# Handwritten Digits with Multilayer Backpropagation Neural Networks

**Shilin Zhu**
Ph.D. student, Computer Science
UCSD
La Jolla, CA
shz338@eng.ucsd.edu

**Yunhui Guo**
Ph.D. student, Computer Science
UCSD
La Jolla, CA
yug185@eng.ucsd.edu

## 1   Abstract

## 2   Classification

### 2.1   Mini-batch gradient descent

In this section, we use mini-batch gradient descent to classify the MNIST dataset. We split the 60000 images in the training set into two parts: the first 50000 images are used to train the model, the last 10000 images are used as validation set to do early stopping. We stop the training procedure once the loss on the validation set goes up and we save the weights that achieves the minimum loss on the validation set. And there are 10000 images in the test set.

We use one hidden layer of 64 nodes, and the mini-batch size is 128. We use a learning rate of 0.01 and sigmoid activation function. We use standard normal distribution to initialize the weights and biases. For the weights, we multiply 0.01 to prevent large initialized values. We report the accuracy and loss on the training set, test se and validation set every batch. The following graphs show the accuracy and loss over each batch on different sets.

## 3   Adding the Tricks of the Trade"

## 4   Experiment with Network Topology

### 4.1   Experiments with differnet hidden units

### 4.2   Doubling the hidden layers

### 4.3   More tricks

In this section, in order to improve the preformance of the network, we consider the following tricks. In our experiments, we found that the network can achieve fast convergence and higher test accuracy with the tricks.

#### 4.3.1   ReLU

We consider using ReLU as the activation function. The ReLU function can be

$$\text{ReLU}(x) = \begin{cases} x & \text{if x} > 0, \\ 0 & \text{otherwise,} \end{cases}$$

The gradient of ReLU can be caculated as,

$$\text{dReLU}(x) = \left\{ \begin{array}{ll} 1 & \text{if x} > 0, \\ 0 & \text{otherwise}, \end{array} \right.$$

The following graphs shows the

### 4.3.2 Leaky ReLU

We consider using leaky ReLU as the activation function. The leaky ReLU function can be

$$\text{LeakyReLU}(x) = \left\{ \begin{array}{ll} x & \text{if x} > 0, \\ 0.01x & \text{otherwise}, \end{array} \right.$$

The gradient of leaky ReLU can be caculated as,

$$\text{dLeakyReLU}(x) = \left\{ \begin{array}{ll} 1 & \text{if x} > 0, \\ 0.01 & \text{otherwise}, \end{array} \right.$$

### 4.3.3 Nesterov momentum

We consider using Nesterov momentum.

### 4.3.4 Xavier initializtion

We consider using Xavier initializtion to initialize the weights.

### 4.3.5 Dropout