

Evaluation of Data-mining techniques for malware detection using system-call Sequence

Shayan Eskandari

Amirreza Soudi

Anomaly Detection

- Modeling the normal behavior of the system
- Flag the patterns that do not conform to the normal model as Anomaly
- This could be done in different layers, our main focus is on system call level

Anomaly Detection & Data-Mining

- Huge Amount of data to be analyzed
- Different learning and evaluating methods
- Machine learning approach to a security problem
- Analysis of system call sequences

•

•

Anubis* Dataset

```
116(desiredaccess: 1048737,filename: "C:\\INSTANTLOVENOTE.EXE",shareaccess: 5,openoptions: 96, nstatus:0)
50(desiredaccess: 983071,sectionname: "C:\\INSTANTLOVENOTE.EXE",filename: "C:\\INSTANTLOVENOTE.EXE", nstatus:0)
119(filename: "HKLM\\SYSTEM\\CURRENTCONTROLSET\\CONTROL\\SESSION MANAGER\\APPCERTDLLS",desiredaccess: 1, nstatus:0)
119(filename: "HKLM\\SYSTEM\\CURRENTCONTROLSET\\CONTROL\\SESSION MANAGER\\APPCOMPATIBILITY",desiredaccess: 1, nstatus:0)
177(keyname: "HKLM\\SYSTEM\\CURRENTCONTROLSET\\CONTROL\\SESSION MANAGER\\APPCOMPATIBILITY",valuenam: "DISABLEAPPCOMPAT",informationclass: 2, nstatus:0)
25(filename: "HKLM\\SYSTEM\\CURRENTCONTROLSET\\CONTROL\\SESSION MANAGER\\APPCOMPATIBILITY", nstatus:0)
179(filename: "C:\\INSTANTLOVENOTE.EXE",fsinformationclass: 4, nstatus:0)
120(desiredaccess: 1179649,objectname: "BASENAMEDOBJECTS\\SHIMCACHEMUTEX", nstatus:0)
125(desiredaccess: 2,objectname: "BASENAMEDOBJECTS\\SHIMSHAREDMEMORY", nstatus:0)
188(mutantname: "BASENAMEDOBJECTS\\SHIMCACHEMUTEX", nstatus:0)
139(filename: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL", nstatus:0)
116(desiredaccess: 1048608,filename: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL",shareaccess: 5,openoptions: 96, nstatus:0)
50(desiredaccess: 14,sectionname: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL",filename: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL", nstatus:0)
25(filename: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL", nstatus:0)
25(filename: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL", nstatus:0)
139(filename: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL", nstatus:0)
116(desiredaccess: 1048608,filename: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL",shareaccess: 5,openoptions: 96, nstatus:0)
50(desiredaccess: 15,sectionname: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL",filename: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL", nstatus:0)
167(sectionname: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL",sectioninformationclass: 1, nstatus:0)
119(filename: "HKLM\\SYSTEM\\CURRENTCONTROLSET\\CONTROL\\SAFEBOOT\\OPTION",desiredaccess: 3, nstatus:0)
