

▼ Prepare Dataset

Download data and model

```
1 !unzip -uq bd.zip -d data
2 !unzip -uq cl.zip -d data
3 !wget https://github.com/csaw-hackml/CSAW-HackML-2020/raw/master/lab3/models/bd_net.h5
```

```
--2021-12-16 23:04:45-- https://github.com/csaw-hackml/CSAW-HackML-2020/raw/master/lab3/models/bd_net.h5
Resolving github.com (github.com)... 140.82.112.3
Connecting to github.com (github.com)|140.82.112.3|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://raw.githubusercontent.com/csaw-hackml/CSAW-HackML-2020/master/lab3/models/bd_net.h5 [following]
--2021-12-16 23:04:45-- https://raw.githubusercontent.com/csaw-hackml/CSAW-HackML-2020/master/lab3/models/bd_net.h5
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.109.133, 185.199.110.133, 185.199.111.133, ...
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.109.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 7275748 (6.9M) [application/octet-stream]
Saving to: 'bd_net.h5.1'

bd_net.h5.1      100%[=====>]   6.94M  --.-KB/s    in 0.09s

2021-12-16 23:04:46 (78.2 MB/s) - 'bd_net.h5.1' saved [7275748/7275748]
```

Load data and model

```
1 import h5py
2 import numpy as np
3 import keras
4 import tensorflow as tf
5 import matplotlib.pyplot as plt
```

```

1 cl_va_data = h5py.File("data/cl/cl/valid.h5", 'r')
2 bd_va_data = h5py.File("data/bd/bd/bd_valid.h5", 'r')
3 cl_ts_data = h5py.File("data/cl/cl/test.h5", 'r')
4 bd_ts_data = h5py.File("data/bd/bd/bd_test.h5", 'r')
5
6 cl_va_x = np.array(cl_va_data['data']).transpose((0, 2, 3, 1))
7 cl_va_y = np.array(cl_va_data['label'])
8 bd_va_x = np.array(bd_va_data['data']).transpose((0, 2, 3, 1))
9 bd_va_y = np.array(bd_va_data['label'])
10 cl_ts_x = np.array(cl_ts_data['data']).transpose((0, 2, 3, 1))
11 cl_ts_y = np.array(cl_ts_data['label'])
12 bd_ts_x = np.array(bd_ts_data['data']).transpose((0, 2, 3, 1))
13 bd_ts_y = np.array(bd_ts_data['label'])
14
15 print(cl_va_x.shape)
16 print(cl_va_y.shape)
17 print(bd_va_x.shape)
18 print(bd_va_y.shape)
19 print(cl_ts_x.shape)
20 print(cl_ts_y.shape)
21 print(bd_ts_x.shape)
22 print(bd_ts_y.shape)

```

```

(11547, 55, 47, 3)
(11547,)
(11547, 55, 47, 3)
(11547,)
(12830, 55, 47, 3)
(12830,)
(12830, 55, 47, 3)
(12830,)

```

```
1 BadNet = keras.models.load_model("bd_net.h5")
```

```
1 # test on original model
```

```

2 cl_pred = np.argmax(BadNet.predict(cl_va_x), axis=1)
3 cl_acc = np.mean(np.equal(cl_pred, cl_va_y))
4 bd_pred = np.argmax(BadNet.predict(bd_va_x), axis=1)
5 bd_acc = np.mean(np.equal(bd_pred, bd_va_y))
6
7 print("Accuracy on the clean dataset", cl_acc)
8 print('Success rate of attact backdoored dataset:', bd_acc)

```

Accuracy on the clean dataset 0.9864899974019226

Success rate of attact backdoored dataset: 1.0

▼ Pruning Defense

Print the backdoored model information to find the last pooling layer and the number of total classes

```
1 BadNet.summary()
```

Model: "model_1"

Layer (type)	Output Shape	Param #	Connected to
input (InputLayer)	[(None, 55, 47, 3)]	0	[]
conv_1 (Conv2D)	(None, 52, 44, 20)	980	['input[0][0]']
pool_1 (MaxPooling2D)	(None, 26, 22, 20)	0	['conv_1[0][0]']
conv_2 (Conv2D)	(None, 24, 20, 40)	7240	['pool_1[0][0]']
pool_2 (MaxPooling2D)	(None, 12, 10, 40)	0	['conv_2[0][0]']
conv_3 (Conv2D)	(None, 10, 8, 60)	21660	['pool_2[0][0]']
pool_3 (MaxPooling2D)	(None, 5, 4, 60)	0	['conv_3[0][0]']
conv_4 (Conv2D)	(None, 4, 3, 80)	19280	['pool_3[0][0]']

