

Learning to Classify PC Malwares

Jingxiang Li Advisor: Yuhong Yang

School of Statistics, UMN

Outline

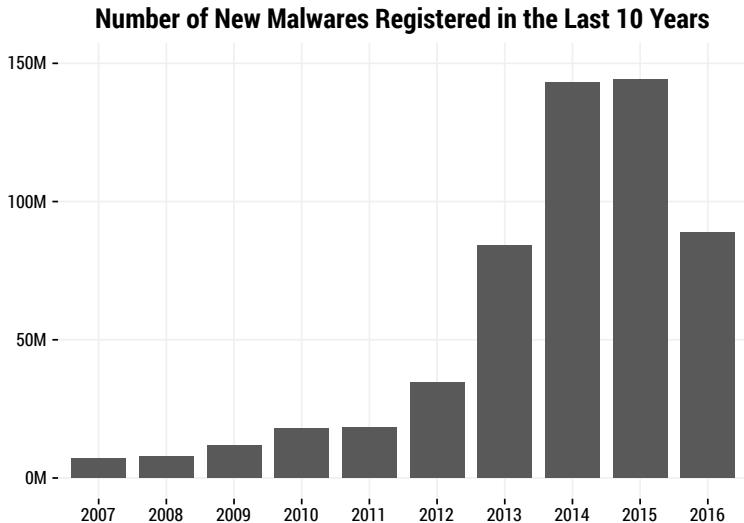
- ▶ Introduction
- ▶ Background
- ▶ Dataset Description
- ▶ Feature Engineering
- ▶ Predictive Modeling
- ▶ Model Evaluation
- ▶ Experiment Result

Introduction

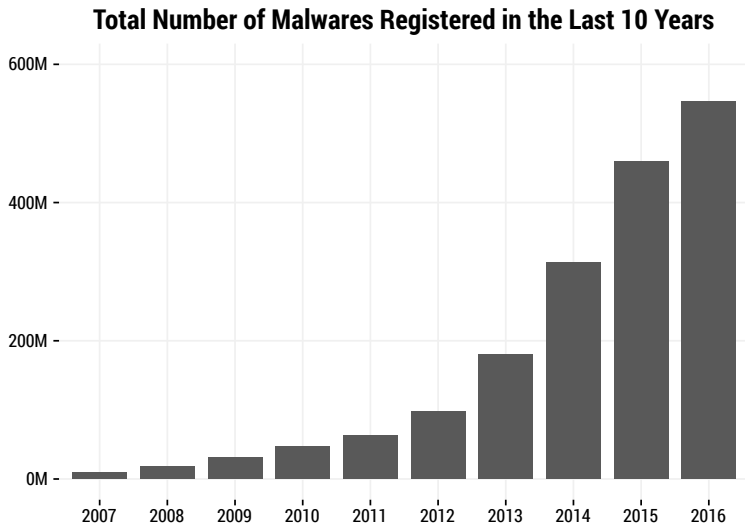
Malware

- ▶ Short for “malicious software”
- ▶ Disrupt computer operations
- ▶ Gather sensitive information
- ▶ Gain access to private computer systems
- ▶ Display unwanted advertising.

Introduction



Introduction



Introduction

Malware Industry

- ▶ Malware industry has become a well organized market involving large amounts of money.
- ▶ “Malware as a Service”
 - ▶ 400,000 Users and Organizations Globally
 - ▶ Anyone can obtain a network controlling 10,000 compromised computers for just \$1,000.

Introduction

Evasion

- ▶ Since the beginning of 2015, a sizable portion of malware utilizes a combination of many techniques designed to avoid detection and analysis.
- ▶ Polymorphism is introduced to malwares. Malicious files belonging to the same family looks like totally different files.
- ▶ Malware classification has become a challenge.

Goal of the Study

Malware Classification

- ▶ Given a malware executable, determine which family it belongs to.
- ▶ Build a statistical learning based framework to automatically classify malware.

