## Week 2: Tackle Overfitting with Data Augmentation

Welcome to this assignment! As in the previous week, you will be using the famous `cats vs dogs` dataset to train a model that can classify images of dogs from images of cats. For this, you will create your own Convolutional Neural Network in Tensorflow and leverage Keras' image preprocessing utilities, more so this time around since Keras provides excellent support for augmenting image data.

You will also need to create the helper functions to move the images around the filesystem as you did last week, so if you need to refresh your memory with the `os` module be sure to take a look a the [docs](#).

Let's get started!

**NOTE:** *To prevent errors from the autograder, please avoid editing or deleting non-graded cells in this notebook . Please only put your solutions in between the `### START CODE HERE` and `### END CODE HERE` code comments, and refrain from adding any new cells.*

```
# grader-required-cell

import os
import zipfile
import random
import shutil
import tensorflow as tf
from tensorflow.keras.preprocessing.image import ImageDataGenerator
from shutil import copyfile
import matplotlib.pyplot as plt
```

Download the dataset from its original source by running the cell below.

Note that the `zip` file that contains the images is unzipped under the `/tmp` directory.

```
# If the URL doesn't work, visit https://www.microsoft.com/en-us/download/confirmation.aspx?id=54765
# And right click on the 'Download Manually' link to get a new URL to the dataset

# Note: This is a very large dataset and will take some time to download

!wget --no-check-certificate \
    "https://download.microsoft.com/download/3/E/1/3E1C3F21-ECDB-4869-8368-6DEBA77B919F/kagglecatsanddogs_5340.zip" \
    -O "/tmp/cats-and-dogs.zip"

local_zip = '/tmp/cats-and-dogs.zip'
zip_ref   = zipfile.ZipFile(local_zip, 'r')
zip_ref.extractall('/tmp')
zip_ref.close()
```

```
    --2023-11-15 03:19:57--  https://download.microsoft.com/download/3/E/1/3E1C3F21-ECDB-4869-8368-6DEBA77B919F/kagglecatsanddogs_5340.z
    Resolving download.microsoft.com (download.microsoft.com)... 23.51.213.51, 2600:1406:5400:485::317f, 2600:1406:5400:4a9::317f
    Connecting to download.microsoft.com (download.microsoft.com)|23.51.213.51|:443... connected.
    HTTP request sent, awaiting response... 200 OK
    Length: 824887076 (787M) [application/octet-stream]
    Saving to: '/tmp/cats-and-dogs.zip'

    /tmp/cats-and-dogs. 100%[===================>] 786.67M  4.40MB/s    in 76s

    2023-11-15 03:21:14 (10.3 MB/s) - '/tmp/cats-and-dogs.zip' saved [824887076/824887076]
```

Now the images are stored within the `/tmp/PetImages` directory. There is a subdirectory for each class, so one for dogs and one for cats.

```
# grader-required-cell

source_path = '/tmp/PetImages'

source_path_dogs = os.path.join(source_path, 'Dog')
source_path_cats = os.path.join(source_path, 'Cat')

# Deletes all non-image files (there are two .db files bundled into the dataset)
!find /tmp/PetImages/ -type f ! -name "*.jpg" -exec rm {} +

# os.listdir returns a list containing all files under the given path
print(f"There are {len(os.listdir(source_path_dogs))} images of dogs.")
print(f"There are {len(os.listdir(source_path_cats))} images of cats.")
```

