

## Exercise 1

In this exercise,  $\sigma_n$  is used to express  $\sigma(\mathbf{w}^T \mathbf{x}_n)$ .

### Exercise 1.1

In this exercise it is required to determine the expression for the negative log likelihood for  $N$  pairs of data  $(x_n, c_n), c_n \in \{0, 1\}$ . It is also reminded that the likelihood loss is  $H(p_n, c_n) = -\mathbb{E}[\log p_n]_{c_n}$ . Finally, it is known that  $p_n$  is in state 1 with probability  $\sigma_n$  and since the setting is binary, it is in state 0 with probability  $1 - \sigma_n$ .

$$\begin{aligned}
 E_{lik}(\mathbf{w}) &= H(p_n, c_n) \\
 &= -\mathbb{E}[\log p_n]_{c_n} \\
 &= -\sum_{n=1}^N \begin{cases} \log(\sigma_n) & c_n = 1 \\ \log(1 - \sigma_n) & c_n = 0 \end{cases} \\
 &= -\sum_{n=1}^N \mathbb{I}[c_n = 1] \log(\sigma_n) + \mathbb{I}[c_n = 0] \log(1 - \sigma_n)
 \end{aligned} \tag{1}$$

The operator  $\mathbb{I}[v]$  is the indicator function: it is equal to 1 if the condition  $v$  specified is true, 0 otherwise. However, the scaled loss is required, in order to normalise it for the number of pairs.

$$E_{lik}(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^N \mathbb{I}[c_n = 1] \log(\sigma_n) + \mathbb{I}[c_n = 0] \log(1 - \sigma_n) \tag{2}$$

### Exercise 1.2

The question requires to determine the convexity of  $E_{lik}$  with respect to  $\mathbf{w}$ . One way to do that is to determine whether the second derivative is ever negative.

$$E_{lik}(\mathbf{w}) = -\frac{1}{N} \sum_{n=1}^n \mathbb{I}[c_n = 1] \log \sigma_n + \mathbb{I}[c_n = 0] \log(1 - \sigma_n) \tag{3}$$

$$\begin{aligned}
 \frac{\partial E}{\partial \mathbf{w}} &= -\frac{1}{N} \sum_{n=1}^n \mathbb{I}[c_n = 1] \frac{x_n \sigma_n (1 - \sigma_n)}{\sigma_n} - \mathbb{I}[c_n = 0] \frac{x_n \sigma_n (1 - \sigma_n)}{1 - \sigma_n} \\
 &= -\frac{1}{N} \sum_{n=1}^n \mathbb{I}[c_n = 1] x_n (1 - \sigma_n) - \mathbb{I}[c_n = 0] x_n \sigma_n \\
 &= -\frac{1}{N} \sum_{n=1}^n x_n (\mathbb{I}[c_n = 1] (1 - \sigma_n) - \mathbb{I}[c_n = 0] \sigma_n)
 \end{aligned} \tag{4}$$

$$\begin{aligned}
 &= -\frac{1}{N} \sum_{n=1}^n x_n (c_n - \sigma_n) \\
 &= -\frac{1}{N} \sum_{n=1}^n x_n c_n - x_n \sigma_n \\
 \frac{\partial^2 E}{\partial \mathbf{w}^2} &= -\frac{1}{N} \sum_{n=1}^n 0 - x_n x_n \sigma_n (1 - \sigma_n) \\
 &= \frac{1}{N} \sum_{n=1}^n x_n^2 \sigma_n (1 - \sigma_n)
 \end{aligned} \tag{5}$$

It can be observed that both  $x_n^2$  and  $\sigma_n(1 - \sigma_n)$  are never negative, which indicates that  $e_{lik}$  is convex with respect to  $\mathbf{w}$ .

### Exercise 1.3

The question requires to determine the convexity of  $E_{sq}$  with respect to  $\mathbf{w}$ . One way to do that is to determine whether the second derivative is ever negative.

$$E_{sq}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^n (c_n - \sigma_n)^2 \quad (6)$$

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{w}} &= \frac{1}{N} \sum_{n=1}^N 2(c_n - \sigma_n)(-1)\sigma_n(1 - \sigma_n)x^n \\ &= -\frac{2}{N} \sum_{n=1}^N (c_n - \sigma_n)\sigma_n(1 - \sigma_n)x^n \\ &= -\frac{2}{N} \sum_{n=1}^N x_n c_n \sigma_n - \sigma_n^2 x_n - c_n \sigma_n^2 x_n + \sigma_n^3 x_n \end{aligned} \quad (7)$$

$$\begin{aligned} \frac{\partial^2 E}{\partial \mathbf{w}^2} &= -\frac{2}{N} \sum_{n=1}^N x_n c_n \sigma_n (1 - \sigma_n) x_n - 2\sigma_n \sigma_n (1 - \sigma_n) x_n x_n - c_n x_n 2\sigma_n \sigma_n (1 - \sigma_n) x_n + 3\sigma_n^2 \sigma_n (1 - \sigma_n) x_n x_n \\ &= -\frac{2}{N} \sum_{n=1}^N c_n x_n^2 \sigma_n (1 - \sigma_n) - 2x_n^2 \sigma_n^2 (1 - \sigma_n) - 2c_n x_n^2 \sigma_n^2 (1 - \sigma_n) + 3x_n^2 \sigma_n^3 (1 - \sigma_n) \\ &= -\frac{2}{N} \sum_{n=1}^N x_n^2 \sigma_n (1 - \sigma_n) (c_n - 2\sigma_n - 2c_n \sigma_n + 3\sigma^2) \\ &= \frac{2}{N} \sum_{n=1}^N x_n^2 \sigma_n (1 - \sigma_n) (2\sigma_n + 2c_n \sigma_n - 3\sigma^2 - c_n) \end{aligned} \quad (8)$$

The factor of interest to determine convexity is then  $(2\sigma_n + 2c_n \sigma_n - 3\sigma^2 - c_n)$ . For  $c_n = 0$ :

$$\begin{aligned} 2\sigma_n - 3\sigma^2 &< 0 \\ 3\sigma^2 - 2\sigma_n &> 0 \\ \sigma_n(3\sigma_n - 2) &> 0 \\ 3\sigma_n - 2 &> 0 \\ \sigma_n &> \frac{2}{3} \end{aligned} \quad (9)$$

For  $c_n = 1$ :

$$\begin{aligned} 2\sigma_n + 2\sigma_n - 3\sigma_n^2 - 1 &< 0 \\ 3\sigma_n^2 - 4\sigma_n + 1 &> 0 \\ \frac{1}{3} &< \sigma_n < 1 \end{aligned} \quad (10)$$

This determines that the second derivative of  $E_{sq}$  with respect to  $\mathbf{w}$  is negative when  $c_n = 0 \wedge \sigma_n \in [\frac{2}{3}, 1]$  and  $c_n = 1 \wedge \sigma_n \in [\frac{1}{3}, 1]$ . Due to this constraints, it cannot be convex.

