

AssignmentReport

February 9, 2024

1 Assignment 1 Report - Group 175

1.1 Group members:

Alfred Indrehus

Andreas Lothe Måkestad

2 Task 1

2.1 task 1a)

Task 1 a)

$$\text{show that: } \frac{dC^n(w)}{dw_i} = -(y^n - \hat{y}^n) x_i^n$$

$$\text{where: } C(w) = \frac{1}{N} \sum_{n=1}^N -(y^n \ln(\hat{y}^n) + (1-y^n) \ln(1-\hat{y}^n))$$

$$\hat{y}^n = \sigma(w_i x_i^n + b) = \sigma(z) = \frac{1}{1+e^{-z}} \quad \left(z_i = w_i x_i^n + b \right)$$

$$C^n(w) = -[y^n \ln(\hat{y}^n) + (1-y^n) \ln(1-\hat{y}^n)]$$

$$\frac{dC^n}{dw_i} = \underbrace{\frac{dC^n}{d\hat{y}^n}}_{(1)} \cdot \underbrace{\frac{d\hat{y}^n}{dz}}_{(2)} \cdot \underbrace{\frac{dz}{dw_i}}_{(3)}$$

$$\boxed{\frac{dC^n}{d\hat{y}^n} = \frac{-y^n}{\hat{y}^n} + \frac{(1-y^n)}{(1-\hat{y}^n)}} \quad (1)$$

$$\frac{d\hat{y}^n}{dz_i} = \frac{d}{dz_i} (1+e^{-z_i})^{-1} \Rightarrow \frac{d\hat{y}^n}{dz_i} = \frac{(e^{-z_i})}{(1+e^{-z_i})^2}$$

$$\frac{d\hat{y}^n}{dz_i} = \frac{(e^{-z_i})}{(1+e^{-z_i})^2}$$

$$\hat{y}^n = \frac{1}{(1+e^z)^2} \Rightarrow \boxed{\frac{d\hat{y}^n}{dz_i} = \hat{y}^n (1-\hat{y}^n)} \quad (2)$$

$$e^{-z} = (1-\hat{y}^n) / \hat{y}^n$$

$$\frac{dz_i}{dw_i} = \frac{d}{dw_i} (w_i x_i + b_i) \Rightarrow \boxed{\frac{dz_i}{dw_i} = x_i} \quad (3)$$

$$\frac{dC^n}{dw_i} = \textcircled{1} \cdot \textcircled{2} \cdot \textcircled{3}$$

$$\begin{aligned} \frac{dC^n}{dw_i} &= \left(\frac{-\gamma^n}{\hat{\gamma}^n} + \frac{(1-\gamma^n)}{(1-\hat{\gamma}^n)} \right) \cdot \hat{\gamma}^n (1-\hat{\gamma}^n) x_i \\ &= \left(\frac{-\gamma^n + \hat{\gamma}^n}{\hat{\gamma}^n (1-\hat{\gamma}^n)} \right) \cdot \hat{\gamma}^n (1-\hat{\gamma}^n) \cdot x_i \end{aligned}$$

$$\frac{dC^n}{dw_i} = (-\gamma^n + \hat{\gamma}^n) \cdot x_i^n$$

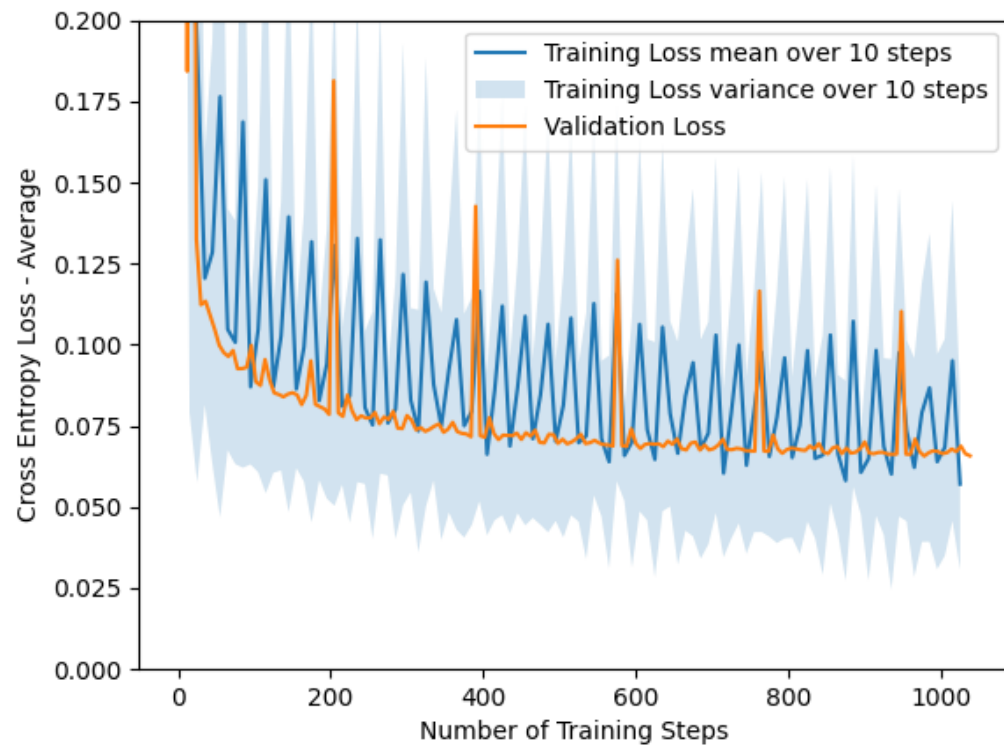
$$\frac{dC^n}{dw_i} = -(\gamma^n - \hat{\gamma}^n) x_i^n$$