119(filename: "HKLM\\SOFTWARE\\POLICIES\\MICROSOFT\\WINDOWS\\SAFER\\CODEIDENTIFIERS",desiredaccess: 1, nstatus:0)
177(keyname: "HKLM\\SOFTWARE\\POLICIES\\MICROSOFT\\WINDOWS\\SAFER\\CODEIDENTIFIERS",valuenam: "TRANSPARENTENABLED",informationclass: 2, nstatus:0)
25(filename: "HKLM\\SOFTWARE\\POLICIES\\MICROSOFT\\WINDOWS\\SAFER\\CODEIDENTIFIERS", nstatus:0)
119(filename: "HKU\\USER\\SOFTWARE\\POLICIES\\MICROSOFT\\WINDOWS\\SAFER\\CODEIDENTIFIERS",desiredaccess: 1, nstatus:0)
25(filename: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL", nstatus:0)
25(filename: "C:\\WINDOWS\\SYSTEM32\\APPHELP.DLL", nstatus:0)
119(filename: "HKLM\\SOFTWARE\\MICROSOFT\\WINDOWS NT\\CURRENTVERSION\\IMAGE FILE EXECUTION OPTIONS\\APPHELP.DLL",desiredaccess: 2147483648, nstatus:0)
37(desiredaccess: 2148532352,filename: "\\SYSTEMROOT\\APPPATCH\\SYSMAIN.SDB",shareaccess: 1,createoptions: 96, nstatus:0)
151(filename: "\\SYSTEMROOT\\APPPATCH\\SYSMAIN.SDB",fileinformationclass: 5, nstatus:0)
50(desiredaccess: 4,sectionname: "\\SYSTEMROOT\\APPPATCH\\SYSMAIN.SDB",filename: "\\SYSTEMROOT\\APPPATCH\\SYSMAIN.SDB", nstatus:0)
151(filename: "\\SYSTEMROOT\\APPPATCH\\SYSMAIN.SDB",fileinformationclass: 5, nstatus:0)
37(desiredaccess: 2148532352,filename: "\\SYSTEMROOT\\APPPATCH\\SYSTEST.SDB",shareaccess: 1,createoptions: 96, nstatus:0)
173(systeminformationclass: 1, nstatus:0)
119(filename: "HKLM\\SYSTEM\\WPA\\TABLETPC",desiredaccess: 257, nstatus:0)
119(filename: "HKLM\\SYSTEM\\WPA\\MEDIACENTER",desiredaccess: 257, nstatus:0)
177(keyname: "HKLM\\SYSTEM\\WPA\\MEDIACENTER",valuenam: "INSTALLED",informationclass: 2, nstatus:0)
25(filename: "HKLM\\SYSTEM\\WPA\\MEDIACENTER", nstatus:0)
37(desiredaccess: 1179926,filename: "\\DEVICE\\NAMEDPIPE\\SHIMVIEWER",shareaccess: 0,createoptions: 0, nstatus:0)
116(desiredaccess: 1048577,filename: "C:\\",shareaccess: 3,openoptions: 16417, nstatus:0)
145(filename: "C:\\",fileinformationclass: 3, nstatus:0)
25(filename: "C:\\", nstatus:0)
139(filename: "C:\\INSTANTLOVENOTE.EXE", nstatus:0)
119(filename: "HKLM\\SOFTWARE\\MICROSOFT\\WINDOWS NT\\CURRENTVERSION\\APPCOMPATFLAGS\\LAYERS",desiredaccess: 2147483904, nstatus:0)
119(filename: "HKU\\USER\\SOFTWARE\\MICROSOFT\\WINDOWS NT\\CURRENTVERSION\\APPCOMPATFLAGS\\LAYERS",desiredaccess: 2147483904, nstatus:0)
119(filename: "HKLM\\SOFTWARE\\MICROSOFT\\WINDOWS NT\\CURRENTVERSION\\APPCOMPATFLAGS\\CUSTOM\\INSTANTLOVENOTE.EXE",desiredaccess: 2147483904, nstatus:0)
139(filename: "C:\\INSTANTLOVENOTE.EXE", nstatus:0)
179(filename: "C:\\INSTANTLOVENOTE.EXE",fsinformationclass: 4, nstatus:0)
151(filename: "C:\\INSTANTLOVENOTE.EXE",fileinformationclass: 4, nstatus:0)
151(filename: "C:\\INSTANTLOVENOTE.EXE",fileinformationclass: 5, nstatus:0)
```