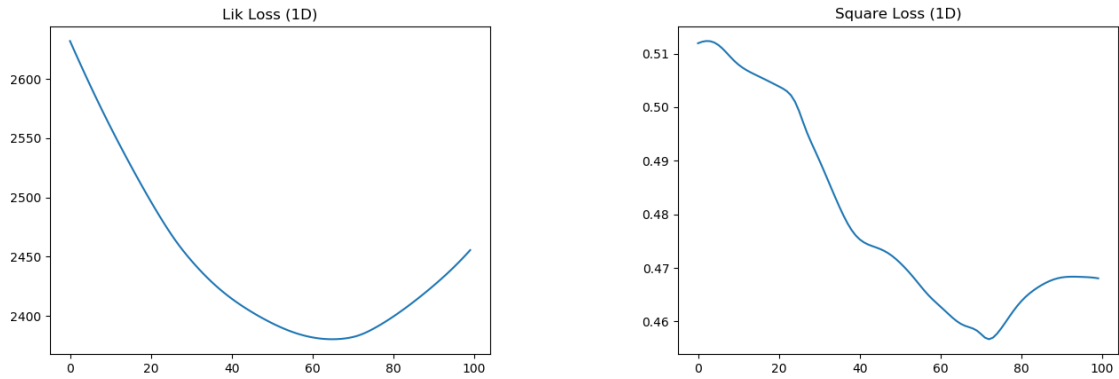


Figure 1: Slices of the 1-dimensional error surface of LikLoss and SquareLoss

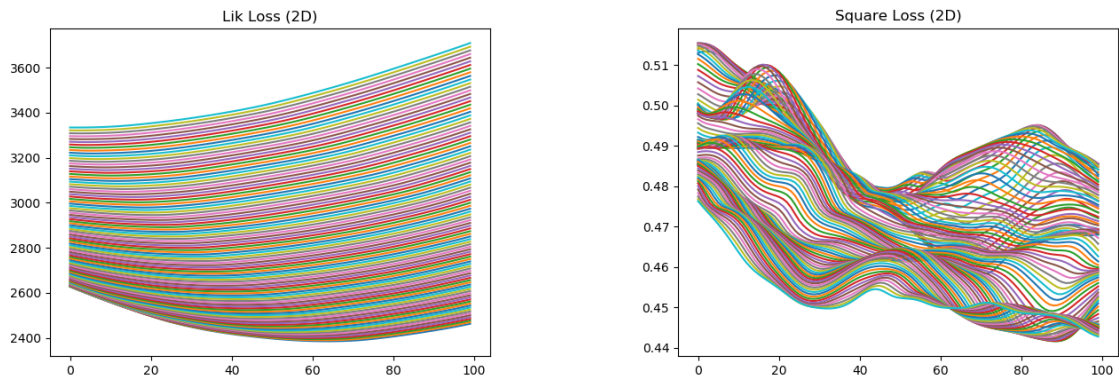


Figure 2: Slices of the 2-dimensional error surface of LikLoss and SquareLoss

### Exercise 1.6 and 1.7

In this exercises, 1-dimensional slices of the error surface of the LogLoss and SquareLoss are plotted to observe whether they are convex or else. While it does not guarantee a definite answer, if a plot shows a function is not convex for the selected interval, then it is not convex. Results are observed in figure 1. Similarly, this is done over two dimensions and the slices are observed in figure 2. For exercise 1.6, the supplied Julia code was used without any change, whereas in exercise 1.7 slight changes have been introduced to account for the additional dimension. The code below is the overall change to the given Julia to plot the four plots.

```
# The values of the three coefficients are selected at random
vec0=10*randn(D)
vec1=10*randn(D)
vec2=10*randn(D)

# The number of steps in the [0, 1] interval to take
I=100

# The values for the two loss functions on 1 dimension
Loss11=zeros(I);
Loss12=zeros(I)
```

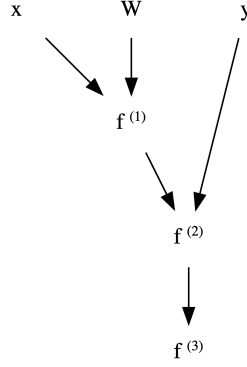


Figure 3: The computation graph of the loss function in Exercise 2.1

```

# The values for the two loss functions on 2 dimensions
Loss21=zeros(I,I);
Loss22=zeros(I,I)

for i=1:I
    # Compute the 1-dimensional error value for both losses
    lambda_i=i/I
    Loss11[i] = SquareLoss(c,NNpred(vec0+lambda_i*vec1,x))
    Loss12[i] = LikLoss(c,vec0+lambda_i*vec1,x)
    for j=1:I
        # Compute the 2-dimensional error value for both losses
        lambda_j=j/I
        Loss21[i,j] = SquareLoss(c,NNpred(vec0+lambda_i*vec1+lambda_j*vec2,x))
        Loss22[i,j] = LikLoss(c,vec0+lambda_i*vec1+lambda_j*vec2,x)
    end
end
end

```

## Exercise 2

### Exercise 2.1

The following loss function is given

$$E = \sum_i \left( y_i - \sum_j W_{ij} x_j \right)^2 \quad (11)$$

and the exercise requires do compute the derivative

$$\frac{\partial E}{\partial W_{ab}} \quad (12)$$

The computation graph is drawn (figure 3) and the functions are determined

$$f_i^{(1)} = \sum_j W_{ij} x_j \quad (13)$$

$$f_i^{(2)} = y_i - f_i^{(1)} \quad (14)$$

$$f^{(3)} = \sum_i \left( f_i^{(2)} \right)^2 \quad (15)$$

The required partial derivatives are calculated.

$$\begin{aligned}
\frac{\partial f_i^{(1)}}{\partial W_{ab}} &= \frac{\partial}{\partial W_{ab}} \sum_j W_{ij} x_j \\
&= \frac{\partial}{\partial W_{ab}} (W_{i0} x_0 + W_{i1} x_1 + \cdots + W_{ib} x_b + \cdots + W_{in} x_n) \\
&= \delta_{ia} x_b
\end{aligned} \tag{16}$$

$$\begin{aligned}
\frac{\partial f_i^{(2)}}{\partial f_j^{(1)}} &= \frac{\partial}{\partial f_j^{(1)}} (y_i - f_i^{(1)}) \\
&= -\delta_{ij}
\end{aligned} \tag{17}$$

$$\begin{aligned}
\frac{\partial f_i^{(3)}}{\partial f_i^{(2)}} &= \frac{\partial}{\partial f_i^{(2)}} \sum_i \left(f_i^{(2)}\right)^2 \\
&= \frac{\partial}{\partial f_i^{(2)}} \left(2f_0^{(2)} + 2f_1^{(2)} + \cdots + 2f_i^{(2)} + \cdots + 2f_n^{(2)}\right) \\
&= 2f_i^{(2)}
\end{aligned} \tag{18}$$

The following reverse propagation schedules are defined.

$$t^{(3)} = 1 \tag{19}$$

$$t_i^{(2)} = \frac{\partial f^{(3)}}{\partial f_i^{(2)}} t^{(3)} \tag{20}$$

$$t_i^{(1)} = \sum_j \frac{\partial f_j^{(2)}}{\partial f_i^{(1)}} t_j^{(2)} \tag{21}$$

$$t_{ab}^W = \sum_i \frac{\partial f_i^{(1)}}{\partial W_{ab}} t_i^{(1)} \tag{22}$$

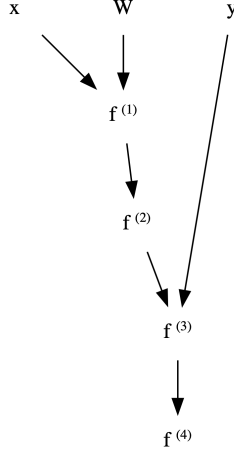


Figure 4: The computation graph of the loss function in Exercise 2.2

From which the required derivative can be calculated.

$$\begin{aligned}
\frac{\partial E}{\partial W_{ab}} &= t_{ab}^W \\
&= \sum_i \frac{\partial f_i^{(1)}}{\partial W_{ab}} t_i^{(1)} \\
&= \sum_i \delta_{ia} x_b t_i^{(1)} \\
&= x_b t_a^{(1)} \\
&= x_b \sum_j \frac{\partial f_j^{(2)}}{\partial f_a^{(1)}} t_j^{(2)} \\
&= x_b \sum_j -\delta_{ja} t_j^{(2)} \\
&= -x_b t_a^{(2)} \\
&= -x_b \frac{\partial f^{(3)}}{\partial f_a^{(2)}} t^{(3)} \\
&= -2x_b f_a^{(2)} \\
&= -2x_b \left( y_a - f_a^{(1)} \right) \\
&= -2x_b \left( y_a - \sum_i W_{ai} x_i \right)
\end{aligned} \tag{23}$$

## Exercise 2.2

An extended loss function (compared to the previous exercise) is given.

$$E = \sum_i \left( y_i - \phi \left( \sum_j W_{ij} x_j \right) \right)^2 \tag{24}$$

and the exercise requires do compute the derivative

$$\frac{\partial E}{\partial W_{ab}} \quad (25)$$

The computation graph is drawn (figure 4) the functions are derived

$$f_i^{(1)} = \sum_j W_{ij} x_j \quad (26)$$

$$f_i^{(2)} = \phi(f_i^{(1)}) \quad (27)$$

$$f_i^{(3)} = y_i - f_i^{(2)} \quad (28)$$

$$f^{(4)} = \sum_i \left( f_i^{(2)} \right)^2 \quad (29)$$

The required partial derivatives are calculated.

$$\frac{\partial f_i^{(1)}}{\partial W_{ab}} = \delta_{ia} x_b \quad (30)$$

$$\frac{\partial f_i^{(2)}}{\partial f_j^{(1)}} = \phi'(f_j^{(1)}) \delta_{ij} \quad (31)$$

$$\frac{\partial f_i^{(3)}}{\partial f_j^{(2)}} = -\delta_{ij} \quad (32)$$

$$\frac{\partial f^{(4)}}{\partial f_i^{(3)}} = 2f_i^{(3)} \quad (33)$$

The following reverse propagation schedules are defined.

