

Machine Learning Final Project

Prof. Yuh-Jye Lee

Kaggle - Digit Recognizer

M10412005 龔鈺傑 (Yu-Chieh Kung)

0034980 王澤浩 (Ze-Hao Wang)

2016 / 01 / 18

National Taiwan University of Science and Technology

Kaggle - Digit Recognizer

Competition link: <https://www.kaggle.com/c/digit-recognizer/>

Dataset: MNIST (Mixed National Institute of Standards and Technology database)

The goal in this competition is to take an image of a handwritten single digit, and determine what that digit is.

Our best submission

Leaderboard: 65th / 965 total

Error rate: 0.757%

63	↓5	andrei	0.99257	16	Wed, 09 Dec 2015 09:31:23 (-7.5h)
64	↓5	indy256	0.99257	18	Sun, 03 Jan 2016 07:46:36
65	new	NTUST 1041 ML - Group 8	0.99243	10	Tue, 12 Jan 2016 10:00:52
66	↓6	pklfz	0.99229	5	Wed, 09 Dec 2015 14:28:48
67	↓6	wanglingdeemo	0.99229	20	Mon, 14 Dec 2015 02:45:41

Our project source code and submission data:

<https://github.com/grass0916/NTUST-1041-ML-final-project>

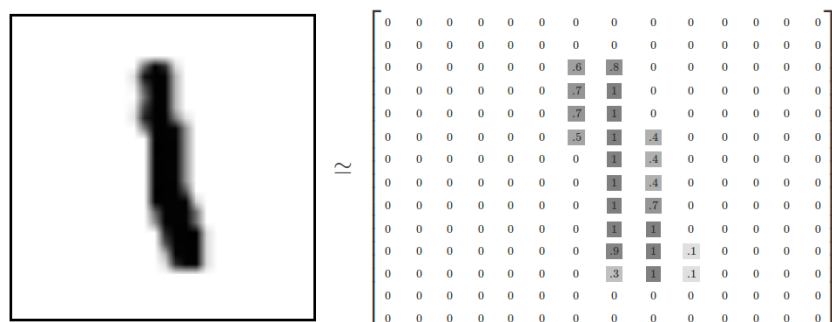
About MNIST Dataset



Pic-1. MNIST Examples. (source: myselph.de)

Each image is $28 * 28$ pixels. We can interpret this image as a big array of numbers like Pic-1. We can flatten this array into $28 * 28 = 784$ numbers, this MNIST images are a bunch of points in a 784-dimensional vector space.

The data files `train.csv` and `test.csv` contain gray-scale images of hand-drawn digits, from zero through nine. Each pixel has a single pixel-value associated with it, indicating the lightness or darkness of that pixel, with higher numbers meaning darker. This pixel-value is an integer between 0 and 255, inclusive.



Pic-2. MNIST matrix. (source: tensorflow.org)

The training data set, `train.csv`, has 785 columns. The first column, called **"label"**, is the digit that was drawn by the user. The rest of the

columns contain the pixel-values of the associated image.

The test data set, `test.csv`, is the same as the training set, except that it does not contain the **"label"** column.

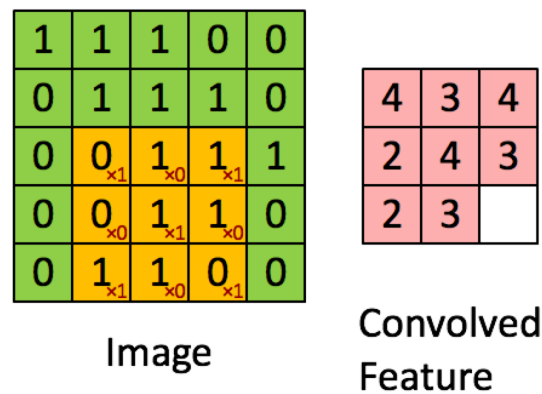
Related work

Convolutional Neural Network has wide applications in image and video recognition, recommender systems and natural language processing. CNN consist of multiple layers of small neuron collections which process portions of the input image, called receptive fields. [1] [8] Our datasets, the MNIST, there are plenty of references. There are quite great results on the MNIST by CNN.