* anubis.iseclab.org

Pre-Processing

- Data Cleaning
- Extracting the system calls sequence
- Define a window size and separating the sequences

```
139,116,151,224,25,116,151,224,25,116,151,224,25,116,151,224,25,116,151,224,25,116,151,224,25,37,179,183,25,116,50,119,119,177,25,179,120,125,188,
139,116,50,25,25,139,116,50,167,119,119,177,25,119,25,25,119,37,151,50,151,37,173,119,119,177,25,37,116,145,25,139,119,119,119,139,179,151,151,188,
,25,119,119,177,177,25,119,119,177,25,119,71,119,177,177,25,71,25,119,71,119,177,177,177,177,25,71,119,177,177,177,177,25,71,119,177,177,177,177,2
5,71,119,177,177,177,177,25,71,119,177,177,177,177,25,71,25,119,119,119,119,119,119,119,119,119,119,119,119,119,119,119,119,119,119,119,119,119,11
9,119,119,119,119,119,119,119,177,25,119,119,177,25,179,151,139,151,50,25,119,119,25,177,25,119,177,25,119,167,119,173,228,154,116,154,53,25,257,1
19,177,25,119,173,116,179,84,139,125,25,173,173,50,25,173,119,177,25,125,25,125,25,125,167,25,125,25,125,125,125,139,139,116,50,167,119,119,177,25
,119,25,25,125,25,125,25,125,25,125,25,125,25,125,25,119,177,25,119,119,173,119,119,177,25,119,177,25,119,51,119,119,119,177,177,25,
119,177,25,119,119,119,173,125,25,119,116,173,173,119,177,25,119,177,25,173,173,173,173,119,177,177,177,25,119,177,177,25,119,119,119,119,119,173,
43,173,119,119,119,177,25,139,51,139,116,50,25,25,173,139,116,50,25,25,139,116,50,167,25,25,119,139,139,116,50,25,25,50,119,119,177,25,43,43,43,43
,43,119,119,177,177,25,139,119,177,25,43,125,188,173,173,119,177,25,139,125,139,139,116,50,167,25,25,119,173,119,119,177,177,177,119,50,139,37
,151,179,224,151,224,224,224,173,50,51,50,151,100,265,100,265,100,265,100,265,100,265,37,151,224,183,224,183,25,139,188,188,125,188,25,188,188,188
```

Data Cleaning

```
1  #!/usr/bin/env python
2
3  import fileinput
4  import os
5
6
7  file = "1grams-seq_Process-" #name of the raw trace files without the last number
8  path = "/Users/sbeta/Desktop/PProject/ISSTA12-mal-data/anubis-good" #path to the raw trace files
9  W=5 #window size
10 K=1 #shift size
11
12 def parsefile(file): #parse the files to extract the sequence files
13     with open(file) as f:
14         for line in f:
15             Syscall=line.split("")
16             print Syscall[0]
17             if Syscall[0].isdigit():
18                 seq = open(file + "-seq.txt", "a")
19                 seq.write(Syscall[0])
20                 seq.write(",")
21
22
23 def window(seq): # to extract windows of size W with the shift size of K from the sequence files (-seq)
24     with open(seq) as s:
25         for line in s:
26             sequence=line.split(",")
27             end=len(sequence)-W
28             for i in xrange(0,end,K):
29                 for j in xrange(i,W+i):
30                     print "K: ", K, " end: ", end, " seq: ", sequence[j]
31                     norm = open("NormDB.txt", "a")
32                     norm.write(sequence[j])
33                     norm.write(",")
34                 norm.write("\n")
35
36
37
38
39 os.chdir(path)
40 for i in xrange(1,36):
41     parsefile(file + str(i)) # To extract sequences from the raw files, the i would be the number in the file name
42     window(file + str(i) + "-seq.txt") # to extract windows of size W with the shift size of K from the sequence files (-seq)
43
44
```


Window's Size Open Question

- Which window size would be the best size to find the anomaly sequences?
 - For this we tried 5,6,7,8 window sizes and compared the results

s1,s2,s3,s4,s5,flag
139,37,179,183,25,0
37,179,183,25,119,119,119,119,119,177,0
179,183,25,119,119,119,119,177,25,0
183,25,119,119,119,177,25,179,1
25,116,50,119,177,25,179,151,0
116,50,119,177,25,179,151,139,0
50,119,177,25,179,151,139,116,0
119,177,25,179,151,139,116,145,0
177,25,179,151,139,116,145,25,0
25,179,151,139,116,145,25,151,0
179,120,116,145,25,151,50,0
120,125,145,25,151,50,25,0
125,188,145,25,151,50,25,119,0
188,139,25,151,50,25,119,119,0
139,116,50,25,25,139,0
116,50,25,25,139,0
50,25,25,139,116,0
25,25,139,116,50,0