Q.E.D

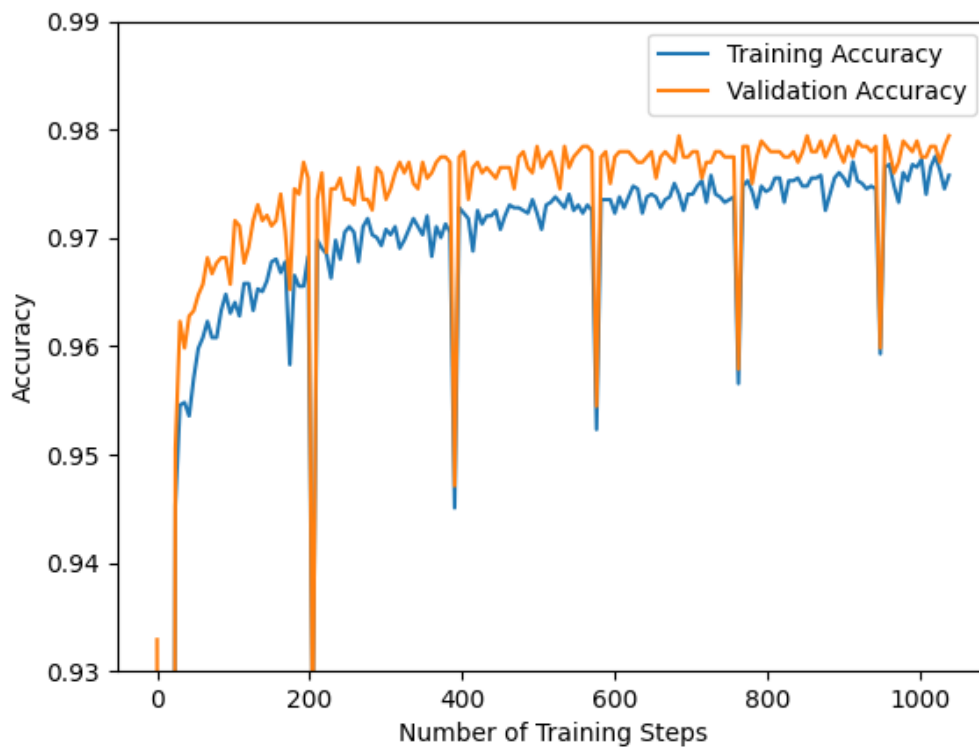
2.2 task 1a)

