NEXT-GENERATION OF EXPLOIT KIT DETECTION BY BUILDING SIMULATED OBFUSCATORS

Tongbo Luo & Xing Jin Palo Alto Networks



Agenda

- Background and Motivation
- Challenges
- Our Approach and Lessons We learned
- Result and Our Observation





www.jsDarwin.com

Similarity-based Detection on malicious Exploit Kit JavaScript

Large Number of Samples Required



Exploit Kit Obfuscator





Obfuscator Reverse-Engineering

Obfuscation Engine (variant)

=

"window [" + \$[eval_ctr] + "] [" + \$v[0] +"(" v[1] + "," + \vee [2] + ")" + "]"

> **EVAL** Template (version)



Challenges

1. Code Complexity

```
<script>
function pDA(AUT) { cm }
ZzF ="lnFloa";
function vzLh(Hpv){gEao="oi"+""+"n",a0C=window,SLVmg="e"+"x",
bUdR="R"+"e"+""+"q", FMwVS="o"+""+"u"+""+"te"+"rH
",HjEn="u"+"s"+""+"e"+"r",Wyc="b"+"o",Wbf="e"+"R"+""+"x";var gVa=new
aOC[bUdR+'Exp']('MSIE (\d+\.\d+);');var xngS=navigator[HjEn+'Agent'];var
prC=qVa[(Wbf+'ec').replace("R","")](xnqS);var
$mrh=!prC,ryG;$mrh=!$mrh;if($mrh){ryG=prC[1]}GcSX='.*'+'>(.'+'*?)'+'<\/p';</pre>
if( com ) {GcSX='.'+ '?>('+'.*'+'?)'+'<\/p'}else { com } VdZW =</pre>
["r"+"e"+"p"+"l"+"a"+"c"+"e"]["j"+gEao]();var UkP=document[Wyc+'dy'],
yh$f=UkP[FMwVS+'TML'],wvnV="", El=new a0C[bUdR+'Exp'](ZzF+Hpv+GcSX,'gi');wvnV=
El[SLVmg+'ec'] (yh$f);wvnV=wvnV [1];var WMK =
wvnV[VdZW](/[-]/q,"/")[VdZW](/[_]/q,"+"); function asafas(a){return a}function
GKcx(s) { m }if(!pDA(GKcx(WMK)))alert("Error");}vzLh("dNTvXuVoyq");
vzLh("xWgaeJKYcclAu");</script>
```

Hundred lines of code Random variables





2. Data Complexity



Big data set (~20000 Samples over 2 year period) Mixed versions and variants



Overview of Our Approach





JavaScript Normalization



- Ignore Superficial Obfuscation (e.g. Randomized variable names)



Statistics

Total Number of samples we collected



networks



Statistics on Normalized Samples



Clustering

- Goal: Cluster Samples based on their obfuscator
- Observation: Similar Structure → Generated by Similar Obfuscator
- Define Similarity.

 $SimScore(norm1, norm2) = 1 - \frac{EditDist(norm1, norm2)}{max(len(norm1), len(norm2))}$

- Similarity Score [0, 1]
- 1 = Identical, 0 = Different



Hierarchical vs Flat

- Flat Model: (K-Means)
 - Easy and Efficient
 - Drawbacks
- Require predefined K as input
 - K IN Number of Obfuscator Version
 - Hard to Predict, lack of knowledge
 - Will be Changed over time.



Agglomerative (Bottom-Up) Hierarchical Clustering



What is the proper Threshold to Identify Obfuscator Version/Variant?

[0.4 – 0.5]

[0.78 ~ 0.85]

Identify Obfuscator Version Identify Obfuscator Variants

Threshold vs K



Dendrogram





Reproduce Obfuscator for Each Cluster



Why This Research







Boost Samples Set Improve Detection Rate Better Understanding on Obfuscator



Life Cycle of Nuclear Exploit Kit Obfuscator



During **December 2014**, a new version of Nuclear Pack emerged. ... new version will **completely replace the old** version.

http://community.WebSense.com/blogs/securitylabs/arc hive/2015/01/15/evolution-of-an-exploit-kit-nuclear-pack.aspx



Life Cycle of Angler Exploit Kit Obfuscator



- Extremely Prevalent
- Aggressive tactics for evading detection
- Only few Obfuscator Versions / Variants.

of Domains Deployed

Version 1: 519 Version 2: 89 Version 3: 753 Version 4: 68 Version 5: 234



Evolution of Variants

Nuclear Obfuscator Version 2







- A new Angle to Explore Exploit-Kit
- The novel method to boost sample set and improve detection rate by reproducing obfuscator.
- The Evolution of Obfuscator in the wild



https://github.com/irobert-tluo/rebuild_obfuscator