s1,s2,s3,s4,s5,s6,s7,s8,flag
139,116,151,224,25,116,151,224,0
119,119,119,119,119,119,119,177,0
119,119,119,119,119,119,177,25,0
119,119,119,119,119,177,25,179,1
119,119,119,119,177,25,179,151,1
119,119,119,177,25,179,151,139,1
119,119,177,25,179,151,139,116,0
119,177,25,179,151,139,116,145,0
177,25,179,151,139,116,145,25,0
25,179,151,139,116,145,25,151,0
179,151,139,116,145,25,151,50,0
151,139,116,145,25,151,50,25,0
139,116,145,25,151,50,25,119,0
116,145,25,151,50,25,119,119,0
145,25,151,50,25,119,119,25,0
224,25,116,151,224,25,116,151,0
25,116,151,224,25,116,151,224,0
116,151,224,25,116,151,224,25,0

Different Models

- Boost
 - boosting algorithms consist of iteratively learning weak classifiers with respect to a distribution and adding them to a final strong classifier.
- Random Forest
 - Each decision tree is constructed by using a random subset of the training data.
- SVM
- Linear Model
 - a family of model-based learning approaches that assume the output can be expressed as a linear algebraic relation with the input attributes
- Neural Network
 - a learning algorithm that is inspired by the structure and functional aspects of biological neural networks

R Code generated by Rattle

```
# Load the data.

crs$dataset <- read.csv("file:///C:/Users/umroot/Desktop/6180 Project/INSE6180/Real-Data/5W/malware-5W-25-33.csv",
"?" ), strip.white=TRUE, encoding="UTF-8")

# Build the training/validate/test datasets.

set.seed(crv$seed)
crs$nobs <- nrow(crs$dataset) # 11999 observations
crs$sample <- crs$train <- sample(nrow(crs$dataset), 0.7*crs$nobs) # 8399 observations
crs$validate <- sample(setdiff(seq_len(nrow(crs$dataset)), crs$train), 0.15*crs$nobs) # 1799 observations
crs$test <- setdiff(setdiff(seq_len(nrow(crs$dataset)), crs$train), crs$validate) # 1801 observations

# The 'kernlab' package provides the 'ksvm' function.

require(kernlab, quietly=TRUE)

# Build a Support Vector Machine model.

set.seed(crv$seed)
crs$ksvm <- ksvm(as.factor(flag) ~ .,
  data=crs$dataset[crs$train,c(crs$input, crs$target)],
  kernel="rbfdot",
  prob.model=TRUE)

# Generate a textual view of the SVM model.

crs$ksvm

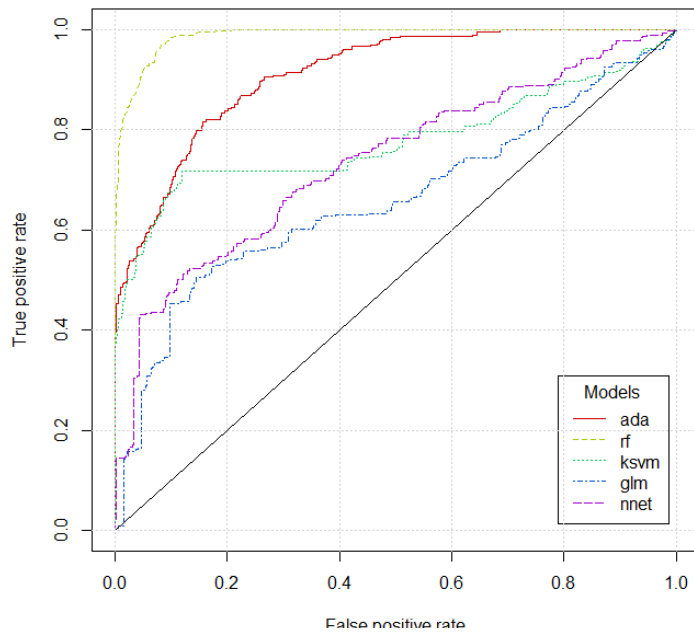
require(ROCR, quietly=TRUE)

# Generate an ROC Curve for the rpart model on malware-5W-25-33.csv [validate].

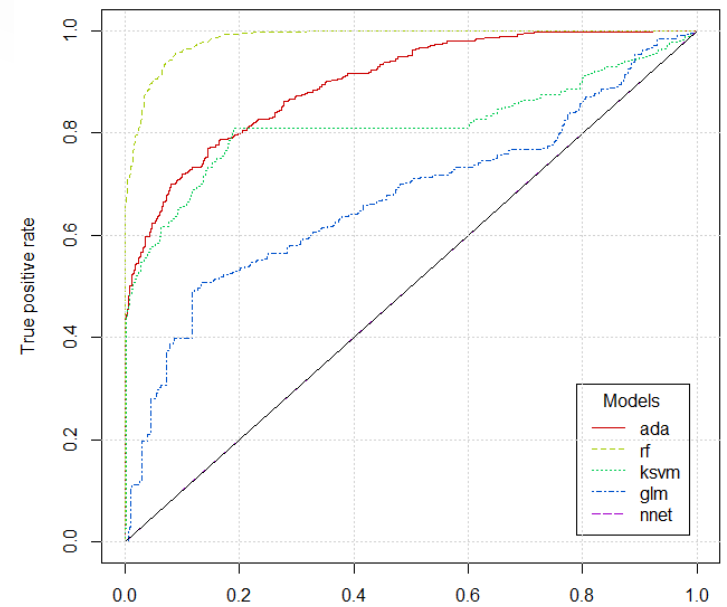
crs$pr <- predict(crs$rpart, newdata=crs$dataset[crs$validate, c(crs$input, crs$target)]),[2]
```