**Expected Output:**

```
  There are 12500 images of dogs.
  There are 12500 images of cats.
```

You will need a directory for cats-v-dogs, and subdirectories for training and validation. These in turn will need subdirectories for 'cats' and 'dogs'. To accomplish this, complete the `create_train_val_dirs` below:

```
# grader-required-cell

# Define root directory
root_dir = '/tmp/cats-v-dogs'

# Empty directory to prevent FileExistsError is the function is run several times
if os.path.exists(root_dir):
  shutil.rmtree(root_dir)

# GRADED FUNCTION: create_train_val_dirs
def create_train_val_dirs(root_path):
  """
  Creates directories for the train and test sets

  Args:
    root_path (string) - the base directory path to create subdirectories from

  Returns:
    None
  """

  ### START CODE HERE

  # HINT:
  # Use os.makedirs to create your directories with intermediate subdirectories
  # Don't hardcode the paths. Use os.path.join to append the new directories to the root_path parameter

  path = os.path.join(root_dir, "training")
  os.makedirs(path)
  path_1 = os.path.join(path, "cats")
  os.makedirs(path_1)
  path_2 = os.path.join(path, "dogs")
  os.makedirs(path_2)
  path = os.path.join(root_dir, "validation")
  os.makedirs(path)
  path_3 = os.path.join(path, "cats")
  os.makedirs(path_3)
  path_4 = os.path.join(path, "dogs")
  os.makedirs(path_4)


  ### END CODE HERE


try:
  create_train_val_dirs(root_path=root_dir)
except FileExistsError:
  print("You should not be seeing this since the upper directory is removed beforehand")
```

```
# grader-required-cell

# Test your create_train_val_dirs function

for rootdir, dirs, files in os.walk(root_dir):
    for subdir in dirs:
        print(os.path.join(rootdir, subdir))

    /tmp/cats-v-dogs/validation
    /tmp/cats-v-dogs/training
    /tmp/cats-v-dogs/validation/cats
    /tmp/cats-v-dogs/validation/dogs
    /tmp/cats-v-dogs/training/cats
    /tmp/cats-v-dogs/training/dogs
```

**Expected Output (directory order might vary):**

Code the `split_data` function which takes in the following arguments:

- SOURCE_DIR: directory containing the files

- TRAINING_DIR: directory that a portion of the files will be copied to (will be used for training)

- VALIDATION_DIR: directory that a portion of the files will be copied to (will be used for validation)

- SPLIT_SIZE: determines the portion of images used for training.

The files should be randomized, so that the training set is a random sample of the files, and the validation set is made up of the remaining files.

For example, if `SOURCE_DIR` is `PetImages/Cat`, and `SPLIT_SIZE` is .9 then 90% of the images in `PetImages/Cat` will be copied to the `TRAINING_DIR` directory and 10% of the images will be copied to the `VALIDATION_DIR` directory.

All images should be checked before the copy, so if they have a zero file length, they will be omitted from the copying process. If this is the case then your function should print out a message such as `"filename is zero length, so ignoring."`. **You should perform this check before the split so that only non-zero images are considered when doing the actual split.**

Hints:

- `os.listdir(DIRECTORY)` returns a list with the contents of that directory.

- `os.path.getsize(PATH)` returns the size of the file

- `copyfile(source, destination)` copies a file from source to destination

- `random.sample(list, len(list))` shuffles a list

```python
# grader-required-cell

# GRADED FUNCTION: split_data
def split_data(SOURCE_DIR, TRAINING_DIR, VALIDATION_DIR, SPLIT_SIZE):
  """
  Splits the data into train and test sets

  Args:
    SOURCE_DIR (string): directory path containing the images
    TRAINING_DIR (string): directory path to be used for training
    VALIDATION_DIR (string): directory path to be used for validation
    SPLIT_SIZE (float): proportion of the dataset to be used for training

  Returns:
    None
  """
  ### START CODE HERE
  files = []
  for filename in os.listdir(SOURCE_DIR):
    file = SOURCE_DIR + filename
    if os.path.getsize(file) > 0:
      files.append(filename)
    else:
      print(filename + ' is zero length, so ignoring.')

    training_length = int(len(files) * SPLIT_SIZE)
    testing_length = int(len(files) - training_length)
    shuffled_set = random.sample(files, len(files))
    training_set = shuffled_set[0:training_length]
    testing_set = shuffled_set[-testing_length:]

  for filename in training_set:
    src_file = SOURCE_DIR + filename
    dest_file = TRAINING_DIR + filename
    copyfile(src_file, dest_file)

  for filename in testing_set:
    src_file = SOURCE_DIR + filename
    dest_file = VALIDATION_DIR + filename
    copyfile(src_file, dest_file)


  ### END CODE HERE
```

```python
# grader-required-cell

# Test your split_data function
```

```python
# Define paths
CAT_SOURCE_DIR = "/tmp/PetImages/Cat/"
DOG_SOURCE_DIR = "/tmp/PetImages/Dog/"

TRAINING_DIR = "/tmp/cats-v-dogs/training/"
VALIDATION_DIR = "/tmp/cats-v-dogs/validation/"

TRAINING_CATS_DIR = os.path.join(TRAINING_DIR, "cats/")
VALIDATION_CATS_DIR = os.path.join(VALIDATION_DIR, "cats/")

TRAINING_DOGS_DIR = os.path.join(TRAINING_DIR, "dogs/")
VALIDATION_DOGS_DIR = os.path.join(VALIDATION_DIR, "dogs/")

# Empty directories in case you run this cell multiple times
if len(os.listdir(TRAINING_CATS_DIR)) > 0:
  for file in os.scandir(TRAINING_CATS_DIR):
    os.remove(file.path)
if len(os.listdir(TRAINING_DOGS_DIR)) > 0:
  for file in os.scandir(TRAINING_DOGS_DIR):
    os.remove(file.path)
if len(os.listdir(VALIDATION_CATS_DIR)) > 0:
  for file in os.scandir(VALIDATION_CATS_DIR):
    os.remove(file.path)
if len(os.listdir(VALIDATION_DOGS_DIR)) > 0:
  for file in os.scandir(VALIDATION_DOGS_DIR):
    os.remove(file.path)

# Define proportion of images used for training
split_size = .9

# Run the function
# NOTE: Messages about zero length images should be printed out
split_data(CAT_SOURCE_DIR, TRAINING_CATS_DIR, VALIDATION_CATS_DIR, split_size)
split_data(DOG_SOURCE_DIR, TRAINING_DOGS_DIR, VALIDATION_DOGS_DIR, split_size)

# Your function should perform copies rather than moving images so original directories should contain unchanged images
print(f"\n\nOriginal cat's directory has {len(os.listdir(CAT_SOURCE_DIR))} images")
print(f"Original dog's directory has {len(os.listdir(DOG_SOURCE_DIR))} images\n")

# Training and validation splits. Check that the number of images matches the expected output.
print(f"There are {len(os.listdir(TRAINING_CATS_DIR))} images of cats for training")
print(f"There are {len(os.listdir(TRAINING_DOGS_DIR))} images of dogs for training")
print(f"There are {len(os.listdir(VALIDATION_CATS_DIR))} images of cats for validation")
print(f"There are {len(os.listdir(VALIDATION_DOGS_DIR))} images of dogs for validation")
```

```
    666.jpg is zero length, so ignoring.
    11702.jpg is zero length, so ignoring.


    Original cat's directory has 12501 images
    Original dog's directory has 12501 images

    There are 11250 images of cats for training
    There are 11250 images of dogs for training
    There are 1250 images of cats for validation
    There are 1250 images of dogs for validation
```

