

MNIST 참고자료

19.04.01

NLP LAB,
Department of Computer Engineering,
Kyung Hee University.

Index

- What is MNIST ?
- How to download MNIST data
- Data description

What is MNIST ? (1/1)

- MNIST

- Modified National Institute of Standards and Technology database
- The digits have been size-normalized and centered in a fixed-size image (28 x 28)
- taken from NIST's training dataset
- An extended dataset similar to MNIST called EMNIST has been published in 2017



Download MNIST data (1/2)

- Website

<http://yann.lecun.com/exdb/mnist/>

THE MNIST DATABASE of handwritten digits

Yann LeCun, Courant Institute, NYU

Corinna Cortes, Google Labs, New York

Christopher J.C. Burges, Microsoft Research, Redmond

The MNIST database of handwritten digits, available from this page, has a training set of 60,000 examples, and a test set of 10,000 examples. It is a subset of a larger set available from NIST. The digits have been size-normalized and centered in a fixed-size image.

It is a good database for people who want to try learning techniques and pattern recognition methods on real-world data while spending minimal efforts on preprocessing and formatting.

Four files are available on this site:

train-images-idx3-ubyte.gz: training set images (9912422 bytes)

train-labels-idx1-ubyte.gz: training set labels (28881 bytes)

t10k-images-idx3-ubyte.gz: test set images (1648877 bytes)

t10k-labels-idx1-ubyte.gz: test set labels (4542 bytes)

⇒ Click

Download MNIST data (2/2)

- Code (using by Python 2.x or 3.x)

```
> pip install tensorflow
```

```
>>> from tensorflow.examples.tutorials.mnist import input_data
```

```
>>> input_data.read_data_sets("YOUR_FOLDER/", one_hot=True)
```

```
1  # -*- coding: utf-8 -*-
2
3  # MNIST 데이터를 다운로드 한다.
4  from tensorflow.examples.tutorials.mnist import input_data
5  mnist = input_data.read_data_sets("MNIST_data/", one_hot=True)
6
```

Data description (1/)

- These files are not in any standard image format
- The data is stored in a very simple file format designed for storing vectors and multidimensional matrices.
- All the integers in the files are stored in the MSB first (high endian) format used by most non-Intel processors.
- Users of Intel processors and other low-endian machines must flip the bytes of the header.

Data description (1/)

- The training set contains 60000 examples, and the test set 10000 examples.
- The first 5000 examples of the test set are taken from the original NIST training set. The last 5000 are taken from the original NIST test set. The first 5000 are cleaner and easier than the last 5000.
- Result 60000 examples
= 55000 train data + 5000 validation data

Data description (1/)

File	Purpose
<code>train-images-idx3-ubyte.gz</code>	training set images - 55000 training images, 5000 validation images
<code>train-labels-idx1-ubyte.gz</code>	training set labels matching the images
<code>t10k-images-idx3-ubyte.gz</code>	test set images - 10000 images
<code>t10k-labels-idx1-ubyte.gz</code>	test set labels matching the images

DataSet Object

The underlying code will download, unpack, and reshape images and labels for the following datasets:

Dataset	Purpose
<code>data_sets.train</code>	55000 images and labels, for primary training.
<code>data_sets.validation</code>	5000 images and labels, for iterative validation of training accuracy.
<code>data_sets.test</code>	10000 images and labels, for final testing of trained accuracy.

Data description (1/)

TRAINING SET LABEL FILE (train-labels-idx1-ubyte):

[offset]	[type]	[value]	[description]
0000	32 bit integer	0x00000801(2049)	magic number (MSB first)
0004	32 bit integer	60000	number of items
0008	unsigned byte	??	label
0009	unsigned byte	??	label
.....			
xxxx	unsigned byte	??	label

The labels values are 0 to 9.

TRAINING SET IMAGE FILE (train-images-idx3-ubyte):

[offset]	[type]	[value]	[description]
0000	32 bit integer	0x00000803(2051)	magic number
0004	32 bit integer	60000	number of images
0008	32 bit integer	28	number of rows
0012	32 bit integer	28	number of columns
0016	unsigned byte	??	pixel
0017	unsigned byte	??	pixel
.....			
xxxx	unsigned byte	??	pixel

Pixels are organized row-wise. Pixel values are 0 to 255. 0 means background (white), 255 means foreground (black).

Data description (1/)

- Image example

train-images.idx3-ubyte x																								
									00	01	02	03	04	05	06	07	08	09	0a	0b	0c	0d	0e	0f
1	0000	0803	0000	ea60	0000	001c	0000	001c	00 00 08 03	00 00	ea 60	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
2	0000	0000	0000	0000	0000	0000	0000	0000	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
3	0000	0000	0000	0000	0000	0000	0000	0000	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
4	0000	0000	0000	0000	0000	0000	0000	0000	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
5	0000	0000	0000	0000	0000	0000	0000	0000	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
6	0000	0000	0000	0000	0000	0000	0000	0000	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
7	0000	0000	0000	0000	0000	0000	0000	0000	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
8	0000	0000	0000	0000	0000	0000	0000	0000	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
9	0000	0000	0000	0000	0000	0000	0000	0000	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
10	0000	0000	0000	0000	0000	0000	0000	0000	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
11	0000	0000	0000	0000	0312	1212	7e88	af1a	00 0e	01 9a	fd 5a	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
12	a6ff	f77f	0000	0000	0000	0000	0000	0000	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00
13	1e24	5e9a	aa fd	fd fd	fd fd	e1 ac	fd f2	c340	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00	00 00

