# AssignmentReport-Group187

February 20, 2024

# 1 Assignment 1 Report

This is an outline for your report to ease the amount of work required to create your report. Jupyter notebook supports markdown, and I recommend you to check out this cheat sheet. If you are not familiar with markdown.

Before delivery, **remember to convert this file to PDF**. You can do it in two ways: 1. Print the webpage (ctrl+P or cmd+P) 2. Export with latex. This is somewhat more difficult, but you'll get somehwat of a "prettier" PDF. Go to File -> Download as -> PDF via LaTeX. You might have to install nbconvert and pandoc through conda; `conda install nbconvert pandoc`.

# 2 Task 1

## 2.1 task 1a)

Fill in task 1a image of hand-written notes which are easy to read, or latex equations here

## Task 1a

Input → Hidden layer →

**eq.2:** $\quad w_{ji} := w_{ji} - \alpha \dfrac{\partial C}{\partial w_{ji}}$

↓ Chain-rule

$$\frac{\partial C}{\partial w_{ji}} = \frac{\partial C}{\partial a_j} \cdot \frac{\partial a_j}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_{ji}}$$

with $a_j = f(z_j)$
Sigmoid function:

using $z_j = \sum_i w_{ji} x_i$ :
$$\frac{\partial z_j}{\partial w_{ji}} = x_i$$

$$\frac{\partial a_j}{\partial z_j} = f'(z_j) = f(z_j) \cdot (1 - f(z_j))$$

↑

$\dfrac{d}{dz}\big(f(z)\big) = \dfrac{d}{dz}\Big(\dfrac{1}{1+e^{-z}}\Big) = \dfrac{e^{-z}}{(1+e^{-z})^2} = \dfrac{1+e^{-z}-1}{(1+e^{-z})^2} = \dfrac{1}{1+e^{-z}} \cdot \dfrac{1+e^{-z}-1}{1+e^{-z}} = \dfrac{1}{1+e^{-z}}\Big(1 - \dfrac{1}{1+e^{-z}}\Big) = f(z)(1-f(z))$

$$\frac{\partial C}{\partial w_{ji}} = \frac{\partial C}{\partial a_j} \cdot f'(z_j) \cdot x_i$$

backward-prop! (Output→hidden)

$$\frac{\partial C}{\partial a_j} = \sum_u \frac{\partial C}{\partial z_u} \cdot \frac{\partial z_u}{\partial a_j} = \sum_u \delta_u w_{uj}$$

$= \delta_u$ $\quad = w_{uj}$
$= -(y_u - \hat{y}_u)$ $\quad$ since $z_u = \sum_j w_{uj} a_j$

$$\frac{\partial C}{\partial w_{ji}} = \sum_u w_{uj} \delta_u \cdot f'(z_j) \cdot x_i$$

$= \delta_j$ , q.e.d.

$$\frac{\partial C}{\partial w_{ji}} = \delta_j \, x_i$$

Into eq.2:

$$w_{ji} := w_{ji} - \alpha \frac{\partial C}{\partial w_{ji}} = w_{ji} - \alpha \delta_j x_i$$

$= $ **eq.3** , q.e.d.

## 2.2 task 1b)

**Task 1b)**

Output → Hidden layer:

$$w_{kj} := w_{kj} - \alpha \frac{\partial C}{\partial w_j'}$$

$$\underbrace{\delta_k a_j' = \frac{\partial C}{\partial z_k} a_j'}_{} = -(y_k - \hat{y}_k) a_j'$$

$$\rightarrow w_{kj} := w_{kj} - \alpha (\hat{y}_k - y_k) a_j^T$$

→ Vectorized:

$$(k \times J)$$
$$\overbrace{\delta^{(2)}}$$

$$\boxed{W^{(2)} := W^{(2)} - \alpha (\hat{Y} - Y) \cdot A^T}$$

$$(k \times J) \quad (k \times 1) \quad (J \times 1)$$
$$\qquad \quad N \qquad \quad N$$

+ averaged over all
samples / batch
(N)

Hidden layer → Input layer :

$$w_{ji} := w_{ji} - \alpha \delta_j x_i \quad , \quad \delta_j = f'(z_j) \cdot \sum_k w_{kj} \cdot \overbrace{\delta_k}^{-(y_k - \hat{y}_k)}$$

$$w_{ji} := w_{ji} - \alpha f'(z_j) \underbrace{\sum_k w_{kj} \cdot (\hat{y}_k - y_k)}_{} \cdot x_i^T$$

$$(J \times I) \quad (J \times 1) \quad (k \times J) \ (k \times 1) \quad (1 \times I)$$

$$\underbrace{w_{kj}^T \cdot (\hat{y}_k - y_k)}_{} = (J \times 1)$$

$$\underbrace{f'(z_j) \odot (w_{kj}^T \cdot (\hat{y}_k - y_k))}_{} = (J \times 1)$$

