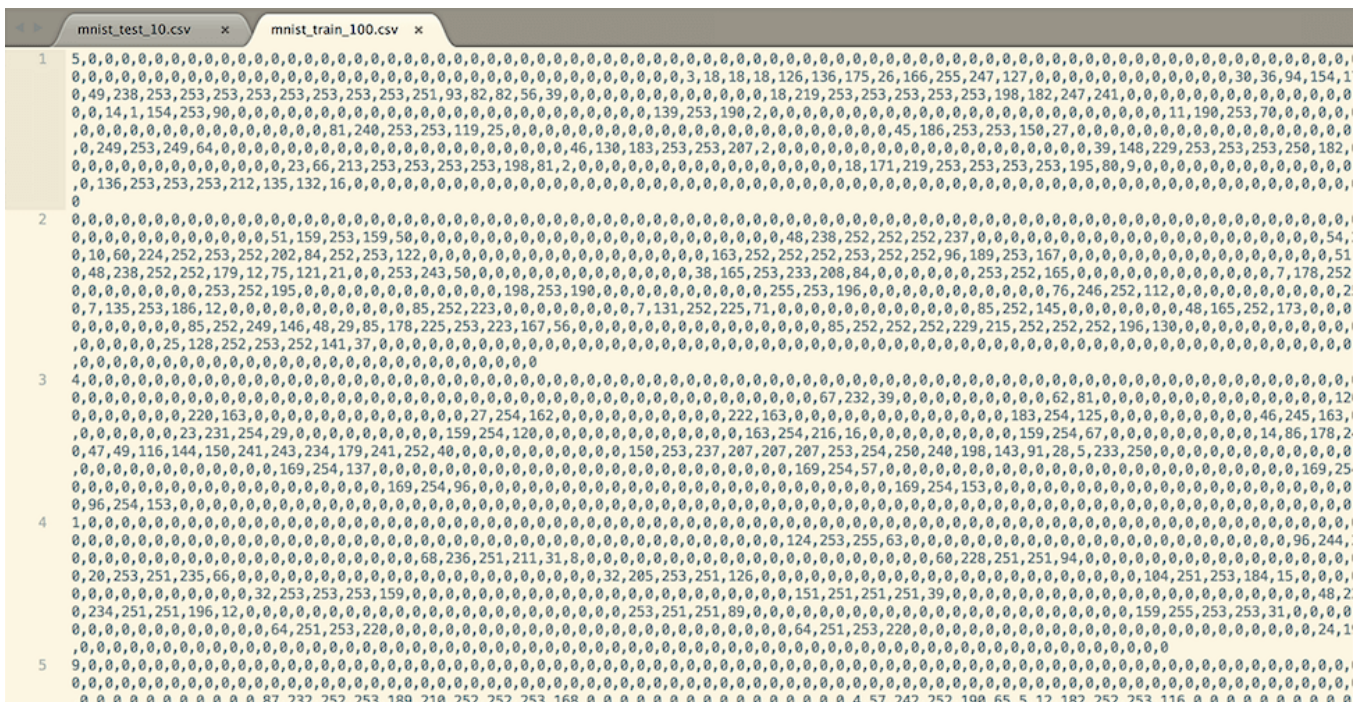# A Quick Look at the Data Files

Quick Look at the files of MNIST handwritten digits database to get the idea.

Let's take a peek at these files. The following shows a section of the MNIST test set loaded into a text editor.



Whoah! That looks like something went wrong! Like one of those movies from the 80s where a computer gets hacked.

Actually, all is well. The text editor is showing long lines of text. Those lines consist of numbers, separated by commas. That is easy enough to see. The lines are quite long, so they wrap around a few times. Helpfully this text editor shows the real line numbers in the margin, and we can see four whole lines of data, and part of the fifth one.

The content of these records, or lines of text, is easy to understand:

- The first value is the *label*, that is, the actual digit that the handwriting is supposed to represent, such as a 7 or a 9. This is the answer the neural network is trying to learn to get right.

- The subsequent values, all comma separated, are the *pixel* values of the handwritten digit. The size of the pixel array is 28 by 28, so there are 784 values after the label. Count them if you really want!

So that first record represents the number 5 as shown by the first value, and the rest of the text on that line is the pixel values for someone's handwritten number 5. The second record represents a handwritten 0, the third represents 4, the fourth record is 1, and the fifth represents 9. You can pick any line from the MNIST data files, and the first number will tell you the label for the following image data.

But it is hard to see how that long list of 784 values makes up a picture of someone's handwritten number 5. We should plot those numbers as an image to confirm that they really are the color values of the handwritten number.

Before we dive in and do that, we should download a smaller subset of the MNIST dataset. The MNIST data files are pretty big and working with a smaller subset is helpful because it means we can experiment, trial and develop our code without being slowed down by a large data set slowing our computers down. Once we've settled on an algorithm and code we're happy with; we can use the full data set.

The following are the links to a smaller subset of the MNIST dataset, also in CSV format:

- Ten records from the MNIST test data set - https://raw.githubusercontent.com/makeyourownneuralnetwork/makeyourownneuralnetwork/master/mnist_dataset/mnist_test_10.csv

- 100 records from the MNIST training dataset - https://raw.githubusercontent.com/makeyourownneuralnetwork/makeyourownneuralnetwork/master/mnist_dataset/mnist_train_100.csv