**Expected Output:**

```
 666.jpg is zero length, so ignoring.
 11702.jpg is zero length, so ignoring.



 Original cat's directory has 12500 images
 Original dog's directory has 12500 images


 There are 11249 images of cats for training
 There are 11249 images of dogs for training
 There are 1250 images of cats for validation
 There are 1250 images of dogs for validation
```

Now that you have successfully organized the data in a way that can be easily fed to Keras' `ImageDataGenerator`, it is time for you to code the generators that will yield batches of images, both for training and validation. For this, complete the `train_val_generators` function below.

Something important to note is that the images in this dataset come in a variety of resolutions. Luckily, the `flow_from_directory` method allows you to standarize this by defining a tuple called `target_size` that will be used to convert each image to this target resolution. **For this exercise use a `target_size` of (150, 150).**

```
# grader-required-cell

# GRADED FUNCTION: train_val_generators
def train_val_generators(TRAINING_DIR, VALIDATION_DIR):
  """
  Creates the training and validation data generators

  Args:
    TRAINING_DIR (string): directory path containing the training images
    VALIDATION_DIR (string): directory path containing the testing/validation images

  Returns:
    train_generator, validation_generator - tuple containing the generators
  """
  ### START CODE HERE

  # Instantiate the ImageDataGenerator class (don't forget to set the arguments to augment the images)
  train_datagen = ImageDataGenerator(rescale=1.0/255.,
                                     rotation_range=40,
                                     width_shift_range=0.2,
                                     height_shift_range=0.2,
                                     shear_range=0.2,
                                     zoom_range=0.2,
                                     horizontal_flip=True,
                                     fill_mode='nearest')

  # Pass in the appropriate arguments to the flow_from_directory method
  train_generator = train_datagen.flow_from_directory(directory=TRAINING_DIR,
                                                      batch_size=128,
                                                      class_mode='binary',
                                                      target_size=(150, 150))

  # Instantiate the ImageDataGenerator class (don't forget to set the rescale argument)
  validation_datagen = ImageDataGenerator(rescale=1.0/255.)

  # Pass in the appropriate arguments to the flow_from_directory method
  validation_generator = validation_datagen.flow_from_directory(directory=VALIDATION_DIR,
                                                               batch_size=16,
                                                               class_mode='binary',
                                                               target_size=(150, 150))
  ### END CODE HERE
  return train_generator, validation_generator
```

```
# grader-required-cell

# Test your generators
train_generator, validation_generator = train_val_generators(TRAINING_DIR, VALIDATION_DIR)
```

```
    Found 22499 images belonging to 2 classes.
    Found 2499 images belonging to 2 classes.
```

**Expected Output:**

```
 Found 22498 images belonging to 2 classes.
 Found 2500 images belonging to 2 classes.
```

One last step before training is to define the architecture of the model that will be trained.