$$t^{(4)} = 1 \quad (34)$$

$$t_i^{(3)} = \frac{\partial f^{(4)}}{\partial f_i^{(3)}} t^{(4)} \quad (35)$$

$$t_i^{(2)} = \sum_j \frac{\partial f_j^{(3)}}{\partial f_i^{(2)}} t_j^{(3)} \quad (36)$$

$$t_i^{(1)} = \sum_j \frac{\partial f_j^{(2)}}{\partial f_i^{(1)}} t_j^{(2)} \quad (37)$$

$$t_{ab}^W = \sum_i \frac{\partial f_i^{(1)}}{\partial W_{ab}} t_i^{(1)} \quad (38)$$

From which the required derivative can be calculated.

$$\begin{aligned}
\frac{\partial E}{\partial W_{ab}} &= t_{ab}^W \\
&= \sum_i \frac{\partial f_i^{(1)}}{\partial W_{ab}} t_i^{(1)} \\
&= \sum_i \delta_{ia} x_b t_i^{(1)} \\
&= x_b t_a^{(1)} \\
&= x_b \sum_j \frac{\partial f_j^{(2)}}{\partial f_a^{(1)}} t_j^{(2)} \\
&= x_b \sum_j \phi'(f_j^{(1)}) \delta_{aj} t_j^{(2)} \\
&= x_b \phi'(f_a^{(1)}) t_a^{(2)} \\
&= x_b \phi'(f_a^{(1)}) \sum_j \frac{\partial f_j^{(3)}}{\partial f_a^{(2)}} t_j^{(3)} \\
&= x_b \phi'(f_a^{(1)}) \sum_j -\delta_{ja} t_j^{(3)} \\
&= -x_b \phi'(f_a^{(1)}) t_a^{(3)} \\
&= -x_b \phi'(f_a^{(1)}) \frac{\partial f^{(4)}}{\partial f_a^{(3)}} t^{(4)} \\
&= -2f_a^{(3)} x_b \phi'(f_a^{(1)})
\end{aligned} \tag{39}$$

This result can be proved by differentiating traditionally by applying the chain rule.

$$\frac{dE}{dW_{ab}} = \sum_i \left( 2f_i^{(3)} \right) (-1) \left( \phi' \left( f_i^{(1)} \right) \right) (x_b \delta_{ia}) = -2f_a^{(3)} \phi' \left( f_a^{(1)} \right) x_b \tag{40}$$

### Exercise 2.3

The following loss function is given.

$$E = \sum_i \left( y_i - \phi_2 \left( \sum_k W_{ik}^{(2)} \phi_1 \left( \sum_j W_{kj}^{(1)} x_j \right) \right) \right)^2 \tag{41}$$

The computation graph is drawn (figure 5) The functions are defined.

$$f_k^{(1)} = \sum_j W_{kj}^{(1)} x_j \tag{42}$$

$$f_k^{(2)} = \phi_1(f_k^{(1)}) \tag{43}$$

$$f_i^{(3)} = \sum_k W_{ik}^{(2)} f_k^{(2)} \tag{44}$$

$$f_i^{(4)} = \phi_2(f_i^{(3)}) \tag{45}$$

$$f_i^{(5)} = y_i - f_i^{(4)} \tag{46}$$

$$f^{(6)} = E = \sum_i \left( f_i^{(5)} \right)^2 \tag{47}$$



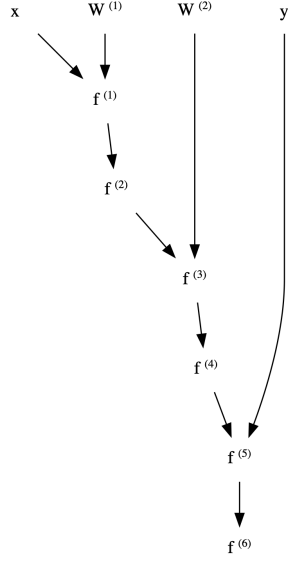


Figure 5: The computation graph of the loss function in Exercise 2.3

The partial derivatives are calculated.

$$\frac{\partial f_i^{(1)}}{\partial W_{ab}^{(1)}} = x_b \delta_{ia} \quad (48)$$

$$\frac{\partial f_i^{(2)}}{\partial f_j^{(1)}} = \phi_1' \left( f_i^{(1)} \right) \delta_{ij} \quad (49)$$

$$\frac{\partial f_i^{(3)}}{\partial f_j^{(2)}} = W_{ij}^{(2)} \quad (50)$$

$$\frac{\partial f_i^{(3)}}{\partial W_{ab}^{(2)}} = f_b^{(2)} \delta_{ia} \quad (51)$$

$$\frac{\partial f_i^{(4)}}{\partial f_j^{(3)}} = \phi_2' \left( f_i^{(3)} \right) \delta_{ij} \quad (52)$$

$$\frac{\partial f_i^{(5)}}{\partial f_j^{(4)}} = -\delta_{ij} \quad (53)$$

$$\frac{\partial f^{(6)}}{\partial f_i^{(5)}} = 2f_i^{(5)} \quad (54)$$

The propagation schedule is determined and messages calculated.

$$t^{(6)} = 1 \quad (55)$$

$$t_i^{(5)} = \frac{\partial f^{(6)}}{\partial f_i^{(5)}} t^{(6)} = 2f_i^{(5)} t^{(6)} \quad (56)$$

$$t_i^{(4)} = \sum_j \frac{\partial f_j^{(5)}}{\partial f_i^{(4)}} t_j^{(5)} = \sum_j -\delta_{ji} t_j^{(5)} = -t_i^{(5)} \quad (57)$$

$$t_i^{(3)} = \sum_j \frac{\partial f_j^{(4)}}{\partial f_i^{(3)}} t_j^{(4)} = \sum_j \phi'_2 \left( f_j^{(3)} \right) \delta_{ji} t_j^{(4)} = \phi'_2 \left( f_i^{(3)} \right) t_i^{(4)} \quad (58)$$

$$t_i^{(2)} = \sum_j \frac{\partial f_j^{(3)}}{\partial f_i^{(2)}} t_j^{(3)} = \sum_j W_{ji}^{(2)} t_j^{(3)} \quad (59)$$

$$t_i^{(1)} = \sum_j \frac{\partial f_j^{(2)}}{\partial f_i^{(1)}} t_j^{(2)} = \sum_j \phi'_1 \left( f_j^{(1)} \right) \delta_{ji} t_j^{(2)} = \phi'_1 \left( f_i^{(1)} \right) t_i^{(2)} \quad (60)$$

$$t_{ab}^{W^{(1)}} = \sum_i \frac{\partial f_i^{(1)}}{\partial W_{ab}^{(1)}} t_i^{(1)} = \sum_i x_b \delta_{ia} t_i^{(1)} = x_b t_a^{(1)} \quad (61)$$

$$t_{ab}^{W^{(2)}} = \sum_i \frac{\partial f_i^{(3)}}{\partial W_{ab}^{(2)}} t_i^{(1)} = \sum_i f_b^{(2)} \delta_{ia} t_i^{(3)} = f_b^{(2)} t_a^{(3)} \quad (62)$$

The required derivatives are calculated by joining messages.

$$\begin{aligned} \frac{dE}{dW_{ab}^{(1)}} &= t_{ab}^{W^{(1)}} \\ &= x_b t_a^{(1)} \\ &= x_b \phi'_1 \left( f_a^{(1)} \right) t_a^{(2)} \\ &= x_b \phi'_1 \left( f_a^{(1)} \right) \sum_j W_{ja}^{(2)} t_j^{(3)} \\ &= x_b \phi'_1 \left( f_a^{(1)} \right) \sum_j W_{ja}^{(2)} \phi'_2 \left( f_j^{(3)} \right) t_j^{(4)} \\ &= -x_b \phi'_1 \left( f_a^{(1)} \right) \sum_j W_{ja}^{(2)} \phi'_2 \left( f_j^{(3)} \right) t_j^{(5)} \\ &= -2x_b \phi'_1 \left( f_a^{(1)} \right) \sum_j W_{ja}^{(2)} \phi'_2 \left( f_j^{(3)} \right) f_j^{(5)} \end{aligned} \quad (63)$$

$$\begin{aligned} \frac{dE}{dW_{ab}^{(2)}} &= t_{ab}^{W^{(2)}} \\ &= f_b^{(2)} t_a^{(3)} \\ &= f_b^{(2)} \phi'_2 \left( f_a^{(3)} \right) t_a^{(4)} \\ &= -f_b^{(2)} \phi'_2 \left( f_a^{(3)} \right) t_a^{(5)} \\ &= -2f_b^{(2)} \phi'_2 \left( f_a^{(3)} \right) f_a^{(5)} \end{aligned} \quad (64)$$

