

File
 Edit
 View
 Run
 Kernel
 Git
 Tabs
 Settings
 Help

DL0321EN-2-1-Data-Prepa

Code

 git
 Run as Pipeline

Python

IBM Developer SKILLS NETWORK

Data Preparation

Objective

In this lab, you will learn how to load images and manipulate them for training using Keras ImageDataGenerator.

Table of Contents

1. Download Data
 2. Import Libraries and Packages
 3. Construct an ImageDataGenerator Instance
 4. Visualize Batches of Images
 5. Questions

Download Data

For your convenience, I have placed the data on a server which you can retrieve easily using the **wget** command. So let's run the following line of code to get the data. Given the large size of the image dataset, it might take some time depending on your internet speed.

```
[ ]: ## get the data
wget https://s3-us-east-1.amazonaws.com/cognitive-class-datasets/concrete_data_week2.zip
...

And now if you check the left directory pane, you should see the zipped file concrete_data_week2.zip appear. So, let's go ahead and unzip the file to access the images. Given the large number of images in the dataset, this might take a couple of minutes, so please be patient, and wait until the code finishes running.
```

```
[ ]: !unzip concrete_data_week2.zip
...

Now, you should see the folder concrete_data_week2 appear in the left pane. If you open this folder by double-clicking on it, you will find that it contains two folders: Positive and Negative. These are the same folders that we saw in the lab in the previous module of this course, where Negative is the negative class and it represents the concrete images with no cracks. Positive on the other hand is the positive class and represents the concrete images with cracks.


Important Note: There are thousands and thousands of images in each folder, so please don't attempt to double click on the Negative and Positive folders. This may consume all of your memory and you may end up with a 50* error. So please DO NOT DO IT.



### Import Libraries and Packages



Before we proceed, let's import the libraries and packages that we will need to complete the rest of this lab.



```
[]: import os
import numpy as np
import matplotlib.pyplot as plt

import keras
from keras.preprocessing.image import ImageDataGenerator
...
```



You can check the content of ./concrete_data_week2 by running the following:



```
[]: !ls ./concrete_data_week2
```



or the following:



```
[]: os.listdir('concrete_data_week2')
```



### Construct an ImageDataGenerator Instance



In this section, you will learn how to define a Keras ImageDataGenerator instance and use it to load and manipulate data for building a deep learning model.



Before we proceed, let's define a variable that represents the path to the folder containing our data which is concrete_data_week2 in this case.



```
[]: dataset_dir = './concrete_data_week2'
```



Keras ImageDataGenerator requires images be arranged in a certain folder hierarchy, where the main directory would contain folders equal to the number of classes in your problem. Since in this case we are trying to build a classifier of two classes, then our main directory, which is concrete_data_week2, should contain two folders, one for each class. This has already been done for you as the negative images are in one folder and the positive images are in another folder.



Let's go ahead and define an instance of the Keras ImageDataGenerator.


```

Support/Feedback

Standard ImageDataGenerator

You can define a standard one like this, where you are simply using the ImageDataGenerator to train your model in batches.

```
[ ]: # instantiate your image data generator
data_generator = ImageDataGenerator()
```

Next, you use the `flow_from_directory` methods to loop through the images in batches. In this method, you pass the directory where the images reside, the size of each batch, `batch_size`, and since batches are sampled randomly, then you can also specify a random seed, `seed`, if you would like to reproduce the batch sampling. In case you would like to resize your images, then you can using the `target_size` argument to accomplish that.

```
[ ]: image_generator = data_generator.flow_from_directory(
    dataset_dir,
    batch_size=4,
    class_mode='categorical',
    seed=24
)
```

What is great about this method, is it prints a summary of it found in the directory passed. Here, it found 40,000 images in total belonging to 2 classes.

Now, to access the batches, you use the `next` method as follows:

```
[ ]: first_batch = image_generator.next()
first_batch
***
```

As you can see, this returned the images along with their labels. Therefore, the following returns the images only,

```
[ ]: first_batch_images = image_generator.next()[0]
first_batch_images
***
```

and the following returns the labels only.

```
[ ]: first_batch_labels = image_generator.next()[1]
first_batch_labels
***
```

Custom ImageDataGenerator

Did you know? IBM Watson Studio lets you build and deploy an AI solution, using the best of open source and IBM software and giving your team a single environment to work in. [Learn more here.](#)

You can also specify some transforms, like scaling, rotations, and flips, that you would like applied to the images when you define an ImageDataGenerator object. Say you want to normalize your images, then you can define your ImageDataGenerator instance as follows:

```
[ ]: # instantiate your image data generator
data_generator = ImageDataGenerator(
    rescale=1./255
)
***
```

And then you proceed with defining your `image_generator` using the `flow_from_directory` method, just like before.

```
[ ]: image_generator = data_generator.flow_from_directory(
    dataset_dir,
    batch_size=4,
    class_mode='categorical',
    seed=24
)
***
```

However, now we explore the first batch using the `next` method,

```
[ ]: first_batch = image_generator.next()
first_batch
***
```

we find that the values are not integer values anymore, but scaled resolution since the original number are divided by 255.

You can learn more about the Keras ImageDataGeneration class [here](#).

Visualize Batches of Images

Let write some code to visualize a batch. We will use subplots in order to make visualizing the images easier.

Recall that we can access our batch images as follows:

```
first_batch_images = image_generator.next()[0] # first batch
```

```
second_batch_images = image_generator.next()[0] # second batch
```

and so on.

```
[ ]: fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(20, 10))...# define your figure and axes

ind = 0
for ax1 in axs:
    for ax2 in ax1:
        image_data = first_batch_images[ind]
        ax2.imshow(image_data)
        ind += 1

fig.suptitle('First Batch of Concrete Images')..
plt.show()
***
```

Remember that batches are sampled randomly from the data. In our first batch, we ended up with two negative image and two positive images.

Important Note: Because of a bug with the `imshow` function in Matplotlib, if you are plotting the unscaled RGB images, then the `image_data` in the code above would be like this:

```
image_data = first_batch_images[ind].astype(np.uint8)
```

where the image_data is cast to uint8 before you call the `imshow` function.

Questions

Question: Create a plot to visualize the images in the third batch.

```
[ ]: ## You can use this cell to type your code to answer the above question
```

Question: How many images from each class are in the fourth batch?

```
[ ]: ## You can use this cell to type your code to answer the above question
```

Question: Create a plot to visualize the second image in the fifth batch.

```
[ ]: ## You can use this cell to type your code to answer the above question
```

Question: How many images from each class are in the fifth batch?

```
[ ]: ## You can use this cell to type your code to answer the above question
```

Make sure to answer the above questions as the quiz in this module is heavily based on them.

Thank you for completing this lab!

This notebook was created by Alex Aklson. I hope you found this lab interesting and educational.

This notebook is part of a course on **Coursera** called *AI Capstone Project with Deep Learning*. If you accessed this notebook outside the course, you can take this course online by clicking [here](#).

Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2020-09-18	2.0	Shubham	Migrated Lab to Markdown and added to course repo in GitLab

Copyright © 2020 [IBM Developer Skills Network](#). This notebook and its source code are released under the terms of the [MIT License](#).