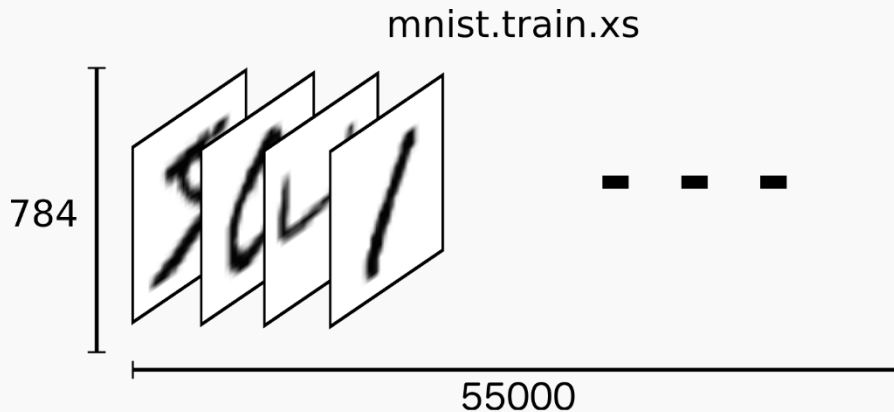
Data description (1/)

- Label example

	train-images.idx3-ubyte	train-labels.idx1-ubyte
1	0000 0801 0000 ea60 0500 0401 0902 0103	0000000000 00 00 08 01 00 00 ea 60 05 00 04 01 09 02 01 03
2	0104 0305 0306 0107 0208 0609 0400 0901	0000000010 01 04 03 05 03 06 01 07 02 08 06 09 04 00 09 01
3	0102 0403 0207 0308 0609 0005 0600 0706	0000000020 01 02 04 03 02 07 03 08 06 09 00 05 06 00 07 06
4	0108 0709 0309 0805 0903 0300 0704 0908	0000000030 01 08 07 09 03 09 08 05 09 03 03 00 07 04 09 08
5	0009 0401 0404 0600 0405 0601 0000 0107	0000000040 00 09 04 01 04 04 06 00 04 05 06 01 00 00 01 07
6	0106 0300 0201 0107 0900 0206 0708 0309	0000000050 01 06 03 00 02 01 01 07 09 00 02 06 07 08 03 09
7	0004 0607 0406 0800 0708 0301 0507 0107	0000000060 00 04 06 07 04 06 08 00 07 08 03 01 05 07 01 07
8	0101 0603 0002 0903 0101 0004 0902 0000	0000000070 01 01 06 03 00 02 09 03 01 01 00 04 09 02 00 00
9	0200 0207 0108 0604 0106 0304 0509 0103	0000000080 02 00 02 07 01 08 06 04 01 06 03 04 05 09 01 03
10	0308 0504 0707 0402 0805 0806 0703 0406	0000000090 03 08 05 04 07 07 04 02 08 05 08 06 07 03 04 06
11	0109 0906 0003 0702 0802 0904 0406 0409	00000000a0 01 09 09 06 00 03 07 02 08 02 09 04 04 06 04 09
12	0700 0902 0905 0105 0901 0203 0203 0509	00000000b0 07 00 09 02 09 05 01 05 09 01 02 03 02 03 05 09
13	0107 0602 0802 0205 0007 0409 0708 0302	00000000c0 01 07 06 02 08 02 02 05 00 07 04 09 07 08 03 02
14	0101 0803 0601 0003 0100 0001 0702 0703	00000000d0 01 01 08 03 06 01 00 03 01 00 00 01 07 02 07 03

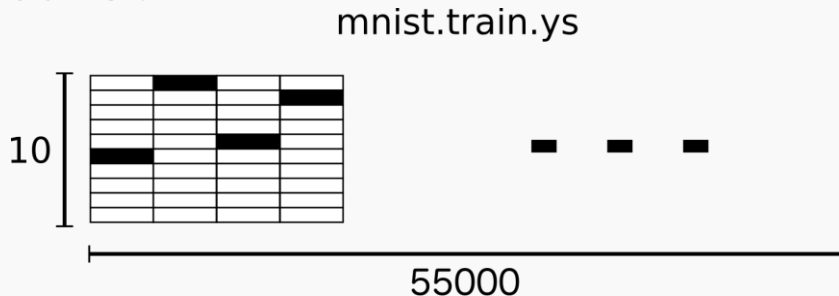
Data description (1/)

- We can flatten this array into a vector of $28 \times 28 = 784$ numbers.
- The result is that `mnist.train.images` is a tensor (an n-dimensional array) with shape (55000L, 784L).



Data description (1/)

- Each image in MNIST has a corresponding label, a number between 0 and 9 representing the digit drawn in the image.
- A one-hot vector is a vector which is 0 in most dimensions, and 1 in a single dimension.
- Consequently, `mnist.train.labels` is a (55000L, 10L) array of floats.



- Sample Code

[illegible]

Natural Language Processing Laboratory

Q & A