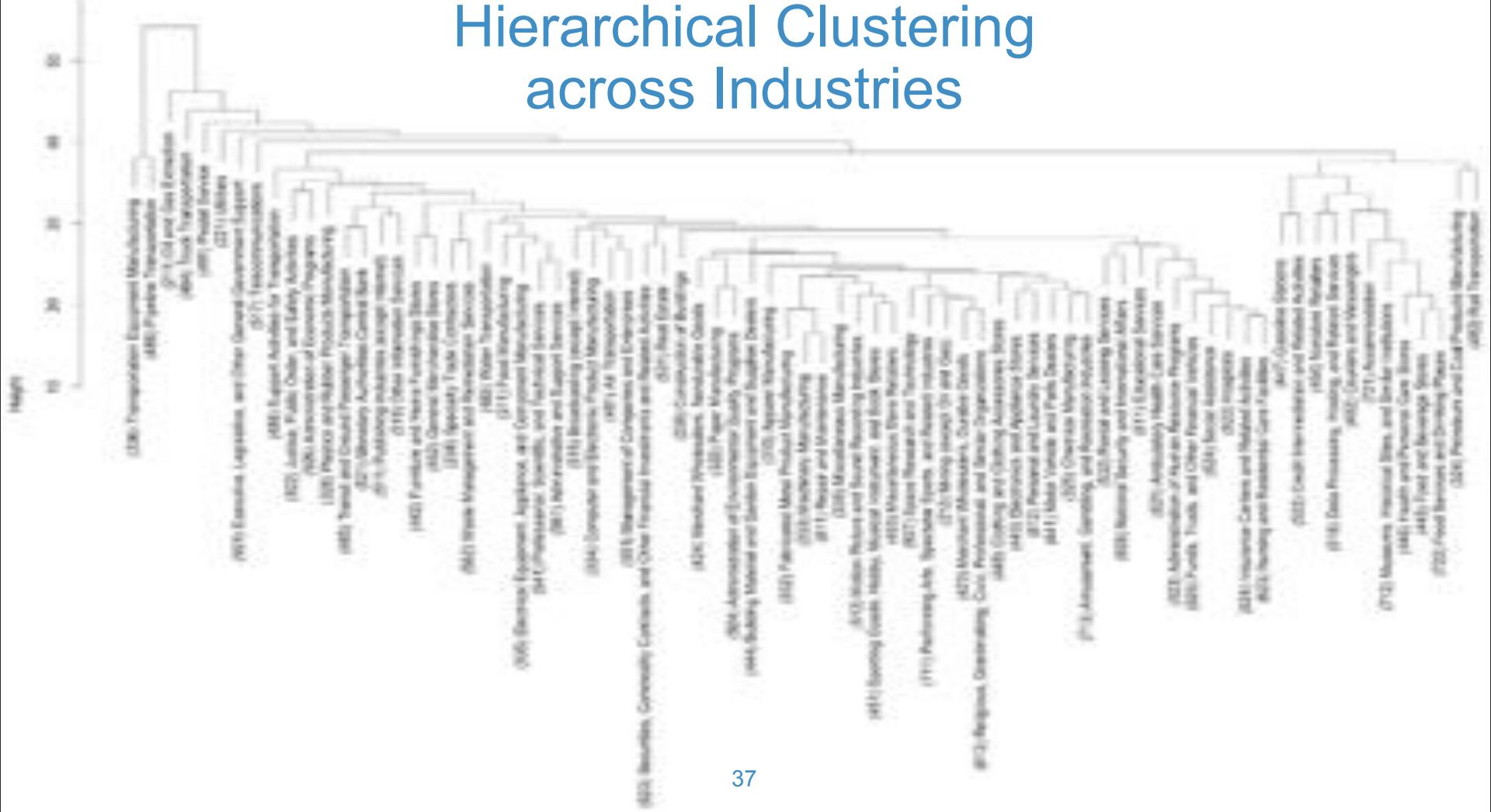
Research Question:

What industries are comparable (similar threat

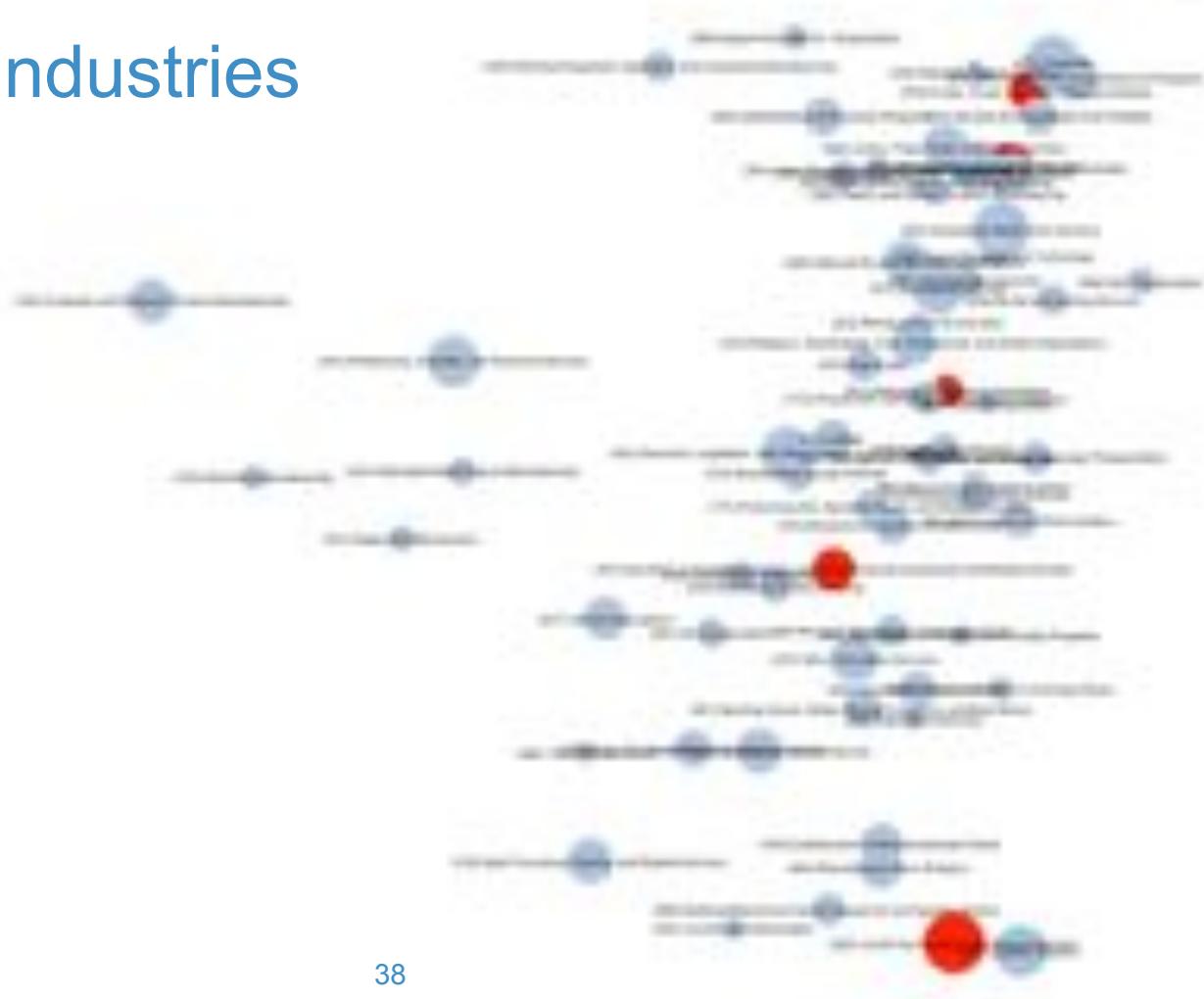


Need to answer this before,
“How am I compared to my peers?”

Hierarchical Clustering across Industries



Clustering on Industries



RSA CONFERENCE 2014

FEBRUARY 24 - 28 | MOSCONE CENTER | SAN FRANCISCO

