Bài 32 - Kĩ thuật tensorflow Dataset

09 Apr 2020 - phamdinhkhanh

Menu

- 1. Vai trò của tensorflow Dataset
- · 2. Đinh nghĩa generator
- 3. Các cách khởi tạo một Dataset
 - 3.1. In Memory Dataset
 - 3.2. Generator Dataset
 - 3.2.1. Sử dụng ImageGenerator
 - 3.2.2. Customize ImageGenerator
- 4. Tổng kết
- 5. Tài liệu tham khảo

1. Vai trò của tensorflow Dataset

Chắc hẳn các bạn từng thắc mắc vì sao trong deep learning các bộ dữ liệu bigdata có kích thước rất lớn mà các máy tính có RAM nhỏ hơn vẫn có thể huấn luyện được?

Xuất phát từ lý do đó, bài này mình sẽ lý giải các cách thức dữ liệu có thể được truyền vào mô hình để huấn luyện theo cách tiếp cận dễ hiểu nhất. Các bạn chuyên gia và giàu kinh nghiệm huấn luyện mô hình có thể bỏ qua bài viết này vì nó khá cơ bản.

Vì sao có thể truyền các bộ dữ liệu lớn vào mô hình huấn luyện?

Các bộ dữ liệu deep learning thường có kích thước rất lớn. Trong quá trình huấn luyện các model deep learning chúng ta không thể truyền toàn bộ dữ liệu vào mô hình cùng một lúc bởi dữ liệu thường có kích thước lớn hơn RAM máy tính. Xuất phát từ lý do này, các framework deep learning đều hỗ trợ các hàm huấn luyện mô hình theo generator. Dữ liệu sẽ không được khởi tạo ngay toàn bộ từ đầu mà sẽ huấn luyện đến đâu sẽ được khởi tạo đến đó theo từng phần nhỏ gọi là batch.

Tùy theo định dạng dữ liệu là text, image, data frame, numpy array,... mà chúng ta sẽ sử dụng những module tạo dữ liệu huấn luyện khác nhau.

Vậy thì với từng kiểu dữ liệu khác nhau sẽ có phương pháp xử lý như thế nào để đưa vào huấn luyện mô hình? Có những kĩ thuật khởi tạo dataset trong tensorflow nào? Bài viết này mình sẽ giới thiệu tới các bạn.

2. Định nghĩa generator

generator có thể coi là một người vay nợ, được quyền sử dụng tiền của người khác mà không trả ngay. Nếu chúng ta coi tiền là dữ liệu thì ta có thể hình dung generator sẽ sử dụng và biến đổi dữ liệu như cách người vay nợ sử dụng tiền vào các mục đích của mình. Tuy nhiên dữ liệu sau biến đổi không được trả về như các hàm return thông thường của python.

Để đơn giản hóa mình lấy ví dụ một hàm tính lãi suất phải trả theo năm như sau:

Giả sử một người vay n món nợ với cùng lãi suất là 1%/tháng. Để tính lãi suất phải trả c $\frac{1}{100}$ pác khoản vay chúng ta có thể sử dụng vòng for và tính để tính kết quả trong 1 lần.

Note: Bạn đọc có thể mở google colab để cùng thực hành tensorflow Dataset - khanh blog (https://colab.research.google.com/drive/1mVwq7Py4Rv2MCDp1lOD8mQ1FXQWJDLlp)

```
1
       import numpy as np
2
       from datetime import datetime
3
       def _interest_rate(month):
4
5
         return (1+0.01)**month - 1
6
7
8
       periods = [1, 3, 6, 9, 12]
       scales = [_interest_rate(month) for month in periods]
9
10
       print('scales of origin balance: ', scales)
      scales of origin balance: [0.010000000000000, 0.030301000000000133
1
```

