

---

# Using Neural Networks to Effectively Classify Hand-Written Digits of the MNIST Dataset

---

Chetan Gandotra  
UC San Diego  
9500 Gilman Drive  
cgandotr@ucsd.edu

## Abstract

1 This report discusses the second programming assignment of our course CSE 253:  
2 Neural Networks and Pattern Recognition, its solutions and the inferences we drew.  
3 A variable layer neural network was implemented from the scratch and observations  
4 were made based on various parameters and before/after adding certain trades of  
5 tricks as discussed in Yann LeCun's famous paper "Efficient BackProp". The  
6 data-set used was the famous MNIST data-set and a ten-way classification was  
7 performed on it. A test data-set accuracy in excess of 97% was achieved using  
8 various mechanisms and tricks, which is almost at par with the accuracy reported by  
9 LeCun on his website. This is the individual part of the programming assignment.

## 10 1 Linear Regression and Sum-of-squared Error

### 11 1.1 Introduction

12 In linear regression problems, we try to fit the data generated from:

$$t = h(x) + \epsilon \quad (1)$$

13 where  $x$  is a  $K$ -dimensional vector,  $h(x)$  is a deterministic function of  $x$ , where  $x$  includes the bias term  
14  $x_0$  and  $\epsilon$  is random noise that has a Gaussian probability distribution with zero mean and variance  $\sigma^2$ ,  
15 i.e.,  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . Now, consider  $y$  to be of form

$$y = \sum_{i=1}^K w_i \cdot x_i \quad (2)$$

16 We need to prove that finding the optimal parameter  $w$  for the above linear regression problem on  
17 the dataset  $D = (x^1, t^1), (x^2, t^2), \dots, (x^N, t^N)$  is equal to finding the  $w^*$  that minimizes the sum of  
18 squared error (SSE):

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{n=1}^N (t^n - y^n)^2 \quad (3)$$

### 19 1.2 Methodology

20 In this section, we will prove that finding the optimal parameter  $w$  for the above linear regression  
21 problem on the dataset  $D = (x^1, t^1), (x^2, t^2), \dots, (x^N, t^N)$  is equal to finding the  $w^*$  that minimizes  
22 the sum of squared error (SSE):

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{n=1}^N (t^n - y^n)^2 \quad (4)$$

23 To begin with, we look at (1). We are saying that the actual label is equal to the predicted label, plus  
 24 some noise  $\epsilon$ . Now, we are given that the noise  $\epsilon$  has a Gaussian probability distribution with zero  
 25 mean and variance  $\sigma^2$ , i.e.,  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ .

26 Modifying (1), we can write:

$$\epsilon = h(x) - t \quad (5)$$

27 Also, the Gaussian distribution for  $\epsilon$  can be written as:

$$p(\epsilon) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \quad (6)$$

28 Substituting value of  $\epsilon$  from (5) and putting  $h(x) = y$  (which corresponds to the model we are using),  
 29 we get:

$$p(\epsilon) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{(y - t)^2}{2\sigma^2}\right) \quad (7)$$

For least error, the probability of predicting a label must be maximized, when provided with the corresponding input features  $X$ . This is repeated for all training examples from 1 to  $N$  and hence, we end up taking the multiplication of all probabilities. In a nutshell, we can maximize the probability of correct prediction by finding weights  $w$  such that the following expression is maximized:

$$\max_w \prod_{n=1}^N \frac{1}{\sigma(2\pi)^{1/2}} \exp\left(-\frac{(t_n - y_n)^2}{2\sigma^2}\right)$$

The above maximization objective follows from (2), which contains  $w$ . Removing the constant terms,

$$= \max_w \prod_{n=1}^N \exp(-(t_n - y_n)^2)$$

Taking log and converting maximization objective to minimize by discarding the negative sign, we get:

$$\begin{aligned} &= \min_w \sum_{n=1}^N (t_n - y_n)^2 \\ &= \operatorname{argmin}_w \sum_{n=1}^N (t_n - y_n)^2 \end{aligned}$$

30 which is essentially the same as equation (4).

### 31 1.3 Results

Using the derivation described above, we have proved that finding the optimal parameter  $w$  for the above linear regression problem on the dataset  $D = (x^1, t^1), (x^2, t^2), \dots, (x^N, t^N)$  is equal to finding the  $w^*$  that minimizes the sum of squared error (SSE):

$$w^* = \operatorname{argmin}_w \sum_{n=1}^N (t^n - y^n)^2$$

### 32 1.4 Discussion

33 Adjusting the weight vector to get maximum probability of finding the correct label is equivalent to  
 34 minimizing the sum of squared error. This is true because in both the cases we are trying to reduce the  
 35 number of mistakes by reducing the distance between actual and predicted values. A good classifier  
 36 must have less sum of squared error or maximum probability of the label being correct, and for these  
 37 two statements being equivalent makes sense.

