

CS 189 Homework 6

Edwin Xie, John Du, Daniel Zezula

April 21, 2014

External Sources

We used CUDA to parallelize the training process.

Single layer neural network

(i) Stochastic gradient descent for loss functions:

We want to find a $n_{in} \times n_{out}$ $\frac{\partial J}{\partial \mathbf{W}}$ where $(\frac{\partial J}{\partial \mathbf{W}})_{ij}$ is taken over W_{ij} .

Define W_k such that $W = [W_1 \ W_2 \ \dots \ W_{out}]$.

Then $\frac{\partial J}{\partial \mathbf{W}} = [\frac{\partial J}{\partial \mathbf{W}_1} \ \frac{\partial J}{\partial \mathbf{W}_2} \ \dots \ \frac{\partial J}{\partial \mathbf{W}_{out}}]$.

We also want to find a $n_{out} \times 1$ $\frac{\partial J}{\partial \mathbf{b}}$, where $\frac{\partial J}{\partial \mathbf{b}} = [\frac{\partial J}{\partial b_1} \ \frac{\partial J}{\partial b_2} \ \dots \ \frac{\partial J}{\partial b_{out}}]^T$.

The stochastic gradient descent updates are then

$$\mathbf{W} = \mathbf{W} - \eta \frac{\partial J}{\partial \mathbf{W}}$$

$$\mathbf{b} = \mathbf{b} - \eta \frac{\partial J}{\partial \mathbf{b}}$$

for every training point x .

Derivations for mean squared error:

$$J = \frac{1}{2} \sum_{k=1}^{n_{out}} (t_k - y_k)^2$$

$$J = \frac{1}{2} (\mathbf{t} - \mathbf{y})^T (\mathbf{t} - \mathbf{y})$$

$$\frac{\partial J}{\partial \mathbf{W}_{\mathbf{k}}} = -(t_k - y_k) \frac{\partial y_k}{\partial \mathbf{W}_{\mathbf{k}}}$$

$$\frac{\partial y_k}{\partial \mathbf{W}_{\mathbf{k}}} = y_k(1 - y_k) \frac{\partial}{\partial \mathbf{W}_{\mathbf{k}}} (\mathbf{W}_{\mathbf{k}}^T \mathbf{x} + b_k)$$

$$\frac{\partial y_k}{\partial \mathbf{W}_{\mathbf{k}}} = y_k(1 - y_k) \mathbf{x}$$

$$\frac{\partial J}{\partial \mathbf{W}_{\mathbf{k}}} = -(t_k - y_k)y_k(1 - y_k)\mathbf{x}$$

$$\frac{\partial J}{\partial \mathbf{b}_{\mathbf{k}}} = -(t_k - y_k)\frac{\partial y_k}{\partial \mathbf{b}_{\mathbf{k}}}$$

$$\frac{\partial y_k}{\partial \mathbf{b}_{\mathbf{k}}} = y_k(1 - y_k)$$

$$\frac{\partial J}{\partial \mathbf{b}_{\mathbf{k}}} = -(t_k - y_k)y_k(1 - y_k)$$

Derivation for cross-entropy error:

$$J = - \sum_{k=1}^{n_{out}} [t_k \ln y_k + (1 - t_k) \ln (1 - y_k)]$$

$$\frac{\partial J}{\partial \mathbf{W}_{\mathbf{k}}} = \frac{\partial J}{\partial y_k} \frac{\partial y_k}{\partial \mathbf{W}_{\mathbf{k}}}$$

$$\frac{\partial J}{\partial y_k} = -\left(\frac{t_k}{y_k} - \frac{1 - t_k}{1 - y_k}\right) = \frac{y_k - t_k}{y_k(1 - y_k)}$$

$$\frac{\partial y_k}{\partial \mathbf{W}_{\mathbf{k}}} = y_k(1 - y_k)\mathbf{x}$$

$$\frac{\partial J}{\partial \mathbf{W}_{\mathbf{k}}} = (y_k - t_k)\mathbf{x}$$

$$\frac{\partial J}{\partial \mathbf{b}_{\mathbf{k}}} = (y_k - t_k)$$

Derivation for vectorizing, Forward Propagation

For one training example:

$$y_k = \sigma\left(\sum_{j=1}^{n_{features}} \mathbf{W}_{jk}\mathbf{x}_j + \mathbf{b}_k\right)$$

where $1 \leq k \leq n_{out}$ and $1 \leq j \leq n_{features}$

$$y_k = \sigma(\mathbf{W}_{\mathbf{k}}^T \mathbf{x} + \mathbf{b}_k)$$

For element-wise sigmoid:

$$\mathbf{y} = \sigma(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

above we have the vector dimmensions of:

$$\mathbf{y} = n_{out} \times 1, \mathbf{b} = n_{out} \times 1, \mathbf{W}_k = 1 \times n_{features}, \mathbf{x} = 1 \times n_{features}, \mathbf{W} = n_{features} \times n_{out}$$