Nếu sử dụng generator thì chúng ta chỉ việc thay return bằng yield.

```
def _gen_interest_rate(month):
    yield (1+0.01)**month - 1

periods = [1, 3, 6, 9, 12]
    scales = [_gen_interest_rate(month) for month in periods]
    print('scales of origin balance: ', scales)

scales of origin balance: [<generator object _gen_interest_rate at 0x]</pre>
```

Ta thấy generator sẽ không trả về kết quả ngay mà chỉ tạo sẵn các ô nhớ lưu hàm generator mô tả cách tính lãi suất. Do đó chúng ta sẽ không tốn chi phí thời gian để thực hiện các phép tính. Thực tế là chúng ta đang nợ máy tính kết quả trả về. Chỉ khi nào được chủ nợ gọi tên bằng cách kích hoạt trong hàm <code>next()</code> thì mới tính kết quả.

```
1     [next(_gen_interest_rate(0.01, n)) for n in periods]
1     [1.01,
2      1.0303010000000001,
3      1.0615201506010001,
4      1.0936852726843609,
5      1.1268250301319698]
```

Chúng ta có thể thấy generator có lợi thế là:

- Không sinh toàn bộ dữ liệu cùng một lúc, do đó sẽ nâng cao hiệu suất vì sử dụng ít bộ nhớ hơn.
- Không phải chờ toàn bộ các vòng lặp được xử lý xong thì mới xử lý tiếp nên tiết kiệm thời gian tính toán.

Top Đó chính là lý do generator chính là giải pháp được lựa chọn cho huấn luyện mô hình deep learning với dữ liệu lớn.

3. Các cách khởi tạo một Dataset

Dataset là một class của tensorflow được sử dụng để wrap dữ liệu trước khi truyền vào mô hình để huấn luyện. Bạn hình dung dữ liệu của bạn có input là ma trận X và output là Y. Ban đầu X và Y chỉ là các dữ liệu thô định dạng numpy. Tất nhiên chúng ta có thể truyền trực tiếp chúng vào hàm fit() của mô hình. Nhưng để kiểm soát được X và Y chẳng hạn như fit vào với batch size bằng bao nhiêu? có shuffle dữ liệu hay không thì chúng ta nên wrap chúng trong tf.Dataset.

Có 2 phương pháp chính để khởi tạo một tf.Dataset trong tensorflow:

- In memory Dataset: Khởi tạo các dataset ngay từ đầu và dữ liệu được lưu trữ trên memory.
- Generator Dataset: Dữ liệu được sinh ra theo từng batch và xen kẽ với quá trình huấn luyện từ các hàm khởi tạo generator.

Phương pháp In memory Dataset sẽ phù hợp với các bộ dữ liệu kích thước nhỏ mà RAM có thể load được. Quá trình huấn luyện theo cách này thì nhanh hơn so với phương pháp Generator Dataset vì dữ liệu đã được chuẩn bị sẵn mà không tốn thời gian chờ khởi tạo batch. Tuy nhiên dễ xảy ra out of memory trong quá trình huấn luyện.

Theo cách Generator Dataset chúng ta sẽ qui định cách mà dữ liệu được tạo ra như thế nào thông qua một hàm generator. Quá trình huấn luyện đến đâu sẽ tạo batch đến đó. Do đó các bộ dữ liệu big data có thể được load theo từng batch sao cho kích thước vừa được dung lượng RAM. Theo cách huấn luyện này chúng ta có thể huấn luyện được các bộ dữ liệu có kích thước lớn hơn nhiều so với RAM bằng cách chia nhỏ chúng theo batch. Đồng thời có thể áp dụng thêm các step preprocessing data trước khi dữ liệu được đưa vào huấn luyện. Do đó đây thường là phương pháp được ưa chuộng khi huấn luyện các model deep learning.

