
Machine Learning for Malware Analysis

Mike Slawinski
Data Scientist

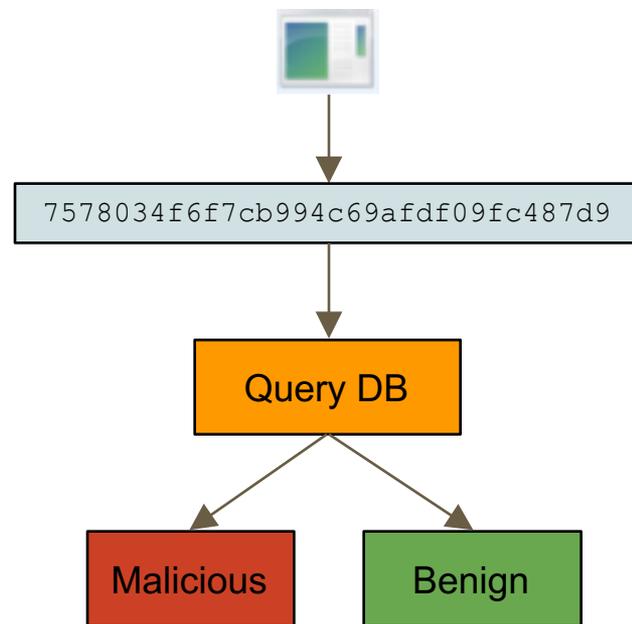


Introduction - What is Malware?

- Software intended to cause harm or inflict damage on computer systems
- Many different kinds:
 - Viruses
 - Trojans
 - Worms
 - Adware/Spyware
 - Ransomware
 - Rootkits
 - Backdoors
 - Botnets
 - ...

Malware Detection - Hashing

- Simplest method:
 - Compute a fingerprint of the sample (MD5, SHA1, SHA256, ...)
- Check for existence of hash in a database of known malicious hashes
- If the hash exists, the file is malicious
- Fast and simple
- Requires work to keep up the database



Malware Detection - Signatures

Look for specific strings, byte sequences, ... in the file.

If attributes match, the file is likely the piece of malware in question

Signature Example

```
93 rule Stuxnet_Malware_3
94 {
95
96     meta:
97         description = "Stuxnet Sample - file -WTR4141.tmp"
98         author = "Florian Roth"
99         reference = "Internal Research"
100        date = "2016-07-09"
101        hash1 = "6bcf88251c876ef00b2f32cf97456a3e306c2a263d487b0a50216c6e3cc07c6a"
102        hash2 = "70f8789b03e38d07584f57581363afa848dd5c3a197f2483c6dfa4f3e7f78b9b"
103
104    strings:
105        $x1 = "SHELL32.DLL.ASLR." fullword wide
106        $s1 = "-WTR4141.tmp" fullword wide
107        $s2 = "-WTR4132.tmp" fullword wide
108        $s3 = "totalcmd.exe" fullword wide
109        $s4 = "wincmd.exe" fullword wide
110        $s5 = "http://www.realtek.com@" fullword ascii
111        $s6 = "{%08x-%08x-%08x-%08x}" fullword wide
112
113    condition:
114        ( uint16(0) == 0x5a4d and filesize < 150KB and ( $x1 or 3 of ($s*) ) ) or ( 5 of them )
115 }
```

Problems with Signatures

- Can be thought of as an overfit classifier
- No generalization capability to novel threats
- Requires reverse engineers to write new signatures
- Signature may be trivially bypassed by the malware author