Background

How are Executables Generated

```
gcc -o helloworld helloworld.c
```

1. Write code in a programming language
2. Compiler generates assembly code version
3. Assembler converts assembly code into binary object file
4. Linker merges object files and libraries together and generates the executable

Background

Representations of an Executable

- ▶ Original code in programming language
- ▶ Assembly code generated by compiler
- ▶ Binary object file generated by assembler and linker

Background

Representations of an Executable

- ▶ Original code in programming language ✗
- ▶ Assembly code generated by compiler ✓
- ▶ Binary object file generated by assembler and linker ✓

Assembly and Binary code can be obtained by Interactive Disassembler (IDA)

Background

Example of Assembly Code

```
.text:00401090 8B 44 24 10    mov     eax, [esp+10h]
.text:00401094 8B 4C 24 0C    mov     ecx, [esp+0Ch]
.text:00401098 8B 54 24 08    mov     edx, [esp+8]
.text:0040109C 56            push    esi
.text:0040109D 8B 74 24 08    mov     esi, [esp+8]
```

.....

```
.idata:005265CC    ; Imports from WS2_32.dll
.idata:005265CC    ;
.idata:005265D0    ; int __stdcall WSACleanup()
.idata:005265D0 ?? ?? ?? ?? extrn WSACleanup:dword ;
                                CODE XREF: _WinMain@16_0+
.idata:005265D0    ; DATA XREF: _WinMain@16_0+
```

Background

Information in Assembly Code

- ▶ Code segments (e.g. .text, .idata) ✓
- ▶ Memory address of the line of code (e.g. 00401090)
- ▶ Binary representation of the code (e.g. 8B 44 24 10) ✓
- ▶ Operation code (Opcode), single instruction executed by the CPU (e.g. mov, push) ✓
- ▶ Operands (e.g. eax, [esp+10h])
- ▶ Dynamic link libraries imported to the program (e.g. WS2_32.dll)
- ▶ Function called by the program (e.g. WSACleanup)

Dataset Description

Dataset

- ▶ From Microsoft
- ▶ 10,826 malwares, all labeled
- ▶ 8 malware families
- ▶ 185 GB
- ▶ 2 text files for each malware variant

Table 1: Malware Families

Family	Count
Ramnit	1541
Lollipop	2478
Kelihos_ver3	2942
Vundo	475
Tracur	751
Kelihos_ver1	398
Obfuscator.ACY	1228
Gatak	1013

Dataset

```
4380 .text:0040E328      jmp     loc_40E32C
4381 .text:0040E328      loc_40E328:          ; CODE XREF: .text:0040E336(0)
4382 .text:0040E328 33 C0      xor     eax, eax
4383 .text:0040E32A EB 72      jmp     short loc_40E39E
4384 .text:0040E32C      ; -----
4385 .text:0040E32C
4386 .text:0040E32C      loc_40E32C:          ; CODE XREF: .text:0040E326(0)
4387 .text:0040E32C      push    5
4388 .text:0040E32E E8 10 A2 FF FF      call    __mtinitlocknum
4389 .text:0040E333 59      pop     ecx
4390 .text:0040E334 85 C0      test    eax, eax
4391 .text:0040E336 74 F0      jz      short loc_40E328
4392 .text:0040E338 6A 05      push    5
4393 .text:0040E33A E8 C7 A2 FF FF      call    __lock
4394 .text:0040E33F 59      pop     ecx
4395 .text:0040E340 89 75 FC      mov     [ebp-4], esi
4396 .text:0040E343 89 3D E8 52 52 00      mov     dword_5252E8, edi
4397 .text:0040E349 8B 45 18      mov     eax, [ebp+18h]
4398 .text:0040E34C A3 EC 52 52 00      mov     dword_5252EC, eax
4399 .text:0040E351 89 35 F8 52 52 00      mov     dword_5252F8, esi
4400 .text:0040E357 89 35 F0 52 52 00      mov     dword_5252F0, esi
4401 .text:0040E35D 89 35 F4 52 52 00      mov     dword_5252F4, esi
4402 .text:0040E363 FF 75 20      push    dword ptr [ebp+20h]
4403 .text:0040E366 FF 75 1C      push    dword ptr [ebp+1Ch]
4404 .text:0040E369 FF 75 10      push    dword ptr [ebp+10h]
4405 .text:0040E36C FF 75 0C      push    dword ptr [ebp+0Ch]
4406 .text:0040E36F FF 75 08      push    dword ptr [ebp+8]
4407 .text:0040E372 8D 4D 8C      lea     ecx, [ebp-74h]
4408 .text:0040E375 E8 96 BD FF FF      call    ??0UnDecorator@@@PADPBDHP6APADJ@ZK@Z ; UnDecorator::UnDecorator(char
*,char const *,int,char * (*)(long),ulong)
4409 .text:0040E37A 8D 4D 8C      lea     ecx, [ebp-74h]
4410 .text:0040E37D E8 76 FC FF FF      call    ??0UnDecorator@@@QAE@PADXPZ ; UnDecorator::operator char *(void)
4411 .text:0040E382 89 45 E4      mov     [ebp+1Ch], eax
4412 .text:0040E385 E8 52 52 00      mov     ecx, offset dword_5252E8
4413 .text:0040E38A E8 EE B2 FF FF      call    unknown_libname.4 ; Microsoft VisualC 2-11/net runtime
4414 .text:0040E38F C7 45 FC FE FF FF      mov     dword ptr [ebp-4], 0FFFFFFFh
4415 .text:0040E396 E8 00 00 00 00      call    sub_40E3A4
4416 .text:0040E39B 8B 45 E4      mov     eax, [ebp+1Ch]
4417 .text:0040E39E      loc_40E39E:          ; CODE XREF: .text:0040E32A(0)
4418 .text:0040E39E
4419 .text:0040E39E E8 66 9E FF FF      call    __SEH_epilog4
4420 .text:0040E3A3 C3      ret     0
```