3.1. In Memory Dataset

Bên dưới chúng ta sẽ thử nghiệm khởi tạo một dtaset trên tensorflow theo phuwong pháp In Memory Dataset . Bộ dữ liệu được lựa chọn là chữ số viết tay mnist với kích thước của tập train và validation lần lượt là 60000 và 10000 .

```
1
      %tensorflow_version 2.x
2
3
      from google.colab import drive
4
      import os
5
      drive.mount("/content/gdrive")
6
7
      path = 'gdrive/My Drive/Colab Notebooks/TensorflowData'
      os.chdir(path)
8
9
      os.listdir()
1
      Enter your authorization code:
      . . . . . . . . . .
2
3
      Mounted at /content/gdrive
4
      ['Dog-Cat-Classifier', 'TensorflowDataset.ipynb']
5
```

```
from tensorflow.keras.datasets import mnist
1
2
3
     (X_train, y_train), (X_test, y_test) = mnist.load_data()
4
     print(X_train.shape)
5
     print(X_test.shape)
     print(y_train.shape)
6
7
     print(y_test.shape)
1
     Downloading data from https://storage.googleapis.com/tensorflow/tf-kera
     2
3
     (60000, 28, 28)
     (10000, 28, 28)
4
     (60000,)
5
6
     (10000,)
```

Như vậy các dữ liệu train và test của bộ dữ liệu mnist đã được load vào bộ nhớ. Tiếp theo chúng ta sẽ khởi tạo Dataset cho những dữ liệu in memory này bằng hàm

tf.data.Dataset.from_tensor_slices(). Hàm này sẽ khai báo dữ liệu đầu vào cho mô hình huấn luyện.

```
import tensorflow as tf
train_dataset = tf.data.Dataset.from_tensor_slices((X_train, y_train))
valid_dataset = tf.data.Dataset.from_tensor_slices((X_test, y_test))
```

Khi đó chúng ta đã có thể fit vào mô hình huấn luyện các dữ liệu được truyền vào tf.Dataset là (X_train, y_train).

Chúng ta cũng có thể áp dụng các phép biến đổi bằng các hàm như Dataset.map() hoặc Dataset.batch() để biến đổi dữ liệu trước khi fit vào model. Các bạn xem thêm tại tf.Dataset (https://www.tensorflow.org/api_docs/python/tf/data/Dataset). Chẳng hạn trước khi truyền batch vào huấn luyện tôi sẽ thực hiện chuẩn hóa batch theo phân phối chuẩn.

```
import numpy as np
1
       from tensorflow.keras.backend import std, mean
2
3
       from tensorflow.math import reduce_std, reduce_mean
4
       def _normalize(X_batch, y_batch):
5
6
7
         X_batch: matrix digit images, shape batch_size x 28 x 28
         y_batch: labels of digit.
8
9
10
         X_batch = tf.cast(X_batch, dtype = tf.float32)
         # Padding về 2 chiều các giá trị 0 để được shape là 32 x 32
11
12
         pad = tf.constant([[0, 0], [2, 2], [2, 2]])
         X_batch = tf.pad(X_batch, paddings=pad, mode='CONSTANT', constant_value
13
14
         X_batch = tf.expand_dims(X_batch, axis=-1)
         mean = reduce_mean(X_batch)
15
16
         std = reduce_std(X_batch)
17
         X_{norm} = (X_{batch-mean})/std
18
         return X_norm, y_batch
19
       train_dataset = train_dataset.batch(32).map(_normalize)
20
21
       valid_dataset = valid_dataset.batch(32).map(_normalize)
```

train_dataset và valid_dataset lần lượt thực hiện các bước xử lý dữ liệu sau:

- Hàm .batch(32): Trích xuất ra từ list (X_train, y_train) các batch_size có kích thước là 32.
- Hàm .map(_normalize): Mapping đầu vào là các batch (X_batch, y_batch) kích thước 32 vào hàm số _normalize(). Kết quả trả về là giá trị đã chuẩn hóa theo batch của X_batch và y_batch. Dữ liệu này sẽ được sử dụng để huấn luyện model.

