```python
# Copyright 2016 The TensorFlow Authors. All Rights Reserved.
#
# Licensed under the Apache License, Version 2.0 (the "License");
# you may not use this file except in compliance with the License.
# You may obtain a copy of the License at
#
#     http://www.apache.org/licenses/LICENSE-2.0
#
# Unless required by applicable law or agreed to in writing, software
# distributed under the License is distributed on an "AS IS" BASIS,
# WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
# See the License for the specific language governing permissions and
# limitations under the License.
# ==============================================================================

"""Functions for downloading and reading MNIST data."""

from __future__ import absolute_import
from __future__ import division
from __future__ import print_function

import gzip

import numpy
from six.moves import xrange  # pylint: disable=redefined-builtin

from tensorflow.contrib.learn.python.learn.datasets import base
from tensorflow.python.framework import dtypes
from tensorflow.python.framework import random_seed

# CVDF mirror of http://yann.lecun.com/exdb/mnist/
m = 1500
SOURCE_URL = 'http://storage.googleapis.com/cvdf-datasets/mnist/'


def _read32(bytestream):
    dt = numpy.dtype(numpy.uint32).newbyteorder('>')
    return numpy.frombuffer(bytestream.read(4), dtype=dt)[0]


def extract_images(f):
    """Extract the images into a 4D uint8 numpy array [index, y, x, depth].

  Args:
    f: A file object that can be passed into a gzip reader.

  Returns:
    data: A 4D uint8 numpy array [index, y, x, depth].

  Raises:
    ValueError: If the bytestream does not start with 2051.
```

```python
  """
  #print('Extracting', f.name)
  with gzip.GzipFile(fileobj=f) as bytestream:
    magic = _read32(bytestream)
    if magic != 2051:
      raise ValueError('Invalid magic number %d in MNIST image file: %s' %
                       (magic, f.name))
    num_images = _read32(bytestream)
    rows = _read32(bytestream)
    cols = _read32(bytestream)
    buf = bytestream.read(rows * cols * num_images)
    data = numpy.frombuffer(buf, dtype=numpy.uint8)
    data = data.reshape(num_images, rows, cols, 1)
    return data


def dense_to_one_hot(labels_dense, num_classes):
  """Convert class labels from scalars to one-hot vectors."""
  num_labels = labels_dense.shape[0]
  index_offset = numpy.arange(num_labels) * num_classes
  labels_one_hot = numpy.zeros((num_labels, num_classes))
  labels_one_hot.flat[index_offset + labels_dense.ravel()] = 1
  return labels_one_hot


def extract_labels(f, one_hot=False, num_classes=10):
  """Extract the labels into a 1D uint8 numpy array [index].

  Args:
    f: A file object that can be passed into a gzip reader.
    one_hot: Does one hot encoding for the result.
    num_classes: Number of classes for the one hot encoding.

  Returns:
    labels: a 1D uint8 numpy array.

  Raises:
    ValueError: If the bystream doesn't start with 2049.
  """
  #print('Extracting', f.name)
  with gzip.GzipFile(fileobj=f) as bytestream:
    magic = _read32(bytestream)
    if magic != 2049:
      raise ValueError('Invalid magic number %d in MNIST label file: %s' %
                       (magic, f.name))
    num_items = _read32(bytestream)
    buf = bytestream.read(num_items)
    labels = numpy.frombuffer(buf, dtype=numpy.uint8)
    if one_hot:
      return dense_to_one_hot(labels, num_classes)
    return labels
```

```python
106  class DataSet ( object ):
107
108    def __init__ ( self ,
109                 images ,
110                 labels ,
111                 fake_data = False ,
112                 one_hot = False ,
113                 dtype = dtypes . float32 ,
114                 reshape = True ,
115                 seed = None ):
116      """ Construct a DataSet .
117    one_hot arg is used only if fake_data is true . 'dtype ' can be either
118    'uint8 ' to leave the input as '[0, 255] ', or 'float32 ' to rescale into
119    '[0, 1] '. Seed arg provides for convenient deterministic testing .
120    """
121      seed1 , seed2 = random_seed . get_seed ( seed )
122      # If op level seed is not set , use whatever graph level seed is returned
123      numpy . random . seed ( seed1 if seed is None else seed2 )
124      dtype = dtypes . as_dtype ( dtype ). base_dtype
125      if dtype not in ( dtypes . uint8 , dtypes . float32 ):
126        raise TypeError ( 'Invalid image dtype %r, expected uint8 or float32 ' %
127                       dtype )
128      if fake_data :
129        self . _num_examples = 10000
130        self . one_hot = one_hot
131      else :
132        assert images . shape [0] == labels . shape [0] , (
133            'images . shape : %s labels . shape : %s' % ( images . shape , labels . shape ))
134        self . _num_examples = images . shape [0]
135
136        # Convert shape from [num examples , rows , columns , depth]
137        # to [num examples , rows * columns ] ( assuming depth == 1)
138        if reshape :
139          assert images . shape [3] == 1
140          images = images . reshape ( images . shape [0] ,
141                                    images . shape [1] * images . shape [2])
142        if dtype == dtypes . float32 :
143          # Convert from [0, 255] -> [0.0 , 1.0].