3 Task 2

3.1 Task 2b)



3.2 Task 2c)

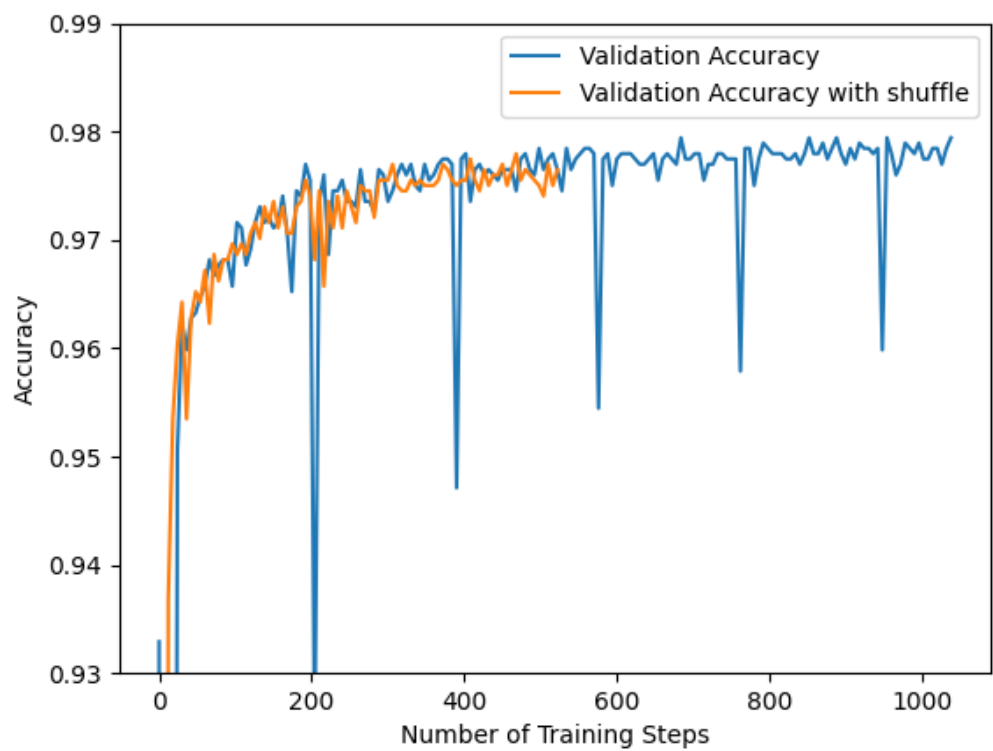


3.3 Task 2d)

Early stopping kicks in after 33 epochs.

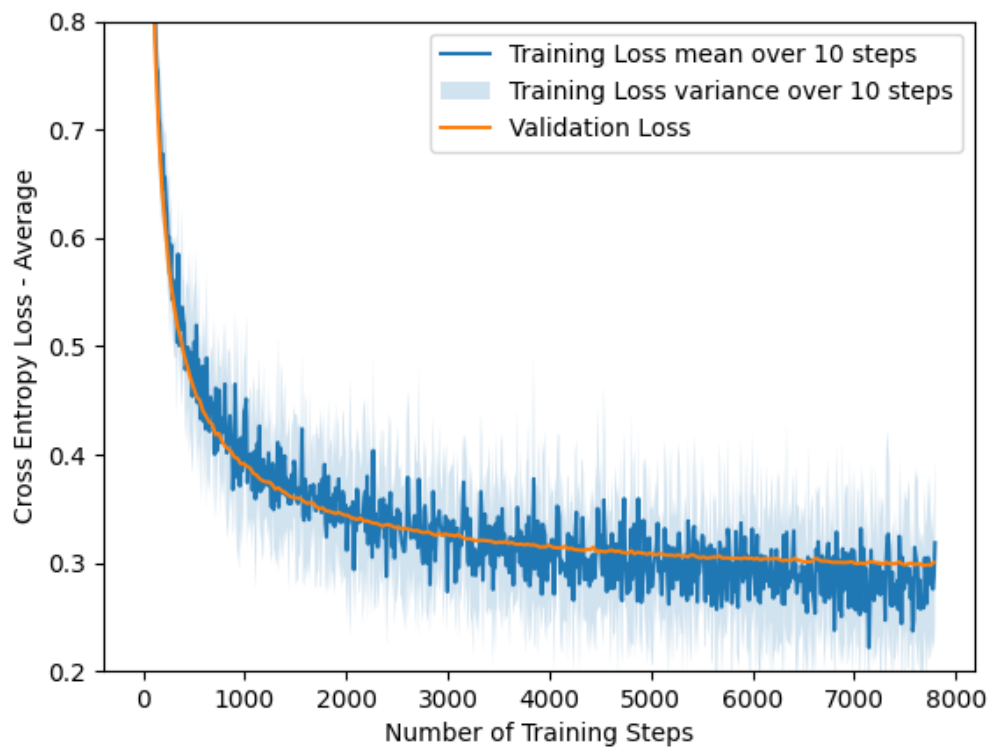
3.4 Task 2e)

The graph has notably fewer spikes and is in general more stable. This happens because shuffling the dataset makes it less prone to overfitting, causing the model to not be as sensitive to the data it is trained on. If for instance the data contains certain orders or patterns, shuffling the dataset will make the model avoid learning these patterns.

