

Effective Malware Detection based on Behavior and Data Features

Zhiwu Xu, Cheng Wen, Shengchao Qin, and Zhong Ming

College of Computer Science and Software Engineering, Shenzhen University, China



Approach

• Experiments

Malware

- Malicious software:
 - Computer viruses, worms, Trojan horses, ransomware, spyware, adware, scareware, and other intrusive codes
- Recent report from McAfee:
 - More than 650 million malware samples detected in Q1, 2017, in which more than 30 million ones are new.

Signature-based method

- To compare with the known signatures,
 - Comodo, McAfee, Kaspersky, Kingsoft, and Symantec
- Can be easily evaded by the evasion techniques
 packing, variable-renaming, and polymorphism.

Heuristic-based method

- To identity malicious patterns though either static analysis or dynamic analysis
- However, heavy-weight or Inefficient

Machine learning approaches

- Most of existing work focus on behaviour features, without data information
 binary codes, opcodes and API calls
- Can be easily evaded
 - previously-unseen behaviors
 - obfuscate

Introduction

Approach

• Experiments

Our approaches

- Based on machine learning
- Consider both the behaviour information and the data information.
- Consider the time-split samples and obfuscated samples

Framework



- Decompilation
- Information Extraction
- Feature Selection and representation

Decompilation

- Information Extraction
- Feature Selection and representation



text:00401000	?Bubble_Sort@@YAXPA	ihh@Z pi	roc near	; CODE	XREF:	_main+47 j p	
text:00401000							
text:00401000	temp = d	lword p	tr -0Ch				
text:00401000	j = d	lword p	tr -8				
text:00401000	i = d	lword p	tr -4				
text:00401000	num = d	lword p	tr 8				
text:00401000	n = d	lword p	tr OCh				
text:00401000							
text:00401000	pus	sh el	bp				
text:00401001	mov	v el	bp, esp				
text:00401003	sub) e:	sp, 4Ch				
text:00401006	pus	sh el	bx				
text:00401007	pus	sh e:	si				
text:00401008	pus	sh e	di				
text:00401009	mov	יו י	ebp+i], Ø				
text:00401010	յան) sl	hort loc_40101B				
text:00401012	3						
text:00401012							
text:00401012	loc_401012:			; CODE	XREF:	Bubble_Sort(int	*, <mark>int</mark>):loc_401098 j j
text:00401012	mov	v e	ax, [ebp+i]				
text:00401015	add	i e	ax, 1				
text:00401018	mov	י ני	ebp+i], eax				
text:0040101B							
text:0040101B	loc_40101B:			; CODE	XREF:	Bubble_Sort(<mark>int</mark>	*, <mark>int</mark>)+10îj
text:0040101B	mov	y e	ax, [ebp+i]				
text:0040101E	cmp) e	ax, [ebp+n]				
text:00401021	jge	e sl	hort loc_40109D				
text:00401023	mov	יו י	ebp+j], Ø				
text:0040102A	յան) sl	hort loc_401035				
text:0040102C	;						
L							

ASM codes

Decompilation

Information Extraction

Feature Selection and representation



- Decompilation
- Information Extraction
- Feature Selection and representation

Selection:

Term Frequency and Inverse Document Frequency (TF-IDF)

Framework



Classifier

Classifier Training

An executable *e* can be represented as a vector *x*. D_0 represent the available dataset with known categories. Our training problem is to find a classifier $C: X \rightarrow [0,1]$ such that

$$min\sum_{(x,c)\in D_0}d(C(x)-c)$$

Malware Detection

➢ Given an executable *e* and its vector representation, the goal of the detection is to find *c* such that min d(C(x) − c) Introduction

Approach

• Experiments

Experiments

Malware dataset (11376 samples)BIG 2015 ChallengetheZoo aka Malware DB

Benign dataset (8003 samples)QIHU 360 software

(with the total size of 250 GB)

Cross Validation Experiments 10-fold cross validation

Table 2.	Results	of Different	Methods
----------	---------	--------------	---------