These results are proved by differentiating in the traditional way with the chain rule.

$$\begin{aligned} \frac{dE}{dW_{ab}^{(1)}} &= \sum_i \left( 2f_i^{(5)} \right) (-1) \left( \phi'_2 \left( f_i^{(3)} \right) \right) \sum_k W_{ik}^{(2)} \left( \phi'_1 \left( f_k^{(1)} \right) \right) (x_b \delta_{ka}) \\ &= -2x_b \phi'_1 \left( f_a^{(1)} \right) \sum_i f_i^{(5)} \phi'_2 \left( f_i^{(3)} \right) W_{ia}^{(2)} \end{aligned} \quad (65)$$

$$\begin{aligned} \frac{dE}{dW_{ab}^{(2)}} &= \sum_i \left( 2f_i^{(5)} \right) (-1) \left( \phi'_2 \left( f_i^{(3)} \right) \right) f_b^{(2)} \delta_{ia} \\ &= -2f_a^{(5)} \phi'_2 \left( f_a^{(3)} \right) f_b^{(2)} \end{aligned} \quad (66)$$

## Exercise 2.4

The loss functions observed in the previous exercises compute the loss for a given unit of data. This question, instead, requires to extend this scheme for all the units of data. In other words, rather than computing the loss for a single pair  $(x, y)$ , it requires to do so for the  $n$  pairs  $(x^n, y^n)$ . Consequently, the loss function analysed in the previous exercise is extended by introducing a sum for each pair and the actual loss sees each occurrence  $y$  and  $x$  being replaced by  $y^n$  and  $x^n$ , respectively. This is equivalent to computing the loss seen in the previous exercise for each unit of data and summing up the individual losses.

$$E = \sum_n \sum_i \left( y_i^n - \phi_2 \left( \sum_k W_{ik}^{(2)} \phi_1 \left( \sum_j W_{kj}^{(1)} x_j^n \right) \right) \right)^2 \quad (67)$$

Consequently, the individual functions are redefined to account for this change: all function previously defined will take an additional parameter  $n$  which is passed down by the newly introduced function  $f^{(7)} = E = \sum_n f_n^{(6)}$ . Partial derivatives which are both passing down the unit index are updated accordingly, that is, a bind  $\delta_{mn}$  is introduced for each derivative to ensure that within the summation for  $n$  derivation is being carried out with respect to the same unit of data; the partial derivative  $\frac{\partial f^{(7)}}{\partial f_n^{(6)}} = 1$  is added. The schedule is updated to reflect the introduction of the new parameter. Eventually, the two required derivatives are eventually calculated.

$$\frac{dE}{dW_{ab}^{(1)}} = -2 \sum_n x_b^n \phi_1' \left( f_{na}^{(1)} \right) \sum_i f_{ni}^{(5)} \phi_2' \left( f_{ni}^{(3)} \right) W_{ia}^{(2)} \quad (68)$$

$$\frac{dE}{dW_{ab}^{(2)}} = -2 \sum_n f_{na}^{(5)} \phi_2' \left( f_{na}^{(3)} \right) f_{nb}^{(2)} \quad (69)$$

## Exercise 2.5

This exercise is similar to exercise 2.3 in that the two weights matrices are identical, such that  $W^{(1)} = W^{(2)} = W$ . Consequently, the loss function becomes

$$E = \sum_i \left( y_i - \phi_2 \left( \sum_k W_{ik} \phi_1 \left( \sum_j W_{kj} x_j \right) \right) \right)^2 \quad (70)$$

The computation graph is drawn (figure 6) the individual functions are updated.

$$f_k^{(1)} = \sum_j W_{kj} x_j \quad (71)$$

$$f_k^{(2)} = \phi_1(f_k^{(1)}) \quad (72)$$

$$f_i^{(3)} = \sum_k W_{ik} f_k^{(2)} \quad (73)$$

$$f_i^{(4)} = \phi_2(f_i^{(3)}) \quad (74)$$

$$f_i^{(5)} = y_i - f_i^{(4)} \quad (75)$$

$$f^{(6)} = E = \sum_i \left( f_i^{(5)} \right)^2 \quad (76)$$

Partial derivatives are updated as well.

$$\frac{\partial f_i^{(1)}}{\partial W_{ab}} = x_b \delta_{ia} \quad (77)$$

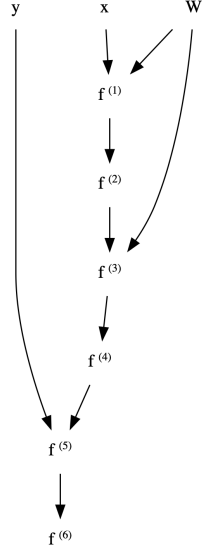


Figure 6: The computation graph of the loss function in Exercise 2.5

$$\frac{\partial f_i^{(2)}}{\partial f_j^{(1)}} = \phi'_1 \left( f_i^{(1)} \right) \delta_{ij} \quad (78)$$

$$\frac{\partial f_i^{(3)}}{\partial f_j^{(2)}} = W_{ij} \quad (79)$$

$$\frac{\partial f_i^{(3)}}{\partial W_{ab}} = f_b^{(2)} \delta_{ia} \quad (80)$$

$$\frac{\partial f_i^{(4)}}{\partial f_j^{(3)}} = \phi'_2 \left( f_i^{(3)} \right) \delta_{ij} \quad (81)$$

$$\frac{\partial f_i^{(5)}}{\partial f_j^{(4)}} = -\delta_{ij} \quad (82)$$

$$\frac{\partial f^{(6)}}{\partial f_i^{(5)}} = 2f_i^{(5)} \quad (83)$$

Reverse schedules and messages are updated too.

$$t^{(6)} = 1 \quad (84)$$

$$t_i^{(5)} = \frac{\partial f^{(6)}}{\partial f_i^{(5)}} t^{(6)} = 2f_i^{(5)} t^{(6)} \quad (85)$$

$$t_i^{(4)} = \sum_j \frac{\partial f_j^{(5)}}{\partial f_i^{(4)}} t_j^{(5)} = \sum_j -\delta_{ji} t_j^{(5)} = -t_i^{(5)} \quad (86)$$

$$t_i^{(3)} = \sum_j \frac{\partial f_j^{(4)}}{\partial f_i^{(3)}} t_j^{(4)} = \sum_j \phi'_2 \left( f_j^{(3)} \right) \delta_{ji} t_j^{(4)} = \phi'_2 \left( f_i^{(3)} \right) t_i^{(4)} \quad (87)$$

$$t_i^{(2)} = \sum_j \frac{\partial f_j^{(3)}}{\partial f_i^{(2)}} t_j^{(3)} = \sum_j W_{ji} t_j^{(3)} \quad (88)$$

$$t_i^{(1)} = \sum_j \frac{\partial f_j^{(2)}}{\partial f_i^{(1)}} t_j^{(2)} = \sum_j \phi'_1 \left( f_j^{(1)} \right) \delta_{ji} t_j^{(2)} = \phi'_1 \left( f_i^{(1)} \right) t_i^{(2)} \quad (89)$$

$$t_{ab}^W = \sum_i \frac{\partial f_i^{(1)}}{\partial W_{ab}} t_i^{(1)} + \sum_j \frac{\partial f_j^{(3)}}{\partial W_{ab}} t_j^{(3)} \quad (90)$$

The main change is observed in  $t_{ab}^W$  as it is updated to account for both *paths*, as seen in the computation graph. The required derivative is then calculated with the new schedule.

$$\begin{aligned} \frac{dE}{dW_{ab}} &= t_{ab}^W \\ &= \sum_i \frac{\partial f_i^{(1)}}{\partial W_{ab}} t_i^{(1)} + \sum_k \frac{\partial f_k^{(3)}}{\partial W_{ab}} t_k^{(3)} \\ &= \left( -2x_b \phi'_1 \left( f_a^{(1)} \right) \sum_j W_{ja}^{(2)} \phi'_2 \left( f_j^{(3)} \right) f_j^{(5)} \right) + \left( -2f_a^{(5)} \phi'_2 \left( f_a^{(3)} \right) f_b^{(2)} \right) \\ &= -2 \left( f_a^{(5)} \phi'_2 \left( f_a^{(3)} \right) f_b^{(2)} + x_b \phi'_1 \left( f_a^{(1)} \right) \sum_j W_{ja}^{(2)} \phi'_2 \left( f_j^{(3)} \right) f_j^{(5)} \right) \end{aligned} \quad (91)$$

