



Selecting Features to Classify Malware

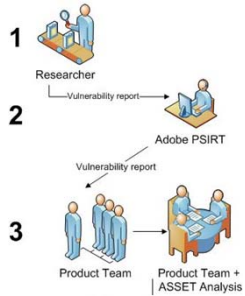
Karthik Raman | Security Researcher, Adobe Secure Software Engineering Team (ASSET)



© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

About Us

- Adobe PSIRT = Adobe Product Security Incident Response Team
- PSIRT is part of ASSET, the Adobe Secure Software Engineering Team



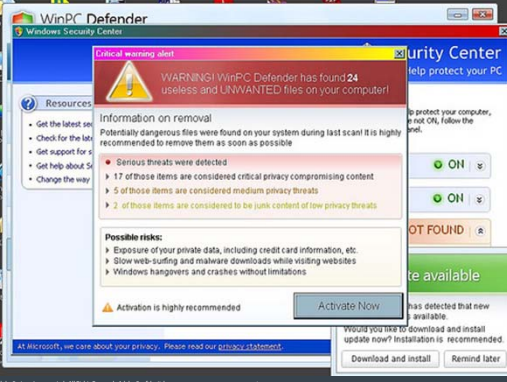
© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

What Adobe PSIRT Does, contd.

- Work with product teams to create fixes
- Work with researchers to verify fixes
- Publish bulletins
- Drive Adobe's involvement in MAPP

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

Did Malware Ever Infect your Computer(s)?



© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

Agenda

- Part I: What is the Malware Menace?
 - "How did I just get infected?"
- Part II: Using Machine Learning For Malware Classification

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

The Story of Mass Malware

- Regular Web site compromised
 - Whistleblowing Site Cryptome.org Infected With Drive-by Exploits**
 - By Lucian Constantin, IDG News
 - Cryptome.org, a website dedicated to disclosing confidential information, was compromised last week and was used to infect PCs running Internet Explorer through drive-by exploits.
- Malicious site visited because of Search Engine Optimization (SEO)

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

Malware Testing: Quality Assurance

SHA256: 7e3669a58bb7830e55e7d20b54bcf3b8053bd6e07f0c1655e247260f88c99e

SHA1: d25d94b2b1d5991f3beac2d049f00436dd1692

MD5: 66d4d07bc10a2db402fc4b69621580c6

File size: 129.9 KB (133065 bytes)

File name: 66d4d07bc10a2db402fc4b69621580c6

File type: Win32 EXE

Detection ratio: **28 / 42**

Analysis date: 2012-02-07 15:05:10 UTC (1 week, 1 day ago)

Malware Testing: Quality Assurance

Detection ratio: **28 / 42**

Malware Obfuscation: Zeus/Zbot

Information
Builder

Builder

Config and loader building

Source config file:
C:\zbot\config.txt

Output

Loading succeeded!
Loading config from file 'C:\zbot\config.txt'...
Loading succeeded!
Building bot file...
botnet== default --
timer_config=360000ms, 6000ms
timer_loggs=60000ms, 60000ms
timer_stats=120000ms, 60000ms
url_config=http://localhost/test3/bot/cfg.bin
url_compid=http://localhost/test3/bot/php/encryption_key=OK
Build succeeded!

Malware Obfuscation: Packers in the House

PolyPack:
An Automated Online Packing Service for Optimal Antivirus Evasion

Jon Oberheide, Michael Bailey, Farnam Jahanian
Electrical Engineering and Computer Science Department
University of Michigan, Ann Arbor, MI 48109
{jonojono, mibailey, farnam}@umich.edu

We show that PolyPack provides 258% more effective evasion of antivirus engines than using an average packer and out-evades the best evaluated packer (Themida) for over 40% of the binary samples.

Automation Cycle

Obfuscation, Testing, Release

Detection ratio: 28 / 42

```

[... obfuscated code snippet ...]