Huấn luyện và kiểm định model

```
from tensorflow.keras.applications import MobileNet
1
2
       from tensorflow.keras.models import Sequential
3
       from tensorflow.keras.layers import Dense, Flatten
       from tensorflow.keras.optimizers import Adam
4
5
       base_extractor = MobileNet(input_shape = (32, 32, 1), include_top = Face
6
7
       flat = Flatten()
8
       den = Dense(10, activation='softmax')
       model = Sequential([base_extractor,
9
10
                           flat,
                           den])
11
12
       model.summary()
```

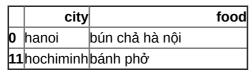
```
Model: "sequential_5"
1
2
3
    Layer (type)
                      Output Shape
                                      Param #
4
    ______
5
    mobilenet_1.00_32 (Model)
                      (None, 1, 1, 1024)
                                      3228288
6
7
    flatten_5 (Flatten)
                      (None, 1024)
                                      0
8
9
    dense_11 (Dense)
                                      10250
                      (None, 10)
10
    ______
    Total params: 3,238,538
11
12
    Trainable params: 3,216,650
13
    Non-trainable params: 21,888
14
1
   model.compile(Adam(), loss='sparse_categorical_crossentropy', metrics :
2
   model.fit(train_dataset,
3
         validation_dataset = valid_dataset,
4
          epochs = 5)
1
   Epoch 1/5
   2
3
   Epoch 2/5
   4
5
   Epoch 3/5
   6
```

3.2. Generator Dataset

Theo cách khởi tạo từ generator chúng ta sẽ không phải ghi nhớ toàn bộ dữ liệu vào RAM. Thay vào đó có thể tao dữ liêu trong quá trình huấn luyên ở mỗi lượt fit từng batch.

Giả sử bên dưới chúng ta có tên các món ăn được chia thành hai nhóm thuộc các địa phương 'hà nội' và 'hồ chí minh'. Chúng ta sẽ khởi tạo data generator để sinh dữ liệu cho mô hình phân loại món ăn theo địa phương.

```
import pandas as pd
1
2
      hanoi = ['bún chả hà nội', 'chả cá lã vọng hà nội', 'cháo lòng hà nội'
3
      hochiminh = ['bánh canh sài gòn', 'hủ tiếu nam vang sài gòn', 'hủ tiếu
4
      city = ['hanoi'] * len(hanoi) + ['hochiminh'] * len(hochiminh)
5
6
      corpus = hanoi+hochiminh
7
      data = pd.DataFrame({'city': city, 'food': corpus})
8
9
      data.sample(5)
```



```
1
       class Voc(object):
2
         def __init__(self, corpus):
3
           self.corpus = corpus
4
           self.dictionary = {'unk': 0}
5
           self._initialize_dict(corpus)
6
7
         def _add_dict_sentence(self, sentence):
           words = sentence.split(' ')
8
           for word in words:
9
             if word not in self.dictionary.keys():
10
               max_indice = max(self.dictionary.values())
11
12
               self.dictionary[word] = (max_indice + 1)
13
14
         def _initialize_dict(self, sentences):
15
           for sentence in sentences:
             self._add_dict_sentence(sentence)
16
17
         def _tokenize(self, sentence):
18
19
           words = sentence.split(' ')
20
           token_seq = [self.dictionary[word] for word in words]
21
           return np.array(token_seq)
22
23
       voc = Voc(corpus = corpus)
```

corpus là list toàn bộ tên các món ăn. Class Voc có tác dụng khởi tạo index từ điển cho toàn bộ corpus (bộ văn bản).