Complete the `create_model` function below which should return a Keras' `Sequential` model.

Aside from defining the architecture of the model, you should also compile it so make sure to use a `loss` function that is compatible with the `class_mode` you defined in the previous exercise, which should also be compatible with the output of your network. You can tell if they aren't compatible if you get an error during training.

**Note that you should use at least 3 convolution layers to achieve the desired performance.**

```
# grader-required-cell

# GRADED FUNCTION: create_model
def create_model():
  # DEFINE A KERAS MODEL TO CLASSIFY CATS V DOGS
  # USE AT LEAST 3 CONVOLUTION LAYERS

  ### START CODE HERE

  model = tf.keras.models.Sequential([
```

```
    tf.keras.layers.Conv2D(16, (3,3), activation='relu', input_shape=(150, 150, 3)),
    tf.keras.layers.MaxPooling2D(2, 2),
    tf.keras.layers.Conv2D(32, (3,3), activation='relu'),
    tf.keras.layers.MaxPooling2D(2,2),
    tf.keras.layers.Conv2D(64, (3,3), activation='relu'),
    tf.keras.layers.MaxPooling2D(2,2),
    tf.keras.layers.Conv2D(64, (3,3), activation='relu'),
    tf.keras.layers.MaxPooling2D(2,2),
    tf.keras.layers.Flatten(),
    # 512 neuron hidden layer
    tf.keras.layers.Dense(512, activation='relu'),
    tf.keras.layers.Dense(128, activation='relu'),
    # Only 1 output neuron. It will contain a value from 0-1 where 0 for 1 class ('horses') and 1 for the other ('humans')
    tf.keras.layers.Dense(1, activation='sigmoid')
  ])


  model.compile(optimizer=tf.keras.optimizers.RMSprop(learning_rate=0.001),
                loss='binary_crossentropy',
                metrics=['accuracy'])

  ### END CODE HERE

  return model
```

Now it is time to train your model!

Note: You can ignore the `UserWarning: Possibly corrupt EXIF data.` warnings.

```
# Get the untrained model
model = create_model()

# Train the model
# Note that this may take some time.
history = model.fit(train_generator,
                    epochs=15,
                    verbose=1,
                    validation_data=validation_generator)

    Epoch 1/15
    176/176 [==============================] - ETA: 0s - loss: 0.6892 - accuracy: 0.5477/usr/local/lib/python3.10/dist-packages/PIL/Tif
      warnings.warn(str(msg))
    176/176 [==============================] - 180s 958ms/step - loss: 0.6892 - accuracy: 0.5477 - val_loss: 0.6696 - val_accuracy: 0.56
    Epoch 2/15
    176/176 [==============================] - 167s 945ms/step - loss: 0.6443 - accuracy: 0.6260 - val_loss: 0.6341 - val_accuracy: 0.63
    Epoch 3/15
    176/176 [==============================] - 170s 969ms/step - loss: 0.6076 - accuracy: 0.6695 - val_loss: 0.7729 - val_accuracy: 0.66
    Epoch 4/15
    176/176 [==============================] - 167s 949ms/step - loss: 0.5867 - accuracy: 0.6867 - val_loss: 0.5381 - val_accuracy: 0.73
    Epoch 5/15
    176/176 [==============================] - 167s 949ms/step - loss: 0.5614 - accuracy: 0.7092 - val_loss: 0.4802 - val_accuracy: 0.77
    Epoch 6/15
    176/176 [==============================] - 170s 966ms/step - loss: 0.5422 - accuracy: 0.7244 - val_loss: 0.4944 - val_accuracy: 0.76
    Epoch 7/15
    176/176 [==============================] - 168s 957ms/step - loss: 0.5195 - accuracy: 0.7439 - val_loss: 0.4592 - val_accuracy: 0.78
    Epoch 8/15
    176/176 [==============================] - 170s 962ms/step - loss: 0.5114 - accuracy: 0.7470 - val_loss: 0.5518 - val_accuracy: 0.71
    Epoch 9/15
    176/176 [==============================] - 167s 947ms/step - loss: 0.4930 - accuracy: 0.7589 - val_loss: 0.4374 - val_accuracy: 0.79
    Epoch 10/15
    176/176 [==============================] - 167s 949ms/step - loss: 0.4796 - accuracy: 0.7684 - val_loss: 0.4653 - val_accuracy: 0.78
    Epoch 11/15
    176/176 [==============================] - 166s 942ms/step - loss: 0.4636 - accuracy: 0.7791 - val_loss: 0.4540 - val_accuracy: 0.77
    Epoch 12/15
    176/176 [==============================] - 168s 956ms/step - loss: 0.4512 - accuracy: 0.7881 - val_loss: 0.3825 - val_accuracy: 0.82
    Epoch 13/15
    176/176 [==============================] - 166s 942ms/step - loss: 0.4303 - accuracy: 0.7992 - val_loss: 0.3802 - val_accuracy: 0.82
    Epoch 14/15
    176/176 [==============================] - 169s 958ms/step - loss: 0.4196 - accuracy: 0.8049 - val_loss: 0.3891 - val_accuracy: 0.82
    Epoch 15/15
    176/176 [==============================] - 165s 938ms/step - loss: 0.4028 - accuracy: 0.8155 - val_loss: 0.3388 - val_accuracy: 0.85
```