```

What Users Suffer

Privacy Protection
Full PC Scan

Scan is being performed

File Name	Malware Name
C:\Windows\System32\... (FAKE)	Infected: WS2/ChM/Pwn/PROXY/Server
C:\Windows\System32\... (FAKE)	Infected: Mal/Genetic_A/Trojan/Agent
C:\Windows\System32\... (FAKE)	Infected: WS2/ChM/Pwn/PROXY/Server
C:\Windows\System32\... (FAKE)	Infected: Mal/Genetic_A/Trojan/Agent

What Users Suffer

The screenshot shows the Windows Security Center interface. It features a 'Resources' sidebar on the left and a main area titled 'Security essentials'. Under 'Security essentials', there are three sections: 'Firewall' (set to OFF), 'Automatic updates' (set to ON), and 'Virus Protection' (set to OFF). A large, semi-transparent red watermark with the word 'FAKE' is overlaid diagonally across the center of the image.

What Users Suffer

The screenshot shows the Internet Security Guard interface. It displays a 'Quick scanning' window with the text 'Scanning... Please wait' and 'Fast scan of the most typical places where viruses store their files. Scans do not take much time and is good for everyday use.' Below this, it lists 'Threats Found: 11' and provides details for a scan: 'Scan started: 4:30:35 PM', 'Scan duration: 00:00', and 'Objects scanned: 0'. A large, semi-transparent red watermark with the word 'FAKE' is overlaid diagonally across the center of the image.

What Dated AV Really Means

A large, bright green thought bubble is centered on a dark background. Inside the bubble, the text reads: 'Malware SDLC Outpaces Antivirus SDLC'. The bubble has a tail pointing towards the bottom left.

Making AV Current



- Automate everything
- Published research discusses
 - Static detection
 - Dynamic detection
 - Cloud detection
- What else?

Got Machine Learning?

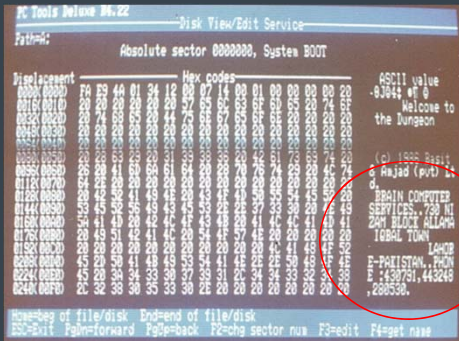
A detailed, close-up image of a metallic, humanoid robot. The robot has a complex, segmented body with various joints and mechanical details. It appears to be a character from a science fiction or action movie.

What is a Virus?

- Fred Cohen's definition
 - A program that can 'infect' other programs by modifying them to include a possibly evolved copy of itself
- Peter Szor's definition
 - A program that recursively and explicitly copies a possibly evolved copy of itself

Down (Computer) Memory Lane



Blasted Worms



A Trojan Horse



Trojan Horse Malware



Agenda

- Part I: What is the Malware Menace?
 - "How did I just get infected?"
- Part II: Using Machine Learning For Malware Classification

Scoping of Research

- Classification of Polymorphic Malware
 - Multiple variants
 - Do not infect other programs
- Examples
 - Backdoors
 - Downloaders
 - Remote Administration Tools
- Infectors and packers out of scope

Why is Polymorphic Malware Important?

Trojans Make Up 80 Percent Of All New Malware

China has the most infected PCs in the world, and 6 million new pieces of malware appeared in Q1 2012, new PandaLabs report says

May 08, 2012 | 11:46 PM |

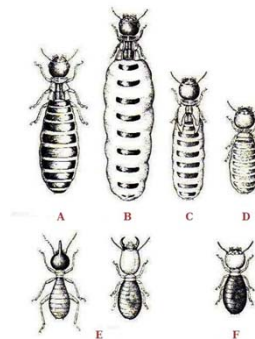
By Kelly Jackson Higgins
Dark Reading

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

31



Polymorphism in Biology



© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

32



Spot the Polymorphic Cylons



© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

33



After Classification

- Clustering
- Detection
- Cleaning for infected files
- Deletion

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

34



Strategies for Polymorphic Malware Classification

- Rieck et al.: mine malware behavior on sandboxed system
 - Machine learning approach
- Karim et al., Venable et al.: search for similar malware
 - Machine learning/Search engine approach

