Designing a backdoor detector for BadNets trained on the YouTube Face dataset using the pruning defense.

```
# All necessary imports
import os
import tarfile
import requests
import re
import sys
import warnings
warnings.filterwarnings('ignore')
import h5py
import numpy as np
import tensorflow as tf
from tensorflow import keras
from keras import backend as K
from keras.models import Model
import matplotlib.pyplot as plt
from mpl_toolkits.axes_grid1.inset_locator import inset_axes
import matplotlib.font_manager as font_manager
import cv2
Define function to load the data
# Load data
def data_loader(filepath):
    data = h5py.File(filepath, 'r')
    x_data = np.array(data['data'])
    y_data = np.array(data['label'])
    x_{data} = x_{data.transpose((0,2,3,1))}
    return x_data, y_data
```

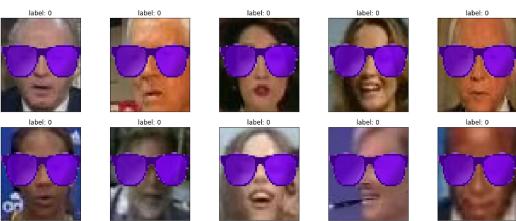
Follow instructions under <u>Data Section</u> to download the datasets.

We will be using the clean validation data (valid.h5) from cl folder to design the defense and clean test data (test.h5 from cl folder) and sunglasses poisoned test data (bd\_test.h5 from bd folder) to evaluate the models.

```
## To-do ##
# After downloading the datasets, provide corresponding filepaths below
clean_data_valid_filename = "data/cl/valid.h5"
clean_data_test_filename = "data/cl/test.h5"
poisoned_data_test_filename = "data/bd/bd_test.h5"
Read the data:
cl_x_valid, cl_y_valid = data_loader(clean_data_valid_filename)
cl_x_test, cl_y_test = data_loader(clean_data_test_filename)
bd_x_test, bd_y_test = data_loader(poisoned_data_test_filename)
Visualizing the clean test data
# Plot some images from the validation set (see https://mrdatascience.com/how-to-plot-mnist-digits-using-matplotlib/)
num = 10
np.random.seed(45)
randIdx = [np.random.randint(10000) for i in range(num)]
num row = 2
num_col = 5# plot images
fig, axes = plt.subplots(num_row, num_col, figsize=(3*num_col,3*num_row))
for i in range(num):
    ax = axes[i//num_col, i%num_col]
   ax.imshow(cl_x_test[randIdx[i]].astype('uint8'))
    ax.set_title('label: {:.0f}'.format(cl_y_test[randIdx[i]]))
```

## Visualizing the sunglasses poisioned test data

```
# Plot some images from the validation set (see https://mrdatascience.com/how-to-plot-mnist-digits-using-matplotlib/)
num = 10
np.random.seed(45)
randIdx = [np.random.randint(10000) for i in range(num)]
num\_row = 2
num_col = 5# plot images
fig, axes = plt.subplots(num_row, num_col, figsize=(3*num_col,3*num_row))
for i in range(num):
   ax = axes[i//num_col, i%num_col]
   ax.imshow(bd_x_test[randIdx[i]].astype('uint8'))
   ax.set_title('label: {:.0f}'.format(bd_y_test[randIdx[i]]))
   ax.set_xticks([])
   ax.set_yticks([])
plt.tight_layout()
plt.show()
                                 label: 0
                                                                                                  label: 0
```



Load the backdoored model.