flatten_1 (Flatten)	(None, 1200)	0	['pool_3[0][0]']
flatten_2 (Flatten)	(None, 960)	0	['conv_4[0][0]']
fc_1 (Dense)	(None, 160)	192160	['flatten_1[0][0]']
fc_2 (Dense)	(None, 160)	153760	['flatten_2[0][0]']
add_1 (Add)	(None, 160)	0	['fc_1[0][0]', 'fc_2[0][0]']
activation_1 (Activation)	(None, 160)	0	['add_1[0][0]']
output (Dense)	(None, 1283)	206563	['activation_1[0][0]']

```

=====
Total params: 601,643
Trainable params: 601,643
Non-trainable params: 0
=====

```

Get the activation list with increading order

```

1 # split the whole BadNet into 2 subNet by the last pooling layer
2 subNet1 = keras.Model(inputs = BadNet.input, outputs = BadNet.layers[6].output)
3 subNet2 = keras.Model(inputs = BadNet.layers[7].input, outputs = BadNet.output)
4 # get the output value after the last pooling layer
5 activations = subNet1.predict(cl_va_x)
6 print(activations.shape)
7 # calculate the average values
8 avg_activations = np.mean(activations, axis=0)
9 # the indexes in increasing order
10 ordered_indexes = np.unravel_index(np.argsort(avg_activations, axis=None), avg_activations.shape)
11 print(len(ordered_indexes))

```

```

(11547, 5, 4, 60)
3

```

Repair the bad net by pruning neurals

```

1 # subNet1 --RepairedSubNet--> subNet2
2 def RepairedSubNet(custom_mask):
3     mask = tf.Variable(custom_mask, trainable=False, dtype=tf.float32)
4     masked = keras.layers.Lambda(lambda x: x * mask)(subNet1.output)
5     return keras.Model(inputs=subNet1.output, outputs=subNet2(masked))
6
7 # prune the neural or not with the value of 0 or 1
8 def PruningDefense(drop_rate):
9     prune_mask = np.ones(activations[0].shape)
10    num_neurals = activations[0].shape[0] * activations[0].shape[1] * activations[0].shape[2]
11    for i in range(num_neurals):
12        # prune one neural, by setting the mask to be 0
13        prune_mask[ordered_indexes[0][i], ordered_indexes[1][i], ordered_indexes[2][i]] = 0
14        if i <= 900:
15            continue # to save time
16        repairedSubNet = RepairedSubNet(prune_mask)
17        pruned_pred = np.argmax(repairedSubNet.predict(activations), axis=1)
18        repaired_acc = np.mean(np.equal(pruned_pred, cl_va_y))
19        # stop once acheiving the drop rate
20        if repaired_acc <= cl_acc - drop_rate:
21            print("i:", i, " repaired_acc:", repaired_acc, " drop:", cl_acc - repaired_acc)
22            break
23    return repairedSubNet, prune_mask

```

```

1 # record each mask
2 path = "models/repairedNet"
3 drop_rate = [0.02, 0.04, 0.1] # [0.02, 0.04, 0.1]
4 repairedSubNets = []
5 repairedSubNet = []
6 masks = []
7 mask = []
8 for rate in drop_rate:

```

```

9     print(rate)
10    repairedSubNet, mask = PruningDefense(rate)
11    #repairedSubNet.save(path + str(rate * 100) + ".h5")
12    repairedSubNets.append(repairedSubNet)
13    masks.append(mask)

0.02
i: 928  repaired_acc: 0.9660517883432926  drop: 0.02043820905862992
0.04
i: 964  repaired_acc: 0.946306399930718  drop: 0.04018359747120459
0.1
i: 1023  repaired_acc: 0.8852515804970988  drop: 0.10123841690482371

```

Test on test_dataset

```

1 def Predict(repairedNet, ts_x, ts_y):
2     original_pred = np.argmax(BadNet.predict(ts_x), axis=1)
3     pooling_activation = subNet1.predict(ts_x)
4     repaired_pred = np.argmax(repairedNet.predict(pooling_activation), axis=1)
5     ts_pred = []
6     for o_pred, r_pred in zip(original_pred, repaired_pred):
7         if o_pred == r_pred:
8             ts_pred.append(o_pred) # correct class
9         else:
10             ts_pred.append(1283) # class N+1
11     return ts_pred