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

35



Strategies for Polymorphic Malware Detection

- Karim et al.: build malware phylogeny
 - Bioinformatics approach
- Karim et al.: use *n-perms* to build malware phylogeny
 - Machine learning/Search engine approach
- Kruegel et al.: fingerprinting malware using CFGs
 - Structural similarity approach
- Vinod P. et al.: analyze CFG and Basic Blocks
 - Machine learning/Search engine approach
- Various academics: normalize code
 - Formal methods approach

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

36

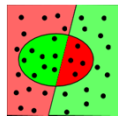


Strategies for Polymorphic Malware Clustering

- Jang, Brumley; Wicherski: fingerprint malware
 - Structural similarity approach
- Gurrutxaga et al.: apply distance algorithms
 - Structural similarity/Machine learning approach
- Bayer et al.: derive behavioral profile (ANUBIS)
 - Machine learning approach

Applying Machine Learning (ML)


- Steps:
 1. Extract features
 2. Train models using ML algorithms
 3. Use models as classifiers
 4. Use models to classify unknown files as 0 or 1



- Started with 600 features

What are the Features?

- EXE and DLL are PE file formats



Microsoft Portable Executable and Common Object File Format Specification

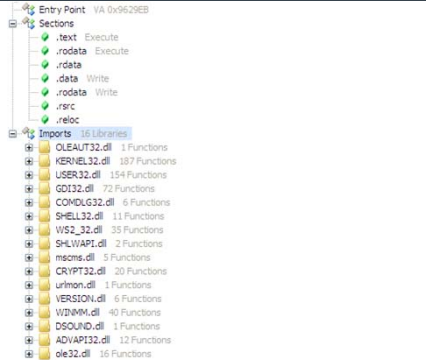
Revision 8.2 – September 21, 2010

What is the PE Format?

Structures

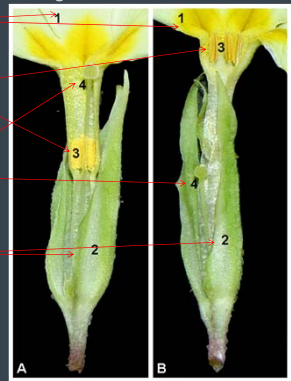
FILETIME	struct(8)
GUID	struct(16)
IMAGE_DATA_DIRECTORY	struct(8)
VirtualAddress	
Size	
IMAGE_DOS_HEADER	struct(6-9)
IMAGE_FILE_HEADER	struct(20)
Machine	
NumberOfSections	
TimeDateStamp	
PointerToSymbolTable	
NumberOfSymbols	
SizeOfOptionalHeader	
Characteristics	
IMAGE_NT_HEADERS	struct(2-8)
IMAGE_OPTIONAL_HEADER	struct(22-4)
IMAGE_SECTION_HEADER	struct(40)
SYSTEMTIME	struct(16)

What is the PE Format?



- Imports: 16 Libraries
 - OLEAUT32.dll 1 Functions
 - KERNEL32.dll 187 Functions
 - USER32.dll 154 Functions
 - GDI32.dll 72 Functions
 - COMDLG32.dll 6 Functions
 - SHELL32.dll 11 Functions
 - WIS2_32.dll 35 Functions
 - SHLWAPI.dll 2 Functions
 - mscms.dll 5 Functions
 - CRYPT32.dll 20 Functions
 - urlmon.dll 1 Functions
 - VERSION.dll 6 Functions
 - WINMM.dll 40 Functions
 - DSOUND.dll 1 Functions
 - ADVAPI32.dll 12 Functions
 - ole32.dll 16 Functions

Features Illustrated using Primula



- Corolla
- Stamen
- Pistil
- Calix

Why Fewer Features?