The backdoor model and its weights can be found here

```
## To-do ##
# First create clones of the original badnet model (by providing the model filepath below)
# The result of repairing B_clone will be B_prime
B = keras.models.load_model("model/bd_net.h5")
B.load_weights("model/bd_weights.h5")
B_clone = keras.models.load_model("model/bd_net.h5")
B_clone.load_weights("model/bd_weights.h5")
Output of the original badnet accuracy on the validation data:
# Get the original badnet model's (B) accuracy on the validation data
cl_label_p = np.argmax(B(cl_x_valid), axis=1)
clean_accuracy = np.mean(np.equal(cl_label_p, cl_y_valid)) * 100
print("Clean validation accuracy before pruning {0:3.6f}".format(clean_accuracy))
K.clear_session()
     Clean validation accuracy before pruning 98.649000
Write code to implement pruning defense
## To-do ##
# Redefine model to output right after the last pooling layer ("pool_3")
intermediate_model = Model(inputs=B.inputs, outputs=B.get_layer('pool_3').output)
# Get feature map for last pooling layer ("pool_3") using the clean validation data and intermediate_model
feature maps cl = intermediate model(cl x valid)
# print(feature_maps_cl.shape)
# Get average activation value of each channel in last pooling layer ("pool 3")
averageActivationsCl = np.mean(feature_maps_cl, axis = (0, 1, 2))
# print(averageActivationsCl)
# print(averageActivationsCl.shape)
# Store the indices of average activation values (averageActivationsCl) in increasing order
idxToPrune = np.argsort(averageActivationsCl)
# Get the conv_4 layer weights and biases from the original network that will be used for prunning
# Hint: Use the get_weights() method (https://stackoverflow.com/questions/43715047/how-do-i-get-the-weights-of-a-layer-in-keras)
lastConvLayerWeights = B.get_layer('conv_3').get_weights()[0]
lastConvLayerBiases = B.get_layer('conv_3').get_weights()[1]
for chIdx in idxToPrune:
    # Prune one channel at a time
    # Hint: Replace all values in channel 'chIdx' of lastConvLayerWeights and lastConvLayerBiases with 0
   # lastConvLayerWeights[:, :, chIdx, :] = 0;
   lastConvLayerWeights[:, :, :, chIdx] = 0;
   lastConvLayerBiases[chIdx] = 0;
   # Update weights and biases of B_clone
    # Hint: Use the set_weights() method
   conv3_layer = [layer for layer in B_clone.layers if layer.name=='conv_3'][0]
    conv3_layer.set_weights([lastConvLayerWeights, lastConvLayerBiases])
   \# Evaluate the updated model's (B_clone) clean validation accuracy
    cl_label_p_valid = np.argmax(B_clone(cl_x_valid), axis=1)
    clean_accuracy_valid = np.mean(np.equal(cl_label_p_valid, cl_y_valid)) * 100
    # print(clean accuracy valid)
    # print(clean_accuracy - clean_accuracy_valid)
    # If drop in clean_accuracy_valid is just greater (or equal to) than the desired threshold compared to clean_accuracy, then save B_clone
    if clean_accuracy - clean_accuracy_valid >= 0.1 * clean_accuracy:
        # Save B clone as B prime and break
        B_prime = B_clone
        B_prime.save("model/B_prime_net_threshold_0.1.h5")
```

```
B_prime.save_weights("model/B_prime_weights_threshold_0.1.h5")
        break
    elif clean_accuracy - clean_accuracy_valid >= 0.04 * clean_accuracy:
        B_prime = B_clone
        B prime.save("model/B prime net threshold 0.04.h5")
        B_prime.save_weights("model/B_prime_weights_threshold_0.04.h5")
    elif clean_accuracy - clean_accuracy_valid >= 0.02 * clean_accuracy:
        B_prime = B_clone
        B_prime.save("model/B_prime_net_threshold_0.02.h5")
        B_prime.save_weights("model/B_prime_weights_threshold_0.02.h5")
# averageActivationsCl = np.mean(feature_maps_cl, axis = (0, 1, 2))
# print(averageActivationsCl)
# print(averageActivationsCl.shape)
# idxToPrune = np.argsort(averageActivationsCl)
# print(idxToPrune)
# thirdConvLayerWeights = B.get_layer('conv_3').get_weights()[0]
# thirdConvLayerBiases = B.get_layer('conv_3').get_weights()[1]
# lastConvLayerWeights = B.get_layer('conv_4').get_weights()[0]
# lastConvLayerBiases = B.get_layer('conv_4').get_weights()[1]
# print(thirdConvLayerWeights.shape)
# print(thirdConvLayerBiases.shape)
# print(lastConvLayerWeights.shape)
# print(lastConvLayerBiases.shape)
# lastConvLayerWeights[:, :, chIdx, :] = 0;
# conv4_layer = [layer for layer in B_clone.layers if layer.name=='conv_4'][0]
# conv4_layer.set_weights([lastConvLayerWeights, lastConvLayerBiases])
# print(lastConvLayerWeights[0, 0, chIdx, 0])
```

