# Malware Classification

Andrea Mambretti

mbr@ccs.neu.edu

CS 7775

# Problem?

- Hundreds of malware samples discovered every day

- Reverse engineering not always possible

- Manual categorization is

  - Time consuming

  - Expensive

  - Subject to human mistakes

# Malware

- Software that is intended to damage or disable computers and computer systems

- Divided in families (e.g Zeus, S...

- Each family conta... ...centralized, updated versi...

- **Famil... ...istics**

IDEAS?

CLASSIFICATION

# Yeah...what about the data though?

- Good malware collections are hard to find

- Encryption and packing do not help in the task

- Virus Total and similar?

    - Hard to distinguish new and old samples

    - Miss-classification can affect the ground truth

# Microsoft Malware Classification Challenge

- 9 Families
  - Ramnit
  - Lollipop
  - Kelihos_ver3
  - Simda

```
00401000 E8 0B 00 00 00 E9 16 00 00 00 90 90 90 90 90 90
00401010 B9 25 2B 56 00 FF 25 80 23 41 00 90 90 90 90 90
00401020 68 30 10 40 00 E8 70 03 01 00 59 C3 90 90 90 90
00401030 B9 25 2B 56 00 FF 25 74 23 41 00 90 90 90 90 90
00401040 E8 0B 00 00 00 E9 16 00 00 00 90 90 90 90 90 90
00401050 B9 24 2B 56 00 FF 25 78 23 41 00 90 90 90 90 90
00401060 68 70 10 40 00 E8 30 03 01 00 59 C3 90 90 90 90
00401070 B9 24 2B 56 00 FF 25 7C 23 41 00 90 90 90 90 90
00401080 B0 9F C2 04 00 90 90 90 90 90 90 90 90 90 90 90
00401090 83 05 C4 24 56 00 20 8B 15 A4 24 56 00 85 D2 56
```

```
text:00401050                              ; =============== S U B R O U T I N E ========================================
text:00401050
text:00401050
text:00401050                              sub_401050      proc near               ; CODE XREF: .text:00401040p
text:00401050 B9 24 2B 56 00                             mov     ecx, offset unk_562B24
text:00401055 FF 25 78 23 41 00                          jmp     ds:??0_Winit@std@@QAE@XZ ; std::_Winit::_Winit(void)
text:00401055                              sub_401050      endp
text:00401055
text:00401055                              ; ----------------------------------------------------------------------------
text:0040105B 90 90 90 90 90                             align 10h
text:00401060
text:00401060                              loc_401060:                             ; CODE XREF: .text:00401045j
text:00401060 68 70 10 40 00                             push    offset loc_401070
text:00401065 E8 30 03 01 00                             call    _atexit
text:0040106A 59                                         pop     ecx
text:0040106B C3                                         retn
text:0040106B
text:0040106C 90 90 90 90                  ; ----------------------------------------------------------------------------
text:00401070                                             align 10h
```

```
004011B0 00 00 A3 E4 23 56 00 EB 0B 29 35 84 24 56 00 BF
```

- IDA pro disassembled

# Have we got features?

- Size of the bytes file
- Strings within the file
- Number of Basic Block
- Sections (e.g. .text, .idata, .rdata)
- # of calls
- # of mov
- # of jmps (e.g. jge, jmp, je etc.)
- # of pop
- # of push
- # of xor
- # of sub
- # of add

Dynamic features non-available due to missing PE header

# Multi-class classifiers

- Support Vector Machines



- Random Forest



- Neural Networks

Total features vector of

**410318** elements


Number of samples of

**10868** elements

# SVM Results

Best configuration found   {'kernel': 'rbf', 'C': 100, 'gamma': 0.0001}

Test size 33%                        Train size 66%

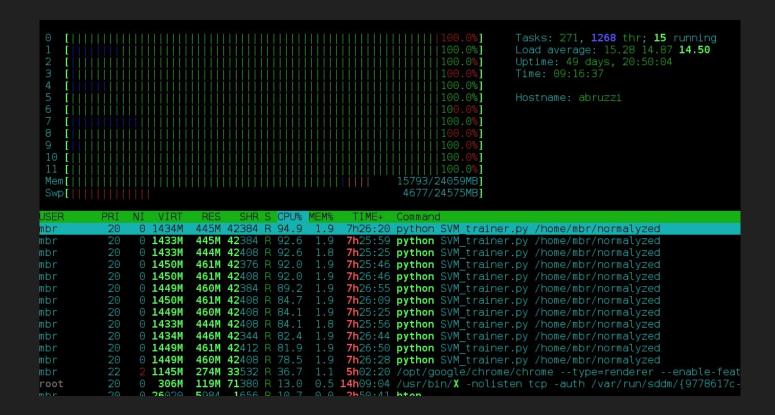|             | precision | recall | f1-score | support |
|-------------|-----------|--------|----------|---------|
| 1           | 0.98      | 0.99   | 0.98     | 518     |
| 2           | 1.00      | 1.00   | 1.00     | 844     |
| 3           | 1.00      | 1.00   | 1.00     | 956     |
| 4           | 0.90      | 0.98   | 0.94     | 145     |
| 5           | 1.00      | 0.87   | 0.93     | 15      |
| 6           | 0.98      | 0.99   | 0.98     | 248     |
| 7           | 1.00      | 0.97   | 0.99     | 139     |
| 8           | 0.98      | 0.96   | 0.97     | 400     |
| 9           | 0.99      | 1.00   | 1.00     | 322     |
| avg / total | 0.99      | 0.99   | 0.99     | 3587    |

# Random Forest Results

Best configuration found   {'bootstrap': False, 'min_samples_leaf': 1, 'min_samples_split': 2, 'criterion': 'gini', 'max_features': 3, 'max_depth': None}

Test size 33%                          Train size 66%

|         | precision | recall | f1-score | support |
|---------|-----------|--------|----------|---------|
| 1       | 0.93      | 0.98   | 0.96     | 518     |
| 2       | 0.99      | 0.99   | 0.99     | 844     |
| 3       | 1.00      | 0.99   | 1.00     | 956     |
| 4       | 0.90      | 0.98   | 0.94     | 145     |
| 5       | 0.90      | 0.60   | 0.72     | 15      |
| 6       | 1.00      | 0.99   | 0.99     | 248     |
| 7       | 0.98      | 0.96   | 0.97     | 139     |
| 8       | 0.98      | 0.92   | 0.95     | 400     |
| 9       | 0.99      | 1.00   | 0.99     | 322     |
| avg / total | 0.98  | 0.98   | 0.98     | 3587    |

# Neural Network Results

# Related Works

- Learning and Classification of Malware Behaviour  (DIMVA 2008)
    - ~10K samples, SVM, 14 Families, Dynamic features
- Lines of malicious code: Insights into the Malicious Software Industry (ACSAC 12)
    - 11 Families, Dynamic features, Dormant Functionalities

- Say no to overfitting (winner of Microsoft competition)
    - Semi-supervised learning, n-grams, bytes visualization

# Conclusion

- So far the best algorithm is SVM for this kind of classification
    - If it will ever terminate we will see if NN are better in this setting

- The feature selection seems reflect real characteristics of the malware families

- Acknowledge to sklearn, my SSD disk and multiprocessing :)

# The end...questions?