Why are fewer features better than more features?

```

    graph TD
      Input[Input] --> Algorithm((Algorithm))
      Algorithm --> YES[YES]
      Algorithm --> NO[NO]
  
```

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

Less is More

- Irrelevant features negatively affect learning
- Using fewer features...
 - Improves algorithm performance
 - Represents problem better
 - Lets user focus on important variables

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

Related Work

- Mining n-grams (Siddiqui et al.) → 94% accuracy
- Multiple algorithms (Schultz et al.) → 97.76% accuracy
- Multiple algorithms, 189 features (Shafiq et al.) → 99% accuracy
- Association mining (Ye et al.) → 92% accuracy
- SVM on program strings (Ye et al.) → 93.8% accuracy
- Key Questions
 - Which features were used and why?
 - What are the minimum features for good classification?

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

Contributions

- Excellent classification using **seven** features
- Another layer to existing antivirus technology
- Still need:
 - Unpackers and deobfuscators
 - Clustering, detection, cleaning, deletion, etc.

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

System Diagram

```

    graph LR
      Dataset[Dataset] --> Parser[Parser]
      Parser --> Classifier[Classifier]
      Classifier --> Evaluator[Evaluator]
      Classifier --> Model[Model]
      Evaluator --> Model
      Model --> Evaluation[Evaluation]
  
```

PE Parser:
pedump tool

pefile is a Python module to read and work with PE (Portable Executable) files

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.


The Haystack (Dataset)

- 100,000 pieces of malware
- 16,000 clean programs
- 645 initial features
 - Structures in PE file format
 - Some calculated features
- See M. Pietrek's "An In-Depth Look into the Win32 Portable Executable File Format"
 - <http://msdn.microsoft.com/en-us/magazine/cc301805.aspx>

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

Classifier and Evaluator: Say Hello to WEKA

Machine Learning Toolkit
<http://www.cs.waikato.ac.nz/ml/weka/>



Scriptable!

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 49


Preliminary Results

- Six numeric machine-learning algorithms
- Experiment I with 645 & Experiment II with 100 features

Check the Classification

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 50

Wait a Minute



© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 51


What Pretty Features You Have

Feature	Accuracy
DebugSize	0.9234
DebugRVA	0.9224
ImageVersion	0.8898
OperatingSystemVersion	0.8850
SizeOfStackReserve	0.8837

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 52

Reduced Feature Set Selection

- Which PE structure does a feature belong to?
- Created seven buckets

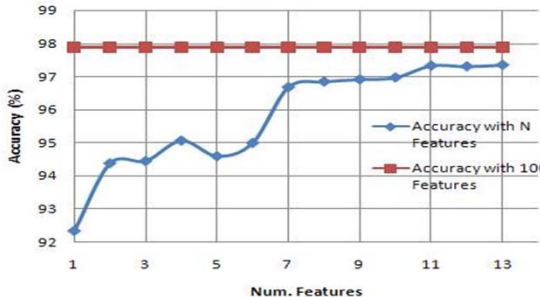


- Algorithm - Start with bucket 1
 - Run ML algorithms on current feature set
 - Add next best feature, modulo 7, to feature set
 - Return to step 1.