Classifier	PPV	TPR	FPR	ACY	1.0		NOC	
KNN K=1	0.9609	0.9830	0.0281	0.9764		-	5	
KNN K=3	0.9572	0.9798	0.0307	0.9736				
KNN $K=5$	0.9556	0.9763	0.0319	0.9715	0.8 -		5 //	
KNN K=7	0.9556	0.9763	0.0319	0.9715				
DT criterion='gini'	0.9658	0.9797	0.0243	0.9773	0.00			
DT criterion='entropy'	0.9685	0.9787	0.0223	0.9781	e Rate			
RF $(n=10, gini)$	0.9678	0.9787	0.0228	0.9778	Positiv			
RF $(n=10, entropy)$	0.9692	0.9800	0.0218	0.9788	an_ 			
Gaussian Naive Bayes	0.9381	0.1990	0.0092	0.6638				
Multinomial Naive Bayes	0.9494	0.6590	0.0247	0.8446			Random Forest(entropy) AUC = 0).9959
Bernoulli Naive Bayes	0.8726	0.6831	0.0701	0.8279	0.2 -		Decision tree(entropy) AUC = 0.9 Bernoulli NB AUC = 0.9005	838
SVM kernel='linear'	0.9581	0.9896	0.0304	0.9778			Gaussian NB AUC = 0.5931 — Multinomia NB AUC = 0.9710	
SVM kernel='rbf'	0.9686	0.9119	0.0207	0.9514		and the second	SGD classifier AUC = 0.9948 SVM (kernel = rbf) AUC = 0.9905	
SVM kernel='sigmoid'	0.9671	0.7308	0.0232	0.8580	0.0		SVM (kernel = sigmoid) AUC = 0.	8528
SGD Classifier	0.9582	0.9862	0.0302	0.9765		0.0	6.2 0.4 0.6 0.8 False Positive Rate	1.0

DOC

Runtime performance

Classifier	Training Time (s)	Testing Time (s)		
KNN (k = 1)	0 + (16.477)	178.789		
KNN (k = 3)	0 + (16.369)	199.474		
KNN (k = 5)	0 + (16.517)	207.052		
KNN (k = 7)	0 + (16.238)	210.557		
DT (criterion = 'gini')	23.442	0.067		
DT (criterion = 'entropy')	13.485	0.066		
RF (n = 10, gini)	4.115	0.086		
RF (n = 10, entropy)	3.791	0.077		
Gaussian Naïve Bayes	3.093	0.480		
Multinomial Naïve Bayes	1.535	0.035		
Bernouli Naïve Bayes	1.826	0.828		
SVM (kernel = 'linear')	150.022	14.494		
SVM (kernel = 'rbf')	799.310	50.196		
SVM (kernel = 'sigmoid')	1303.607	130.178		
SGD Classifier	22.569	0.048		

250GB, 15.6 hours, Decompile0.22s/MB182GB, 10.5 hours, Extract features0.20s/MB

Feature Experiment



Time-Split Experiment

We use some fresh malware samples, which were collected dated from January 2017 to July 2017, from the DAS MALWERK website.

Date	Num	RF Classifier	Accuracy
July, 2017	12	6	0.500
June, 2017	61	29	0.475
May, 2017	87	76	0.874
April, 2017	31	28	0.903
March, 2017	48	42	0.875
February, 2017	69	61	0.884
January, 2017	56	53	0.946
Total	364	295	0.810

Obfuscation Experiments

Obfuscation tools: Obfuscator

• Change code execution flow

Obfuscation tools: Unest

- rewriting digital changes equivalently
- confusing the output string
- pushing the target code segment into the stack and jumping to it to confuse the target code
- obfuscating the static libraries

 Table 6. Results of Obfuscated Malware Samples

Tools	Number	RF Classifier	Accuracy
Obfuscator	50	50	100%
Unest	60	60	100%

Introduction

Approach

• Experiments

- Machine learning methods based on the opcodes, data types and system libraries.
- •Carried out some interesting experiments.
- •Capable of detecting some fresh malware
- •Has a resistance to some obfuscation techniques



That 's all. Thank you very much!