This is confirmed by differentiating traditionally by applying the chain rule.

$$\begin{aligned} \frac{dE}{dW_{ab}} &= \left( \sum_i \left( 2f_i^{(5)} \right) (-1) \left( \phi'_2 \left( f_i^{(3)} \right) \right) \sum_k W_{ik} \left( \phi'_1 \left( f_k^{(1)} \right) \right) (x_b \delta_{ka}) \right) + \\ &+ \left( \sum_j \left( 2f_j^{(5)} \right) (-1) \left( \phi'_2 \left( f_j^{(3)} \right) \right) f_b^{(2)} \delta_{ia} \right) \\ &= -2f_a^{(5)} \phi'_2 \left( f_a^{(3)} \right) f_b^{(2)} - 2x_b \phi'_1 \left( f_a^{(1)} \right) \sum_i f_i^{(5)} \phi'_2 \left( f_i^{(3)} \right) W_{ia}^{(2)} \\ &= -2 \left( f_a^{(5)} \phi'_2 \left( f_a^{(3)} \right) f_b^{(2)} + x_b \phi'_1 \left( f_a^{(1)} \right) \sum_i f_i^{(5)} \phi'_2 \left( f_i^{(3)} \right) W_{ia}^{(2)} \right) \end{aligned} \quad (92)$$

### Exercise 3

A neural network with a single hidden layer is given, representing an autoencoder for the  $28 \times 28$  MNIST images.

$$f_i(x) = \sigma \left( b_i^{(2)} + \sum_{k=1}^H W_{ik}^{(2)} \phi \left( b_k^{(1)} + \sum_{j=1}^{784} W_{kj}^{(1)} x_j \right) \right) \quad (93)$$

The error function  $E$  is defined for a generic loss function  $L$ .

$$E = \frac{1}{N} \sum_n \sum_i^{784} L(x_i^n, f_i(x^n)) \quad (94)$$

### Exercises 3.2 and 3.3

The individual functions are defined and the related computation graph is drawn (figure 7).

$$f_{nk}^{(1)} = \sum_{j=1}^{784} W_{kj}^{(1)} x_j^n \quad (95)$$

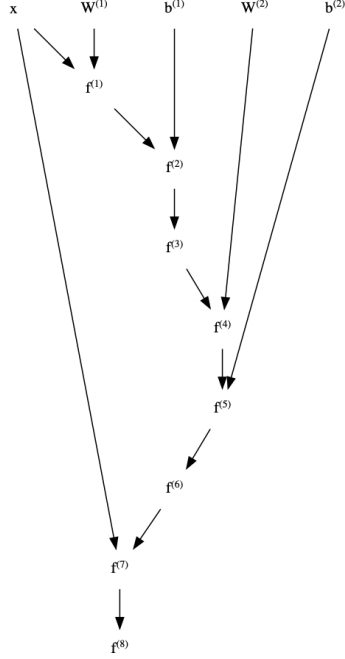


Figure 7: The computation graph of the loss function in Exercise 3.2

$$f_{nk}^{(2)} = b_k^{(1)} + f_{nk}^{(1)} \quad (96)$$

$$f_{nk}^{(3)} = \phi \left( f_{nk}^{(2)} \right) \quad (97)$$

$$f_{ni}^{(4)} = \sum_{k=1}^H W_{ik}^{(2)} f_{nk}^{(3)} \quad (98)$$

$$f_{ni}^{(5)} = b_i^{(2)} + f_{ni}^{(4)} \quad (99)$$

$$f_{ni}^{(6)} = \sigma \left( f_{ni}^{(5)} \right) \quad (100)$$

$$f_{ni}^{(7)} = L \left( x_i^n, f_{ni}^{(6)} \right) \quad (101)$$

$$f^{(8)} = \frac{1}{N} \sum_n \sum_i f_{ni}^{(7)} \quad (102)$$

The partial derivatives are calculated.

$$\frac{\partial f_{ni}^{(1)}}{\partial W_{ab}^{(1)}} = \delta_{ia} x_b \quad (103)$$

$$\frac{\partial f_{ni}^{(2)}}{\partial f_{mj}^{(1)}} = \delta_{ij} \delta_{mn} \quad (104)$$

$$\frac{\partial f_{ni}^{(2)}}{\partial b_a^{(1)}} = \delta_{ia} \quad (105)$$

$$\frac{\partial f_{ni}^{(3)}}{\partial f_{mj}^{(2)}} = \phi'(f_{ni}^{(2)}) \delta_{ij} \delta_{mn} \quad (106)$$

$$\frac{\partial f_{ni}^{(4)}}{\partial f_{mj}^{(3)}} = W_{ij}^{(2)} \delta_{mn} \quad (107)$$

$$\frac{\partial f_{ni}^{(4)}}{\partial W_{ab}^{(2)}} = \delta_{ia} f_{nb}^{(3)} \quad (108)$$

$$\frac{\partial f_{ni}^{(5)}}{\partial f_{mj}^{(4)}} = \delta_{ij} \delta_{mn} \quad (109)$$

$$\frac{\partial f_{ni}^{(5)}}{\partial b_a^{(2)}} = \delta_{ia} \quad (110)$$

$$\frac{\partial f_{ni}^{(6)}}{\partial f_{mj}^{(5)}} = \sigma' \left( f_{ni}^{(5)} \right) \delta_{ij} \delta_{mn} \quad (111)$$

$$\frac{\partial f_{ni}^{(7)}}{\partial f_{mj}^{(6)}} = L' \left( x_i^n, f_{ni}^{(6)} \right) \delta_{ij} \delta_{mn} \quad (112)$$

$$\frac{\partial f_{ni}^{(8)}}{\partial f_{ni}^{(7)}} = \frac{1}{N} \quad (113)$$

The reverse schedule is determined along with the messages.

$$t^{(8)} = 1 \quad (114)$$

$$t_{ni}^{(7)} = \frac{\partial f^{(8)}}{\partial f_n^{(7)} i} t^{(8)} = \frac{1}{N} t^{(8)} \quad (115)$$

$$t_{ni}^{(6)} = \sum_m \sum_j \frac{\partial f_{mj}^{(7)}}{\partial f_n^{(6)} i} t_{mj}^{(7)} = \sum_m \sum_j L' \left( x_j^m, f_{mj}^{(6)} \right) \delta_{ij} \delta_{mn} t_{mj}^{(7)} = L' \left( x_i^n, f_{ni}^{(6)} \right) t_{ni}^{(7)} \quad (116)$$

$$t_{ni}^{(5)} = \sum_m \sum_j \frac{\partial f_{mj}^{(6)}}{\partial f_n^{(5)} i} t_{mj}^{(6)} = \sum_m \sum_j \sigma' \left( f_{mj}^{(5)} \right) \delta_{ij} \delta_{mn} t_{mj}^{(6)} = \sigma' \left( f_{ni}^{(5)} \right) t_{ni}^{(6)} \quad (117)$$

$$t_a^{b(2)} = \sum_m \sum_j \frac{\partial f_{mj}^{(5)}}{\partial b_a^{(2)}} t_{mj}^{(5)} = \sum_m \sum_j \delta_{ja} t_{mj}^{(5)} = \sum_m t_{ma}^{(5)} \quad (118)$$

$$t_{ni}^{(4)} = \sum_m \sum_j \frac{\partial f_{mj}^{(5)}}{\partial f_n^{(4)} i} t_{mj}^{(5)} = \sum_m \sum_j \delta_{ij} \delta_{mn} t_{mj}^{(5)} = t_{ni}^{(5)} \quad (119)$$

$$t_{ab}^{W(2)} = \sum_n \sum_i \frac{\partial f_{ni}^{(4)}}{\partial W_{ab}^{(2)}} t_{ni}^{(4)} = \sum_n \sum_i \delta_{ia} f_{nb}^{(3)} t_{ni}^{(4)} = \sum_n f_{nb}^{(3)} t_{na}^{(4)} \quad (120)$$

$$t_{ni}^{(3)} = \sum_m \sum_j \frac{\partial f_{mj}^{(4)}}{\partial f_n^{(3)} i} t_{mj}^{(4)} = \sum_m \sum_j W_{ji}^{(2)} \delta_{mn} t_{mj}^{(4)} = \sum_j W_{ji}^{(2)} t_{nj}^{(4)} \quad (121)$$

$$t_{ni}^{(2)} = \sum_m \sum_j \frac{\partial f_{mj}^{(3)}}{\partial f_n^{(2)} i} t_{mj}^{(3)} = \sum_m \sum_j \phi' \left( f_{mj}^{(2)} \right) \delta_{ij} \delta_{mn} t_{mj}^{(3)} = \phi' \left( f_{ni}^{(2)} \right) t_{ni}^{(3)} \quad (122)$$

$$t_{ni}^{(1)} = \sum_m \sum_j \frac{\partial f_{mj}^{(2)}}{\partial f_n^{(1)} i} t_{mj}^{(2)} = \sum_m \sum_j \delta_{ij} \delta_{mn} t_{mj}^{(2)} = t_{ni}^{(2)} \quad (123)$$

$$t_a^{b(1)} = \sum_m \sum_j \frac{\partial f_{mj}^{(2)}}{\partial b_a^{(1)}} t_{mj}^{(2)} = \sum_m \sum_j \delta_{ja} t_{mj}^{(2)} = \sum_m t_{ma}^{(2)} \quad (124)$$

$$t_{ab}^{W(1)} = \sum_m \sum_j \frac{\partial f_{mj}^{(1)}}{\partial W_{ab}^{(1)}} t_{mj}^{(1)} = \sum_m \sum_j x_b \delta_{ja} t_{mj}^{(2)} = \sum_m x_b t_{ma}^{(1)} \quad (125)$$