Now we need to combine the models into a repaired goodnet G that outputs the correct class if the test input is clean and class N+1 if the input is backdoored. One way to do it is to "subclass" the models in Keras:

```
#https://stackoverflow.com/questions/64983112/keras-vertical-ensemble-model-with-condition-in-between
class G(tf.keras.Model):
   def __init__(self, B, B_prime):
        super(G, self).__init__()
       self.B = B
       self.B_prime = B_prime
   def predict(self,data):
       y = np.argmax(self.B(data), axis=1)
       y_prime = np.argmax(self.B_prime(data), axis=1)
       tmpRes = np.array([y[i] if y[i] == y\_prime[i] else 1283 for i in range(y.shape[0])])
       res = np.zeros((y.shape[0],1284))
       res[np.arange(tmpRes.size),tmpRes] = 1
        return res
   # For small amount of inputs that fit in one batch, directly using call() is recommended for faster execution,
   \# e.g., model(x), or model(x), training=False) is faster then model.predict(x) and do not result in
   # memory leaks (see for more details https://www.tensorflow.org/api_docs/python/tf/keras/Model#predict)
   def call(self,data):
       y = np.argmax(self.B(data), axis=1)
       y_prime = np.argmax(self.B_prime(data), axis=1)
       tmpRes = np.array([y[i] if y[i] == y\_prime[i] else 1283 for i in range(y.shape[0])])
        res = np.zeros((y.shape[0],1284))
        res[np.arange(tmpRes.size),tmpRes] = 1
        return res
```

However, Keras prevents from saving this kind of subclassed model as HDF5 file since it is not serializable. However, we still can use this architecture for model evaluation.

