

[Home \(../..../index.php\)](#)

>

[Deep Learning \(../..../topics.php?my_id=4\)](#)

>

[TensorFlow \(../..../posts.php?my_id=10\)](#)

Machine Learning Guru (../..

[../index.php\)](#)

Machine Learning and Computer Vision
tutorials using open source packages.

Sections

[Introduction](#)[List images and their labels](#)[Create a TFRecords file](#)[Read the TFRecords file](#)

How to write into and read from a TFRecords file in TensorFlow

In this post, you will learn how to save a large amount of data (images) into a single TFRecords format file and load it batch-wise to train your network in tensorflow.

Introduction

In the previous post ([../hdf5/hdf5.html](#)) we explained the benefits of saving a large dataset in a single HDF5 file. In this post we will learn how to convert our data into the Tensorflow standard format, called TFRecords. When we are training a deep network, we have two options to feed the data into out Tensorflow program: loading the data using pure python code at each step and feed it into a computaion graph or use an input pipeline which takes a list of filenames (any supported format), shuffle them (optional), create a file queue, read, and decode the data. However, TFRecords is the recommended file format for Tensorflow.

In this post, we load, resize and save all the images inside the train folder of the well-known Dogs vs. Cats (<https://www.kaggle.com/c/dogs-vs-cats>) data set into a single TFRecords file and then load and plot a couple of them as samples. To follow the rest of this post you need to download the train part of the Dogs vs. Cats data set.

List images and their labels

First, we need to list all images and label them. We give each cat image a label = 0 and each dog image a label = 1. The following code list all images, give them proper labels, and then shuffle the data. We also divide the data set into three train (%60), validation (%20), and test parts (%20).

List images and label them

```
1. from random import shuffle
2. import glob
3. shuffle_data = True # shuffle the addresses before saving
4. cat_dog_train_path = 'Cat vs Dog/train/*.jpg'
5.
6. # read addresses and labels from the 'train' folder
7. addrs = glob.glob(cat_dog_train_path)
8. labels = [0 if 'cat' in addr else 1 for addr in addrs] # 0 = Cat, 1 = Dog
9.
10. # to shuffle data
11. if shuffle_data:
12.     c = list(zip(addrs, labels))
13.     shuffle(c)
14.     addrs, labels = zip(*c)
15.
16. # Divide the hata into 60% train, 20% validation, and 20% test
17. train_addrs = addrs[0:int(0.6*len(addrs))]
18. train_labels = labels[0:int(0.6*len(labels))]
19.
20. val_addrs = addrs[int(0.6*len(addrs)):int(0.8*len(addrs))]
21. val_labels = labels[int(0.6*len(addrs)):int(0.8*len(addrs))]
22.
23. test_addrs = addrs[int(0.8*len(addrs)):]
24. test_labels = labels[int(0.8*len(labels)):]
```

Create a TFRecords file

First we need to load the image and convert it to the data type (float32 in this example) in which we want to save the data into a TFRecords file. Let's write a function which take an image address, load, resize, and return the image in proper data type:

A function to Load images

```
1. def load_image(addr):
2.     # read an image and resize to (224, 224)
3.     # cv2 load images as BGR, convert it to RGB
4.     img = cv2.imread(addr)
5.     img = cv2.resize(img, (224, 224), interpolation=cv2.INTER_CUBIC)
6.     img = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)
7.     img = img.astype(np.float32)
8.     return img
```

Before we can store the data into a TFRecords file, we should stuff it in a protocol buffer called **Example**. Then, we serialize the protocol buffer to a string and write it to a TFRecords file. Example protocol buffer contains **Features**. Feature is a protocol to describe the data and could have three types: bytes, float, and int64. In summary, to store your data you need to follow these steps:

- Open a TFRecords file using `tf.python_io.TFRecordWriter`
- Convert your data into the proper data type of the feature using `tf.train.Int64List`, `tf.train.BytesList`, Or `tf.train.FloatList`
- Create a feature using `tf.train.Feature` and pass the converted data to it
- Create an Example protocol buffer using `tf.train.Example` and pass the feature to it
- Serialize the Example to string using `example.SerializeToString()`
- Write the serialized example to TFRecords file using `writer.write`

We are going to use the following two functions to create features (Functions are from this [Tensorflow Tutorial \(https://github.com/tensorflow/tensorflow/blob/master/tensorflow/examples/how_tos/reading_data/convert_to_records.py\)](https://github.com/tensorflow/tensorflow/blob/master/tensorflow/examples/how_tos/reading_data/convert_to_records.py))

Convert data to features

```
1. def _int64_feature(value):
2.     return tf.train.Feature(int64_list=tf.train.Int64List(value=[value]))
3.
4.
5. def _bytes_feature(value):
6.     return tf.train.Feature(bytes_list=tf.train.BytesList(value=[value]))
```

Now let's store the training data to the TFRecords file.