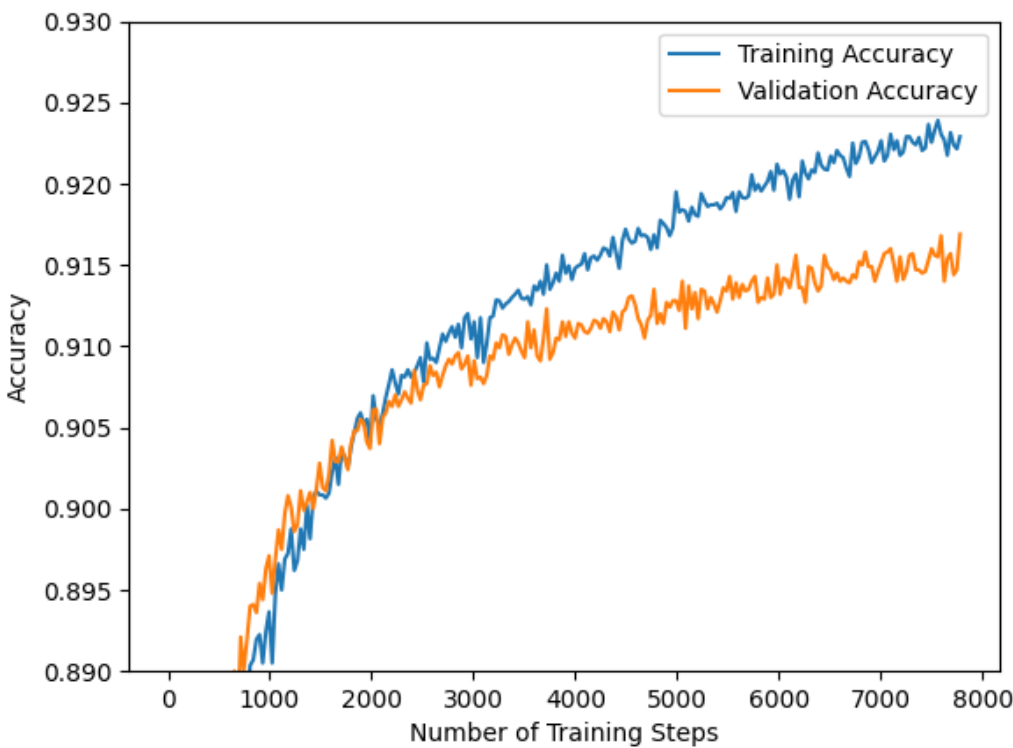


4 Task 3

4.1 Task 3b)



4.2 Task 3c)



4.3 Task 3d)

The difference in accuracy on the training data and the validation data seems to start diverging at around 3000 steps. Here, the training accuracy continues to increase with the same slope, while the accuracy on the validation accuracy starts to slow down, but not by much. The difference between the two could indicate slight overfitting, however, with the difference being so small, one could argue it is a good fit.

5 Task 4

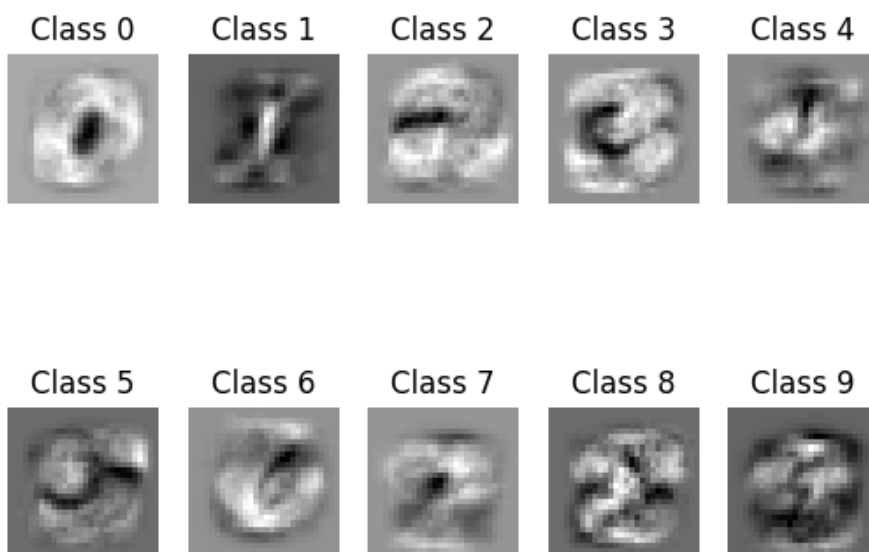
5.1 Task 4a)