144          images = images . astype ( numpy . float32 )
145          images = numpy . multiply ( images , 1.0 / 255.0)
146      self . _images = images
147      self . _labels = labels
148      self . _epochs_completed = 0
149      self . _index_in_epoch = 0
150
151    @property
152    def images ( self ):
153      return self . _images
154
155    @property
156    def labels ( self ):
157      return self . _labels
158
```

```python
     @property
     def num_examples(self):
       return self._num_examples

     @property
     def epochs_completed(self):
       return self._epochs_completed

     def next_batch(self, batch_size, fake_data=False, shuffle=True):
       """Return the next 'batch_size' examples from this data set."""
       if fake_data:
         fake_image = [1] * 784
         if self.one_hot:
           fake_label = [1] + [0] * 9
         else:
           fake_label = 0
         return [fake_image for _ in xrange(batch_size)], [
             fake_label for _ in xrange(batch_size)
         ]
       start = self._index_in_epoch
       # Shuffle for the first epoch
       if self._epochs_completed == 0 and start == 0 and shuffle:
         perm0 = numpy.arange(self._num_examples)
         numpy.random.shuffle(perm0)
         self._images = self.images[perm0]
         self._labels = self.labels[perm0]
       # Go to the next epoch
       if start + batch_size > self._num_examples:
         # Finished epoch
         self._epochs_completed += 1
         # Get the rest examples in this epoch
         rest_num_examples = self._num_examples - start
         images_rest_part = self._images[start:self._num_examples]
         labels_rest_part = self._labels[start:self._num_examples]
         # Shuffle the data
         if shuffle:
           perm = numpy.arange(self._num_examples)
           numpy.random.shuffle(perm)
           self._images = self.images[perm]
           self._labels = self.labels[perm]
         # Start next epoch
         start = 0
         self._index_in_epoch = batch_size - rest_num_examples
         end = self._index_in_epoch
         images_new_part = self._images[start:end]
         labels_new_part = self._labels[start:end]
         return numpy.concatenate((images_rest_part, images_new_part), axis=0) , numpy.co
       else:
         self._index_in_epoch += batch_size
         end = self._index_in_epoch
         return self._images[start:end], self._labels[start:end]
```

```
212 │ def read_data_sets ( train_dir ,
213 │                       fake_data=False ,
214 │                       one_hot=False ,
215 │                       dtype=dtypes.float32 ,
216 │                       reshape=True ,
217 │                       validation_size=int(m/10) ,
218 │                       seed=None ):
219 │   if fake_data:
220 │
221 │     def fake ():
222 │       return DataSet (
223 │             [], [], fake_data=True , one_hot=one_hot , dtype=dtype , seed=seed )
224 │
225 │     train = fake ()
226 │     validation = fake ()
227 │     test = fake ()
228 │     return base.Datasets ( train=train , validation=validation , test=test )
229 │
230 │   TRAIN_IMAGES = 'train-images-idx3-ubyte.gz'
231 │   TRAIN_LABELS = 'train-labels-idx1-ubyte.gz'
232 │   TEST_IMAGES = 't10k-images-idx3-ubyte.gz'
233 │   TEST_LABELS = 't10k-labels-idx1-ubyte.gz'
234 │
235 │   local_file = base.maybe_download ( TRAIN_IMAGES , train_dir ,
236 │                                      SOURCE_URL + TRAIN_IMAGES )
237 │   with open ( local_file , 'rb' ) as f:
238 │     train_images = extract_images ( f )
239 │     train_images = train_images [ range(m) ,:,:,:]
240 │     #print ( train_images.shape )
241 │
242 │   local_file = base.maybe_download ( TRAIN_LABELS , train_dir ,
243 │                                      SOURCE_URL + TRAIN_LABELS )
244 │   with open ( local_file , 'rb' ) as f:
245 │     train_labels = extract_labels ( f , one_hot=one_hot )
246 │     train_labels = train_labels [ range(m) ]
247 │     #print ( train_labels.shape )
248 │
249 │   local_file = base.maybe_download ( TEST_IMAGES , train_dir ,
250 │                                      SOURCE_URL + TEST_IMAGES )
251 │   with open ( local_file , 'rb' ) as f:
252 │     test_images = extract_images ( f )
253 │
254 │   local_file = base.maybe_download ( TEST_LABELS , train_dir ,
255 │                                      SOURCE_URL + TEST_LABELS )
256 │   with open ( local_file , 'rb' ) as f:
257 │     test_labels = extract_labels ( f , one_hot=one_hot )
258 │
259 │   if not 0 <= validation_size <= len ( train_images ):
260 │     raise ValueError (
261 │         'Validation␣size␣should␣be␣between␣0␣and␣{}.␣Received:␣{}.'
262 │         .format ( len ( train_images ), validation_size ))
263 │
264 │   validation_images = train_images [: validation_size ]
```

```
265    validation_labels = train_labels[:validation_size]
266    train_images = train_images[validation_size:]
267    train_labels = train_labels[validation_size:]
268
269
270    options = dict(dtype=dtype, reshape=reshape, seed=seed)
271
272    train = DataSet(train_images, train_labels, **options)
273    validation = DataSet(validation_images, validation_labels, **options)
274    test = DataSet(test_images, test_labels, **options)
275
276    return base.Datasets(train=train, validation=validation, test=test)
277
278
279 def load_mnist(train_dir='MNIST-data'):
280    return read_data_sets(train_dir)
281
282 if __name__ == "__main__":
283     #print("test")
284     load_mnist()
```