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 53

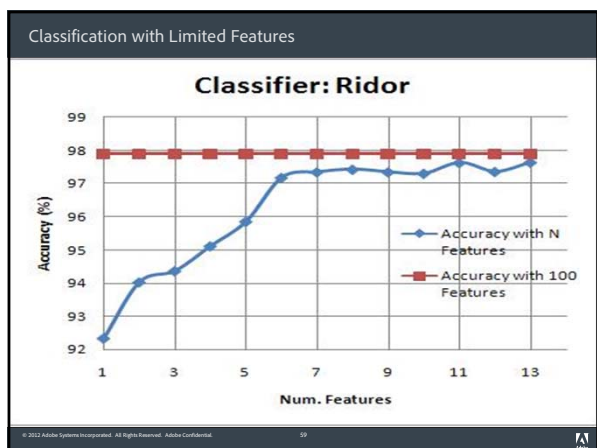
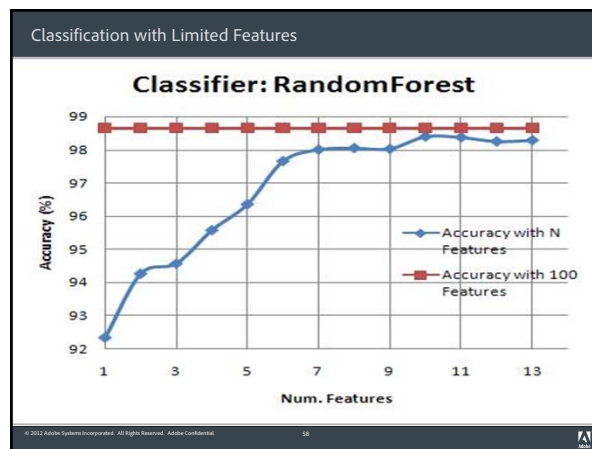
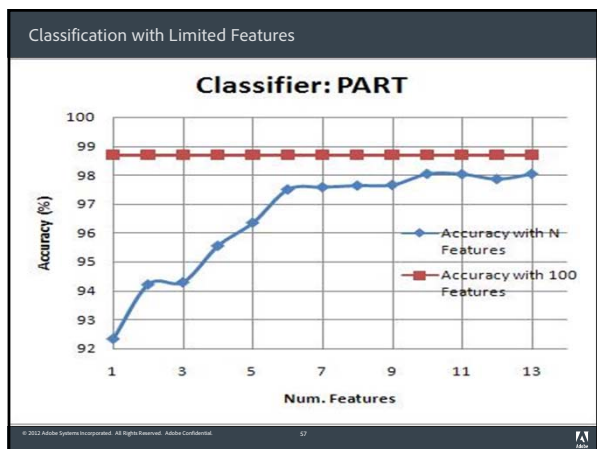
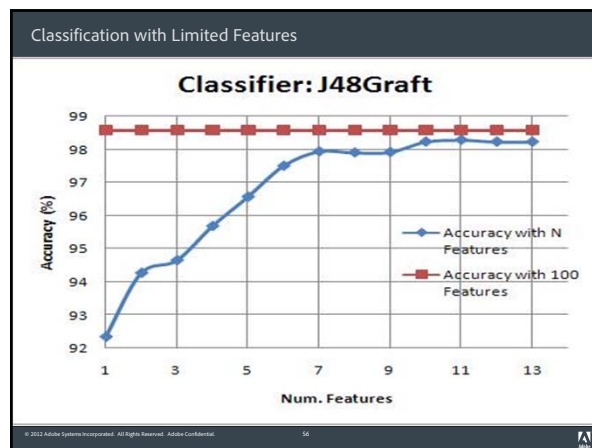
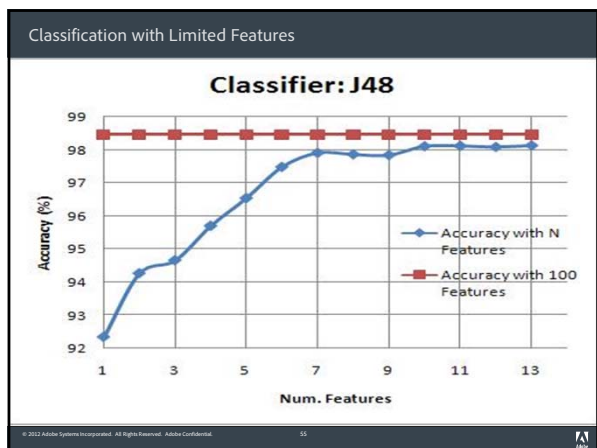
Classification with Limited Features

Classifier: IBk



Num. Features	Accuracy with N Features (%)	Accuracy with 100 Features (%)
1	92.5	98.0
2	94.5	98.0
3	94.5	98.0
4	95.0	98.0
5	94.5	98.0
6	95.0	98.0
7	96.8	98.0
8	96.8	98.0
9	96.8	98.0
10	97.0	98.0
11	97.2	98.0
12	97.2	98.0
13	97.5	98.0