Write data into a TFRecords file

```
1. train_filename = 'train.tfrecords' # address to save the TFRecords file
2.
3. # open the TFRecords file
4. writer = tf.python_io.TFRecordWriter(train_filename)
5.
6. for i in range(len(train_addrs)):
7.     # print how many images are saved every 1000 images
8.     if not i % 1000:
9.         print 'Train data: {}/{}'.format(i, len(train_addrs))
10.        sys.stdout.flush()
11.
12.    # Load the image
13.    img = load_image(train_addrs[i])
14.
15.    label = train_labels[i]
16.
17.    # Create a feature
18.    feature = {'train/label': _int64_feature(label),
19.              'train/image': _bytes_feature(tf.compat.as_bytes(img.tostring()))}
20.
21.    # Create an example protocol buffer
```

```

22.     example = tf.train.Example(features=tf.train.Features(feature=feature))
23.
24.     # Serialize to string and write on the file
25.     writer.write(example.SerializeToString())
26.
27. writer.close()
28. sys.stdout.flush()

```

and finally we close the file using: `writer.close()`. Similarly we write the validation and test data to two other TFRecords files.

Write validation and test data into a TFRecords file

```

1. # open the TFRecords file
2. val_filename = 'val.tfrecords' # address to save the TFRecords file
3. writer = tf.python_io.TFRecordWriter(val_filename)
4.
5. for i in range(len(val_addrs)):
6.     # print how many images are saved every 1000 images
7.     if not i % 1000:
8.         print 'Val data: {}/{}'.format(i, len(val_addrs))
9.         sys.stdout.flush()
10.
11.     # Load the image
12.     img = load_image(val_addrs[i])
13.
14.     label = val_labels[i]
15.
16.     # Create a feature
17.     feature = {'val/label': _int64_feature(label),
18.               'val/image': _bytes_feature(tf.compat.as_bytes(img.tostring()))}
19.
20.     # Create an example protocol buffer
21.     example = tf.train.Example(features=tf.train.Features(feature=feature))
22.
23.     # Serialize to string and write on the file
24.     writer.write(example.SerializeToString())
25.
26. writer.close()
27. sys.stdout.flush()
28.
29. # open the TFRecords file
30. test_filename = 'test.tfrecords' # address to save the TFRecords file
31. writer = tf.python_io.TFRecordWriter(test_filename)
32.
33. for i in range(len(test_addrs)):
34.     # print how many images are saved every 1000 images
35.     if not i % 1000:
36.         print 'Test data: {}/{}'.format(i, len(test_addrs))
37.         sys.stdout.flush()
38.
39.     # Load the image
40.     img = load_image(test_addrs[i])
41.
42.     label = test_labels[i]
43.
44.     # Create a feature
45.     feature = {'test/label': _int64_feature(label),
46.               'test/image': _bytes_feature(tf.compat.as_bytes(img.tostring()))}
47.
48.     # Create an example protocol buffer
49.     example = tf.train.Example(features=tf.train.Features(feature=feature))
50.
51.     # Serialize to string and write on the file
52.     writer.write(example.SerializeToString())
53.
54. writer.close()
55. sys.stdout.flush()

```

Read the TFRecords file

It's time to learn how to read data from the TFRecords file. To do so, we load the data from the train data in batches of an arbitrary size and plot images of the 5 batches. We also check the label of each image. To read from files in tensorflow, you need to do the following steps:

- **Create a list of filenames:** In our case we only have a single file `data_path = 'train.tfrecords'`. Therefore, our list is gonna be like this: `[data_path]`
- **Create a queue to hold filenames:** To do so, we use `tf.train.string_input_producer` `tf.train.string_input_producer` function which hold filenames in a FIFO queue. it gets the list of filenames. It also has some optional arguments including `num_epochs` which indicates the number of

epoch you want to load the data and `shuffle` which indicates whether to shuffle the filenames in the list or not. It is set to **True** by default.