1 voc.dictionary

Tiếp theo chúng ta sẽ khởi tạo một random_generator có tác dụng lựa chọn ngẫu nhiên một tên món ăn trong corpus và tokenize chúng.

```
1
       import tensorflow as tf
2
3
       cat_indices = {
           'hanoi': 0,
4
5
           'hochiminh': 1
6
       }
7
8
       def generators():
         i = 0
9
         while True:
10
11
           i = np.random.choice(data.shape[0])
12
           sentence = data.iloc[i, 1]
           x_indice = voc._tokenize(sentence)
13
14
           label = data.iloc[i, 0]
           y_indice = cat_indices[label]
15
16
           yield x_indice, y_indice
17
           i += 1
18
19
       random_generator = tf.data.Dataset.from_generator(
20
           generators,
21
           output_types = (tf.float16, tf.float16),
22
           output_shapes = ((None,), ())
23
       )
24
25
       random_generator
1
      import numpy as np
2
3
      random_generator_batch = random_generator.shuffle(20).padded_batch(20,
      sequence_batch, label = next(iter(random_generator_batch))
4
5
      print(sequence_batch)
6
7
      print(label)
```

```
tf.Tensor(
1
2
       [[ 8.
               9.
                   3.
                        4.
                            Θ.
                   0.
        [10. 11.
                       Θ.
                            Θ.
                                0.1
4
        [17. 18. 19. 20. 15. 16.]
5
        [22. 23. 15. 16.
6
        [13. 14. 15. 16.
                            0.
                                0.]
7
        [13. 23.
                   Θ.
                        Θ.
                            Θ.
                                0.1
8
        [17. 18. 19. 20. 15. 16.]
9
               2.
                   3.
                        4.
        [ 1.
                            Θ.
                                0.]
10
        [10. 11.
                   0.
                        0.
                            0.
11
        [17. 18. 21. 15. 16.
                                0.]
12
        [ 8.
               9.
                   Θ.
                        Θ.
                            Θ.
                                0.]
        [22. 23. 15. 16.
13
                            Θ.
                                0.]
14
        [ 2.
               5.
                   6.
                        7.
                            3.
                                4.]
15
        [13. 14. 15. 16.
                                0.]
16
               2.
                   3.
                        4.
                            Θ.
                                0.]
        [ 1.
17
        [13. 14. 15. 16.
                            Θ.
                                0.]
18
               9.
                   Θ.
                       Θ.
                            Θ.
                                0.]
        [ 8.
19
        [13. 14. 15. 16.
                            Θ.
                                0.]
20
        [13. 23.
                   Θ.
                       Θ.
                            Θ.
                                0.]
21
                            4.
                                0.]], shape=(20, 6), dtype=float16)
        [10. 11. 12.
                       3.
22
       tf.Tensor([0. 0. 1. 1. 1. 1. 0. 0. 1. 0. 1. 0. 1. 0. 1. 0. 1. 0.
```

hàm shuffle(20) có tác dụng trộn lẫn ngẫu nhiên dữ liệu. Sau đó dữ liệu được chia thành những batch có kích thước là 10 và padding giá trị 0 sao cho bằng với độ dài của câu dài nhất bằng hàm padded_batch().

3.2.1. Sử dụng ImageGenerator

ImageGenerator cũng là một dạng data generator được xây dựng trên framework keras và dành riêng cho dữ liệu ảnh.

Đây là một high level function nên cú pháp đơn giản, rất dễ sử dụng nhưng khả năng tùy biến và can thiệp sâu vào dữ liệu kém.

Khi khởi tạo ImageGenerator chúng ta sẽ khai báo các thủ tục preprocessing image trước khi đưa vào huấn luyện. Mình sẽ không quá đi sâu vào các kĩ thuật preprocessing data này. Bạn đọc quan tâm có thể xem thêm tại ImageDataGenerator

(https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/image/ImageDataGenerator).

```
image_gen = tf.keras.preprocessing.image.ImageDataGenerator(
    featurewise_center=True,
    featurewise_std_normalization=True,
    rescale = 1./255,
    rotation_range = 20,
    horizontal_flip = True
}
```

Tiếp theo chúng ta sẽ truyền dữ liệu vào mô hình thông qua một hàm là flow_from_directory().

```
import glob2

root_folder = 'Dog-Cat-Classifier/Data/Train_Data/'
images, labels = next(image_gen.flow_from_directory(root_folder))

found 1399 images belonging to 2 classes.
```