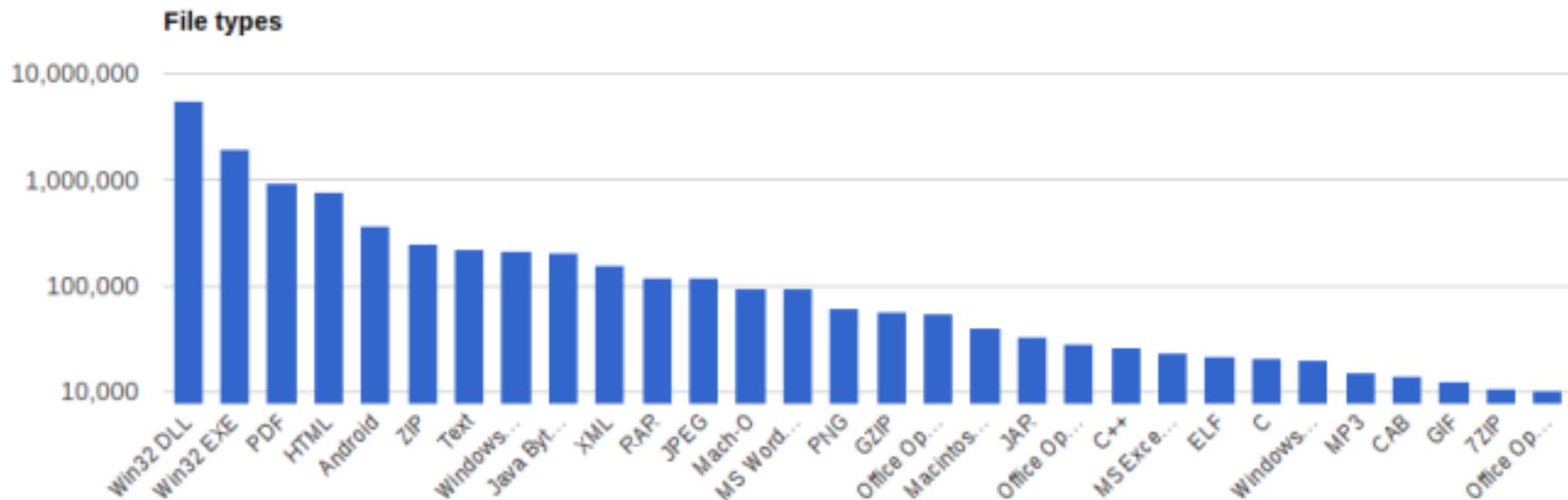
Malware Detection - Behavioral Methods

- Instead of scanning for signatures, examine what the program does when executed
- Very slow - AV must run the program and extract information about what the sample does
- Malicious samples can “run out the clock” on behavior checks

Scaling Malware Detection

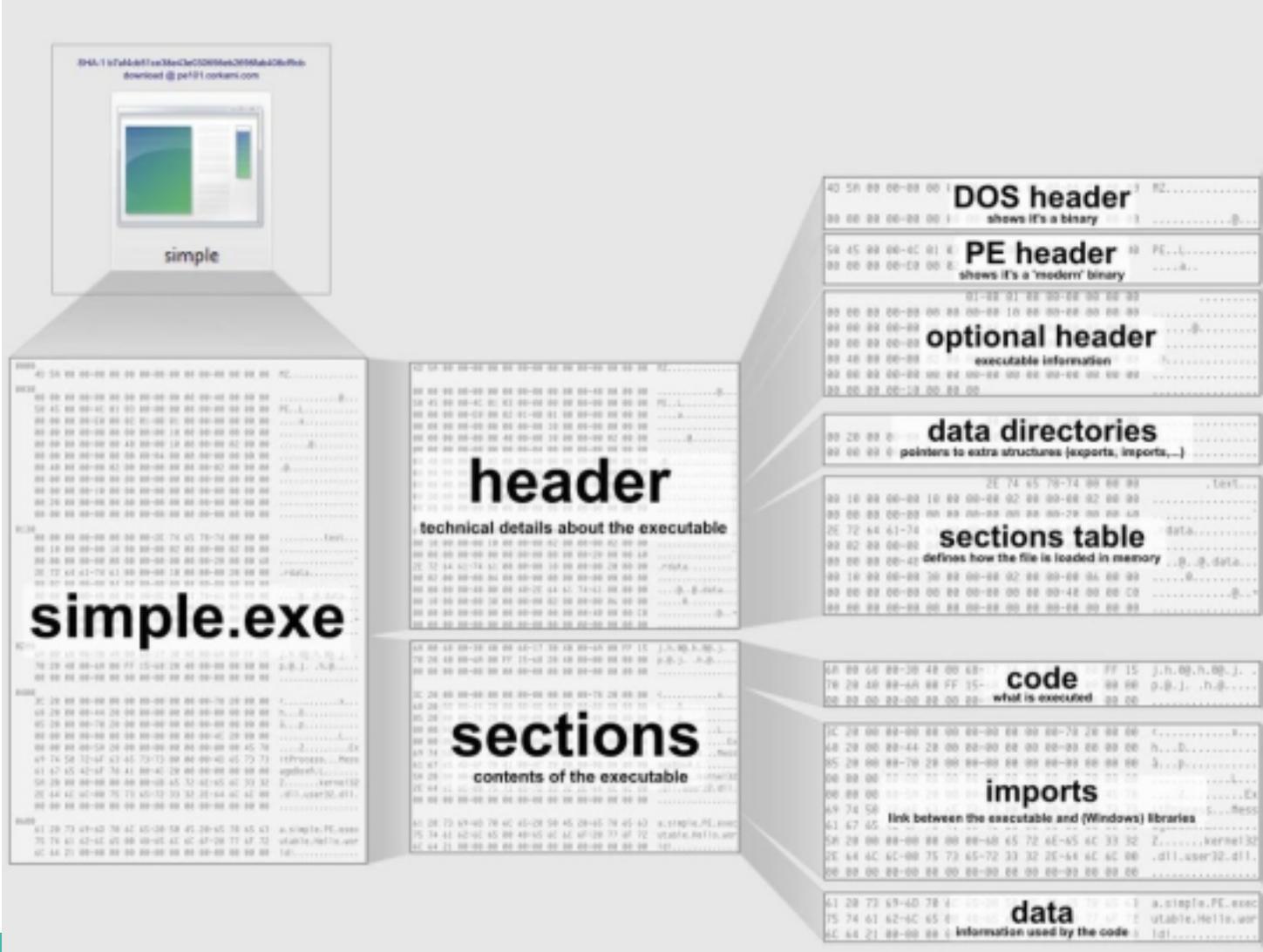
- Previously mentioned approaches have difficulty generalizing to new malware
- New kinds of malware require humans in the loop to reverse-engineer and create new signatures and heuristics for adequate detection
- Can we automate this process with machine learning?

