
Replication Study: Generating Adversarial Malware Examples for Black-Box Attacks Based on GAN

Claire Casanova

Department of Computing Security
Rochester Institute of Technology
Rochester, NY 14623
ccc2876@rit.edu

1

2 1 Summary

3 Recently, machine learning algorithms are being implemented to detect malicious software within
4 systems. Malware authors are learning how to attack these machine learning detection algorithms in
5 order to create malware that is able to bypass detection by the algorithm. However, malware authors
6 do not have access to the inner architecture of the algorithm that they are trying to attack and therefore
7 they can only perform a black box attack. I chose to do a replication study on the paper *Generating*
8 *Adversarial Malware Examples for Black-Box Attacks Based on GAN* that seeks to create adversarial
9 examples that are successfully able to fool a black box machine learning malware detection system. It
10 does this by utilizing a Generative Adversarial Network (GAN) named MalGAN to create malware
11 samples that appear benign and then refit the black box detection system to read these samples as
12 benign. After the black box is retrained, MalGAN is successfully able to trick the black box into
13 believing that the malware samples are benign.

14

15 2 Background

16

17 2.1 Generative Adversarial Network

18 A Generative Adversarial Network (GAN) is comprised of two feed forward neural networks, a
19 generator and a discriminator, that work in competition with each other. The generator's goal is to
20 create adversarial examples and the discriminator's goal is to determine if the examples that it is
21 fed are real or fake. The generator is trained by receiving the output of the discriminator as well as
22 labels that associate the output as real. The discriminator train on fake data by taking the generator's
23 fake data and labeling it as adversarial and trains on real data by feeding in real data with labels that
24 associate it as real examples.

25 2.1.1 Generator in MalGAN

26 The generator is the portion of the architecture that is tasked with creating the adversarial examples. It
27 is fed a binary vector that represents a piece of malware and its corresponding features. The examples
28 that are created by the generator have changed bits of the features within the sample of malware
29 making it appear as though it is a benign. The generator returns a binary vector of malware samples
30 that appear benign to the black box.

31 2.1.2 Discriminator in MalGAN

32 The discriminator is tasked with determining whether the sample it was fed is a malicious program or
33 a benign program. The discriminator is used to train the generator by providing gradient information.
34 The discriminator acts similarly to the black box since the malware authors are unable to determine
35 the specific detection algorithm.

36

37 2.2 Black Box

38 The black box is utilized as the classifier for malware. These are different machine learning algorithms
39 that a machine learning malware detection system would use to determine if a sample it is sent is
40 malware or not. For the purpose of the replication study I chose to examine three different machine
41 learning algorithms as the black in order to compare the output as to which algorithms were more
42 easily fooled than others.

43 2.2.1 Random Forest Classifier

44 A Random Forest Classifier (RF) is a machine learning model that is comprised of many decision
45 trees. A random forest learns from a random sample of the data points that are passed to it. It trains
46 each tree on the subset of these points and as the tree depth goes down the number of features within
47 a node decreases. The random forest makes its final predictions by averaging the predictions of each
48 of the trees.

49 2.2.2 Support Vector Machine

50 A Support Vector Machine (SVM) is a supervised learning algorithm that is utilized for binary
51 classification problems. After training on data, the SVM is able to classify new data into one of
52 the categories. For this replication study, I utilized a linear SVM. This algorithm learns by plotting
53 the data points and clustering them based on their labels and drawing a line that separates the two
54 different classes. This line is used as the decision line and wherever new data falls on the plot it is
55 classified into that pre-determined class.

56 2.2.3 Logistic Regression

57 Logistic Regression (LR) is used as a classification algorithm for a discrete number of outputs.
58 Logistic regression uses the sigmoid function to output a probability that can then be mapped to a
59 binary classification.

60 3 Data

61 The first step I took in replicating this paper was understanding the data sets that were provided in
62 the Github that was associated with the paper. Three data sets were listed *data.npz*, *data1.npz*, and
63 *mydata.npz*. There was no real explanation on what each of these data sets contained or how they
64 were different. By analyzing each of the data sets, I found a flaw within the *mydata.npz* file. Both the
65 benign software and malicious software labels were set to 1 so I decided to rule out this data set for
66 use. *data.npz* was the updated version of *mydata.npz* so I additionally ruled this set out and decided
67 to solely use *data1.npz* as there were no flaws or errors within this set.

68 Within this data set, there were 1,368 samples of malware and 441 samples of benignware. Each of
69 these samples contained 128 features that were the API calls outputted from running the samples
70 through Cuckoo sandbox. The data was then normalized by designated an API call as either 1 or 0
71 indicating that a piece of software either had that API call or did not.

72

73 4 Architecture Specifications

74 The authors of the paper specified that their generator and discriminator were "multi-layer feed-
75 forward neural networks." I decided to start my experiment by examining the results between two
76 hidden layers and three hidden layers within both the generator and discriminator and see if the
77 results varied between the two.

78 Additionally, I decided to test the results between using the sigmoid activation function for the output
79 layer and the tanh function. The tanh function provides stronger gradients over the sigmoid function
80 so I wanted to see if this would result in the adversarial examples being able to fool the black box
81 easier than when using the sigmoid.

82 I tested each of these parameters on all three of the black box algorithms I chose to experiment with
83 and the next section will outline the analysis of the results I achieved.