## 2 Derivation of Delta's, Update Rule and Vectorization

### 2.1 Introduction

In this question, we have been given a three layer neural network with  $J$  units in the hidden layer. We use index  $k$  to represent a node in output layer, index  $j$  to represent a node in hidden layer and index  $i$  to represent a node in the input layer. Additionally, the weight from node  $i$  in the input layer to node  $j$  in the hidden layer is  $w_{ij}$ . Similarly, the weight from node  $j$  in the hidden layer to node  $k$  in the output layer is  $w_{jk}$ .

Firstly, we are required to derive the expression for  $\delta$  for both the units of output layer ( $\delta_k$ ) and the hidden layer ( $\delta_j$ ). The definition of  $\delta$  here is  $\delta_i = -\frac{\partial E}{\partial a_i}$ , where  $a_i$  is the weighted sum of the inputs to unit  $i$ .

Second, we need to derive the update rule for  $w_{ij}$  and  $w_{jk}$  using learning rate  $\eta$ , starting with the gradient descent rule:

$$w_{ij} = w_{ij} - \eta \frac{\partial E}{\partial w_{ij}} \quad (8)$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial a_j} \frac{\partial a_j}{\partial w_{ij}} \quad (9)$$

The derivative should take into account all of the outputs, so:

$$\delta_j = -\frac{\partial E^n}{\partial a_j^n} = \sum_k \frac{\partial E^n}{\partial y_k} \frac{\partial y_k}{\partial a_j} \quad (10)$$

Finally, we have been asked to write these equations in a vectorized form.

### 2.2 Methodology

For the first part for this question, we are required to derive expressions for  $\delta_k$  and  $\delta_j$ , where  $\delta_i = -\frac{\partial E}{\partial a_i}$ . We know that:

$$a_k = \sum_j w_{kj} z_j \quad (11)$$

Also,

$$\delta_k^n = -\frac{\partial E^n}{\partial a_k}$$

$$\Rightarrow \delta_k^n = -\sum_{m=1}^C \frac{\partial E^n}{\partial y_m^n} \cdot \frac{\partial y_m^n}{\partial a_k^n} \quad (12)$$

Now, we know that the error for  $n^{th}$  example -  $E^n$  can be written as:

$$E^n = -\sum_k t_k^n \cdot \ln(y_k^n)$$

Taking derivative with respect to  $y_k$ ,

$$\Rightarrow \frac{\partial E^n}{\partial y_k} = -t_k^n \cdot \frac{1}{y_k^n} \cdot 1$$

$$\Rightarrow \frac{\partial E^n}{\partial y_k} = -\frac{t_k^n}{y_k^n} \quad (13)$$

Also, since

$$y_l^n = \frac{e^{a_l^n}}{\sum_{m=1}^C e^{a_m^n}}$$

58 it will have different derivatives when  $l = k$  and when  $l \neq k$ . Thus, the derivatives  $\frac{\partial y_k^n}{\partial a_k}$  and  $\frac{\partial y_{k'}^n}{\partial a_k}$   
 59 ( $l \neq k$  case) can be written as:

$$\begin{aligned}\frac{\partial y_k^n}{\partial a_k} &= \frac{e^{a_k^n}}{\sum_{m=1}^C e^{a_m^n}} - e^{a_k^n} \left[ \frac{1}{(\sum_{m=1}^C e^{a_m^n})^2} \right] e^{a_k^n} \\ \frac{\partial y_k^n}{\partial a_k} &= y_k^n (1 - y_k^n)\end{aligned}\tag{14}$$

And,

$$\begin{aligned}\frac{\partial y_{k'}^n}{\partial a_k} &= -e^{a_{k'}^n} e^{a_k^n} \frac{1}{(\sum_{m=1}^C e^{a_m^n})^2} \\ \frac{\partial y_{k'}^n}{\partial a_k} &= -(y_{k'}^n)^2\end{aligned}\tag{15}$$

Putting (13), (14) and (15) in (12),

$$\begin{aligned}\delta_k^n &= -\frac{\partial E^n}{\partial a_k^n} = -\sum_{l=1 \& l \neq k}^C \left( \frac{t_l^n}{y_l^n} (y_l^n)^2 \right) - \left( \frac{t_k^n}{y_k^n} (y_k^n (1 - y_k^n)) \right) \\ \delta_k^n &= -\sum_{l=1 \& l \neq k}^C (t_l^n y_k^n) - (t_k^n (1 - y_k^n))\end{aligned}\tag{16}$$

