

Pruning Defense on Backdoored Networks

ECE-GY 9163 ML for Cyber Security Lab 2 Report

Kunal Kashyap (kk4564)

Objective

Design a backdoor detector for BadNets trained on the YouTube Face dataset using the pruning defense discussed in class. The detector will take as input:

1. B , a backdoored neural network classifier with N classes.
2. D_{valid} , a validation dataset of clean, labelled images.

The output is G , a repaired BadNet. G has $N+1$ classes, and given unseen test input, it must:

1. Output the correct class if the test input is clean. The correct class will be in $[1, N]$.
2. Output class $N+1$ if the input is backdoored.

G will be designed using the pruning defense that we discussed in class. That is, we will prune the last pooling layer of BadNet B by removing one channel at a time from that layer. Channels should be removed in decreasing order of average activation values over the entire validation set. Every time we prune a channel, we will measure the new validation accuracy of the new pruned badnet. We will stop pruning once the validation accuracy drops atleast $X\%$ below the original accuracy. This will be our new network B' .

Now, our goodnet G works as follows. For each test input, we will run it through both B and B' . If the classification outputs are the same, i.e., class i , we will output class i . If they differ we will output $N+1$.

Results

We were able to create 3 different repaired BadNets with different levels of accuracy drop (2%, 4% and 10%). Below you can see the plot and table for changes in Clean Validation Accuracy and Attack Success Rate as a function of Fraction of Channels Pruned.

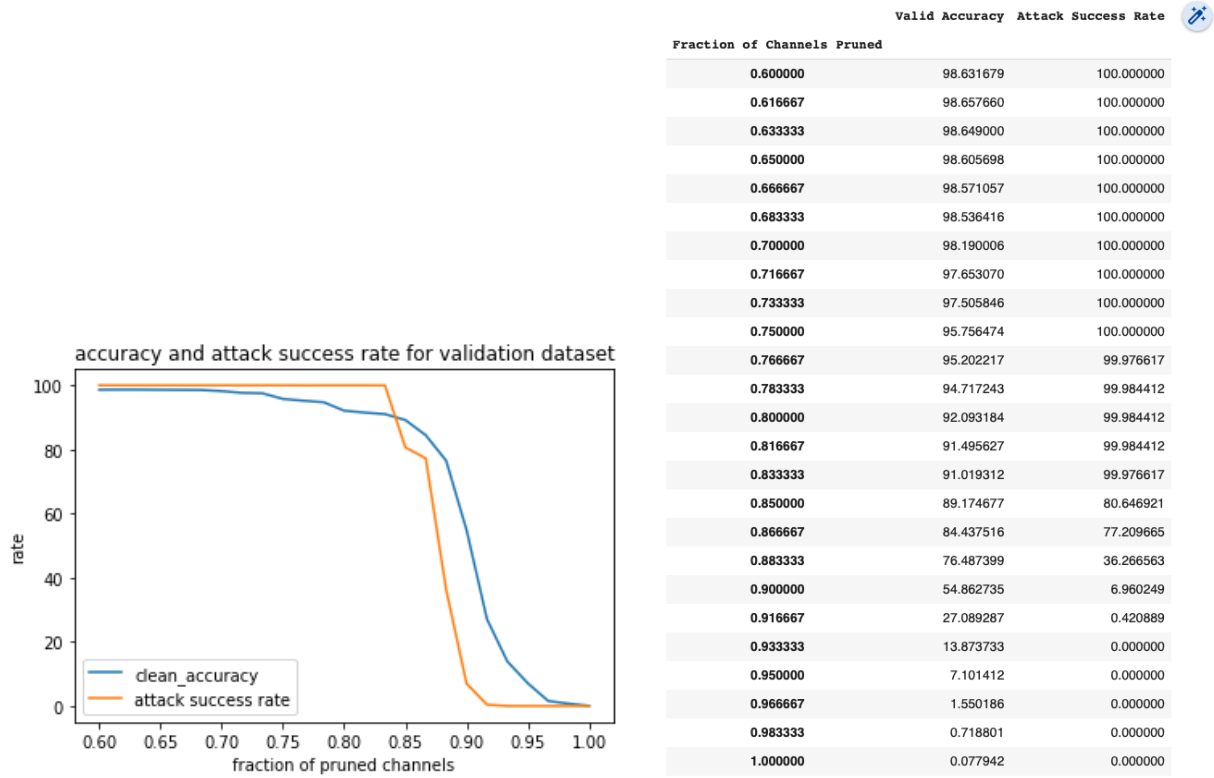


Figure 1: Clean Accuracy and Attack Success Rate as a function of Pruned Channels

In the table below, we see the final test accuracies and attack success rate for the 3 Repaired BadNets:

Model	Clean Test Accuracy (in %)	Attack Success Rate (in %)
RepairedNetX2	95.744	100
RepairedNetX4	92.127	99.984
RepairedNetX10	84.333	77.209

Table 1: Clean Accuracy and Attack Success Rate for final Repaired BadNets

Conclusion

We conclude that pruning is a good technique to defend against backdoored neural networks, but it is not very effective. In the results we can see that the clean accuracies suffer a great loss with minimal loss in the attack success rate.