ROC*

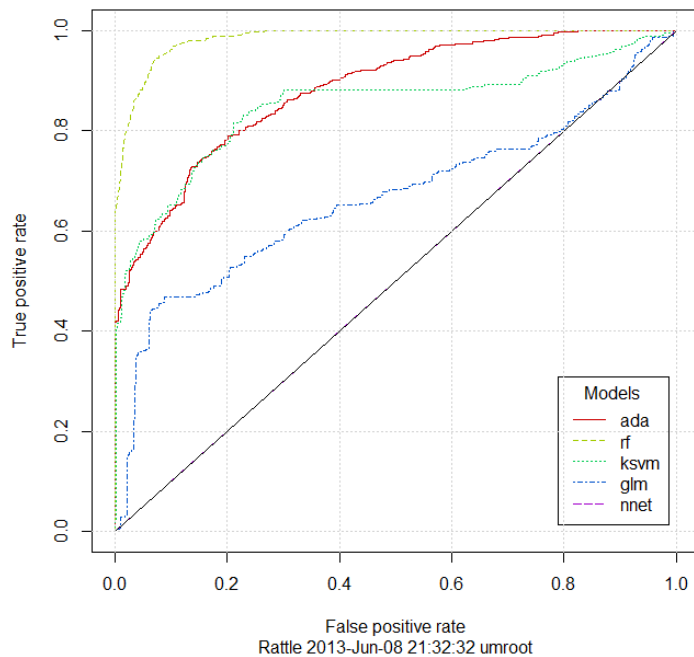
ROC Curve malware-5W-450-590.csv [validate]



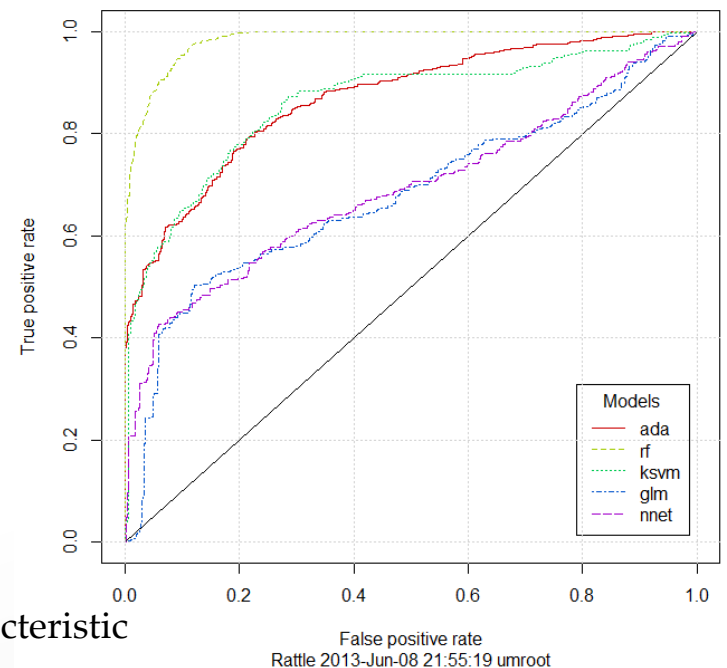
ROC Curve malware-6W-450-590.csv [validate]