The required derivatives are then computed using the schedule and messages.

$$\begin{aligned} \frac{dE}{dW_{ab}^{(1)}} &= t_{ab}^{W(1)} \\ &= \sum_n x_b t_{na}^{(1)} \\ &= \sum_n x_b t_{na}^{(2)} \\ &= \sum_n x_b \phi' \left( f_{na}^{(2)} \right) t_{na}^{(3)} \\ &= \sum_n x_b \phi' \left( f_{na}^{(2)} \right) \sum_i W_{ia}^{(2)} t_{ni}^{(4)} \\ &= \sum_n x_b \phi' \left( f_{na}^{(2)} \right) \sum_i W_{ia}^{(2)} t_{ni}^{(5)} \\ &= \sum_n x_b \phi' \left( f_{na}^{(2)} \right) \sum_i W_{ia}^{(2)} \sigma' \left( f_{ni}^{(5)} \right) t_{ni}^{(6)} \\ &= \sum_n x_b \phi' \left( f_{na}^{(2)} \right) \sum_i W_{ia}^{(2)} \sigma' \left( f_{ni}^{(5)} \right) L' \left( x_i^n, f_{ni}^{(6)} \right) t_{ni}^{(7)} \\ &= \frac{1}{N} \sum_n x_b \phi' \left( f_{na}^{(2)} \right) \sum_i W_{ia}^{(2)} \sigma' \left( f_{ni}^{(5)} \right) L' \left( x_i^n, f_{ni}^{(6)} \right) \end{aligned} \quad (126)$$

$$\begin{aligned} \frac{dE}{db_a^{(1)}} &= t_a^{b(1)} \\ &= \sum_n t_{na}^{(2)} \\ &= \sum_n \phi' \left( f_{na}^{(2)} \right) t_{na}^{(3)} \\ &= \sum_n \phi' \left( f_{na}^{(2)} \right) \sum_i W_{ia}^{(2)} t_{ni}^{(4)} \\ &= \sum_n \phi' \left( f_{na}^{(2)} \right) \sum_i W_{ia}^{(2)} t_{ni}^{(5)} \\ &= \sum_n \phi' \left( f_{na}^{(2)} \right) \sum_i W_{ia}^{(2)} \sigma' \left( f_{ni}^{(5)} \right) L' \left( x_i^n, f_{ni}^{(6)} \right) t_{ni}^{(7)} \\ &= \frac{1}{N} \sum_n \phi' \left( f_{na}^{(2)} \right) \sum_i W_{ia}^{(2)} \sigma' \left( f_{ni}^{(5)} \right) L' \left( x_i^n, f_{ni}^{(6)} \right) \end{aligned} \quad (127)$$



$$\begin{aligned}
\frac{dE}{dW_{ab}^{(2)}} &= t_{ab}^{W^{(2)}} \\
&= \sum_n f_{nb}^{(3)} t_{na}^{(4)} \\
&= \sum_n f_{nb}^{(3)} t_{na}^{(5)} \\
&= \sum_n f_{nb}^{(3)} \sigma' \left( f_{na}^{(5)} \right) t_{na}^{(6)} \\
&= \sum_n f_{nb}^{(3)} \sigma' \left( f_{na}^{(5)} \right) L' \left( x_a^n, f_{na}^{(6)} \right) t_{na}^{(7)} \\
&= \frac{1}{N} \sum_n f_{nb}^{(3)} \sigma' \left( f_{na}^{(5)} \right) L' \left( x_a^n, f_{na}^{(6)} \right)
\end{aligned} \tag{128}$$

$$\begin{aligned}
\frac{dE}{db_a^{(2)}} &= t_a^{b^{(2)}} \\
&= \sum_n t_{na}^{(5)} \\
&= \sum_n \sigma' \left( f_{na}^{(5)} \right) t_{na}^{(6)} \\
&= \sum_n \sigma' \left( f_{na}^{(5)} \right) L' \left( x_a^n, f_{na}^{(6)} \right) t_{na}^{(7)} \\
&= \frac{1}{N} \sum_n \sigma' \left( f_{na}^{(5)} \right) L' \left( x_a^n, f_{na}^{(6)} \right)
\end{aligned} \tag{129}$$

### Exercise 3.4 and 3.5

For this question, Python was used to implement the ADAM optimiser, along with the previously defined AutoDiff schedule. Different transfer functions were tested and the relu was chosen; initial weights are picked randomly from a zero-mean unit-variance normal distribution. Learning rate was set to 0.05, while the three ADAM parameters were tuned as suggested. The error decrease can be observed in figure 8; results of reconstruction with the autoencoder can be observed in figure 9.

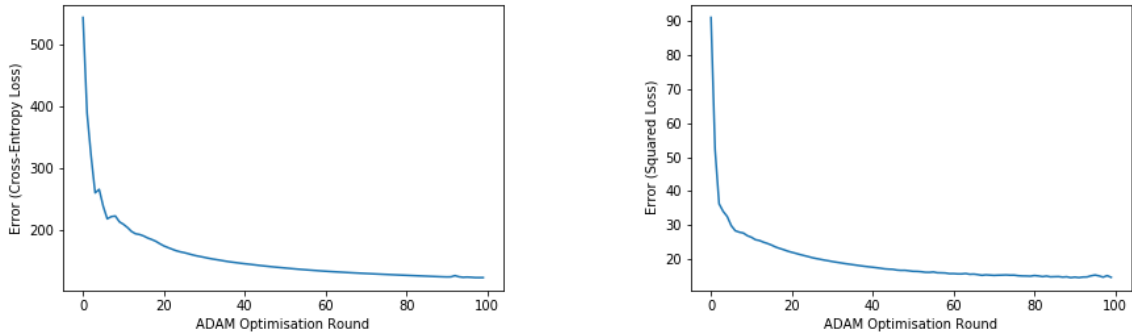


Figure 8: Training error at each ADAM optimisation round with Square Loss (left) and Cross Entropy Loss (right).