- **Define a reader:** For files of TFRecords we need to define a TFRecordReader with `reader = tf.TFRecordReader()`. Now, the reader returns the next record using: `reader.read(filename_queue)`
- **Define a decoder:** A decoder is needed to decode the record read by the reader. In case of using TFRecords files the decoder should be `tf.parse_single_example`. It takes a serialized Example and a dictionary which maps feature keys to FixedLenFeature or VarLenFeature values and returns a dictionary which maps feature keys to Tensor values: `features = tf.parse_single_example(serialized_example, features=feature)`
- **Convert the data from string back to the numbers:** `tf.decode_raw(bytes, out_type)` takes a Tensor of type string and convert it to type `out_type`. However, for labels which have not been converted to string, we just need to cast them using `tf.cast(x, dtype)`
- **Reshape data into its original shape:** You should reshape the data (image) into its original shape before serialization using `image = tf.reshape(image, [224, 224, 3])`
- **Preprocessing:** if you want to do any preprocessing you should do it now.
- **Batching:** Another queue is needed to create batches from the examples. You can create the batch queue using `tf.train.shuffle_batch([image, label], batch_size=10, capacity=30, num_threads=1, min_after_dequeue=10)` where `capacity` is the maximum size of queue, `min_after_dequeue` is the minimum size of queue after dequeue, and `num_threads` is the number of threads enqueueing examples. Using more than one thread, it comes up with a faster reading. The first argument in a list of tensors which you want to create batches from.

Read a TFRecords file

```

1. import tensorflow as tf
2. import numpy as np
3. import matplotlib.pyplot as plt
4.
5. data_path = 'train.tfrecords' # address to save the hdf5 file
6.
7. with tf.Session() as sess:
8.     feature = {'train/image': tf.FixedLenFeature([], tf.string),
9.               'train/label': tf.FixedLenFeature([], tf.int64)}
10.
11.     # Create a list of filenames and pass it to a queue
12.     filename_queue = tf.train.string_input_producer([data_path], num_epochs=1)
13.
14.     # Define a reader and read the next record
15.     reader = tf.TFRecordReader()
16.     _, serialized_example = reader.read(filename_queue)
17.
18.     # Decode the record read by the reader
19.     features = tf.parse_single_example(serialized_example, features=feature)
20.
21.     # Convert the image data from string back to the numbers
22.     image = tf.decode_raw(features['train/image'], tf.float32)
23.
24.     # Cast label data into int32
25.     label = tf.cast(features['train/label'], tf.int32)
26.
27.     # Reshape image data into the original shape
28.     image = tf.reshape(image, [224, 224, 3])
29.
30.     # Any preprocessing here ...
31.
32.     # Creates batches by randomly shuffling tensors
33.     images, labels = tf.train.shuffle_batch([image, label], batch_size=10, capacity=30, num_threads=

```

- **Initialize all global and local variables**
- **Filing the example queue:** Some functions of `tf.train` such as `tf.train.shuffle_batch` add `tf.train.QueueRunner` objects to your graph. Each of these objects hold a list of enqueue op for a queue to run in a thread. Therefore, to fill a queue you need to call `tf.train.start_queue_runners` which starts threads for all the queue runners in the graph. However, to manage these threads you need a `tf.train.Coordinator` to terminate the threads at the proper time.

Everything is ready. Now you can read a batch and plot all batch images and labels. Do not forget to stop the threads (by stopping the coordinator) when you are done with your reading process.

Read a TFRecords file

```
34. # Initialize all global and local variables
35. init_op = tf.group(tf.global_variables_initializer(), tf.local_variables_initializer())
36. sess.run(init_op)
37.
38. # Create a coordinator and run all QueueRunner objects
39. coord = tf.train.Coordinator()
40. threads = tf.train.start_queue_runners(coord=coord)
41.
42. for batch_index in range(5):
43.     img, lbl = sess.run([images, labels])
44.
45.     img = img.astype(np.uint8)
46.
47.     for j in range(6):
48.         plt.subplot(2, 3, j+1)
49.         plt.imshow(img[j, ...])
50.         plt.title('cat' if lbl[j]==0 else 'dog')
51.
52.     plt.show()
53.
54. # Stop the threads
55. coord.request_stop()
56.
57. # Wait for threads to stop
58. coord.join(threads)
59. sess.close()
```

You can download the codes of this post in our Github (TFRecords.html) page.

[Go Top](#)

Related Posts:

 (https://twitter.com/M_L_Guru)  (<https://github.com/Machinelearningguru>)
 (<https://www.linkedin.com/groups/12030461>)

© Machine Learning Guru. All rights reserved

Design: HTML5 UP (<http://html5up.net>)