Hàm flow_from_directory() sẽ có tác dụng đọc các ảnh từ root_folder và lấy ra những thông tin bao gồm ma trận ảnh sau biến đổi và nhãn tương ứng. Cấu trúc cây thư mục của root_folder có dạng như sau:

```
root-folder
sub-folder-class-1
sub-folder-class-2
...
sub-folder-class-C
```

Trong đó bên trong các sub-folder-class-i là list toàn bộ các ảnh thuộc về một class. Hàm flow_from_directory() sẽ tự động xác định các file dữ liệu nào là ảnh để load vào quá trình huấn luyện mô hình.

```
1 !ls {root_folder}

1 cat dog
```

Ở đây trong root folder chúng ta có 2 sub-folders tương ứng với 2 classes là dog, cat.

Tiếp theo ta sẽ khởi tạo một tf.Dataset từ generator thông qua hàm from_generator().

```
image_gen_dataset = tf.data.Dataset.from_generator(
image_gen.flow_from_directory,
args = ([root_folder]),
output_types=(tf.float32, tf.float32),
output_shapes=([32,256,256,3], [32, 1])
```

Trong hàm from_generator() chúng ta phải khai báo bắt buộc định dạng dữ liệu input và output thông qua tham số output_types và output shape thông qua tham số output_shapes.

Như vậy kết quả trả ra sẽ là những batch có kích thước 32 và ảnh có kích thước 256×256 và nhãn tương ứng của ảnh.

3.2.2. Customize ImageGenerator

Giả sử bạn có một bộ dữ liệu ảnh mà kích thước các ảnh là khác biệt nhau. Đồng thời bạn cũng muốn can thiệp sâu hơn vào bức ảnh trước khi đưa vào huấn luyện như giảm nhiễu bằng bộ lọc Gausianblur (https://phamdinhkhanh.github.io/2020/01/06/ImagePreprocessing.html#222-I%C3%A0m-m%E1%BB%9D-%E1%BA%A3nh-image-blurring), rotate ảnh, crop, zoom ảnh p.... Nếu sử dụng các hàm mặc định của image preprocessing trong ImageGenerator thì sẽ gặp hạn

chế đó là bị giới hạn bởi một số phép biến đổi mà hàm này hỗ trợ. Sử dụng high level framework tiện thì rất tiện nhưng khi muốn can thiệp sâu thì rất khó. Muốn can thiệp được sâu vào bên trong các biến đổi chúng ta phải customize lại một chút ImageGenerator.

Cách thức customize như thế nào. Mình sẽ giới thiệu với các bạn qua chương này.

Đầu tiên chúng ta sẽ download bộ dữ liệu Dog & Cat đã được thu nhỏ số lượng ảnh về.

1 !git clone https://github.com/ardamavi/Dog-Cat-Classifier.git

Cloning into 'Dog-Cat-Classifier'... remote: Enumerating objects: 1654, done. remote: Total 1654 (delta 0), reused 0 (delta 0), pack-reused 1654 Receiving objects: 100% (1654/1654), 34.83 MiB | 16.60 MiB/s, done. Resolving deltas: 100% (147/147), done. Checking out files: 100% (1672/1672), done.

