black hat ASIA 2017

MARCH28-31,2017

MARINA BAY SANDS / SINGAPORE

Beyond the blacklists: Detecting malicious URL through machine learning

Hao Dong, Jin Shang, David Yu, Chenghuai Lu

Hillstone Networks



About Us

Hao Dong

- Senior Principle Engineer
 @Hillstonenet.com
- Interested in computer system and network security

Jin Shang, Ph.D.

- Chief Scientist & Fellow
 @Hillstonenet.com
- Marathon Runner In-Training





About Us

David Yu

- Distinguished Engineer @Hillstone Networks
- BS and MS in Electrical & Computer Engineering
- Over 20 years experiences in computer networking and security industry

Chenghuai Lu

- Senior Principle Engineer
 @Hillstonenet.com
- Lead efforts in defending against advanced threats in next generation firewall
- Afternoon Tea Party Ping Pong Champion





Agenda

- Introduction and Problem Overview
- Malicious HTTP Traffic Detection Methods
- Machine Learning Techniques
- The Results
- Beyond the Malicious URL
- Black Hat Sound Bytes
- Q & A



Introduction

- Hillstone is recommended by NSS lab NGFW test with the best values, but today we talk about the research beyond the NGFW
- Network security defense technology evolution



- Intelligent technologies
 - Algorithms: Data Analysis, Machine Learning, Anomaly Behavior, ...
 - □ Modeling: Feature set, customized Model, cloud computing
 - □ Product: False Positive & Negative, Forensic, Mitigation



Introduction

- Our thoughts:
 - □ iNGFW engines: Intelligent + NGFW
 - Deploy in perimeters and inside networks
 - □ Enhance the defense architecture
 - Suspicious event detection using AI: DGA, malicious URL, abnormal network behavior
 - Multi Layer Correlation among iNGFW to improve false positive/negatives

Today, we present one of the best ML work:

malicious URL detection with ML



WHAT: *Problem to Solve*

• To detect malicious URL connections, esp. in **post-breach** advanced threat protection.





Post-breach detection

□Limitation of signature based blacklist

□To know unknowns from knowns

••• There are known knowns. These are things we know that we know. There are known unknowns. That is to say, there are things that we know we don't know. But there are also unknown unknowns. There are things we don't know we don't know.



-- Donald Rumsfeld



HOW: *Explore machine learning technology*

- Predict malicious connection with knowledge learned from existing ones
- Based on URL lexical features
 - □Not by page content
 - □ Focus on solving URL parameters complex issues
 - Catch commonalities and be generic enough to detect variants





WHERE: Model build and deployment

- Network-level detections
- Flexible deployment

□Network perimeter firewall appliance

□Server-centric intranet breach detection

• Constantly model update with incoming threat intelligences, user feedbacks and global correlations.

HTTP detection modules

- User-Agent
- Host Domain Name
- URL Path
- URL Parameters
- HTTP headers patterns

Put URL Under Microscope

Host Domain Name

Known malicious domainsDGA (Domain generation algorithms)

User-Agent

□Known signatures

• URL Path

Inttp://surusegitimmerkezi.com/administrator/components/com_akeeba/akeeba/engine s/proc/mzsystem.php (from ransomware tracker)

URL parameters

Put URL Under Microscope

HTTP Headers Pattern (browser finger print)

- □ Mostly useful to reduce false positive
 - Based on common http header fields
 - > Number of fields and order of fields of common web browsers
- Observations
 - > Malicious connection tend to have less number of header fields
 - Typical web browser traffic have fixed set of headers and orders
- □So we set rules
 - Threshold number of header fields
 - Order of header fields match common web browsers

ACCEPT	text/html,application/xh tml+xml,application/xml ;q=0.9,image/webp,*/*; q=0.8
ACCEPT_ENCODING	gzip, deflate, sdch, br
ACCEPT_LANGUAGE	en-US,en;q=0.8
CONNECTION	keep-alive
DNT	1
HOST	www.whatismybrowser. com
REFERER	https://www.google.co m/
USER_AGENT	Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537
COOKIE	

URL Parameter Features

• Challenges:

Dynamic, huge amount of text strings and keep increasing

- We treat parameters as collection of key/value pairs
 □param_name=param_value is a feature
 □Unlimited amount of strings, means high number of feature dimensions
- The way to reduce dimensions:

□Sorting method vs. hashing

URL Parameter Features

- Sorting method
 - Parameters sorted by hit numbers (total hits, malware hits, etc.)
 - □Select top N as feature set

URL Parameter Features

- Hashing method
 - □Fixed size of feature space
 - Transform a parameter to multiple features to emphasize field weight
 - Expand features using combinations of field value, type and length
 - ➢ Numeric 123456
 - Alphabets abcdefg
 - ➢ NumAlpha abc123
 - Base64 d2VsY29tZSBibGFja2hhdCBhdHRlbmRlZXMh
 - ≻ Etc.

