

Cloudy With a Chance of Breach Forecasting Cyber Security Incidents

Manish Karir

Yang Liu, Armin Sarabi, Jing Zhang, Parinaz Nagzadeh, Michael Bailey, Mingyan Li

Background

Reputation Matters

Security Posture is an important part of a business relationship

. .

Security Posture is the sum total of all factors including people, processes and technology



The Problem

Measuring Security Posture is difficult. Need metrics to assess a network or organization from the outside to determine risk



One-Slide Summary

Use internet scale data measurements and machine learning to do predictive cybersecurity risk analysis of networks.

This information helps:

- Underwriters set premiums for cyber insurance.
- Network managers understand their security posture and their exposure due to partners and vendors.

APPROACH



Data Driven Models and Analysis



Mis-management Symptoms

Datasets used in the study:

- Open Recursive Resolvers
 (openresolverproject.org)
- DNS Source Port Randomization Research DNS test probes
- BGP Misconfiguration Extrapolated from short-lived BGP route announcements
- Untrusted HTTPS Certificates Research SSL scan data
- Open SMTP Mail Relays Research test probes



Malicious Activities Data

Inferred Malicious Activities from RBL Lists:

- SPAMHAUS-XBL, SpamCop
- UCE-PROTECT, SURBL, WPBL, PhishTank, HpHosts
- Darknet Scanners, DSHIELD, OpenBL



Labeling Cyber Security Incidents Datasets

Ground-Truth data used for identifying data breaches:

Reported Data Breach Collections:

- VCDB Veris Community Database (basis for Verizon Data Breach Investigations Report)
- Hackmageddon
- Web Hacking Incidents Database



Building Machine Learning Models

- 250+ features for each "maintainer" or address block including:
- Mis-management symptoms Ratios:
 - Number of invalid certs
 - Number of open relays
- Malicious activities time series:
 - Fraction of address block in RBL
 - 60 day RBL activity trends
- Second order statistics:
 - Frequency, Duration of RBL activity



Model Training and Testing

Machine learning approach

- Random Forest Machine Learning Algorithms – many sub-samples to build a forest of decision trees which are then labeled using the incident or non-victim population data
- 50/50 training-testing split
- Testing Evaluation to demonstrate accuracy
- Training multiple windows, multiple data combination, multiple parameters



RESULTS



CYBERSECURITY INCIDENT PREDICTION RESULTS



RELATIVE IMPORTANCE OF DIFFERENT FEATURES

Feature Category	Normalized Importance
Mismanagement	0.3229
Time Series Data	0.2994
Recent-60 Secondary Features	0.2602
Organization Size	0.0976
Recent-14 Secondary Features	0.02

- Overall mis-management features which are the most directly related with human factors have the largest normalized weight
- This confirms the intuitive understanding that the human element is the most important factor in cyber security

IMPORTANCE OF DATA DIVERSITY



- Mis-management features by themselves are not sufficiently good predictors BUT
- In combination with other features such as malicious activities they add the MOST value.

PREDICTING SOME WELL KNOWN INCIDENTS



- 65% of incidents in blind-test dataset were predicted as 100% chance of breach
- A threshold of 0.85; predicts 92% of breaches

Who could have seen it coming...

- Prediction Models in study built mid-2014
- Anthem Healthcare Prediction Score: 90% Breach Reported Mar 2015
- Penn State Prediction Score: 97% Breach Reported May 2015
- Rutgers Prediction Score: 92% Breach Reported Oct 2014
- Office of Personnel and Management Prediction Score: 90% -Breach Reported April 2015

Conclusion

- It is possible to statistically predict cyber security incidents on the basis of historical incidents and pre-incident security posture data
- Difference between *detection* and *prediction* is key One relies on signatures the other looks at patterns and trends in data that might appear to be un-related
- Security posture is many-dimensional and requires data from many parts of an organization including web applications, network configurations, to DNS
- Protecting against data breaches requires fighting a battle on many fronts and the key almost always is people

Questions?