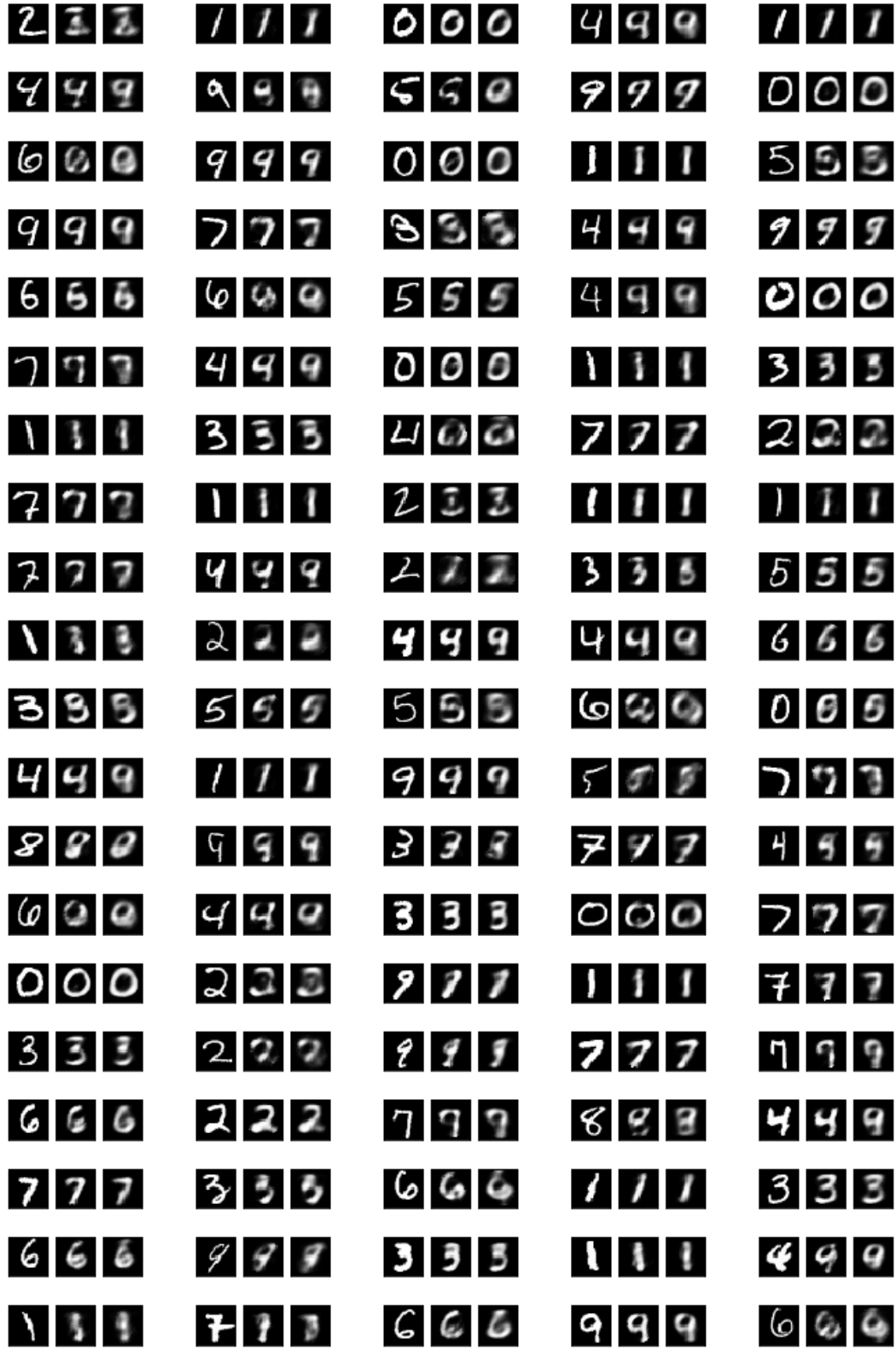


Figure 9: The reconstructions of 100 MNIST digits; each triplet contains (from left to right) the original image, the reconstruction using squared loss and cross-entropy loss.

### Exercise 3.6

This question requires to determine whether the error function is convex with respect to  $W^2$  and  $b^2$  for both the LikLoss and SquareLoss. As usual, convexity is determined by whether the second derivative is always positive. For convenience, throughout this exercise,  $\sigma$  is used to represent  $f_{na}^{(6)} = \sigma(f_{na}^{(5)})$  and  $L$  to represent  $L(x_a^n, f_{na}^{(6)})$ ; the single apex  $'$  and double apex  $''$  represent the first and second derivative of the function they follow with respect to either of the two variables of interest. The loss functions and their derivatives with respect to the second parameter  $y$  (which will be the one of interest in the calculations) are

$$L_{lik}(x, y) = -x \log(y) - (1 - x) \log(1 - y) \quad (130)$$

$$L'_{lik}(x, y) = \frac{y - x}{y(1 - y)} \quad (131)$$

$$L''_{lik}(x, y) = \frac{y^2 + x - 2xy}{y^2(1 - y)^2} \quad (132)$$

$$L_{sq}(x, y) = \frac{1}{2}(x - y)^2 \quad (133)$$

$$L'_{sq}(x, y) = y - x \quad (134)$$

$$L''_{sq}(x, y) = 1 \quad (135)$$

With regard to  $b^{(2)}$ , the first and second derivatives are:

$$\frac{\partial E}{\partial b_a^{(2)}} = \frac{1}{N} \sum_n \sigma' \left( f_{na}^{(5)} \right) L' \left( x_a^n, f_{na}^{(6)} \right) = \frac{1}{N} \sum_n \sigma' L' \quad (136)$$

$$\begin{aligned} \frac{\partial^2 E}{\partial b_a^{(2)2}} &= \frac{1}{N} \sum_n \sigma' L'' \sigma' + \sigma'' L' \\ &= \frac{1}{N} \sum_n \sigma^2 (1 - \sigma)^2 L'' + \sigma(1 - \sigma)(1 - 2\sigma) L' \\ &= \frac{1}{N} \sum_n \sigma(1 - \sigma) [\sigma(1 - \sigma) L'' + (1 - 2\sigma) L'] \end{aligned} \quad (137)$$

It is known that  $\sigma(1 - \sigma)$  is never negative, hence the factor to assess for positivity is  $\sigma(1 - \sigma) L'' + (1 - 2\sigma) L'$ . For the LikLoss:

$$\sigma(1 - \sigma) \frac{\sigma^2 + x - 2\sigma x}{\sigma^2(1 - \sigma)^2} + (1 - 2\sigma) \frac{\sigma - x}{\sigma(1 - \sigma)} = \frac{x + \sigma^2 - 2\sigma x + \sigma - x - 2\sigma^2 + 2\sigma x}{\sigma(1 - \sigma)} = 1 \quad (138)$$

This proves that the second derivative of the error function with respect to  $b^{(2)}$  using the LikLoss is always positive, proving convexity. Similarly, for the SquareLoss:

$$\sigma(1 - \sigma) + (1 - 2\sigma)(\sigma - x) = \sigma - \sigma^2 + \sigma - x - 2\sigma^2 + 2\sigma x = -3\sigma^2 + \sigma(2 + 2x) - x \quad (139)$$

For the error function to be convex, this factor has to be positive for each value of  $x$ , but that is not the case as the function can be negative (e.g. for  $\sigma = 0.75, x = 0$  the factor is equal to  $-0.1875$ ). Due to this inconsistency in the sign of the factor, it is determined that the error function is not convex with respect to  $b^{(2)}$  when the SquareLoss is used.

The process is now repeated for  $W^{(2)}$ . The derivative of the error function with respect to  $W_{ab}^{(2)}$  is:

$$\frac{\partial E}{\partial W_{ab}^{(2)}} = \frac{1}{N} \sum_n f_{nb}^{(3)} \sigma' \left( f_{na}^{(5)} \right) L' \left( x_a^n, f_{na}^{(6)} \right) \quad (140)$$

The factor  $f^{(3)}$  does not depend on  $W^{(2)}$  and can therefore be treated as a constant. The second derivative is:

$$\begin{aligned}
\frac{\partial^2 E}{\partial W_{ab}^{(2)2}} &= \frac{1}{N} \sum_n f_{nb}^{(3)} \left[ \sigma' \left( f_{na}^{(5)} \right) L'' \left( x_a^n, f_{na}^{(6)} \right) \sigma' \left( f_{na}^{(5)} \right) f_{nb}^{(3)} + L' \left( x_a^n, f_{na}^{(6)} \right) \sigma'' \left( f_{na}^{(5)} \right) f_{nb}^{(3)} \right] \\
&= \frac{1}{N} \sum_n \left( f_{nb}^{(3)} \right)^2 \left[ \left( \sigma' \left( f_{na}^{(5)} \right) \right)^2 L'' \left( x_a^n, f_{na}^{(6)} \right) + L' \left( x_a^n, f_{na}^{(6)} \right) \sigma'' \left( f_{na}^{(5)} \right) \right] \\
&= \frac{1}{N} \sum_n \left( f_{nb}^{(3)} \right)^2 \left[ (\sigma')^2 L'' + \sigma'' L' \right] \\
&= \frac{1}{N} \sum_n \left( f_{nb}^{(3)} \right)^2 \left[ \sigma^2 (1 - \sigma)^2 L'' + \sigma (1 - \sigma) (1 - 2\sigma) L' \right] \\
&= \frac{1}{N} \sum_n \left( f_{nb}^{(3)} \right)^2 \sigma (1 - \sigma) \left[ \sigma (1 - \sigma) L'' + (1 - 2\sigma) L' \right]
\end{aligned} \tag{141}$$

The only factor to assess is once again  $\sigma(1 - \sigma)L'' + \sigma(1 - 2\sigma)L'$ . The conclusions on its sign are similar to the ones in the first half of this exercise. This means that for the LikLoss, the error function is concave with respect to  $W^{(2)}$ , but it is not for the SquareLoss.