Type	Classifier	Distortion	Preprocessing	Error rate (%)
Linear classifier	Pairwise linear classifier	None	Deskewing	7.6 ^[9]
K-Nearest Neighbors	K-NN with non-linear deformation (P2DHMDM)	None	Shiftable edges	0.52 ^[14]
Boosted Stumps	Product of stumps on Haar features	None	Haar features	0.87 ^[15]
Non-Linear Classifier	40 PCA + quadratic classifier	None	None	3.3 ^[9]
Support vector machine	Virtual SVM, deg-9 poly, 2-pixel jittered	None	Deskewing	0.56 ^[16]
Neural network	2-layer 784-800-10	None	None	1.6 ^[17]
Neural network	2-layer 784-800-10	elastic distortions	None	0.7 ^[17]
Deep neural network	6-layer 784-2500-2000-1500-1000-500-10	elastic distortions	None	0.35 ^[18]
Convolutional neural network	Committee of 35 conv. net, 1-20-P-40-P-150-10	elastic distortions	Width normalizations	0.23 ^[8]

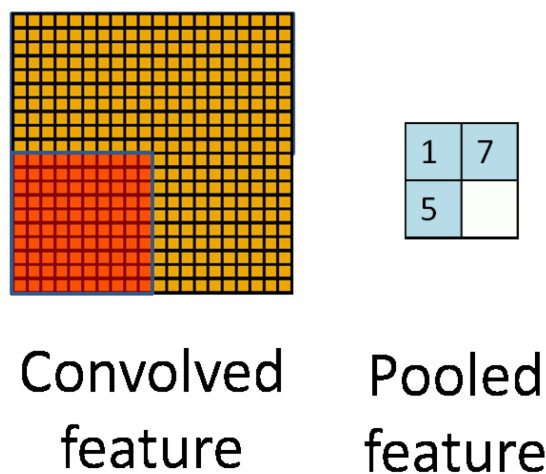
Pic-3. Classifiers performance in MNIST. (source: wikipedia.org)

- **Convolutional:** Convolutional layers consist of a rectangular grid of neurons. It requires that the previous layer also be a rectangular grid of neurons. Each neuron takes inputs from a rectangular section of the previous layer; the weights for this rectangular section are the same for each neuron in the convolutional layer. In other words, the inputs of hidden units in layer **m** are from a subset of units in layer **m-1**, units that have spatially contiguous receptive fields. [1][2][5]



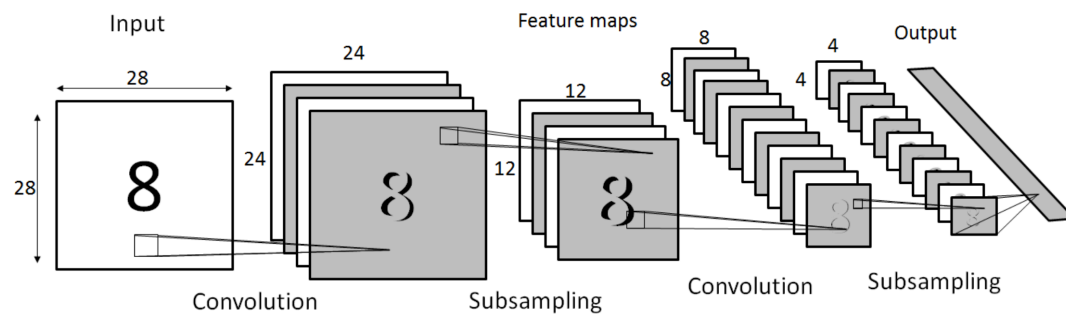
Pic-4. Convolutional schematic. (source: deeplearning.stanford.edu)

- **Pooling:** After each convolutional layer, there may be a pooling layer. The pooling layer takes small rectangular blocks from the convolutional layer and subsamples it to produce a single output from that block. Our pooling layers will always be max-pooling layers. [1][2]



Pic-5. Pooling schematic. (source: deeplearning.stanford.edu)

- **Fully-Connected:** Finally, after several convolutional and max pooling layers, the high-level reasoning in the neural network is done via fully connected layers. A fully connected layer takes all neurons in the previous layer (be it fully connected, pooling, or convolutional) and connects it to every single neuron it has. [2]



Pic-6. LeNet workflow. (source: developer.nvidia.com)

Experiment

Environment: **Mac OS X** El Capitan 10.11.1

Language: **R** version 3.2.3

Toolkit:

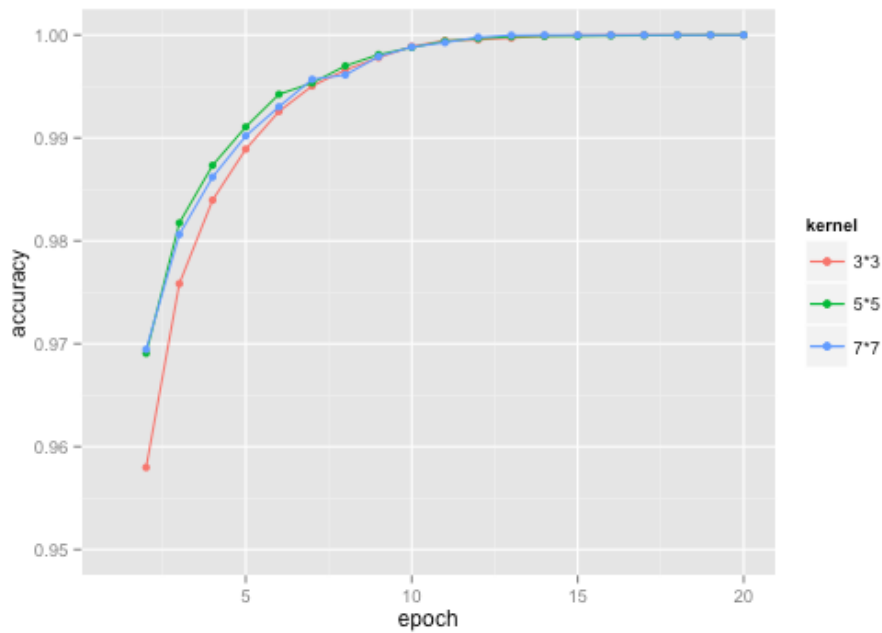
- MXNet 0.5.0 (<http://mxnet.rtd.org>)