```

```

1 cl_accuracy = []
2 defense_acc = []
3 for RepairedNet in repairedSubNets:
4     # test on clean test dataset
5     cl_ts_pred = Predict(RepairedNet, cl_ts_x, cl_ts_y)
6     acc = np.mean(np.equal(cl_ts_pred, cl_ts_y))
7     cl_accuracy.append(acc)
8     print("cl:", acc)
9

```

```

10 # test on doorback test dataset
11 bd_ts_pred = Predict(RepairedNet, bd_ts_x, bd_ts_y)
12 acc = np.mean(np.equal(bd_ts_pred, bd_ts_y))
13 defense_acc.append(acc)
14 print("bd:", acc)

```

```

cl: 0.9658612626656274
bd: 0.9999220576773188
cl: 0.9477786438035853
bd: 0.9846453624318005
cl: 0.8809041309431022
bd: 0.6316445830085736

```

➤ Plot Accuracy and Attack Success Rate

```

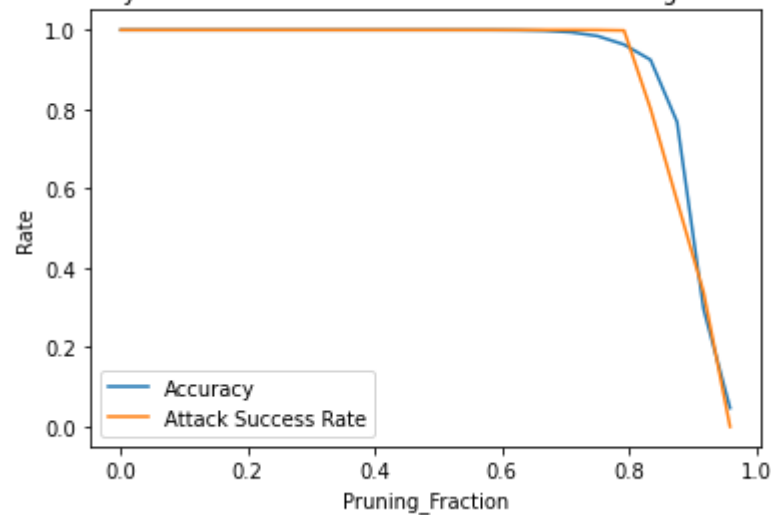
1 acc = []
2 asr = []
3 per = []
4
5 prune_mask = np.ones(activations[0].shape)
6 num_neurals = activations[0].shape[0] * activations[0].shape[1] * activations[0].shape[2]
7 for i in range(num_neurals):
8     prune_mask[ordered_indexes[0][i], ordered_indexes[1][i], ordered_indexes[2][i]] = 0
9     if i % 50 == 0:
10         per.append((i + 1) / num_neurals)
11         repairedSubNet = RepairedSubNet(prune_mask)
12         # the accuracy on clean test data
13         original_pred = np.argmax(BadNet.predict(cl_ts_x), axis=1)
14         pruned_pred = np.argmax(repairedSubNet.predict(subNet1.predict(cl_ts_x)), axis=1)
15         pruned_pred = np.where(original_pred == pruned_pred, original_pred, 1283)
16         repaired_acc = np.mean(np.equal(pruned_pred, original_pred))
17         acc.append(repaired_acc)
18         # the attack success rate on backdoored test data
19         original_pred = np.argmax(BadNet.predict(bd_ts_x), axis=1)
20         bd_pred = np.argmax(repairedSubNet.predict(subNet1.predict(bd_ts_x)), axis=1)

```

```
21     bd_pred = np.where(original_pred == bd_pred, original_pred, 1283)
22     att_suc_rate = np.mean(np.equal(bd_pred, original_pred))
23     asr.append(att_suc_rate)
```

```
1 plt.figure()
2 plt.plot(per, acc, label="Accuracy")
3 plt.plot(per, asr, label="Attack Success Rate")
4 plt.xlabel("Pruning_Fraction")
5 plt.ylabel("Rate")
6 plt.title('Curve for Accuracy and Attack Success Rate with the Percentage of Pruning Channels')
7 plt.legend()
8 plt.savefig("plot.png")
9 plt.show()
```

Curve for Accuracy and Attack Success Rate with the Percentage of Pruning Channels



As shown in the plot, the accuracy on the clean test data and the attack success rate is very high and stable until the percentage comes to a high value around 0.8. After that, both the accuracy and the attack success rate drop quickly, and the accuracy even drops before the attack success rate, which means the pruning defense doesn't work well in this case.