Transforming parameters

One	http URL	is converted	to a sparse vec	tor			
1	1	• •	11	1	1	1	1

Feature Space

Transforming parameters

Feature Space

Machine Learning

- HTTP URL collections
 - From malware network trace PCAPsFrom clean legitimate traffic collections
- Choose learning method
 Supervised vs. Unsupervised
 Pros and Cons

Supervised Machine Learning

Supervised learning is the <u>machine learning</u> task of inferring a function from *labeled training data*.^[1] The <u>training data</u> consist of a set of *training examples*. In supervised learning, each example is a *pair* consisting of an input object (typically a vector) and a desired output value (also called the *supervisory signal*). A supervised learning algorithm analyzes the training data and produces an inferred function, which can be used for mapping new examples. An optimal scenario will allow for the algorithm to correctly determine the class labels for unseen instances. This requires the learning algorithm to generalize from the training data to unseer situations in a "reasonable" way (see <u>inductive bias</u>).

The Free Encyclopedia

Un-supervised Clustering

Cluster analysis or **clustering** is the task of grouping a set of objects in such a way that objects in the same group (called a **cluster**) are more similar (in some sense or another) to each other than to those in other groups (clusters). It is a main task of exploratory <u>data mining</u>, and a common technique for <u>statistical data analysis</u>, used in many fields, including <u>machine learning</u>, pattern recognition, <u>image</u> <u>analysis</u>, <u>information retrieval</u>, <u>bioinformatics</u>, <u>data compression</u>, and <u>computer graphics</u>.

Supervised Machine Learning

Malicious

Adjust models with false positive feedbacks

Hard to change models. Need take feedback and re-learn. Or use whitelist

Easy for incremental update (add and remove)

Add new data and labels, and re-learn the whole model

Supervised ML

Logistic Regression. Small model size, good for binary classification, e.g. black and white
 Random Forrest. Big model size, good result for multi-family classification
 Ground truth labeling of malware family is important. Model can be biased

Unsupervised ML

□Need only malware (black) samples. Good to find common similarities

□Sample noises (mislabeling) can be mitigated

Benign (white) samples are not needed in learning, but for clusters cleanup (pruning)True positive is easy to explain

• We chose semi-supervised clustering method for our work

Clustering of labeled data

A comparison of the clustering algorithms

· · · · · · · · · · · · · · · · · · ·				
Method name	Parameters	Scalability	Usecase	Geometry (metric used)
<u>K-Means</u>	number of clusters	Very large n_samples, medium n_clusters with <u>MiniBatch</u> <u>code</u>	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
Affinity propagation	damping, sample preference	Not scalable with n_samples	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with n_samples	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium n_samples, small n_clusters	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large n_samples and n_clusters	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large n_samples and n_clusters	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large n_samples, medium n_clusters	Non-flat geometry, uneven cluster sizes	Distances between nearest points
Gaussian mixtures	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers
<u>Birch</u>	branching factor, threshold, optional global clusterer.	Large n_clustersand n_samples	Large dataset, outlier removal, data reduction.	Euclidean distance between points

A comparison of the clustering algorithms

Clustering Method

- Coarse-grained grouping by malware families
- Fine-grained clustering within families based on URL parameters
 □Feature vectors
- Clustering algorithm

DSimilarity, Distance (or cost) function $D_{i,j} = f(V_i, V_j)$

Algorithm: K-mean, DBScan

□ Minimal number of URLs in a cluster and distance threshold

Auto removal of noises

Cluster centroid representation (picking an exemplar)

Clustering Method

- Cross-family clusters merge
 - The intuition: code re-use in malware
 - □ For collection of cluster exemplars, do a second-round clustering.
 - □For example, f1c1, f1c2, f2c1, f3c1, f4c1, f4c2, are six clusters from families f1, f2, f3, f4.
 - □Cluster centroids can be further merged to (f1c1, f2c1, f4c2), (f1c2), (f3c1, f4c1) three clusters
 - □ Reduce redundancy

Cross-Family Cluster Merging

Clustering Method

- Cluster Cleanup (pruning)
 - Malware network trace may contain benign legitimate URL connections
 - Clusters are pruned with benign URL parameters
 - □ Reduce false positives

Clusters

Again, the same distance function matching method

White Pruned

Clustering Method

- Cluster Categorization
 - □Known URL categories
 - Apply to clusters and mark similar URLs with the same category type
 - □C&C, remote access, advertisement, etc.