Sigmoid
to every
element

element wise
Multiply

$$\underbrace{\qquad\qquad\qquad}_{} = (J \times I)$$

$$= \delta^{(1)} = (J \times 1)$$
$$\qquad\qquad\qquad N$$

$$\rightarrow \boxed{W^{(1)} := W^{(1)} - \alpha \; f'(z) \odot (W^{(2)} \cdot (\hat{Y} - Y)) \cdot X^T}$$

$$(J \times I) \quad (J \times 1) \quad (J \times k) \ (k \times 1) \quad (1 \times I)$$
$$\qquad\qquad N \qquad\qquad\qquad N \qquad\qquad N$$

Short:
$$\begin{cases} W^{(2)} := W^{(2)} - \alpha \, \delta^{(2)} A^T & , \; \delta^{(2)} = -(Y - \hat{Y}) \\ \quad (k \times J) \quad (k \times 1) \ (1 \times J) & \quad (k \times 1) \ (k \times 1) \\ \qquad\qquad N \qquad N & \qquad N \qquad N \\ \\ W^{(1)} := W^{(1)} - \alpha \, \delta^{(1)} X^T & , \; \delta^{(1)} = f'(z) \odot (W^{(2)T} \cdot \delta^{(2)}) \\ \quad (J \times I) \quad (J \times 1) \ (1 \times I) & \quad (J \times 1) \ (J \times k) \ (k \times 1) \\ \qquad\qquad N \qquad N & \qquad J \qquad N \end{cases}$$