MXNet is a deep learning framework designed for both efficiency and flexibility. It allows you to mix the flavours of deep learning programs together to maximize the efficiency and your productivity.

Filter Shape

At first, we try to tune the convolution kernel size. Size of 3x3 is the most often used. But in this case, size of 5x5 is better than others in the MNIST data. [1]

- The same conditions: activation function (tanh)



Pic-7. Compare different kernel size.

When kernel = 3x3 matrix, **Kaggle submission scored 0.99029.**

When kernel = 5x5 matrix, **Kaggle submission scored 0.99071 (Better).**

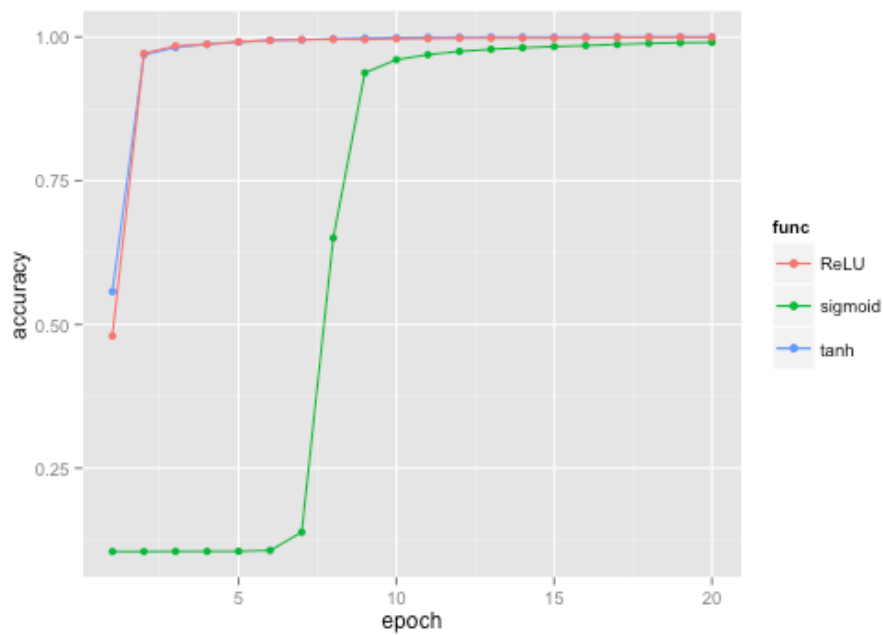
When kernel = 7x7 matrix, **Kaggle submission scored 0.98886.**

Activation functions

Every activation function (or non-linearity) takes a single number and performs a certain fixed mathematical operation on it. [5] In this step, we try to change the activation function to get different outputs.

Specifically, ReLU computes the function $f(x) = \max(0, x)$. In other words, the activation is simply thresholded at zero. [5]

- The same conditions: convolution kernel size (5x5)



Pic-8. Compare different activation function.

When using tanh function, **Kaggle submission scored 0.99071**.

When using sigmoid function, **Kaggle submission scored 0.98500**.

When using ReLU function, **Kaggle submission scored 0.99171 (Better)**

But using 20 epochs still not reaches 1.0 accuracy in ReLU function, we increase the epoch to 30 rounds. Finally, we get **Kaggle submission scored 0.99243 (Best)**.

Max Pooling Shape

Typical values are 2x2 or no max-pooling. [1] But we can't use 1x1 to do our pooling kernel, it costs too large time on iterating.