Tiếp theo ta sẽ khởi tạo một DataGenerator cho bộ dữ liệu ảnh kế thừa class Sequence của keras. Mình sẽ giải thích các phương thức trong DataGenerator này bên dưới.

```
1
       import numpy as np
2
       from tensorflow.keras.utils import Sequence, to_categorical
3
       import cv2
4
5
       class DataGenerator(Sequence):
6
           'Generates data for Keras'
7
           def __init__(self,
8
                         all_filenames,
9
                         labels,
10
                         batch_size,
11
                         index2class,
12
                         input_dim,
13
                         n_channels,
14
                         n_classes=2,
15
                         shuffle=True):
16
               all_filenames: list toàn bộ các filename
17
18
               labels: nhãn của toàn bộ các file
               batch_size: kích thước của 1 batch
19
20
               index2class: index của các class
               input_dim: (width, height) đầu vào của ảnh
21
               n_channels: số lượng channels của ảnh
22
23
               n_classes: số lượng các class
24
               shuffle: có shuffle dữ liệu sau mỗi epoch hay không?
25
26
               self.all_filenames = all_filenames
27
               self.labels = labels
28
               self.batch_size = batch_size
29
               self.index2class = index2class
               self.input_dim = input_dim
30
31
               self.n_channels = n_channels
               self.n_classes = n_classes
32
               self.shuffle = shuffle
33
34
               self.on_epoch_end()
35
36
           def __len__(self):
37
38
               return:
39
                 Trả về số lượng batch/1 epoch
40
               return int(np.floor(len(self.all_filenames) / self.batch_size
41
42
           def __getitem__(self, index):
43
44
45
               params:
46
                 index: index của batch
47
               return:
48
                 X, y cho batch thứ index
49
50
               # Lấy ra indexes của batch thứ index
51
               indexes = self.indexes[index*self.batch_size:(index+1)*self.batch_size:
52
53
               # List all_filenames trong một batch
               all_filenames_temp = [self.all_filenames[k] for k in indexes]
54
55
                                                                       Top
56
               # Khởi tạo data
               X, y = self.__data_generation(all_filenames_temp)
```

```
58
59
               return X, y
60
           def on_epoch_end(self):
61
62
63
               Shuffle dữ liệu khi epochs end hoặc start.
64
               self.indexes = np.arange(len(self.all_filenames))
65
66
               if self.shuffle == True:
67
                   np.random.shuffle(self.indexes)
68
69
           def __data_generation(self, all_filenames_temp):
70
71
               params:
72
                 all_filenames_temp: list các filenames trong 1 batch
73
               return:
74
                 Trả về giá trị cho một batch.
               111
75
               X = np.empty((self.batch_size, *self.input_dim, self.n_channe)
76
77
               y = np.empty((self.batch_size), dtype=int)
78
79
               # Khởi tạo dữ liệu
80
               for i, fn in enumerate(all_filenames_temp):
                   # Đọc file từ folder name
81
82
                   img = cv2.imread(fn)
                   img = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)
83
                   img = cv2.resize(img, self.input_dim)
84
85
                   label = fn.split('/')[3]
                   label = self.index2class[label]
86
87
88
                   X[i,] = img
89
90
                   # Luu class
91
                   y[i] = label
               return X, y
92
```

Một DataGenerator sẽ cần xác định kích thước của một batch, số lượt steps huấn luyện.

Hàm len(): Như chúng ta đã biết, __len__() là một built in function trong python. Bất kì một object nào của python cũng sẽ có hàm __len__() . Đối với Datagenerator thì chúng ta sẽ qui định

$$len = \frac{\text{\# Obs}}{\text{batch size}}$$

Đây chính là số lượng step trong một epoch.

- Hàm __getitem__(): Trong quá trình huấn luyện chúng ta cần phải access vào từng batch trong bộ dữ liệu. Hàm __getitem__() sẽ khởi tạo batch theo thứ tự của batch được truyền vào hàm.
- Hàm on_epoch_end(): Đây là hàm được tự động run mỗi khi một epoch huấn luyện bắt đầu và kết thúc. Tại hàm này chúng ta sẽ xác định các hành động khi bắt đầu hoặc kết thúc một epoch như có shuffle dữ liệu hay không. Điều chỉnh lại tỷ lệ các class tước khi fit vào model,....