Figure 1: .asm File

Dataset

```
1 00401000 56 8D 44 24 08 50 88 F1 E8 1C 18 00 00 C7 06 08
2 00401010 BB 42 00 BB C6 5E C2 04 00 CC CC CC CC CC CC
3 00401020 C7 01 08 BB 42 00 E9 26 1C 00 00 CC CC CC CC CC
4 00401030 56 8B F1 C7 06 08 BB 42 00 E8 13 1C 00 00 F6 44
5 00401040 24 08 01 74 09 56 E8 6C 1E 00 00 83 C4 04 8B C6
6 00401050 5E C2 04 00 CC CC CC CC CC CC CC CC CC CC CC
7 00401060 8B 44 24 08 8A 8B 54 24 04 8B 0A C3 CC CC CC
8 00401070 8D 44 24 04 8D 50 91 8A 08 40 84 C9 75 F9 2B C2
9 00401080 C3 CC CC CC CC CC CC CC CC CC CC CC CC CC
10 00401090 8B 44 24 10 8B 4C 24 0C 8B 54 24 08 56 8B 74 24
11 004010A0 08 50 51 52 56 E8 18 1E 00 00 83 C4 10 8B C6 5E
12 004010B0 C3 CC CC CC CC CC CC CC CC CC CC CC CC CC
13 004010C0 8B 44 24 10 8B 4C 24 0C 8B 54 24 08 56 8B 74 24
14 004010D0 08 50 51 52 56 E8 65 1E 00 00 83 C4 10 8B C6 5E
15 004010E0 C3 CC CC CC CC CC CC CC CC CC CC CC CC CC
16 004010F0 33 C9 C2 10 00 CC CC CC CC CC CC CC CC CC CC
17 00401100 8B 08 00 00 C2 04 00 CC CC CC CC CC CC CC CC
18 00401110 B8 03 00 00 00 C3 CC CC CC CC CC CC CC CC CC
19 00401120 B8 08 00 00 00 C3 CC CC CC CC CC CC CC CC CC
20 00401130 8B 44 24 04 A3 AC 49 52 00 98 FE FF FF FF C2 04
21 00401140 00 CC CC CC CC CC CC CC CC CC CC CC CC CC CC
22 00401150 A1 AC 49 52 00 85 C0 74 16 8B 4C 24 08 8B 54 24
23 00401160 04 51 52 FF D0 C7 05 AC 49 52 00 00 00 00 00 B8
24 00401170 FB FF FF FF C2 08 00 CC CC CC CC CC CC CC CC
25 00401180 6A 04 68 00 10 00 00 68 6E 1C 00 6A 00 FF 15
26 00401190 9C 63 52 00 50 FF 15 C8 63 52 00 8B 4C 24 04 6A
27 004011A0 00 6A 40 68 6E 1C 00 50 89 01 FF 15 C4 63 52
28 004011B0 00 B8 04 00 00 C2 04 00 CC CC CC CC CC CC CC
29 004011C0 CC CC CC CC CC CC CC CC CC CC CC CC CC CC
30 004011D0 8B 44 24 04 A3 AC 49 52 00 98 FF FF FF FF C2 04
31 004011E0 00 CC CC CC CC CC CC CC CC CC CC CC CC CC CC
32 004011F0 A1 AC 49 52 00 85 C0 74 16 8B 4C 24 08 8B 54 24
33 00401200 04 51 52 FF D0 C7 05 AC 49 52 00 00 00 00 00 B8
34 00401210 02 00 00 00 C2 08 00 CC CC CC CC CC CC CC CC
35 00401220 6A 04 68 00 10 00 00 68 6E 3E 18 00 6A 00 FF 15
36 00401230 D0 63 52 00 8B 4C 24 04 6A 00 6A 40 68 6E 3E 18
37 00401240 00 50 89 01 FF 15 9C 63 52 00 50 FF 15 CC 63 52
38 00401250 00 B8 07 00 00 C2 04 00 CC CC CC CC CC CC CC
39 00401260 CC CC CC CC CC CC CC CC CC CC CC CC CC CC
40 00401270 33 C9 C2 04 00 CC CC CC CC CC CC CC CC CC CC
41 00401280 B8 FE FF FF FF C3 CC CC CC CC CC CC CC CC CC
42 00401290 B8 09 00 00 C3 CC CC CC CC CC CC CC CC CC CC
43 004012A0 B8 03 00 00 C3 CC CC CC CC CC CC CC CC CC CC
```