Cluster Centroids

Again, the same distance function matching method

Labeled

Labeled Cluster Centroids

Our Work – Data Source

- Malware Samples:
 - Third party vendor, Hillstone Network security lab
 - □Network trace captured in PCAP when run in sandboxes
- Number of samples:
 - □300,000 initial samples, and 10k updates/week
 - □~4500 malware families (variants are merged)
- Malware samples and their PCAPs are labeled by malware families
 - Gramily variants are grouped under malware family labels.
 - Collect HTTP traces from PCAPs. HTTPs inherit malware family labels.

Our Work – Data Source

- Clustering based on URL parameters
 HTTP GET
- HTTP POST method

□Some suspicious content is stored in http post body

□ Private format/encoding

□Some are simple key: value pairs

Analyze POST body content for certain patterns. And treat them the same as parameters

Our Work – Data Source

HTTP POST Examples

Dhttp://imp.myappz02.com/impression.do

event=loader_start&implementation_id=min.0.0.30&user_id=ef2443da-1705-4437-8280-878105582da1&adprovider=treasure&source=RedUKBeaconRON&page=Workstation

Dhttp://stan.mxp30.com/__dmp__/

data={"msg":";PAYLOAD NOT FOUND","url":"PAYLOAD NOT FOUND","WaitCreateFile":"","MethodTrace":"returnNewPayload","Language":"Chinese","jscrip t":"","jscript":"Notfind","Ino":"","version":"1.7.7","ieversion":"6.0.2900.5512","trace":"U3Rhcn RTZXRVbmhhbmRsZWRFeGNIcHRpb25GaWx0ZXJnZXQgcGF5bG9hZFBBWUxPQUQgTk9UIEZPVU 5E","osname":"Microsoft Windows XP Professional 5.1.2600 x86","av":"","method":"PAYLOAD NOT FOUND"}

Our Work – Process Flow

Experimental Results

- ~9000 clusters. Each cluster is represented as a sparse vector
- 84% detection rate on ~950K malicious URLs
 Some malicious URLs are indistinguishable from benign URLs
- Malware family coverage

□ More than 3000 families of malware PCAPs

□About half have data for clustering

Cluster model covers more than 95% of the 1500 families

• New malware variant detected

□ Malware variant detected in deployment before our partner reported data to us

Example of Detection Result

RiskWare[Downloader]/Win32.Donex

Detection Time: 2016/07/13 11:42:34

Domain: api.down.72zx.com

URI:

/xml/18?winver=6.2&sdsoft=4&webid=18&softid=29621&ver=1.3.1.14&usesnum=1&mac=vYXbi0zPoms%25 3D&filename=office2015%25BC%25A4%25BB%25EE%25B9%25A4%25BE%25DF(KMSpico)+_18@29621.exe& errcode=0&userev=0&rnd=4124

Known URL:

http://api.down.72zx.com/xml/34?winver=5.1&sdsoft=8&webid=34&softid=0&ver=1.3.1.14&usesnum=1&ma c=zPKt%252BjWzpGo%253D&filename=C6F5BAD88E23D89DE798C4D6FFCFA789.BC492044&errcode=0&user ev=0&rnd=9324

Beyond the Malicious URL: Limitations

Sample

□Sandbox execution and evasive technologies

Label accuracy

□Lack large collection of clean benign traffics

Malware

□Non-HTTP

□Non-URL Para. (JSON/Restful API, Post)

Encryption

False Positive/Negative

Intelligent correlation

- The suspicious threat (malicious URL) accuracy
 - In a single HTTP packet, suspicious detection models include HTTP Method, Headers, User-Agent, Path, and host domain name
 - Do the correlation
 - □Spatial correlation: Within the network, check the number of the same incidents in a period of time
 - E.g. more than 20 hosts made the same connection in one hour

The better security arch.

- Multiple defender architecture
 - □NGFW, WAF and other traditional defenders
 - Intelligent: ML and Data Analysis for DGA, Malicious URL and abnormal network behavior

The better security arch.