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 54



More Results

- Six numeric machine-learning algorithms
- Experiment III with 7 features

Check the Classification

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

Results

- Best classifier: RandomForest
 - 98.21% accuracy
 - 6.7% false positive rate
- Why did seven features work so well?
 - Algorithms picked most discriminating features first

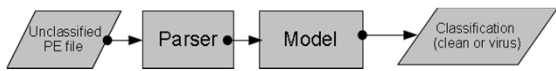
© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 61

Results

- The Seven
 - DebugSize, ImageVersion, latRVA, ExportSize, ResourceSize, VirtualSize2, NumberOfSections
- DebugSize
 - Denotes the size of the debug-directory table
 - Malware vs. clean file discrimination: ...
- ImageVersion
 - Denotes the version of the file
 - Malware vs. clean file discrimination: ...

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 62

How Do I Use That ML Model?



- Desktop antivirus
 - Consolidate signature databases
 - Classifiers in **least aggressive** mode
- Cloud antivirus
 - Quick detection of mass malware variants
 - Classifiers in **more aggressive** mode
- Gateway antivirus
 - Stop worms from spreading
 - Classifiers in **most aggressive** mode

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 63

Tool available at <http://sourceforge.com/adobe/malclassifier>

```

# Program to classify malware into
# 0 = CLEAN
# 1 = DIRTY
# UNKNOWN
##### Results on ~130000 dirty, ~ 16000 clean files:
(False Positives, True Negatives, True Positives, rates
@J48
FP      TN      TP      FN      TP Rate      FP Rate      Accuracy
7683   37171   130302  3451   0.97419871   0.171289071  0.937662018

@J48Graft
FP      TN      TP      FN      TP Rate      FP Rate      Accuracy
6780   38074   129087  4666   0.965114801  0.151157087  0.935915166

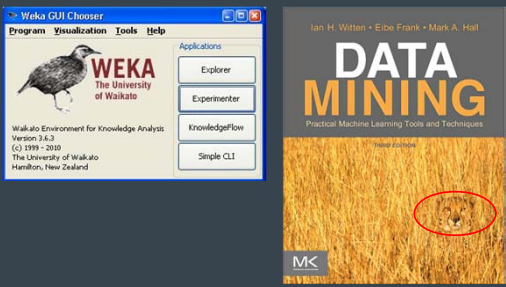
@PART
FP      TN      TP      FN      TP Rate      FP Rate      Accuracy
7074   36492   125060  9412   0.930007734  0.162374329  0.907401791

@Ridor
FP      TN      TP      FN      TP Rate      FP Rate      Accuracy
7390   37935   114194  20930  0.845105237  0.163044677  0.843058149
    
```

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 64

Closing Remarks


Get WEKA (free), get the official book (not free but affordable).



© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 65

Closing Remarks

- The Arms Race
 - "Bad guys can also use machine learning."



- Could ML buy the good guys more time?
- Could self-training ML models strain human analysts less?

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential. 66

Closing Remarks

- The Cost of FPs vs. FNs
 - ML models without tuning can't be used in production
 - Adjust models by adding costs of FPs into probabilities used by algorithms
 - Everyone's calculation is different
- Protecting the User's Privacy
 - What features are you extracting?
 - Is this a development box?
 - Research privacy-preserving data mining