Focus: Windows DLL/EXEs (Portable Executable)



Number of samples submitted to VirusTotal, Jan 29 2017

Portable Executable (PE) Format



Feature Engineering - Static Analysis

- What kinds of features can we extract for PE files?
- Objective: extract features from the EXE without executing anything
- PE-Specific features
 - Information about the structure of the PE file
- Strings
 - Print off all human-readable strings from the binary
- Entropy features
 - Extract information about the predictability of byte sequences
 - Compressed/encrypted data is high entropy
- Disassembly features
 - Get an idea of what kind of code the sample will execute

PE-Specific Features

🔗 FileVersionInfo properties

Copyright	© Microsoft Corporation. All rights reserved.
Product	Microsoft® Windows® Operating System
Original name	NOTEPAD.EXE
Internal name	Notepad
File version	5.1.2600.0 (xpclient.010817-1148)
Description	Notepad

☰ PE header basic information

Target machine	Intel 386 or later processors and compatible processors
Compilation timestamp	2001-08-17 20:52:29
Entry Point	0x00006AE0
Number of sections	3

PE-Specific Features (cont.)

PE sections

Name	Virtual address	Virtual size	Raw size	Entropy	MD5
.text	4096	28018	28160	6.28	ccf25baa681168e6396609387910d90a
.data	32768	7080	1536	1.40	cf692e5fbaebba02c2ad95f4ba0e60be
.rsrc	40960	35144	35328	5.41	c65b2250b6dd3670595004ca95f8f8b3

PE imports

[+] ADVAPI32.dll

[+] COMCTL32.dll

[+] GDI32.dll

[+] KERNEL32.dll

[+] SHELL32.dll

[+] USER32.dll

[+] WINSPOOL.DRV

[+] comdlg32.dll

[+] msvcrt.dll

PE-Specific Features (cont.)

PE imports
[+] ADVAPI32.dll
RegCloseKey
RegSetValueExW
RegQueryValueExA
RegCreateKeyW
RegOpenKeyExA
IsTextUnicode
RegQueryValueExW
[+] COMCTL32.dll
[+] GDI32.dll
[+] KERNEL32.dll
[+] SHELL32.dll
[+] USER32.dll
[+] WINSPOOL.DRV
[+] comdlg32.dll
[+] msvcrt.dll

Feature Engineering - String Features

- Extract contiguous runs of ASCII-printable strings from the binary
- Can see strings used for dialog boxes, user queries, menu items, ...
- Samples trying to obfuscate themselves won't have many strings

```
→ notepad strings Notepad.exe | head -n 25
!This program cannot be run in DOS mode.
Rich
.text
.data
.rsrc
comdlg32.dll
SHELL32.dll
WINSPOOL.DRV
COMCTL32.dll
msvcrt.dll
ADVAPI32.dll
KERNEL32.dll
NTDLL.DLL
GDI32.dll
USER32.dll
@wARAw
j=v?
RegisterPenApp
notepad.chm
hhctrl.ocx
CLSID\{ADB880A6-D8FF-11CF-9377-00AA003B7A11}\InprocServer32
```