Figure 2: .bytes File

Feature Engineering

Source of Features

- ▶ Text features from “.asm” file
- ▶ Text features from “.bytes” file
- ▶ Image features from “.bytes” file
- ▶ Meta data features from both “.asm” file and “.bytes” file

Feature Engineering

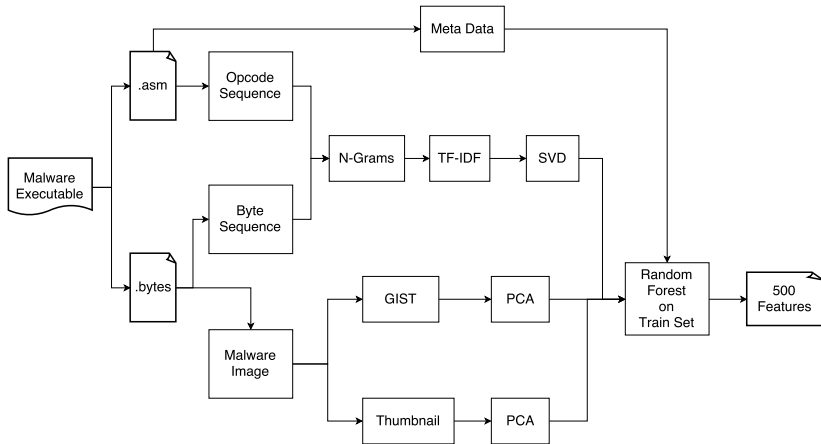


Figure 3: The Feature Engineering Process

Feature Engineering

Text Features from ".asm" File

1. Extract Opcode sequence
2. Count n-grams ($n = 1, 2, 3, 4$) on Opcode sequence
3. TF-IDF to vectorize and re-weight n-grams
4. SVD for Latent Semantic Analysis (LSA) and dimensionality reduction

Text Features from “.asm” File

Extract Opcode sequence

.text:00401390 8B 4C 24 04	mov	ecx, [esp+arg_0]
.text:00401394 B8 1F CD 98 AE	mov	eax, 0AE98CD1Fh
.text:00401399 F7 E1	mul	ecx
.text:0040139B C1 EA 1E	shr	edx, 1Eh
.text:0040139E 69 D2 FA C9 D6 5D	imul	edx, 5DD6C9FAh
.text:004013A4 56	push	esi
.text:004013A5 57	push	edi
.text:004013A6 8B F9	mov	edi, ecx