Load the saved B\_prime model

```
## To-do ##
# Provide B_prime model filepath below

B_prime_002 = keras.models.load_model("model/B_prime_net_threshold_0.02.h5")
B_prime_002.load_weights("model/B_prime_weights_threshold_0.02.h5")
```

```
B_prime_004 = keras.models.load_model("model/B_prime_net_threshold_0.04.h5")
B_prime_004.load_weights("model/B_prime_weights_threshold_0.04.h5")
B_prime_01 = keras.models.load_model("model/B_prime_net_threshold_0.1.h5")
B_prime_01.load_weights("model/B_prime_weights_threshold_0.1.h5")
Check performance of the repaired model on the test data:
cl_label_p = np.argmax(B_prime_002.predict(cl_x_test), axis=1)
clean_accuracy_B_prime_002 = np.mean(np.equal(cl_label_p, cl_y_test))*100
print('Clean Classification accuracy for B_prime_002:', clean_accuracy_B_prime_002)
bd_label_p = np.argmax(B_prime_002.predict(bd_x_test), axis=1)
asr B prime 002 = np.mean(np.equal(bd label p, bd y test))*100
print('Attack Success Rate for B_prime_002:', asr_B_prime_002)
cl label p = np.argmax(B prime 004.predict(cl x test), axis=1)
clean_accuracy_B_prime_004 = np.mean(np.equal(cl_label_p, cl_y_test))*100
print('Clean Classification accuracy for B_prime_004:', clean_accuracy_B_prime_004)
bd_label_p = np.argmax(B_prime_004.predict(bd_x_test), axis=1)
asr_B_prime_004 = np.mean(np.equal(bd_label_p, bd_y_test))*100
print('Attack Success Rate for B_prime_004:', asr_B_prime_004)
cl_label_p = np.argmax(B_prime_01.predict(cl_x_test), axis=1)
clean_accuracy_B_prime_01 = np.mean(np.equal(cl_label_p, cl_y_test))*100
print('Clean Classification accuracy for B_prime_01:', clean_accuracy_B_prime_01)
bd_label_p = np.argmax(B_prime_01.predict(bd_x_test), axis=1)
asr_B_prime_01 = np.mean(np.equal(bd_label_p, bd_y_test))*100
print('Attack Success Rate for B_prime_01:', asr_B_prime_01)
    401/401 [=========== ] - 14s 33ms/step
    Clean Classification accuracy for B_prime_002: 95.0584567420109
    401/401 [========== ] - 13s 32ms/step
    Attack Success Rate for B_prime_002: 99.98441153546376
    401/401 [======== ] - 13s 32ms/step
    Clean Classification accuracy for B_prime_004: 89.84411535463757
    401/401 [========== ] - 18s 45ms/step
    Attack Success Rate for B_prime_004: 80.6469212782541
    401/401 [========== ] - 16s 40ms/step
    Clean Classification accuracy for B_prime_01: 84.54403741231489
    401/401 [========= ] - 13s 33ms/step
    Attack Success Rate for B_prime_01: 77.20966484801247
Check performance of the original model on the test data:
cl_label_p = np.argmax(B.predict(cl_x_test), axis=1)
clean_accuracy_B = np.mean(np.equal(cl_label_p, cl_y_test))*100
print('Clean Classification accuracy for B:', clean_accuracy_B)
bd_label_p = np.argmax(B.predict(bd_x_test), axis=1)
asr_B = np.mean(np.equal(bd_label_p, bd_y_test))*100
print('Attack Success Rate for B:', asr_B)
    401/401 [========= ] - 13s 33ms/step
    Clean Classification accuracy for B: 98.62042088854248
    401/401 [=========== ] - 13s 33ms/step
    Attack Success Rate for B: 100.0
Create repaired network
# Repaired network repaired_net
repaired_net_002 = G(B, B_prime_002)
repaired_net_004 = G(B, B_prime_004)
repaired_net_01 = G(B, B_prime_01)
Check the performance of the repaired_net on the test data
cl_label_p = np.argmax(repaired_net_002(cl_x_test), axis=1)
clean_accuracy_repaired_net_002 = np.mean(np.equal(cl_label_p, cl_y_test))*100
print('Clean Classification accuracy for repaired net when threshold is 0.02:', clean_accuracy_repaired_net_002)
hd label n - nn anomay/nensined net AA3/hd v test) avis-1)
```

```
uu_tauet_p = iip.argiiiax(repatreu_net_ooz(uu_x_test), axts-t)
asr_repaired_net_002 = np.mean(np.equal(bd_label_p, bd_y_test))*100
print('Attack Success Rate for repaired net when threshold is 0.02:', asr repaired net 002)
cl_label_p = np.argmax(repaired_net_004(cl_x_test), axis=1)
clean_accuracy_repaired_net_004 = np.mean(np.equal(cl_label_p, cl_y_test))*100
print('Clean Classification accuracy for repaired net when threshold is 0.04:', clean_accuracy_repaired_net_004)
bd_label_p = np.argmax(repaired_net_004(bd_x_test), axis=1)
asr_repaired_net_004 = np.mean(np.equal(bd_label_p, bd_y_test))*100
print('Attack Success Rate for repaired net when threshold is 0.04:', asr_repaired_net_004)
cl_label_p = np.argmax(repaired_net_01(cl_x_test), axis=1)
clean_accuracy_repaired_net_01 = np.mean(np.equal(cl_label_p, cl_y_test))*100
print('Clean Classification accuracy for repaired net when threshold is 0.1:', clean accuracy repaired net 01)
bd_label_p = np.argmax(repaired_net_01(bd_x_test), axis=1)
asr repaired net 01 = np.mean(np.equal(bd label p, bd y test))*100
print('Attack Success Rate for repaired net when threshold is 0.1:', asr_repaired_net_01)
     Clean Classification accuracy for repaired net when threshold is 0.02: 94.90257209664848
    Attack Success Rate for repaired net when threshold is 0.02: 99.98441153546376
     Clean Classification accuracy for repaired net when threshold is 0.04: 89.68043647700702
    Attack Success Rate for repaired net when threshold is 0.04: 80.6469212782541
    Clean Classification accuracy for repaired net when threshold is 0.1: 84.3335931410756
    Attack Success Rate for repaired net when threshold is 0.1: 77.20966484801247
```