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

67



Further Reading

- M. Siddiqui, M. C. Wang, and J. Lee. Detecting trojans using data mining techniques. In D. M. A. Hussain, A. Q. K. Rajput, B. S. Chowdhry, and Q. Ge, editors, IMTIC, volume 20 of Communications in Computer and Information Science, pages 400-411. Springer, 2008.
- M. G. Schultz, E. Eskin, E. Zadok, and S. J. Stolfo. Data mining methods for detection of new malicious executables. In Proceedings of the 2001 IEEE Symposium on Security and Privacy, pages 38, Washington, DC, USA, 2001. IEEE Computer Society.
- M. Z. Shafiq, S. M. Tabish, F. Mirza, and M. Farooq. Pe-miner: Mining structural information to detect malicious executables in realtime. In Proceedings of the 12th International Symposium on Recent Advances in Intrusion Detection, RAID '09, pages 121-141, Berlin, Heidelberg, 2009. Springer-Verlag.
- Y. Ye, L. Chen, D. Wang, T. Li, Q. Jiang, and M. Zhao. Sbmids: an interpretable string based malware detection system using svm ensemble with bagging. Journal in Computer Virology, 5(4):283-293, 2009.
- Y. Ye, D. Wang, T. Li, and Ye. Imds: Intelligent malware detection system. In Proceedings of ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD 2007), 2007.
- Dan Guido's Exploit Intelligence Project, <http://www.isecpartners.com/storage/docs/presentations/EIP-final.pdf>
- M. Merkel, T. Hoppe, C. Kraetzer, J. Dittman. Statistical Detection of Malicious PE-Executables for Fast Offline Analysis TC11 Conference on Communications and Multimedia Security (CMS 2010)

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

68



Further Reading



TALKS

Mikko Hypponen: Fighting viruses, defending the net

TEDGlobal 2011. Filmed Jul 2011. Posted Jul 2011



http://www.ted.com/talks/mikko_hypponen_fighting_viruses_defending_the_net.html

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

69



References

- Koolkat, <http://www.flickr.com/photos/32936091@N05/3752997536/>
- SANS, <http://isc.sans.edu/diary.html?storyid=4246>
- swankalot, <http://www.flickr.com/photos/swankalot/4335612238/sizes/m/in/photostream/>
- BSOD: http://upload.wikimedia.org/wikipedia/commons/a/a8/Windows_XP_BSOD.png
- AVIRA, http://techblog.avira.com/wp-content/uploads/2010/04/spy_eye.png

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

70



References

- Virustotal, <https://www.virustotal.com/file/7e3669a58bb7830e55e7d2b85a4bcf3b8b53bd6e07cf0c1655e247260f88c59e/analysis/>
- Microsoft, http://www.microsoft.com/security/sir/story/default.aspx#!zbot_works
- Microsoft MPMC, <http://blogs.technet.com/b/mmpc/archive/2012/01/29/when-imitation-isn-t-a-form-of-flattery.aspx>
- PC Magazine, http://www.pcmag.com/slideshow_viewer/0,3253,l%3D205153%26a%3D205149%26p%3D8,00.asp?p=n
- SecurityFocus, <http://www.securityfocus.com/excerpts/2>

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

71



References

- Wikipedia, <http://upload.wikimedia.org/wikipedia/commons/d/da/Brain-virus.jpg>
- Wikipedia, <http://upload.wikimedia.org/wikipedia/commons/8/84/Blaster-virus.jpg>
- darcy m, <http://www.flickr.com/photos/darcym/54086635/>
- darkchacal, <http://www.flickr.com/photos/darkchacal/4252059347/>
- Dark Reading, <http://www.darkreading.com/vulnerability-management/167901026/security/attacks-breaches/240000043/trojans-make-up-80-percent-of-all-new-malware.html>
- Classification, <http://upload.wikimedia.org/wikipedia/commons/d/d1/Binary-classification.svg>

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

72



References

- John Pavelka, <http://www.flickr.com/photos/28705377@N04/4142872268/>
- kmgquidoo, <http://www.flickr.com/photos/38117284@N00/1277420698/>
- LabyrinthX, <http://www.flickr.com/photos/labyrinthx/1955627738/>
- Google Books, http://books.google.com/books/about/Data_Mining.html?id=5FIEAwyn9aoC
- AV Hire Lens, http://www.flickr.com/photos/av_hire_london/5570201239/
- potzuyoko, <http://www.flickr.com/photos/potzuyoko/6549346059/>

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

73



Binary Classification: Cylon or Human?



© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

74



End

QUESTIONS?

kraman@adobe.com

Adobe Malware Classifier:
<http://sourceforge.com/adobe/malclassifier>

© 2012 Adobe Systems Incorporated. All Rights Reserved. Adobe Confidential.

75