Once training has finished, you can run the following cell to check the training and validation accuracy achieved at the end of each epoch.
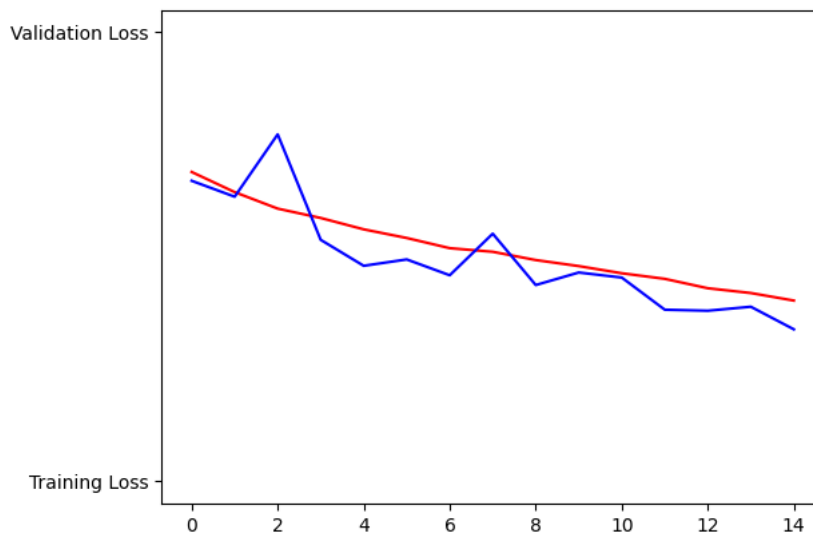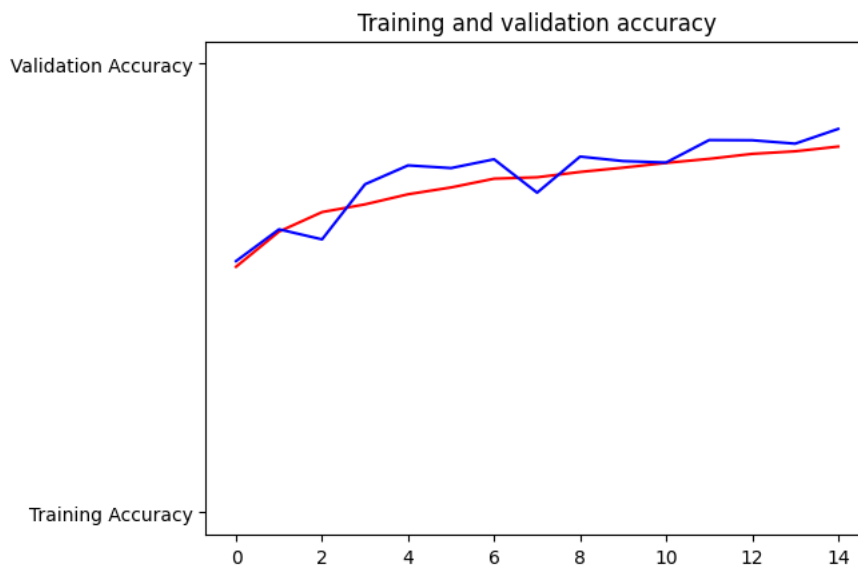
**To pass this assignment, your model should achieve a training and validation accuracy of at least 80% and the final testing accuracy should be either higher than the training one or have a 5% difference at maximum**. If your model didn't achieve these thresholds, try training again with a different model architecture, remember to use at least 3 convolutional layers or try tweaking the image augmentation process.