### Exercise 3.7 and 3.8

In these questions, slices of the error function with the Square Loss first and Cross-Entropy Loss then are plotted to observe whether they are convex or not. In particular, the functions plotted are

$$F_{11}(z_1, z_2) = E(W^{(1)} + z_1 U^{(1)} + z_2^{(2)} U^{(2)}, b^{(1)}, W^{(2)}, b^{(2)}) \tag{142}$$

$$F_{12}(z_1, z_2) = E(W^{(1)} + z_1 U^{(1)}, b^{(1)}, W^{(2)} + z_2 U^{(2)}, b^{(2)}) \tag{143}$$

$$F_{22}(z_1, z_2) = E(W^{(1)}, b^{(1)}, W^{(2)} + z_1 U^{(1)} + z_2^{(2)} U^{(2)}, b^{(2)}) \tag{144}$$

for both loss functions, results are in figures 10, 11 and 12. It can be observed that with cross-entropy loss, the error function appears to be convex, whereas with square loss, the error function is not convex.

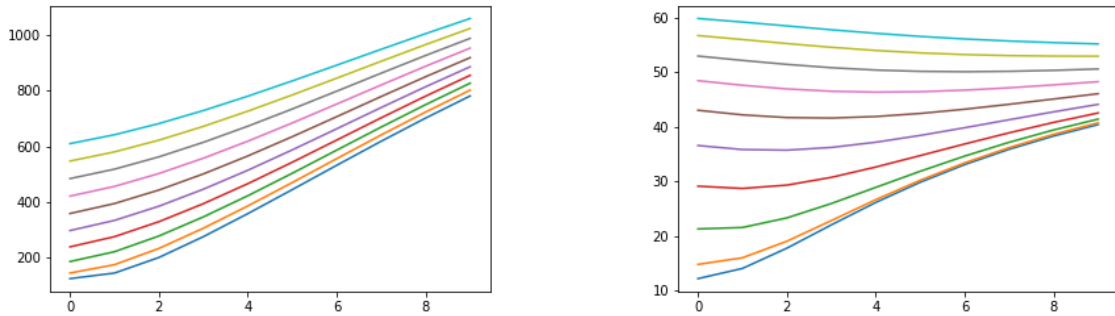


Figure 10:  $F_{1,1}$  for Cross-Entropy Loss (left) and Square Loss (Right)

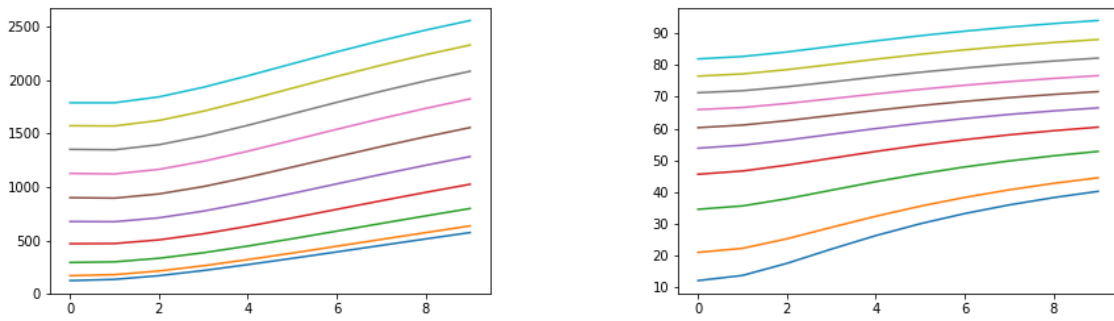


Figure 11:  $F_{1,2}$  for Cross-Entropy Loss (left) and Square Loss (Right)

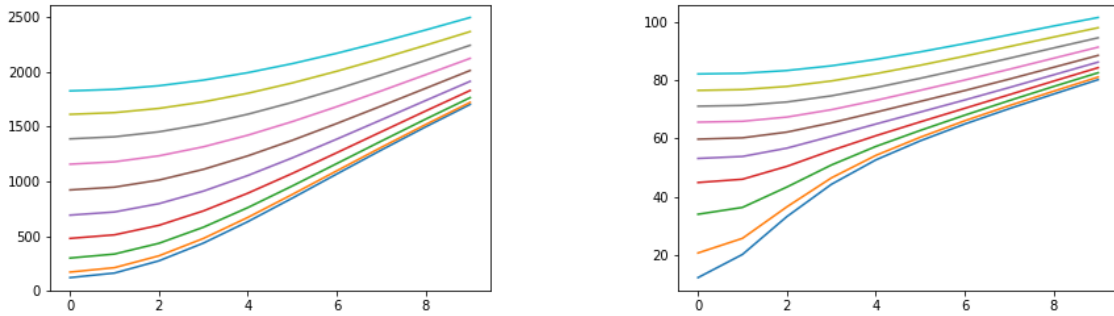


Figure 12:  $F_{2,2}$  for Cross-Entropy Loss (left) and Square Loss (Right)

## Code

### MNIST Autoencoder

The following code was used in Question 3.

```
def sigmoid(x):
    return 1 / (1 + np.exp(0 - x))

def relu(x):
    return np.maximum(x, 0)

def d_relu(x):
    return np.where(x > 0, np.full(np.shape(x), 1), np.full(np.shape(x), 0))

def squared_loss(y, x):
    return (y - x) ** 2 / 2

def d_squared_loss(y, x):
    return -(y - x)

def cross_entropy_loss(x, y):
    return -x*np.log(y) - (1-x)*np.log(1-y)
```

```

def d_cross_entropy_loss(x,y):
    return ((1-x)/(1-y)) - (x/y)

def create_weight(shape):
    stdv = 1. / np.sqrt(shape[-1])
    return np.random.uniform(-stdv, stdv, shape).astype(np.float32)

def autodiff(xs, w1, w2, b1, b2):
    """ Performs automatic differentiation for the required network.

    :param xs: the matrix of features
    :param w1, w2, b1, b2: the matrices and vectors for the pairs of weights and biases
    :return: the matrix representing the loss for each value in xs, the matrices
    representing the total derivative of the error function for each value
    in xs for w1, w2, b1 and b2, the matrix representing the error for each value in xs """
    N = xs.shape[1]
    D = xs.shape[0]
    f1 = w_1 @ xs
    f2 = (b_1 + f1.T).T
    f3 = relu(f2)
    f4 = w_2 @ f3
    f5 = (b_2 + f4.T).T
    f6 = sigmoid(f5)
    f7 = loss(xs, f6)
    f8 = np.sum(f7, axis=0)
    f9 = np.sum(f8)
    f10 = f9 / N
    t10 = 1
    t9 = 1 / N
    t8 = np.full((N), 1 / N)
    t7 = np.full((D, N), 1 / N)
    t6 = d_loss(xs, f6) * t7
    t5 = (f6*(1-f6)) * t6
    t4 = t5
    tb2 = np.sum(t5, axis=1)
    t3 = w_2.T @ t4
    tw2 = (f3 @ t4.T).T
    t2 = d_relu(f2) * t3
    t1 = t2
    tb1 = np.sum(t2, axis=1)
    tw1 = t1 @ xs.T
    return f6, tw1, tw2, tb1, tb2, f10

def adam(weight, grad, k, m_prev, v_prev, b1=0.9, b2=0.999, l=0.0000000001, rate=0.05):
    """ Performs ADAM optimisation.

    :param weight: the list of weights to optimise
    :param grad: the gradient for the error function at the weights
    :param k: the round of optimisation
    :param m_prev, v_prev: the ADAM parameters at the previous round
    :return: the next round of optimisation, the ADAM parameters for the round """
    m = b1 * m_prev + (1 - b1)*grad
    v = b2 * v_prev + (1 - b2)*(grad ** 2)

```

```

    m_hat = m / (1 - (b1 ** k))
    v_hat = v / (1 - (b2 ** k))
    weight -= (rate / (np.sqrt(v_hat) + 1)) * m_hat
    return (k+1, m, v)

w1 = create_weight((30, 784))
w2 = create_weight((784, 30))
b1 = create_weight((30,))
b2 = create_weight((784,))

epochs = 100

past_epochs = []
past_loss_values = []

previous_runnings = [[1,0,0] for i in range(4)]
weights = [w1,w2,b1,b2]

for i in range(epochs):
    res, grad_w1, grad_w2, grad_b1, grad_b2, loss = run(features.T, w1, w2, b1, b2)
    grads = [grad_w1, grad_w2, grad_b1, grad_b2]
    for (j, (weight, previous_running, grad)) in enumerate(zip(weights, previous_runnings, grads)):
        k, m, v = previous_running
        previous_runnings[j] = adam(weight, grad, k, m, v)
    past_epochs.append(i)
    past_loss_values.append(loss)

```