84 I chose to train MalGAN for 200 epochs before retraining the black box and then for 75 epochs after
85 the retrain. Additionally, I chose a batch size of 64 samples when training MalGAN and chose to use
86 the Mean Squared Error as the loss function.

87

88

5 Results

89

90

5.1 Two Hidden Layers using Sigmoid

91 The following tables display the final results from a two hidden layer network utilizing the sigmoid
 92 function for each of the black box algorithms

	Black Box Algorithm	Test TPR Before Retraining	Test TPR After Retraining
93	Random Forest	0.9818	0.9547
	SVM	0.9790	0.20468
	Logistic Reg	0.9918	0.0175

94

95

5.2 Two Hidden Layers using Tanh

96 The following table displays the final results from a two hidden layer network utilizing the tanh
 97 function for each of the black box algorithms

	Black Box Algorithm	Test TPR Before Retraining	Test TPR After Retraining
98	Random Forest	1.0	0.9825
	SVM	0.9927	0.2733
	Logistic Reg	0.9818	0.2325

99

100

5.3 Three Hidden Layers using Sigmoid

101 The following table displays the final results from a three hidden layer network utilizing the sigmoid
 102 function for each of the black box algorithms

	Black Box Algorithm	Test TPR Before Retraining	Test TPR After Retraining
103	Random Forest	1.0	0.9927
	SVM	0.9854	0.2193
	Logistic Reg	1.0	0.0614

104

105

5.4 Three Hidden Layers using Tanh

106 The following table displays the final results from a three hidden layer network utilizing the tanh
 107 function for each of the black box algorithms

	Black Box Algorithm	Test TPR Before Retraining	Test TPR After Retraining
108	Random Forest	1.0	0.9971
	SVM	1.0	0.4576
	Logistic Reg	1.0	0.4795

109 The best results from each testing setup are highlighted and the best result across all of the testing
 110 setups is highlighted in green.

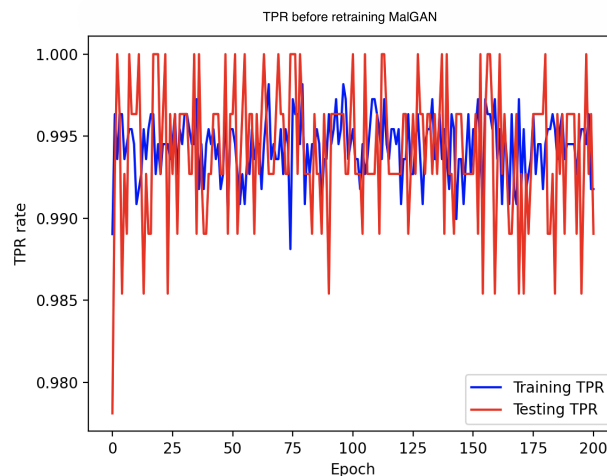
111

112 6 Analysis

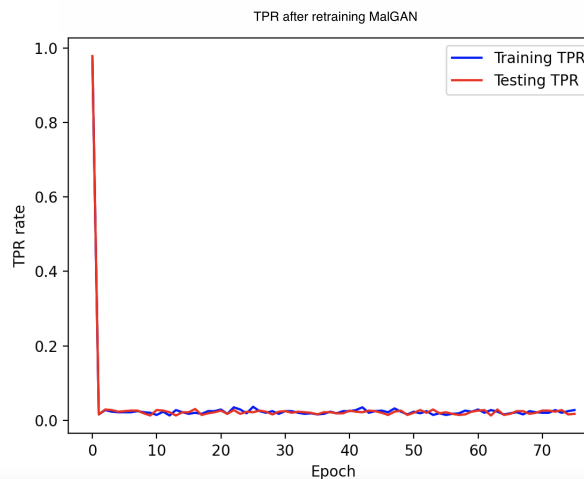
113 Overall, I was unable to reproduce the exact results of the paper. However, I was able to reproduce
114 results where the black box was deceived by the adversarial examples a majority of the time.

115 The most noticeable difference between my results and the results of the paper was in relation the
116 the Random Forest classifier. Even after retraining it with adversarial examples, the Random Forest
117 was not fooled by the adversarial examples and was still able to detect them over the Support Vector
118 Machine and the Logistic Regression algorithms.

119 The logistic regression black box algorithm performed the best over all the other algorithms and
120 it had the lowest true positive rate (TPR) when MalGAN had two hidden layers and used sigmoid
121 activation. Thus it would be the best option for malware authors to use this version of MalGAN when
122 trying to trick malware detection systems into believing that their samples are benign. The below
123 figures outline the Logistic Regression's TPR over time before and after retraining MalGAN:



124



125

126 My created linear Support Vector Machine (SVM) was able to fool the black box in approximately
127 75% of the cases.

128 The best results for the setup of MalGAN was a two hidden layer feed-forward network with sigmoid
129 activation for each the generator and the discriminator. Each of the black boxes were fooled at the
130 highest rate when running against this architecture of MalGAN

131

132 **7 Conclusion**

133 Overall, this project was a rewarding challenge. I was able to learn a lot about different machine
134 learning algorithms and neural network architectures while applying into a real world scenario. This
135 paper along with others that study similar topics open a huge door in the research world as to how
136 malware detection systems that are comprised of machine learning algorithms can be strengthened
137 against adversarial examples.