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

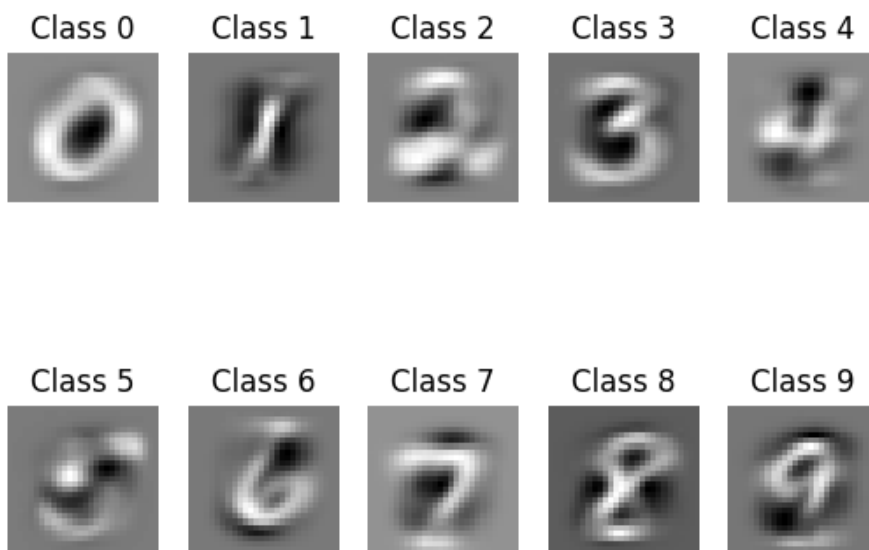
5.2 Task 4b)

The weights with the higher lambda is less noisy because of the increased generalization. This means that the model won't be prone to outliers and it won't become too complex, causing the weights to better reflect the generalized data, since larger weights are penalized.

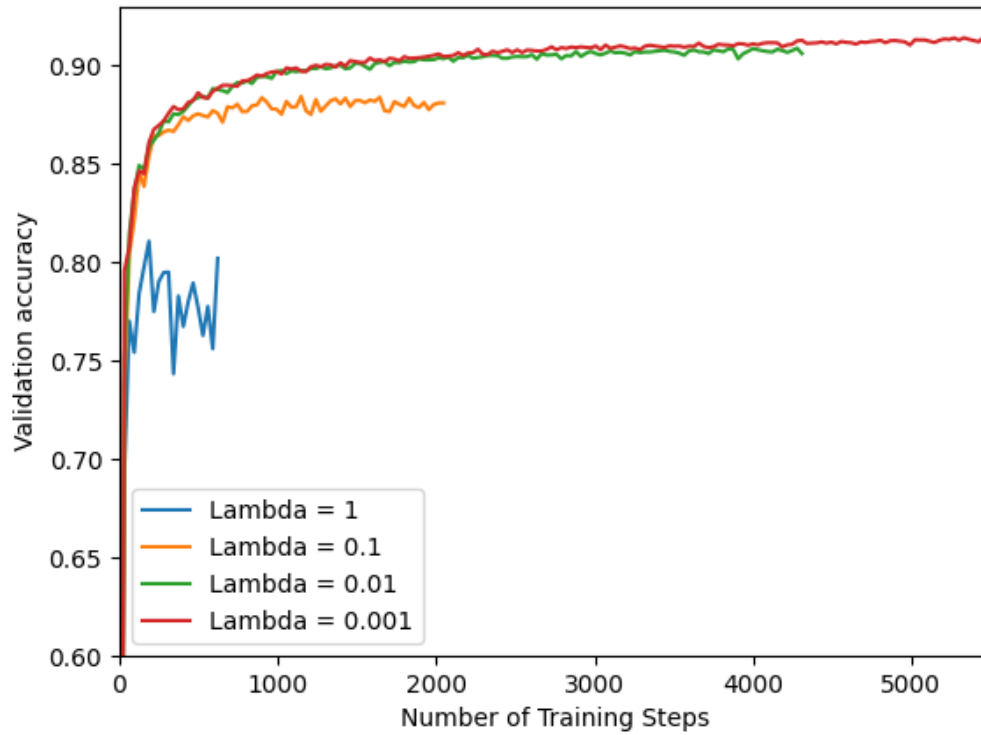
Lambda = 0



Lambda = 1



5.3 Task 4c)



5.4 Task 4d)

The reason the validation accuracy degrades with regularization could be because the model simply becomes too generalized, thus resulting in underfitting.

5.5 Task 4e)

One can see from the graph over the l_2 norm that the length of the the weight vector decreases as λ increases. This makes sense from the definition of L2-regularization, since larger values of λ causes the second term in equation (9) to decrease even more.