ROC Curve malware-7W-450-590.csv [validate]



ROC Curve malware-8W-450-590.csv [validate]



* Receiver operating characteristic

Accuracy

Window Size = 5

Model	Accuracy (%)
Random Forest	98.79
Boost	91.07
SVM	77.71
Linear Model	65.93
Neural Network	73.62

Window Size = 6

Model	Accuracy (%)
Random Forest	98.43
Boost	89.78
SVM	81.36
Linear Model	67.21
Neural Network	50

Window Size = 7

Model	Accuracy (%)
Random Forest	98.38
Boost	88.06
SVM	84.79
Linear Model	66.41
Neural Network	50

Window Size = 8

Model	Accuracy (%)
Random Forest	98.33
Boost	86.79
SVM	85.92
Linear Model	67.86
Neural Network	68.98

Error Matrix $W = 5$

Random Forest

Actual/ Predicted	0	1
0	80	1
1	5	15

Overall Error (%) = 5.16

Boost

Actual/ Predicted	0	1
0	80	0
1	11	9

Overall Error (%) = 10.95

SVM

Actual/ Predicted	0	1
0	80	0
1	12	8

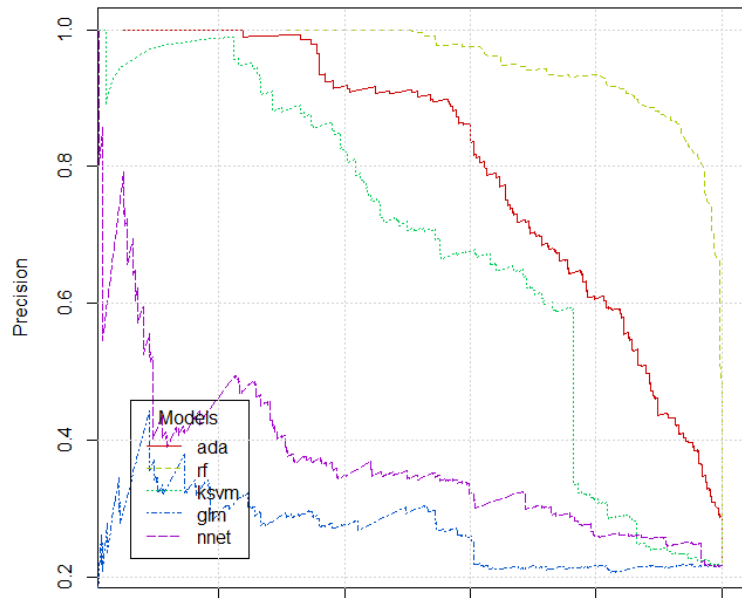
Overall Error (%) = 11.89

Neural Network

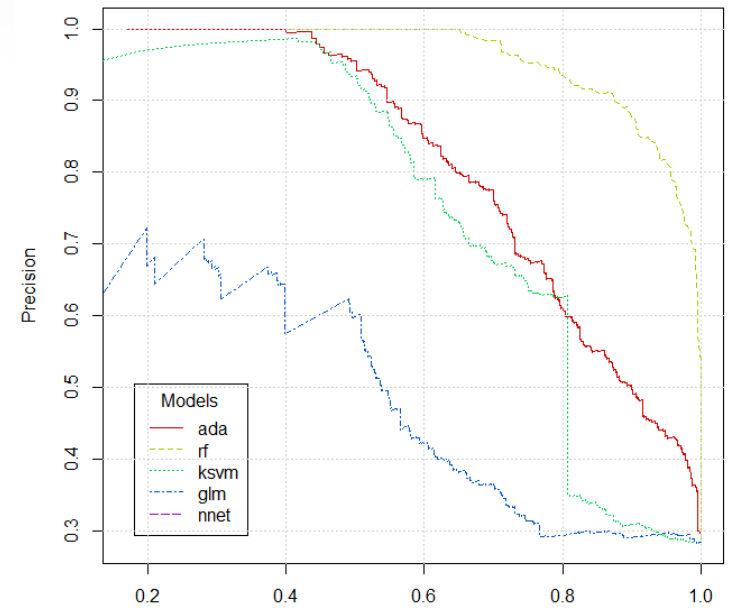
Actual/ Predicted	0	1
0	78	3
1	13	6

Overall Error (%) = 16.12

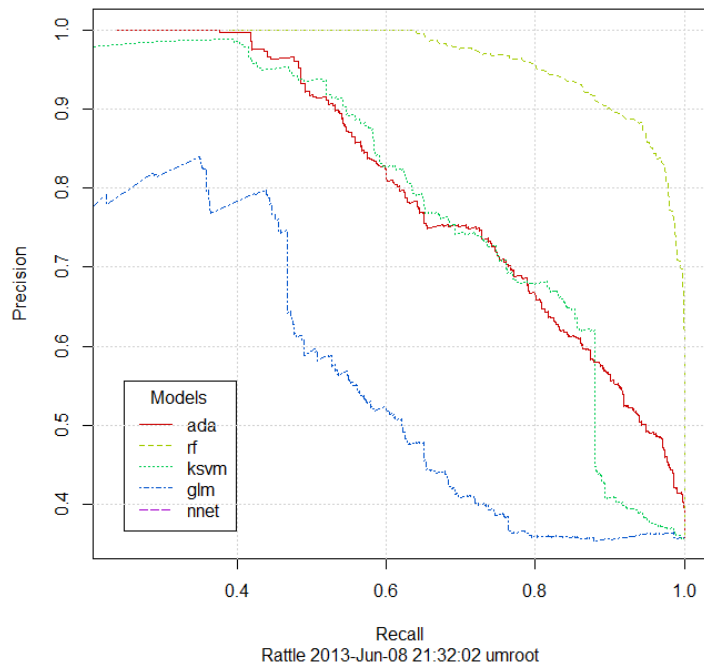
Precision/Recall Plot malware-5W-35-60.csv [validate]



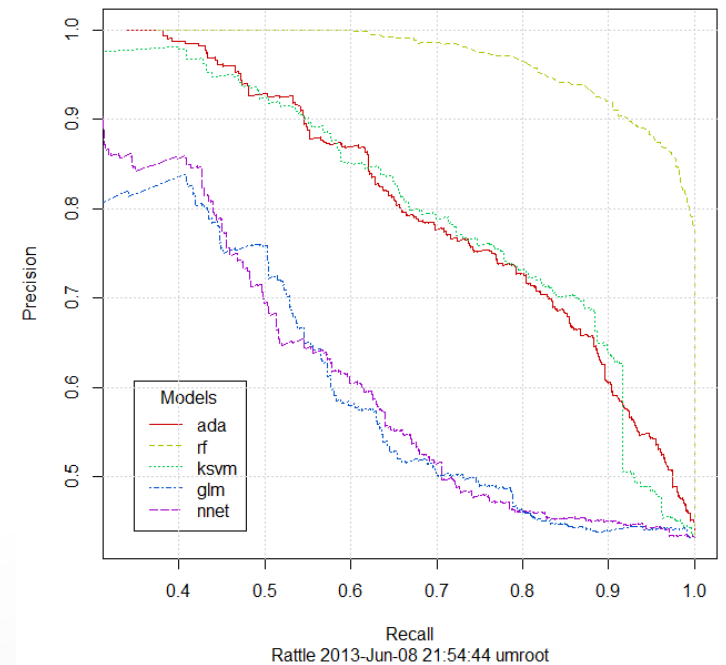
Precision/Recall Plot malware-6W-450-590.csv [validate]



Precision/Recall Plot malware-7W-450-590.csv [validate]



Precision/Recall Plot malware-8W-450-590.csv [validate]



Conclusion

- Which window size has the overall best result?
 - $W = 5$ (highest: 98.79 %)
- Which Model has the overall best result?
 - Random Forest (In the 2nd place Boost)
 - RF has the average of 98% and Boost has the average of 87%
- Improvement of accuracy
 - SVM with increasing the window size
 - From 77% to 85%