

CS 189: Homework 6

Michael Stephen Chen
Kaggle Acct: michaelstchen
SID: 23567341

April 26, 2016

Problem 1

The stochastic gradient descent update for the weights V and W are as follows,

$$V = V - \epsilon \nabla_V J$$

$$W = W - \epsilon \nabla_W J$$

Where J is our loss function. The derivation for $\nabla_V J$ and $\nabla_W J$ when using the **mean squared error** as our loss function is provided below

W Derivation:

$$\begin{aligned} \frac{\partial L}{\partial w_j} &= \frac{\partial L}{\partial z_j} \frac{\partial z_j}{\partial w_j} \\ &= \frac{\partial L}{\partial z_j} \leftarrow z_j(1-z_j)h \\ &= (z_j - y_j)z_j(1-z_j)h \\ \frac{\partial L}{\partial w_j} &= \delta_j h \end{aligned}$$

Annotations: δ_j is a scalar, h is a $(200+1) \times 1$ vector.

V Derivation:

$$\begin{aligned} \frac{\partial L}{\partial v_i} &= \frac{\partial L}{\partial h_i} \frac{\partial h_i}{\partial v_i} \\ &= \frac{\partial L}{\partial h_i} \leftarrow (1-h_i^2)x \\ \frac{\partial L}{\partial v_i} &= \left(\sum_j w_{ji} \delta_j \right) (1-h_i^2)x \end{aligned}$$

Annotations: $\sum_j w_{ji} \delta_j$ is a scalar, x is a $(784+1) \times 1$ vector.

Loss Function and Chain Rule:

$$L = \frac{1}{2} \sum_{j=1}^{10} (y_j - z_j)^2$$

$$\frac{\partial L}{\partial z_j} = z_j - y_j$$

The diagram shows the flow of gradients: $L \rightarrow z \rightarrow g(Wx) \rightarrow h \rightarrow V$ and $h \rightarrow W \rightarrow z$.

Note that $\delta_j = z_j(1-z_j)(z_j - y_j)$ and

The above derivations are done on a per row basis where V_i is the i th row of V and W_j is the j th row of W

The vectorized version is as follows:

$$\begin{aligned} \frac{\partial L}{\partial W} &= \delta h^T \\ &\quad \uparrow \quad \uparrow \\ &\quad 10 \times 1 \quad 1 \times (200+1) \end{aligned}$$

$$\frac{\partial L}{\partial V} = \left[(W^T \delta) (1-h^2) \right] x^T$$

Annotations: $W^T \delta$ is element-wise multiplication (200×1), $(1-h^2)$ is element-wise square (200×1), x^T is a $(784+1) \times 1$ vector.

In the above derivation, I adjusted the indexing notation from what was given (please refer to the starred "Note" above). Also I use the same index j to denote row j of the W matrix as well as to denote the j^{th} element of the vectors z_j and y_j (so $1 \leq j \leq 10$). Likewise I use i to denote the i^{th} row of V as well as the i^{th} column of W (so $1 \leq i \leq 200 \text{ or } 201$). This is the indexing scheme that Prof. Shewchuck used in his lecture notes and I prefer this over the provided scheme because we see the relationship between the elements of the different matrices/arrays.

The derivation for $\nabla_V J$ and $\nabla_W J$ when using the **cross-entropy error** as our loss function is provided below. Note that we only rederive the parts that change (namely δ) because everything else aside from the loss function is the same as before. Just plug the new δ into the previously derived equations for $\nabla_V J$ and $\nabla_W J$.

$$\xrightarrow{z} \left(- \sum_j \left[y_j \ln(z_j) + (1-y_j) \ln(1-z_j) \right] \right) \xrightarrow{L}$$

$$\begin{aligned} \frac{\partial L}{\partial z_j} &= - \left(\frac{y_j}{z_j} + \frac{1-y_j}{1-z_j} \right) \\ &= - \left[\frac{(1-z_j)y_j - z_j(1-y_j)}{z_j(1-z_j)} \right] \\ &= - \left[\frac{y_j - z_j y_j - z_j + z_j y_j}{z_j(1-z_j)} \right] \\ &= \frac{z_j - y_j}{z_j(1-z_j)} \end{aligned}$$

and so,

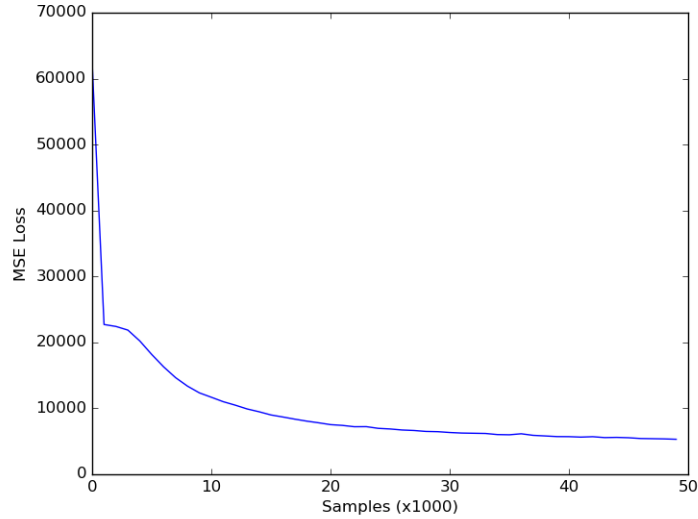
$$\begin{aligned} \delta_j &= z_j(1-z_j) \frac{\partial L}{\partial z_j} \\ &= z_j - y_j \end{aligned} \quad \text{and} \quad \delta = z - y$$