Disassembly Features

- Contains information about what will actually execute
- Disassembly is difficult:
 - Hard to get all of the compiled instructions from a sample
 - x86 instruction set is variable-length
 - Ambiguity about what is executed depending on where one starts interpreting the stream of x86 instructions

```
01001000 <.text>:
1001000: 65 1b dd          gs sbb %ebp,%ebx
1001003: 77 9a             ja 0x1000f9f
1001005: 18 dd           sbb %bl,%ch
1001007: 77 ce           ja 0x1000fd7
1001009: 5f             pop %edi
100100a: dd 77 ca       fnsave -0x36(%edi)
100100d: 60            pusha
100100e: df 77 d7       fbstp -0x29(%edi)
1001011: 23 dd          and %ebp,%ebx
1001013: 77 ea         ja 0x1000fff
1001015: 22 dd          and %ch,%bl
1001017: 77 0b         ja 0x1001024
1001019: 58            pop %eax
100101a: dd 77 00       fnsave 0x0(%edi)
100101d: 00 00         add %al,(%eax)
100101f: 00 0d 77 96 71 00 add %cl,0x719677
1001025: 00 00         add %al,(%eax)
1001027: 00 9a 86 c8 77 b7 add %bl,-0x4888377a(%edx)
100102d: 20 ca         and %cl,%dl
100102f: 77 1d         ja 0x100104e
1001031: 87 c8         xchg %ecx,%eax
1001033: 77 6b         ja 0x10010a0
1001035: 2c c7         sub $0xc7,%al
1001037: 77 1e         ja 0x1001057
1001039: 88 c8         mov %cl,%al
100103b: 77 1d         ja 0x100105a
100103d: 51            push %ecx
100103e: c7           (bad)
100103f: 77 68         ja 0x10010a9
1001041: 6a c7         push $0xffffffffc7
```

Difficulties for Static Analysis

- Polymorphic code
 - Code that can modify itself as it executes
- Packing
 - Samples that compress themselves prior to execution, and decompress themselves while executing
 - Can hide malicious behavior in a compressed blob of bytes
 - Can obscure benign code as well
 - Requires expensive implementation of many unpackers (UPX, ASPack, Mew, Mpress, ...)
- Disassembly
 - Malware authors can intentionally make the disassembly difficult to obtain

Modelling - Malicious versus Benign

- Boils down to a binary classification task
- N: hundreds of millions of samples
- P: millions of highly sparse features ($s=0.9999$)



Modelling - Training on ~600 million samples

- Strong preference for minibatch methods and fast, compact models
- Logistic regression works very well
- Neural networks coupled with dimensionality reduction techniques are the workhorse
- Tend to combine lasso, dimensionality reduction, and neural networks

Filesystems – interesting topological structure

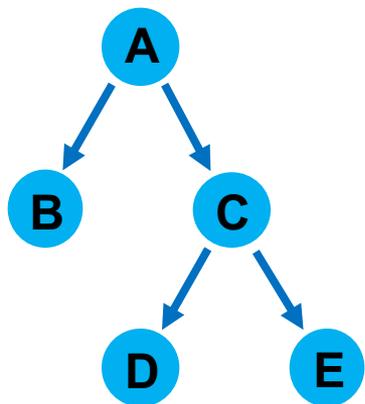
Idea: construct a map which measures the similarity between graphs G and H , which takes into account both the topological differences of the trees **and** the label differences.

$$K: \Gamma \times \Gamma \rightarrow \mathbb{R}$$

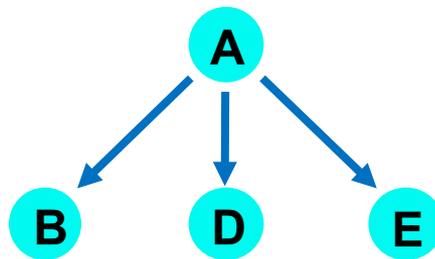
$K(G, H)$ measures the similarity between G and H , taking into account both the topological structure of the trees and their labels.

Upshot: We can measure the similarity between two file systems A and B by measuring the similarity between their labeled tree structure.

Graph Comparison and Vectorization



X



\mathbb{R}

$$\begin{pmatrix} 0 & ab & ac & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & cd & ce \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

X

$$\begin{pmatrix} 0 & ab & ad & ae \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$



\mathbb{R}

Filesystems – interesting topological structure

Can leverage GPU hardware in two ways:

- Kernel computations $K: \Gamma \times \Gamma \rightarrow \mathbb{R}$
- Neural Network training on features derived from these kernels

Upshot: The framing a given problem/procedure in terms of matrix algebra translates to massive computational advantages (GPU).

Selected Hardware

AWS P2 instances - up to
16 NVIDIA K80 GPUs

AWS G3 instance - four
NVIDIA Tesla M60 GPUs



Thank You!

— **Questions?** —