Opcode Seq \Rightarrow mov, mov, mul, shr, imul, push, push, mov

Text Features from “.asm” File

Count n-grams

- ▶ mov, mov, mul, shr, imul, push, push, mov
- ▶ 1-gram \Rightarrow {mov: 3, push: 2, mul: 1, imul: 1, shr: 1}
- ▶ bi-grams \Rightarrow {(mul, shr): 1, (push, push): 1, (push, mov): 1, (imul, push): 1, (mov, mul): 1, (mov, mov): 1, (shr, imul): 1}
- ▶ tri-grams \Rightarrow {(shr, imul, push): 1, (mov, mov, mul): 1, (mov, mul, shr): 1, (mul, shr, imul): 1, (imul, push, push): 1, (push, push, mov): 1}

Text Features from “.asm” File

n-grams

- ▶ 658 1-gram, all
- ▶ 1524 bi-grams, > 150 times in at least one file
- ▶ 5225 tri-grams, > 100 times in at least one file
- ▶ 9130 4-grams, > 100 times in at least one file

Text Features from “.asm” File

TF-IDF

- ▶ Let t be the n-gram term, d be the document, D be the Corpus

$$\text{TF}(t, d) = f_{t,d} \quad \text{IDF}(t, D) = \log \left(1 + \frac{N}{|\{d \in D : t \in D\}|} \right)$$

$$\text{TFIDF}(t, d, D) = \text{TF}(t, d) \cdot (\text{IDF}(t, D) + 1)$$

- ▶ Construct matrix, file as row, n-gram as column, TF-IDF as value
- ▶ Apply L2 normalize for each row (cosine distance)
- ▶ Why IDF? Tokens that occur very frequently in a given corpus are not informative

Text Features from “.asm” File

SVD

- ▶ Latent Semantic Analysis and Dimensionality Reduction
- ▶ $M_{m,n} = U_{m,m} \Sigma_{m,n} V_{n,n}^T$
- ▶ $M_{m,n} \approx U_{m,k} \Sigma_{k,k} V_{n,k}^T$ closest approximation with rank k
- ▶ $U_{m,k}$ the reduced version of M , with k Topic vectors
- ▶ Choose k , use theorem

$$||M||_F^2 = \sum \sigma_i^2$$

- ▶ 95% norm explained is fine

SVD

- ▶ 100 1-gram topics
- ▶ 200 2-gram topics
- ▶ 400 3-gram topics
- ▶ 500 4-gram topics
- ▶ Total from 16537 to 1200

Feature Engineering

Text Features from “.bytes” File

1. Extract Byte sequence
2. Count n-grams ($n = 1, 2, 4$) on Byte sequence
3. TF-IDF to vectorize and re-weight n-grams
4. SVD for Latent Semantic Analysis (LSA) and dimensionality reduction

Feature Engineering

Image Features From ".bytes" File

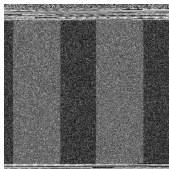
1. Visualize malware image and scale it
2. Obtain feature vector by GIST descriptor and PCA
3. Obtain feature vector by image thumbnail and PCA

Image Features From “.bytes” File

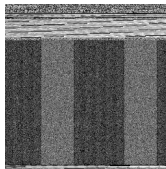
Malware Visualization

- ▶ From .bytes file, instructions are represented sequence of double hexadecimal points
- ▶ 00 to FF in hex is 0 to 255
- ▶ Use this value as a gray scale value
- ▶ Determine the width of image