- Multiple intelligent correlation: correlation and scoring of
 - a cyber kill stage
 - among the cyber kill stages for a Host
 - □ among Hosts and global scope
 - updatable correlation models

Hos	t Name/IP: 10.210.3.189						Risk Level	Certainty
Ope	erating System:							
Acti							Medium	75%
-	ve. 📑 macave							
Zon	e: vpn							
	hain Threats Mitigation							
			Mo	netization				
	Initial Europait	Delivery		C*C	Internal	Lateral	Fufiltration	
	Initial Exploit	Delivery]	C&C	Internal Recon	Lateral Movement	Exfiltration	
	Initial Exploit	Delivery]	C&C	Internal Recon	Lateral Movement	Exfiltration	
	Initial Exploit	Delivery	Severity	C&C	Internal Recon Source	Lateral Movement Destination	Exfiltration Detected at	Status
1	Initial Exploit	Delivery Type Malware - Grayware	Severity	C&C Certainty 100%	Internal Recon Source	Lateral Movement Destination	Exfiltration Detected at 2017/01/26 08:42:10	Status
1 2	Initial Exploit Name The Domain Name of DNS Re The Domain Name of DNS Re	Delivery Type Malware - Grayware Malware - Grayware	Severity High High	 C&C Certainty 100% 100% 	Internal Recon Source 208.201.224.11 8.8.8.8	Lateral Movement Destination & 10.210.3.189 & 10.210.3.189	Exfiltration Detected at 2017/01/26 08:42:10 2017/01/26 08:42:06	Status Detected Detected
1 2 3	Initial Exploit Name The Domain Name of DNS Re The Domain Name of DNS Re High Frequency DNS Query	Delivery Type Malware - Grayware Malware - Grayware Attack - Suspicious	Severity High High High	 C&C Certainty 100% 100% 70% 	Internal Recon Source 208.201.224.11 8.8.8.8 WUEZHANG-SZ(10	Lateral Movement Destination & 10.210.3.189 & 10.210.3.189 = 119.28.48.212	Detected at 2017/01/26 08:42:10 2017/01/26 08:42:06 2016/06/15 04:03:00	Status Detected Detected
1 2 3 4	Initial Exploit Name The Domain Name of DNS Re The Domain Name of DNS Re High Frequency DNS Query Hidden DNS Tunnel	Delivery Type Malware - Grayware Malware - Grayware Attack - Suspicious Attack - Suspicious	Severity High High High High	 C&C Certainty 100% 100% 70% 50% 	Internal Recon Source 208.201.224.11 8.8.8.8 WUEZHANG-SZ(10 3.10.210.3.189	Lateral Movement Destination 	Detected at 2017/01/26 08:42:10 2017/01/26 08:42:06 2016/06/15 04:03:00 2016/05/24 19:40:00	Status Detected Detected Detected Detected
1 2 3 4 5	Initial Exploit Name The Domain Name of DNS Re The Domain Name of DNS Re High Frequency DNS Query Hidden DNS Tunnel Hidden DNS Tunnel	Delivery Type Malware - Grayware Malware - Grayware Attack - Suspicious Attack - Suspicious Attack - Suspicious	Severity High High High High High	 ▼ Certainty 100% 100% 70% 50% 50% 	Internal Recon Source 208.201.224.11 8.8.8.8 WUEZHANG-SZ(10 3.10.210.3.189 3.10.210.3.189	Lateral Movement Destination 0.210.3.189 10.210.3.189 119.28.48.212 101.226.11.38 101.226.11.34	Detected at 2017/01/26 08:42:10 2017/01/26 08:42:06 2016/06/15 04:03:00 2016/05/24 19:40:00 2016/05/24 19:35:00	Status Detected Detected Detected Detected Detected
1 2 3 4 5 6	Initial Exploit Name The Domain Name of DNS Re The Domain Name of DNS Re High Frequency DNS Query Hidden DNS Tunnel Hidden DNS Tunnel Hidden DNS Tunnel Hidden DNS Tunnel	Delivery Type Malware - Grayware Malware - Grayware Attack - Suspicious Attack - Suspicious Attack - Suspicious Attack - Suspicious	Severity High High High High High High	 ▼ Certainty 100% 100% 70% 50% 50% 76% 	Internal Recon Source 208.201.224.11 8.8.8 WUEZHANG-SZ(10 3. 10.210.3.189 3. 10.210.3.189 3. 10.210.3.189	Lateral Movement Destination	Exfiltration Detected at 2017/01/26 08:42:10 2017/01/26 08:42:06 2016/06/15 04:03:00 2016/05/24 19:40:00 2016/05/24 19:35:00 2016/04/14 01:49:00	Status Detected Detected Detected Detected Detected