For batch of m examples we want:

Y has columns of $\mathbf{y}_1 \dots \mathbf{y}_m$

$$Y = \sigma([\mathbf{W}^T \mathbf{x}_1 + \mathbf{b} \dots \mathbf{W}^T \mathbf{x}_m + \mathbf{b}])$$

$$Y = \sigma(\mathbf{W}^T \mathbf{X} + \mathbf{B})$$

\mathbf{B} has m columns equal to \mathbf{b}

$$X = m \times n_{features}$$

Derivation for vectorizing, Back Propagation

For 1 training example, mean squares:

$$\frac{\partial J}{\partial \mathbf{W}_k} = -(t_k - y_k)y_k(1 - y_k)\mathbf{x}$$

$$\delta_k^L = (y_k - t_k)y_k(1 - y_k)$$

For cross entropy:

$$\frac{\partial J}{\partial \mathbf{W}_k} = \frac{\partial J}{\partial y_k} \frac{\partial y_k}{\partial \mathbf{W}_k}$$

$$\frac{\partial J}{\partial \mathbf{W}_k} = (y_k - t_k)\mathbf{x}$$

$$\delta_k^L = (y_k - t_k)$$

In both cases continue with:

$$\frac{\partial J}{\partial \mathbf{W}_k} = \mathbf{x} \times \delta_k^L$$

define:

$$\frac{\partial J}{\partial \mathbf{W}} = [\frac{\partial J}{\partial \mathbf{W}_1} \dots \frac{\partial J}{\partial \mathbf{W}_k}]$$

where $k = 1, \dots n_{features}$ then:

$$\frac{\partial J}{\partial \mathbf{W}} = \mathbf{x} \times \delta^{L^T}$$

For the bias:

$$\frac{\partial J}{\partial \mathbf{b}_k} = \delta_k^L$$

define:

The rows of $\frac{\partial J}{\partial \mathbf{b}}$ are $\frac{\partial J}{\partial b_1} \dots \frac{\partial J}{\partial b_{out}}$

$$\frac{\partial J}{\partial \mathbf{b}} = \delta^L$$

$$\frac{\partial J}{\partial \mathbf{W}_k} = n_{features} \times 1$$

$$\frac{\partial J}{\partial \mathbf{W}} = n_{features} \times n_{out}$$

where the columns of $\frac{\partial J}{\partial \mathbf{W}}$ are:

$$\frac{\partial J}{\partial \mathbf{W}_1} \cdots \frac{\partial J}{\partial \mathbf{W}_k}$$

$$\delta^L = n_{out} \times 1$$

where the columns of Δ^L are $\delta_1^L \dots \delta_m^L$
for batch of m examples:

$$\frac{\partial J}{\partial \mathbf{W}_m} = X_m \delta_m^{L^T}$$

$$X \times \Delta^L = \sum_{i=1}^m (x_i \delta_i^{L^T})$$

$$X \times \Delta^L = \sum_{i=1}^m \left(\frac{\partial J}{\partial \mathbf{W}_i} \right)$$

$$X \times \Delta^L = \frac{\partial J}{\partial \mathbf{W}}$$

Which is what we want for batch gradient descent. Similarly:

$$\frac{\partial J}{\partial \mathbf{b}} = \sum_{i=1}^m \delta_i^L$$

$$\frac{\partial J}{\partial \mathbf{b}} = \Delta^L [1_1 \dots 1_m]$$

Where Δ^L has dimensions $n_{out} \times m$

Multilayer feed forward neural network

(i) Parameter update equations for loss functions at output layer:

$$\frac{\partial J}{\partial \mathbf{W}_{ij}^l} = \frac{\partial J}{\partial \mathbf{y}_j^l} \cdot \frac{\partial \mathbf{y}_j^l}{\partial \mathbf{s}_j^l} \cdot \frac{\partial \mathbf{s}_j^l}{\partial \mathbf{W}_{ij}^l}$$

For mean squares:

$$\frac{\partial J}{\partial \mathbf{y}_j^l} \cdot \frac{\partial \mathbf{y}_j^l}{\partial \mathbf{s}_j^l} = (\mathbf{y}_j^L - t_j) \mathbf{y}_j^L (1 - \mathbf{y}_j^L) = \delta_j^L$$