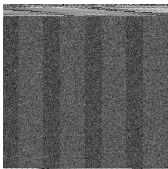
Image Features From “.bytes” File



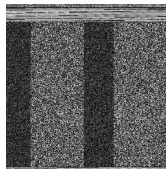
(a) 01IsoiSMh5gxyDYTI4CB



(b) 02zcUmKV16Lya5xqnPGB



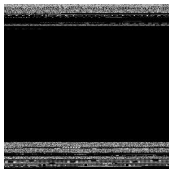
(c) 2ZNUOyAMv67lcYPoetkR



(d) 40ofYnDCI8TKJXwyGmrP

Figure 4: Sample Malware Images for Lollipop

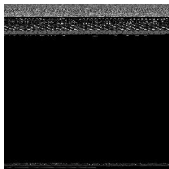
Image Features From “.bytes” File



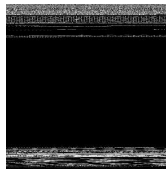
(a) 0AguvpOCcaf2myVDYFGb



(b) 2TaBCxc4msyVU45YDRuOd



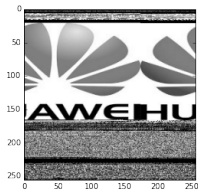
(c) 6jCs7bHULZXW3kKtxNRw



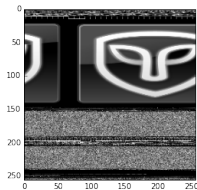
(d) 8lTjbp3rnwtLh104E57v

Figure 5: Sample Malware Images for Obfuscator.ACY

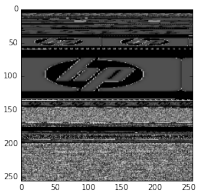
Image Features From ".bytes" File



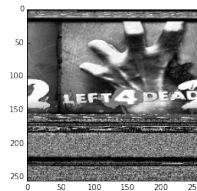
(a) 5764RnJYTirWq1utgdkN



(b) 2DBKbxPnVCyiLzqAHU9c



(c) NxrC3kQvsl5AZ4WtowHD



(d) goYPdUatw3BQEl6iCl0D

Figure 6: Fun Malware Image

Image Features From “.bytes” File

GIST

- ▶ Convolve the image with 20 pre-designed Gabor filters to construct 20 feature maps
- ▶ Divide each feature map into 16 regions (by a 4×4 grid)
- ▶ Average the feature values within each region
- ▶ Concatenate feature values together, length 320
- ▶ PCA for dimension reduction

Image Features From “.bytes” File

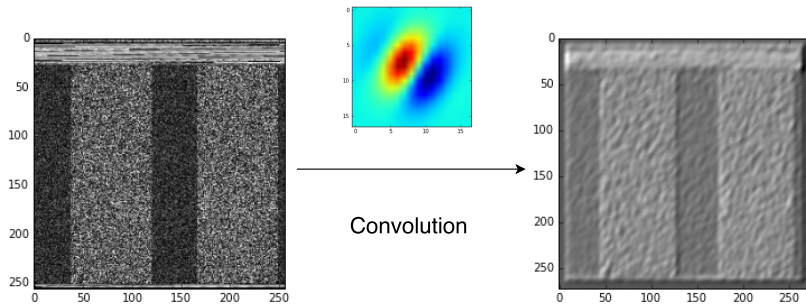


Figure 7: GIST

Image Features From “.bytes” File

Image Thumbnail

- ▶ Rescale the image to 32×32 and flatten it to a vector
- ▶ PCA for dimension reduction

Feature Engineering

Meta Data Features From Both ".asm" File and ".bytes" File

- ▶ Segment count from ".asm" file (.text, .data, .idata), TF-IDF, SVD
- ▶ File size for ".asm" file
- ▶ File size for ".bytes" file

Feature Engineering

Concatenation and Selection

- ▶ Concatenate all features together, more than 4,000
- ▶ Supervised feature selection by training Random Forest on the Training Set, 500 ultimate features