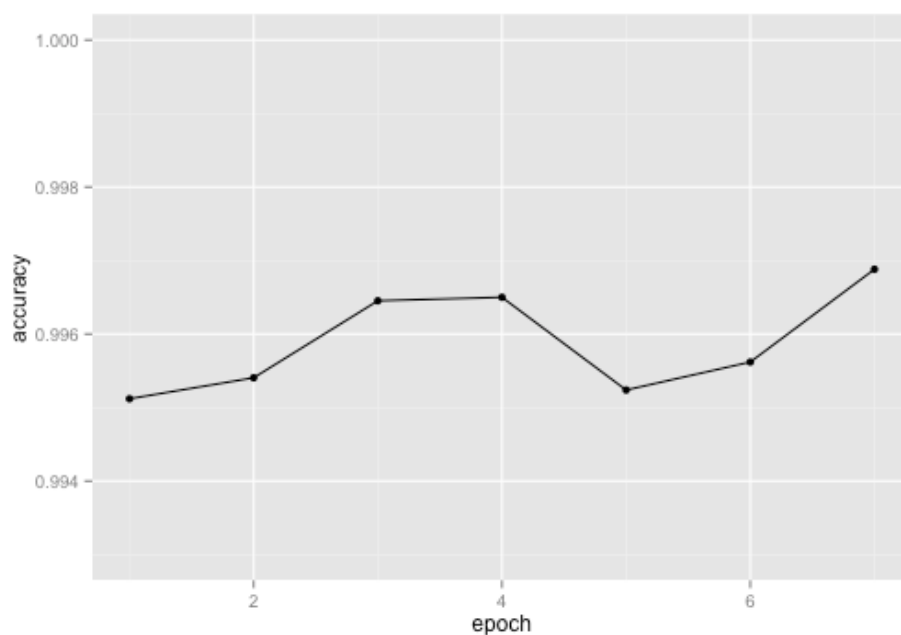
Conclusion

In this final project, we try to use the CNN algorithm (Convolutional Neural Network) that didn't learn in class. The image recognition domain. Researchers optimize the algorithms about image recognition in the past years. There are exist more and more packages, libraries and framework to apply on Machine Learning.

Overfitting

About overfitting problem. Below output is one case in overfitting, we can see the results not all get a better model by increasing epochs. Epoch that 14 to 17 increases, but then decreases.

```
[14] Train-accuracy=0.995119047619049
[15] Train-accuracy=0.995404761904763
[16] Train-accuracy=0.996452380952383
[17] Train-accuracy=0.996500000000002
[18] Train-accuracy=0.995238095238097
[19] Train-accuracy=0.995619047619049
[20] Train-accuracy=0.996880952380954
```



Pic-9. Overfitting problem.

Smoothing kernel and others

Look the image below. We can see some different desired results for images that based on choosed kernel (smooth, blur, enhance or sharpen). [7] Regrettably, **MXNet** doesn't support assign the kernel by user in current version.

Some other kernel examples

1	1	1
1	1	1
1	1	1

Unweighted 3x3 smoothing kernel

0	1	0
1	4	1
0	1	0

Weighted 3x3 smoothing kernel with Gaussian blur

0	-1	0
-1	5	-1
0	-1	0

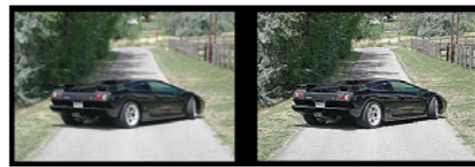
Kernel to make image sharper

-1	-1	-1
-1	9	-1
-1	-1	-1

Intensified sharper image



Gaussian Blur



Sharpened image

Pic-10. Smoothing kernel.

Speed up on computing

Computing is too slow in this project. Average computing time around 4500 - 5000 seconds. We can try to use GPU computing to speed up.

References

[1] Convolutional Neural Networks (LeNet), 2015/12,
<http://deeplearning.net/tutorial/lenet.html>

- [2] *Andrew Gibiansky*, Convolutional Neural Networks, 2015/12,
<http://andrew.gibiansky.com/blog/machine-learning/convolutional-neural-networks/>
- [3] *Michael Nielsen*, Deep learning, 2015/12,
<http://neuralnetworksanddeeplearning.com/chap6.html>
- [4] Feature extraction using convolution, 2015/12,
http://deeplearning.stanford.edu/wiki/index.php/Feature_extraction_using_convolution
- [5] Convolutional Neural Networks for Visual Recognition, 2016/01,
<http://cs231n.github.io/neural-networks-1/>
- [6] *Chen-Yu Lee, Patrick W. Gallagher, Zhuowen Tu*, Generalizing Pooling Functions in Convolutional Neural Networks: Mixed, Gated, and Tree. 2015.
- [7] *Jamie Ludwig*, Image Convolution, 2016/01
http://web.pdx.edu/~jduh/courses/Archive/geog481w07/Students/Ludwig_ImageConvolution.pdf
- [8] Convolutional neural network, 2016/01,
https://en.wikipedia.org/wiki/Convolutional_neural_network

Special Thanks

Kai-Hsiang Hsu, NTU. Thanks for giving information about the better algorithms in image recognition.