60 Following the one hot representation, only one of the labels from  $l = 1$  to  $C$  in  $t^n$  would be 1 and all  
 61 the others would be 0. Therefore, for input term  $u$  where  $t_k^u$  is 1, the derivative becomes:

$$\begin{aligned}\frac{\partial E^u}{\partial a_k^u} &= y_k^u - 1 \\ \Rightarrow \frac{\partial E^u}{\partial a_k^u} &= y_k^u - t_k^u\end{aligned}\tag{17}$$

For input term  $v$  where some  $t_{k'}^v$  is 1 such that  $k' \neq k$ , this term becomes:

$$\begin{aligned}\frac{\partial E^v}{\partial a_k^v} &= y_k^v \\ \Rightarrow \frac{\partial E^v}{\partial a_k^v} &= y_k^v - t_k^v\end{aligned}\tag{18}$$

64 Therefore, combining both the cases together, we can form an expression for  $\delta_k^n$  as:

$$\delta_k^n = -\frac{\partial E^n}{\partial a_k^n} = t_k^n - y_k^n\tag{19}$$

65 Now, we wish to derive an expression for  $\delta_j$ . Note that the derivation below is taken from Professor  
 66 Gary's class notes.

We have,

$$\delta_j^n = -\frac{\partial E^n}{\partial a_j^n}$$

Using chain rule and summation over all labels,

$$\delta_j^n = -\sum_{k=1}^C \frac{\partial E^n}{\partial y_k^n} \frac{\partial y_k^n}{\partial a_j^n}$$

$$\Rightarrow \delta_j^n = - \sum_{k=1}^C \frac{\partial E^n}{\partial y_k^n} \frac{\partial y_k^n}{\partial y_j^n} \frac{\partial y_j^n}{\partial a_j^n}$$

Taking the terms independent of  $k$  outside of summation,

$$\Rightarrow \delta_j^n = - \frac{\partial y_j^n}{\partial a_j^n} \sum_{k=1}^C \frac{\partial E^n}{\partial y_k^n} \frac{\partial y_k^n}{\partial y_j^n}$$

Using chain rule again,

$$\Rightarrow \delta_j^n = - \frac{\partial y_j^n}{\partial a_j^n} \sum_{k=1}^C \frac{\partial E^n}{\partial y_k^n} \frac{\partial y_k^n}{\partial a_k^n} \frac{\partial a_k^n}{\partial y_j^n}$$

67 Combining terms using chain rule, we get:

$$\Rightarrow \delta_j^n = - \frac{\partial y_j^n}{\partial a_j^n} \sum_{k=1}^C \frac{\partial E^n}{\partial a_k^n} \frac{\partial a_k^n}{\partial y_j^n} \quad (20)$$

Since,

$$a_k^n = \sum_j w_{jk} y_j$$

68 we can write the derivative  $\frac{\partial a_k^n}{\partial y_j^n}$  as:

$$\frac{\partial a_k^n}{\partial y_j^n} = w_{jk}$$

69 Substituting in (20):

$$\delta_j^n = - \frac{\partial y_j^n}{\partial a_j^n} \sum_{k=1}^C \delta_k^n w_{jk} \quad (21)$$

Since we are using the sigmoid function as activation, we can say that:

$$\begin{aligned} y_j^n &= \sigma(a_j^n) \\ \Rightarrow \frac{\partial y_j^n}{\partial a_j^n} &= \sigma(a_j^n)(1 - \sigma(a_j^n)) \\ \Rightarrow \frac{\partial y_j^n}{\partial a_j^n} &= y_j^n(1 - y_j^n) \end{aligned}$$

Substituting back in equation (21):

$$\delta_j^n = -y_j^n(1 - y_j^n) \sum_{k=1}^C \delta_k^n w_{jk}$$

which can also be written as: (without substituting value of  $\frac{\partial y_j^n}{\partial a_j^n}$ , as done in class):

$$\delta_j^n = - \frac{\partial y_j^n}{\partial a_j^n} \sum_k \delta_k^n w_{jk}$$

70 For the final part of this problem, we need to derive the update rule for  $w_{ij}$  and  $w_{jk}$ .

Looking at the update rule for  $w_{jk}$  first, we have:

$$w_{jk} = w_{jk} - \eta \frac{\partial E}{\partial w_{jk}}$$

Using chain rule,

$$\Rightarrow w_{jk} = w_{jk} - \eta \frac{\partial E}{\partial a_k} \frac{\partial a_k}{\partial w_{jk}}$$

Since we need to repeat this for all data points,

$$\Rightarrow w_{jk} = w_{jk} - \eta \sum_{n=1}^N \frac{\partial E^n}{\partial a_k^n} \frac{\partial a_k^n}{\partial w_{jk}}$$

71 From the definitions of  $\delta_k$  and  $z_j$ , we can write:

$$\Rightarrow w_{jk} = w_{jk} + \eta \sum_{n=1}^N \delta_k^n z_j^n \quad (22)$$

72 From equations derived above for previous parts of question 2, we have:

$$\Rightarrow w_{jk} = w_{jk} + \eta \sum_{n=1}^N (t_k^n - y_k^n) z_j^n \quad (23)$$