You might wonder why the training threshold to pass this assignment is significantly lower compared to last week's assignment. Image augmentation does help with overfitting but usually this comes at the expense of requiring more training time. To keep the training time

reasonable, the same number of epochs as in the previous assignment are kept.

However, as an optional exercise you are encouraged to try training for more epochs and to achieve really good training and validation accuracies.

```
#-----------------------------------------------------------
# Retrieve a list of list results on training and test data
# sets for each training epoch
#-----------------------------------------------------------
acc=history.history['accuracy']
val_acc=history.history['val_accuracy']
loss=history.history['loss']
val_loss=history.history['val_loss']

epochs=range(len(acc)) # Get number of epochs

#-------------------------------------------------
# Plot training and validation accuracy per epoch
#-------------------------------------------------
plt.plot(epochs, acc, 'r', "Training Accuracy")
plt.plot(epochs, val_acc, 'b', "Validation Accuracy")
plt.title('Training and validation accuracy')
plt.show()
print("")

#-------------------------------------------------
# Plot training and validation loss per epoch
#-------------------------------------------------
plt.plot(epochs, loss, 'r', "Training Loss")
plt.plot(epochs, val_loss, 'b', "Validation Loss")
plt.show()
```





You will probably encounter that the model is overfitting, which means that it is doing a great job at classifying the images in the training set but struggles with new data. This is perfectly fine and you will learn how to mitigate this issue in the upcomming week.

Before closing the assignment, be sure to also download the `history.pkl` file which contains the information of the training history of your model. You can download this file by running the cell below:

```
def download_history():
  import pickle
  from google.colab import files

  with open('history_augmented.pkl', 'wb') as f:
    pickle.dump(history.history, f)

  files.download('history_augmented.pkl')

download_history()
```

## ⌄ Download your notebook for grading

Along with the `history.pkl` file, you will also need to submit your solution notebook for grading. The following code cells will check if this notebook's grader metadata (i.e. hidden data in the notebook needed for grading) is not modified by your workspace. This will ensure that the autograder can evaluate your code properly. Depending on its output, you will either:

- *if the metadata is intact*: Download the current notebook. Click on the File tab on the upper left corner of the screen then click on `Download` -> `Download .ipynb`. You can name it anything you want as long as it is a valid `.ipynb` (jupyter notebook) file.

- *if the metadata is missing*: A new notebook with your solutions will be created on this Colab workspace. It should be downloaded automatically and you can submit that to the grader.

```
# Download metadata checker
!wget -nc https://storage.googleapis.com/tensorflow-1-public/colab_metadata_checker.py
```

```
    --2023-11-15 05:28:01--  https://storage.googleapis.com/tensorflow-1-public/colab_metadata_checker.py
    Resolving storage.googleapis.com (storage.googleapis.com)... 142.251.2.207, 2607:f8b0:4023:c0d::cf
    Connecting to storage.googleapis.com (storage.googleapis.com)|142.251.2.207|:443... connected.
    HTTP request sent, awaiting response... 200 OK
    Length: 1997 (2.0K) [text/x-python-script]
    Saving to: 'colab_metadata_checker.py'

    colab_metadata_chec 100%[===================>]   1.95K  --.-KB/s    in 0s

    2023-11-15 05:28:01 (40.3 MB/s) - 'colab_metadata_checker.py' saved [1997/1997]
```

```
import colab_metadata_checker

# Please see the output of this cell to see which file you need to submit to the grader
colab_metadata_checker.run('C2W2_Assignment_fixed.ipynb')
```

```
    Grader metadata detected! You can download this notebook by clicking `File > Download > Download as .ipynb` and submit it to the gr
```

**Please disregard the following note if the notebook metadata is detected**

*Note: Just in case the download fails for the second point above, you can also do these steps:*

- *Click the Folder icon on the left side of this screen to open the File Manager.*
- *Click the Folder Refresh icon in the File Manager to see the latest files in the workspace. You should see a file ending with a `_fixed.ipynb`.*
- *Right-click on that file to save locally and submit it to the grader.*

**Congratulations on finishing this week's assignment!**

You have successfully implemented a convolutional neural network that classifies images of cats and dogs, along with the helper functions needed to pre-process the images!

**Keep it up!**