Hàm __data_generation(): Hàm này sẽ được gọi trong __getitem__().
 __data_generation() sẽ trực tiếp biến đổi dữ liệu và quyết định các kết quả dữ liệu trả về cho người dùng. Tại hàm này ta có thể thực hiện các phép preprocessing image.

```
1
       import cv2
2
       import glob2
3
4
       dict_labels = {
5
            'dog': 0,
6
           'cat': 1
7
       }
8
9
       root_folder = 'Dog-Cat-Classifier/Data/Train_Data/*/*'
10
       fns = glob2.glob(root_folder)
       print(len(fns))
11
12
13
       image_generator = DataGenerator(
14
           all_filenames = fns,
15
           labels = None,
           batch_size = 32,
16
17
           index2class = dict_labels,
           input_dim = (224, 224),
18
19
           n_{channels} = 3,
20
           n_{classes} = 2,
           shuffle = True
21
22
       )
1
      X, y = image_generator.__getitem__(1)
2
3
      print(X.shape)
4
      print(y.shape)
1
      (32, 224, 224, 3)
2
      (32,)
```

Như vậy ta có thể thấy, tại mỗi lượt huấn luyện model lấy ra một batch có kích thước là 32. Mặc dù ảnh của chúng ta có kích thước khác nhau nhưng đã được resize về chung một kích thước là width \times height = 224 \times 224 .

Chúng ta sẽ thử nghiệm huấn luyện model với generator. Đầu tiên là khởi tạo model.

```
1
       from tensorflow.keras.applications import MobileNet
2
       from tensorflow.keras.models import Sequential
3
       from tensorflow.keras.layers import Dense, Flatten
4
       from tensorflow.keras.optimizers import Adam
5
6
       base_extractor = MobileNet(input_shape = (224, 224, 3), include_top =
7
       flat = Flatten()
8
       den = Dense(1, activation='sigmoid')
9
       model = Sequential([base_extractor,
10
                          flat,
                          den])
11
                                                                      Top
12
       model.summary()
```

```
1
      Downloading data from https://storage.googleapis.com/tensorflow/keras
      17227776/17225924 [============= ] - Os Ous/step
2
      Model: "sequential_3"
3
4
5
      Layer (type)
                               Output Shape
                                                      Param #
      ______
6
7
      mobilenet_1.00_224 (Model)
                               (None, 7, 7, 1024)
                                                      3228864
8
9
      flatten_3 (Flatten)
                               (None, 50176)
                                                      0
10
      dense_9 (Dense)
11
                               (None, 1)
                                                      50177
12
      Total params: 3,279,041
13
14
      Trainable params: 3,257,153
      Non-trainable params: 21,888
15
16
```

Tiếp theo để huấn luyện model chúng ta chỉ cần thay generator vào vị trí của train data trong hàm fit().

```
model.compile(Adam(), loss='binary_crossentropy', metrics = ['accuracy
model.fit(image_generator,
epochs = 5)
```

```
1
  Epoch 1/5
2
  3
  Epoch 2/5
  4
5
  Epoch 3/5
6
  7
  Epoch 4/5
8
  9
  Epoch 5/5
10
  11
12
13
14
15
16
  <tensorflow.python.keras.callbacks.History at 0x7fd031252fd0>
```

Chỉ với khoảng 5 epochs nhưng kết quả đã đạt 98.4% độ chính xác. Đây là một kết quả khá ấn tương.

4. Tổng kết

Như vậy qua bài viết này tôi đã giới thiệu tới các bạn các phương pháp chính để khởi tạo một Dataset trong tensorflow, ưu nhược điểm và trường hợp sử dụng của từng phương pháp.

Khi nắm vững được kiến thức này, các bạn sẽ không còn phải lo lắng nếu phải đối mặt với những bộ dữ liệu rất lớn mà không biết cách truyền vào mô hình huấn luyện.

Chúc các bạn thành công với những mô hình sắp tới. Cuối cùng không thể thiếu là các tài liệu mà tôi đã tham khảo để viết bài này.

5. Tài liệu tham khảo

- 1. tensorflow data tensorflow (https://www.tensorflow.org/guide/data)
- 2. tensorflow ImageDataGenerator tensorflow (https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/image/ImageDataGenerator)
- 3. how to generate data on the fly standford.edu (https://stanford.edu/~shervine/blog/keras-how-to-generate-data-on-the-fly)
- 4. generator python wiki (https://wiki.python.org/moin/Generators)