Feature Engineering

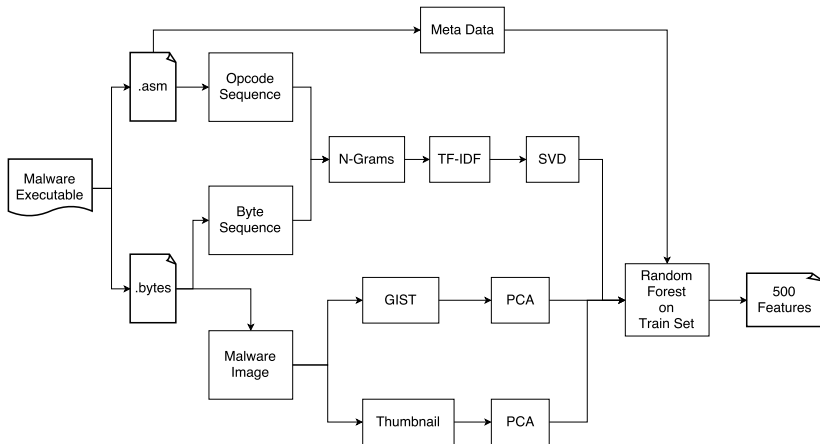


Figure 8: The Feature Engineering Process

Feature Engineering

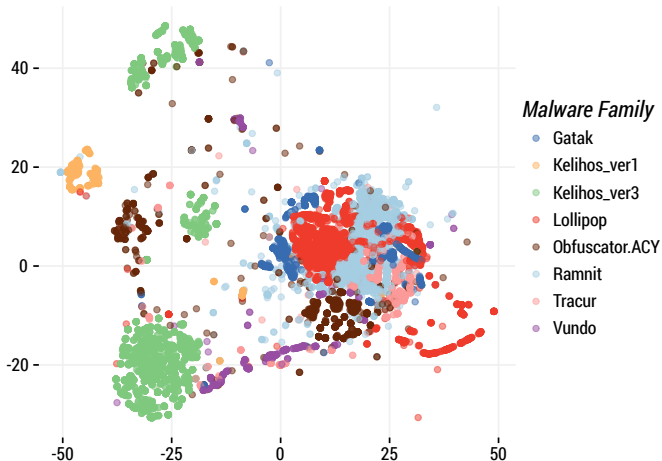


Figure 9: t-SNE Visualization with 500 Selected Features

Predictive Modeling

Candidates

- ▶ Generative Model
 - ▶ Naive Bayes (NB)
 - ▶ Linear Discriminant Analysis (LDA)
 - ▶ Quadratic Discriminant Analysis (QDA)
- ▶ Discriminative model
 - ▶ Multinomial Logistic Regression (Multi-Logreg)
 - ▶ One-vs-Rest Logistic Regression (OvR-Logreg)
 - ▶ Support Vector Machine with Probability Calibration (SVM)
 - ▶ Random Forest Classifier (RF)
 - ▶ Extremely randomized trees (ExtraTrees)
 - ▶ eXtreme Gradient Boosting (XGBoost)

Generative Model

Idea

- ▶ Estimate $P(\vec{x}|y)$ based on different assumptions
- ▶ Make prediction by Bayes' Rule

$$P(y|\vec{x}) = \frac{P(\vec{x}|y)p(y)}{P(\vec{x})}$$

Generative Model

Candidates

- ▶ NB, independent Gaussian: $P(\vec{x}|y) = P(x_1|y)P(x_2|y) \dots$
- ▶ LDA, Multi-Gaussian, same Covariance Matrix across classes
- ▶ QDA, Multi-Gaussian, different Covariance Matrix

Discriminative Model

Idea

- ▶ Directly model $P(y|\vec{x})$
- ▶ Training by Maximum a Posteriori, or Maximum Likelihood

Discriminative Model

Logistic Regression

- ▶ Linear combination of features
- ▶ Softmax function
- ▶ Minimize cross entropy with one-hot encoding
- ▶ Multinomial and One-vs-Rest

Discriminative Model

Support Vector Machine

► Primal form

$$\min_{w,b} \frac{1}{n} \sum_{i=1}^n \zeta_i + \lambda ||w||^2 \quad \text{s.t.} \quad y_i(x_i \cdot w + b) \geq 1 - \zeta_i \text{ and } \zeta_i \geq 0$$

► Dual form

$$\begin{aligned} \max_{a_1, \dots, a_n} f(a_1 \dots a_n) &= \sum_{i=1}^n a_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i a_i (x_i \cdot x_j) y_j a_j \\ \text{s.t.} \quad \sum_{i=1}^n a_i y_i &= 0, \text{ and } 0 \leq a_i \leq \frac{1}{2n\lambda} \end{aligned}$$

Discriminative Model

Support Vector Machine

- ▶ Kernel tricks for non-linear classification (Gaussian kernel)
- ▶ Probability calibration: Bining or softmax fitting

Discriminative Model

Random Forest

- ▶ Ensemble of classification trees
- ▶ Bootstrap sampling
- ▶ Randomly select features for each node split