73 Similarly, we need to work for the update rule for  $w_{ij}$ . Using the generic rule, we have:

$$w_{ij} = w_{ij} - \eta \frac{\partial E}{\partial w_{ij}}$$

Using chain rule,

$$\Rightarrow w_{ij} = w_{ij} - \eta \frac{\partial E}{\partial a_j} \frac{\partial a_j}{\partial w_{ij}}$$

Changing to summation over all input points form, we get:

$$\Rightarrow w_{ij} = w_{ij} - \eta \sum_{n=1}^N \frac{\partial E^n}{\partial a_j^n} \frac{\partial a_j^n}{\partial w_{ij}}$$

74 Substituting the value of  $\delta_j$ , as computed previously,

$$\Rightarrow w_{ij} = w_{ij} + \eta \sum_{n=1}^N \delta_j^n z_i^n \quad (24)$$

$$\Rightarrow w_{ij} = w_{ij} + \eta \sum_{n=1}^N y_j^n (y_j^n - 1) \sum_{k=1}^C \delta_k^n w_{jk} z_i^n$$

$$\Rightarrow w_{ij} = w_{ij} + \eta \sum_{n=1}^N y_j^n (y_j^n - 1) \sum_{k=1}^C (t_k^n - y_k^n) w_{jk} z_i^n \quad (25)$$

75 Finally, we have been asked to write the vector representation of the equations we just derived.

76 Let  $X$  be the input matrix of dimension  $N \times I$ , where  $I$  is the number of input features and  $N$  is the  
 77 number of examples in our data-set. Note that if we add the bias term as well, the dimension of  $X$   
 78 changes to  $N \times (I + 1)$ . We define  $W_{ij}$  as the weights matrix from input layer to hidden layer, having  
 79 dimensions  $(I + 1) \times J$ . Similarly, we have another weights matrix  $W_{jk}$  which includes weights  
 80 from the hidden layer to the output layer. Given the matrix form, a row in  $W$  matrix holds weight  
 81 parameters from that starting unit to all units in the next layer.

Continuing our general representation, we define  $Z^l$  as the output of layer  $l$ . Thus, we can write:

$$Z^1 = X$$

$$A^{l+1} = W^l Z^l$$

$$Z^{l+1} = \sigma(A^{l+1})$$

where  $W^l$  is  $W_{ij}$  when  $l = 0$  and  $W_{jk}$  when  $l = 1$ . Note that we can have used  $\sigma$  or the sigmoid activation function above while it can be generalized for any other activation function as well.

Now, we re-write our update rules from equations (22) and (24) generically using a vectorized form as follows:

$$W^l = W^l + \eta(\delta^{l+1}(Z^l)^T) \quad (26)$$

[As stated above,  $W^l$  is  $W_{ij}$  when  $l = 0$  and  $W_{jk}$  when  $l = 1$ ] where

$$\delta^l = (W^l)^T \delta^{l+1} .* \sigma'(Z^l) \quad (27)$$

and  $\delta$  for last layer (output) is  $\delta^k$ , written as:

$$\delta^k = T - Y \quad (28)$$

In (27), the operator  $.*$  means element wise product. In (28),  $T$  is the actual output whereas  $Y$  is the output predicted by our model, and is obtained using the  $Z$  value of last layer.

## 2.3 Results

We derived the expression for  $\delta_k^n$  as:

$$\delta_k^n = t_k^n - y_k^n$$

and the expression for  $\delta_j^n$  as:

$$\delta_j^n = -\frac{\partial y_j^n}{\partial a_j^n} \sum_k \delta_k w_{jk}$$

The update rule for  $w_{ij}$  is given as:

$$w_{ij} = w_{ij} + \eta \sum_{n=1}^N y_j^n (y_j^n - 1) \sum_{k=1}^C (t_k^n - y_k^n) w_{jk} z_i^n$$

and that for  $w_{jk}$  is written as:

$$w_{jk} = w_{jk} + \eta \sum_{n=1}^N (t_k^n - y_k^n) z_j^n$$

The vectorized form of update rule can be written as:

$$W^l = W^l + \eta(\delta^{l+1}(Z^l)^T)$$

where

$$\delta^l = (W^l)^T \delta^{l+1} .* \sigma'(Z^l)$$

and

$$\delta^k = T - Y$$

## 2.4 Discussion

All the results above seem to be self explanatory and are derived mathematically using known facts. The expressions above are hence genuine and will be used by us for the remaining of this report. The computation is much faster when we update all  $w_{ij}$ s and  $w_{jk}$ s at the same time using matrix multiplications rather than **for** loops.