For cross entropy:

$$\frac{\partial J}{\mathbf{y}_j^L} \cdot \frac{\partial \mathbf{y}_j^L}{\mathbf{s}_j^L} = \frac{\mathbf{y}_j^L - t_k}{\mathbf{y}_j^L(1 - \mathbf{y}_j^L)} \mathbf{y}_j^L(1 - \mathbf{y}_j^L) = (\mathbf{y}_j^L - t_j) = \delta_j^L$$

Continue in either case with:

$$\frac{\partial \mathbf{s}_j^L}{\mathbf{W}_{ij}^L} = \mathbf{y}_i^{L-1}$$

Inductive update for the hidden layers:

$$\begin{aligned} \delta_i^{L-1} &= \frac{\partial J}{\partial \mathbf{s}_i^{L-1}} = \sum_{j=1}^{d(l)} \frac{\partial J}{\partial \mathbf{s}_j^L} \cdot \frac{\partial \mathbf{s}_j^L}{\partial \mathbf{y}_i^{L-1}} \cdot \frac{\partial \mathbf{y}_i^{L-1}}{\partial \mathbf{s}_i^{L-1}} \\ &= \sum_{j=1}^{d(l)} (\delta_j^L \cdot \mathbf{W}_{ij}^L \cdot (1 - \mathbf{y}_i^{L-1^2})) \\ &= (1 - \mathbf{y}_i^{L-1^2}) \sum_{j=1}^{d(l)} (\delta_j^L \cdot \mathbf{W}_{ij}^L) \end{aligned}$$

For mean squares:

$$\delta^L = (\mathbf{y}^L - \mathbf{t}) \circ (\mathbf{y}^L - \mathbf{y}^{L^2})$$

For cross entropy:

$$\delta^L = (\mathbf{y}^L - \mathbf{t})$$

Where \circ is the Hadamard product.

$$\delta^{L-1} = (\mathbf{1} - \mathbf{y}^{L-1^2}) \circ (\mathbf{W}^L \delta^L)$$

\mathbf{y} has dimensions $d^{l-1} \times n$
 $\delta^{\mathbf{L}^T}$ has dimensions $d^l \times n$

$$\frac{\partial J}{\mathbf{W}^L} = \mathbf{y}^{L-1} \cdot \delta^{L^T}$$

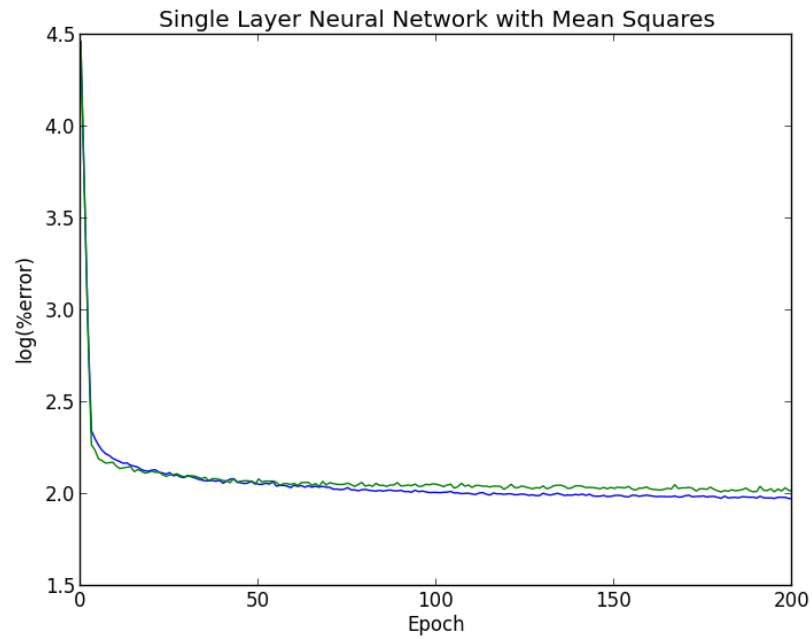
$$\frac{\partial J}{\mathbf{b}^L} = \delta^L$$

$$\mathbf{y}^0 = \mathbf{x}$$

$$\delta^{L-1} = \frac{\partial \mathbf{y}^{L-1}}{\partial \mathbf{s}^{L-1}} \mathbf{W}^L \delta^L$$

Results

We implemented a single layer neural network with both mean squares and cross entropy as error functions. We also implemented a two hidden layer multi layer neural network with both mean squares and cross entropy loss functions. The graphs of these can be seen below, where the green line is the test set, and blue line is the training set. The lowest error we received is with the multi layer neural network with the cross entropy error function, which resulted in a



1.26