Extremely Randomized Trees

- ▶ Randomly draw threshold for each node split

Discriminative Model

eXtreme Gradient Boosting

- ▶ Minimize a loss function (Logloss in this case)
- ▶ Additively minimize the loss function by adding classification trees
- ▶ The first and second order information is used in building each classification tree

Evaluation

Data Split

- ▶ 5% train set (542), 95% test set (10,248) by randomized stratified sampling.
- ▶ Enlarge the difference among classifiers
- ▶ More stable estimation of predictive metrics
- ▶ Cross-validation for hyper-parameter tuning
- ▶ Repeat the procedure for 10 times

Evaluation

Metrics

- ▶ Multi-class Logloss

$$\text{Logloss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(p_{ij})$$

- ▶ Multi-class Accuracy

$$\text{Accuracy} = 1 - \frac{\text{\#misclassified}}{N}$$

Evaluation

Metrics

- Recall, Precision and F_1 score for each malware family

$$\text{Recall} = \text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{Precision} = 1 - \text{FDR} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

Result

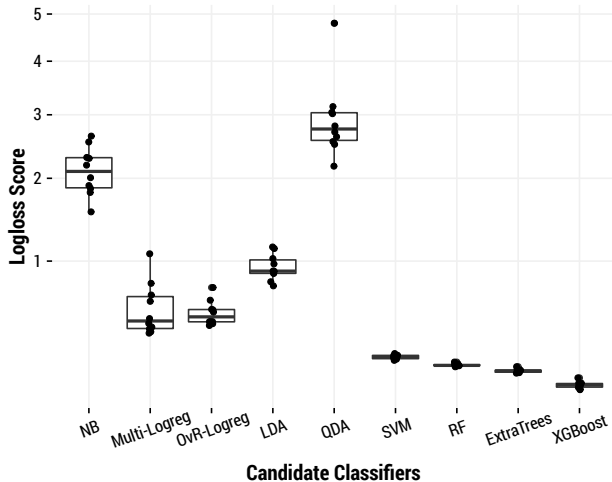


Figure 10: Model Evaluation: Multi-Class Logloss

Result

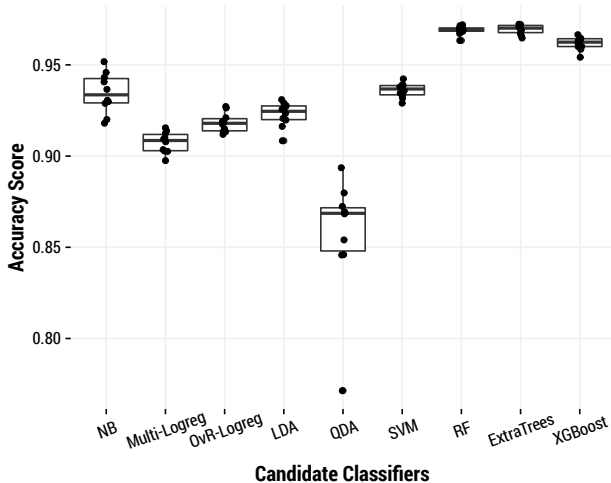


Figure 11: Model Evaluation: Multi-Class Accuracy

Result

Table 2: Confusion Matrix for Extremely Randomized Trees

1422	28	0	0	0	0	14	0
24	2329	0	0	1	0	0	0
0	9	2785	0	1	0	0	0
3	24	0	414	2	0	6	2
7	32	0	2	671	0	1	0
4	8	0	0	0	364	2	0
80	16	0	2	4	0	1065	0
2	6	0	0	2	0	12	940

Result

Table 3: Classification Report for Extremely Randomized Trees

Family	Rrecision	Recall	F_1 Score	Support
Ramnit	0.92	0.97	0.94	1464
Lollipop	0.96	0.99	0.97	2354
Kelihos_ver3	1.00	1.00	1.00	2795
Vundo	0.99	0.92	0.95	451
Tracur	0.99	0.94	0.96	713
Kelihos_ver1	1.00	0.96	0.98	378
Obfuscator.ACY	0.96	0.91	0.93	1167
Gatak	1.00	0.98	0.99	962
Avg / Total	0.97	0.97	0.97	10284

Question?

Thank You!