3

# 3   Task 2

# 4   Task 2a)

### 4.0.1   mean: 33.55, std: 78.88 (rounded values)
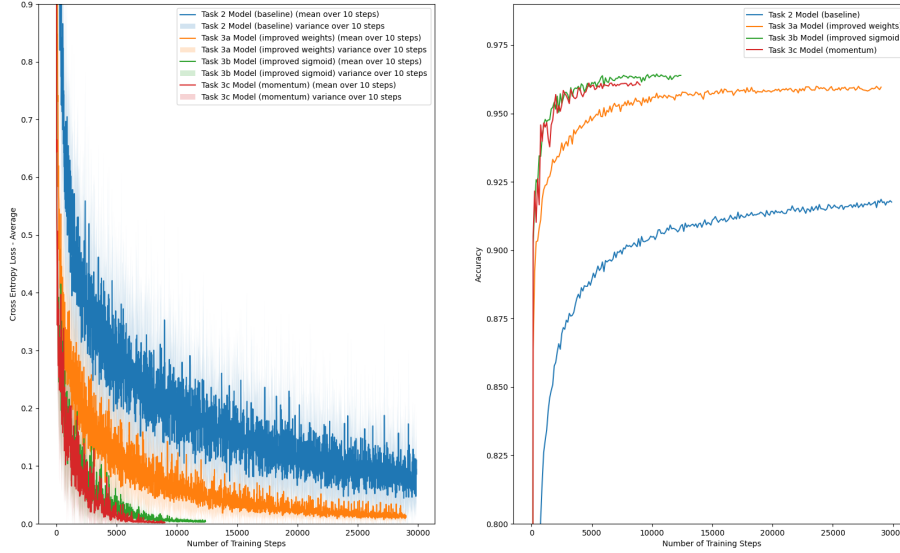
## 4.1   Task 2c)



## 4.2   Task 2d)

Total number of parameters (weights + biases) = 784 * 64 + 64 + 64 * 10 = 50880

# 5 Task 3



Each addition led to faster convergence speed and also generally a lower loss value. The only addition that stood out in an opposite way was momentum: it improved the convergence speed and loss, albeit on the cost of accuracy. On three models the accuracy is pretty high, this could indicate overfitting and hence lower generalization performance.
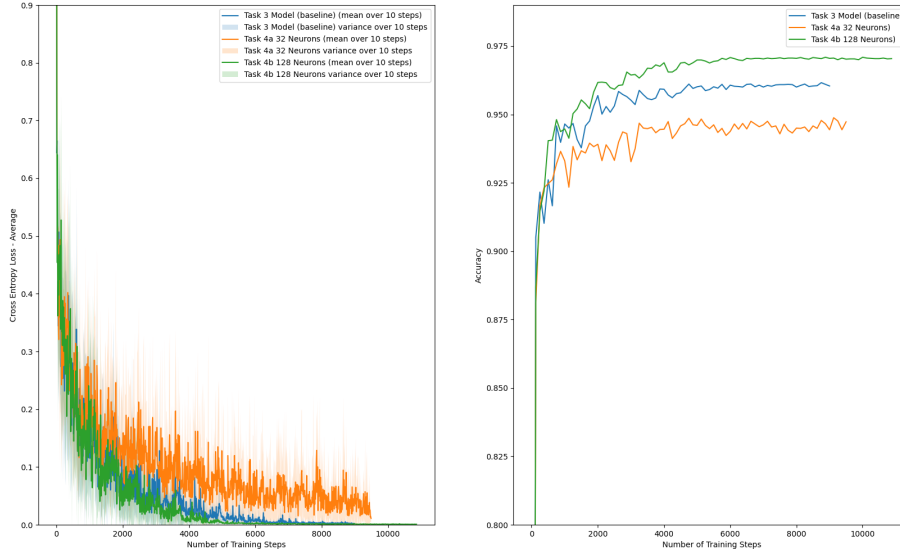
# 6 Task 4

## 6.1 Task 4a)

Setting the number of hidden units to 32 leads to a lower accuracy, slower convergence speed and a higher loss value.

## 6.2 Task 4b)

Setting the number of hidden units to 128 on the other hand leads to a very high accuracy, which could even intigate overfitting. However, it improves the loss value at the start of the training, but then performs worse than the baseline with 64 hidden units.
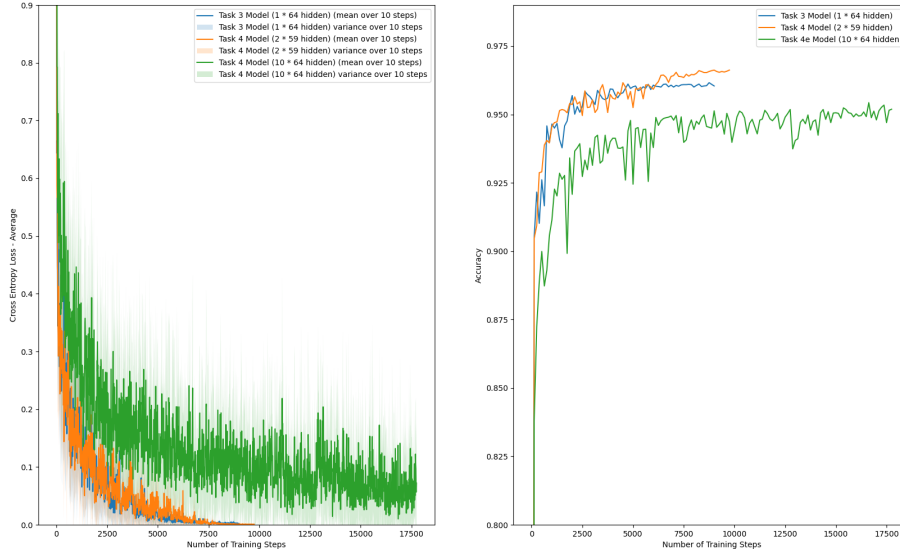
## 6.3  Task 4d

Task 3 model:

- Inputs: 785 neurons (including bias)
- Hidden layer 1: 64 neurons
- Output layer: 10 neurons
- total number of parameters: 785 * 64 (number of connections from input to hidden layer, including bias) + 64*10 (number of connections from hidden layer to output layer) = **50880**

New model used:

- Inputs: 785 neurons
- Hidden layer 1: 59 neurons
- Hidden layer 2: 59 neurons
- Output layer: 10 neurons
- total number of parameters: 785 * 59 (number of connections from input to hidden layer 1, including bias) + 59 * 59 (number of connections from hidden layer 1 to hidden layer 2) + 59*10 (number of connections from hidden layer to output layer) = **50386**

The network with more layers has a similar convergence speed as the baseline from task 3. However, it has a better accuracy.

## 6.4 Task 4e)

The model has a worse convergence speed, loss and accuracy. This effect is called the problem of vanishing gradients, because the network has too many layers ("is too deep") the gradients get decrease over time and vanish which worsens the performance.