Problem 2

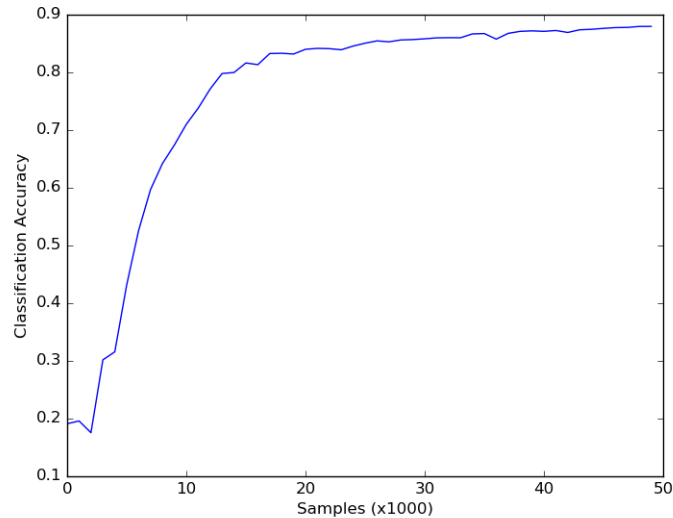
a. Mean-Squared Error

- All of the weights were initialized randomly, drawing from a normal distribution with a mean of 0 and a standard deviation of 0.01. A constant learning rate of 0.01 was used. Training was stopped after one epoch (50000 images drawn randomly with replacement).

- The training accuracy was **0.89872**. The validation accuracy was **0.8814**.
- The total training time was **839.56** seconds (14 minutes).
- Below is a plot of the batch MSE loss over the training period calculated for every 1000 samples trained on

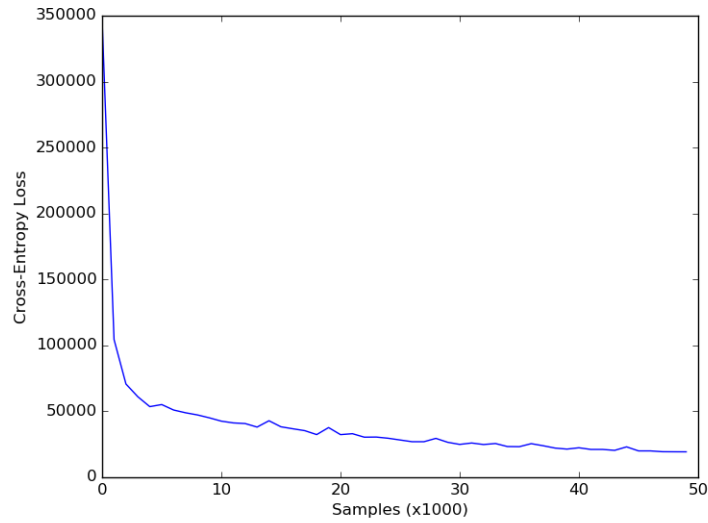


Below is a plot of the training accuracy over the training period.

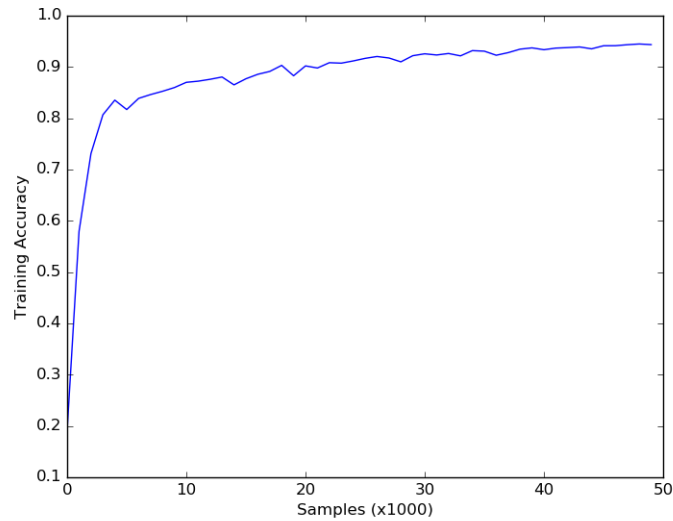


b. Cross-Entropy Error

- All of the weights were initialized randomly, drawing from a normal distribution with a mean of 0 and a standard deviation of 0.01. A constant learning rate of 0.01 was used. Training was stopped after one epoch (50000 images drawn randomly with replacement).
- The final training accuracy was **0.9536**. The validation accuracy was **0.9412**.
- The total training time was **811.48** seconds (13.5 minutes).
- Below is a plot of the batch cross-entropy loss over the training period calculated for every 1000 samples trained on



Below is a plot of the training accuracy over the training period.



- c. From our results we see that the cross-entropy loss function performs considerably better than the mean-squared error. The difference in performance can be explained by examining their corresponding gradients with respect to the weight vectors, and in particular looking at the δ term that I define in my derivation above. For the mean-squared error,

$$\delta = [z(1 - z)](z - y)$$

and for the cross-entropy error,

$$\delta = z - y$$

Where z are our neural network outputs and y are the true labels. We see that the mse δ contains (inside the brackets) the derivative of the sigmoid function, which approaches 0 when our outputs z are close to 0 or 1. Therefore δ and consequently the gradients also approach 0 as z gets close to 0 and 1. For these cases, learning is slow because our gradients are small. In contrast, the δ for the cross-entropy error does not suffer from this “learning slowdown” (I got the term from

<http://neuralnetworksanddeeplearning.com/chap3.html> which provides a nice description of this phenomenon) and is strictly proportional to the error, which is what we would like because the more wrong we are, the larger our gradient and the more we learn.

- d. My best Kaggle score (after training for 10 epochs with 50000 images per epoch) was **0.97860**