1 2 3 4 5 6 7	Initial Exploit Name The Domain Name of DNS Re The Domain Name of DNS Re High Frequency DNS Query Hidden DNS Tunnel Hidden DNS Tunnel High Frequency DNS Query Suspicious Encrypted Channel	Delivery Type Malware - Grayware Malware - Grayware Attack - Suspicious Attack - Suspicious Attack - Suspicious Attack - Suspicious Malware - Riskware	Severity High High High High High High	 ▼ Certainty 100% 100% 50% 50% 76% 27% 	Internal Recon Source 208.201.224.11 8.8.8.8 VUEZHANG-SZ(10 10.210.3.189 3.10.210.3.189 3.10.210.3.189 10.210.3.189 140.205.152.166	Lateral Movement Destination	Exfiltration Detected at 2017/01/26 08:42:10 2017/01/26 08:42:06 2016/06/15 04:03:00 2016/05/24 19:40:00 2016/05/24 19:35:00 2016/04/14 01:49:00 2016/08/03 02:02:00	Status Detected Detected Detected Detected Detected Detected
1 2 3 4 5 6 7 8	Initial Exploit Name The Domain Name of DNS Re The Domain Name of DNS Re High Frequency DNS Query Hidden DNS Tunnel Hidden DNS Tunnel High Frequency DNS Query Suspicious Encrypted Channel The TTL of DNS Response Is 0	Delivery Type Malware - Grayware Malware - Grayware Attack - Suspicious Attack - Suspicious Attack - Suspicious Attack - Suspicious Malware - Riskware Malware - Grayware	Severity High High High High High High Medum Medum	 ▼ Certainty 100% 100% 70% 50% 50% 76% 27% 90% 	Internal Recon Source 208.201.224.11 8.8.8.8 YUEZHANG-SZ(10 10.210.3.189 3.10.210.210.3.189 3.10.210.210.200,200,200,200,2	Lateral Movement Destination 0.210.3.189 10.210.3.189 119.28.48.212 101.226.11.38 101.226.11.34 208.201.224.11 XUEZHANG-SZ(10 XUEZHANG-SZ(10 XISONG-PC(10.21	Detected at 2017/01/26 08:42:10 2017/01/26 08:42:06 2016/06/15 04:03:00 2016/05/24 19:40:00 2016/05/24 19:35:00 2016/05/24 19:35:00 2016/04/14 01:49:00 2016/08/03 02:02:00 2016/08/03 02:02:00	Status Detected Detected Detected Detected Detected Detected Detected
1 2 3 4 5 6 7 8	Initial Exploit Name The Domain Name of DNS Re The Domain Name of DNS Re High Frequency DNS Query Hidden DNS Tunnel Hidden DNS Tunnel High Frequency DNS Query Suspicious Encrypted Channel The TTL of DNS Response Is 0	Delivery Type Malware - Grayware Malware - Grayware Attack - Suspicious Attack - Suspicious Attack - Suspicious Attack - Suspicious Malware - Riskware Malware - Grayware	Severity High High High High High High Medium	 Cac Certainty 100% 100% 50% 50% 50% 76% 27% 90% 	Internal Recon Source 208.201.224.11 8.8.8.8 VUEZHANG-SZ(10 10.210.3.189 10.210.3.189 10.210.3.189 10.210.3.189 10.210.3.189 10.210.3.189 10.210.3.189 10.210.3.189 10.210.3.189 10.210.3.189	Lateral Movement Destination 	Detected at 2017/01/26 08:42:10 2017/01/26 08:42:06 2016/06/15 04:03:00 2016/05/24 19:40:00 2016/05/24 19:35:00 2016/05/24 19:35:00 2016/04/14 01:49:00 2016/08/03 02:02:00 2016/07/13 22:17:21	Status Detected Detected Detected Detected Detected Detected Detected

Black Hat Sound Bytes

- Analysis of HTTP traffic has benefit in detecting compromised hosts.
- A novel way to extract URL parameter feature dimensions, and an unique method to transform infinite dimensions into limited feature space.
- Semi-supervised clustering method with filtering, post-processing and correlation demonstrates strength in precision, model size, variance coverage.

Thank you!