# AssignmentReport-Group1

February 6, 2021

## 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)

1. a) $\dfrac{\partial c^n(w)}{\partial w_i} = \dfrac{\partial c^n(w)}{\partial f(x^n)} \dfrac{\partial f(x^n)}{\partial w_i}$

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

$\hat{y}^n = f(x^n)$

$c^n(w) = -(y^n \ln(f(x^n)) + (1-y^n)\ln(1-f(x^n)))$

$\dfrac{\partial c^n(w)}{\partial f(x^n)} = -\left(\dfrac{y^n}{f(x^n)} + (1-y^n)\left(-\dfrac{1}{1-f(x^n)}\right)\right)$

$\qquad\qquad = -\dfrac{y^n}{f(x^n)} + \dfrac{1-y^n}{1-f(x^n)}$

$\dfrac{\partial f(x^n)}{\partial w_i} = x_i^n \, f(x^n)(1-f(x^n))$

$\dfrac{\partial c^n(w)}{\partial w_i} = \left(\dfrac{y^n}{f(x^n)} + \dfrac{1-y^n}{1-f(x^n)}\right) x_i^n \, f(x^n)(1-f(x^n))$

$\qquad\qquad = (-y^n(1-f(x^n)) + (1-y^n)f(x^n))\, x_i^n$

$\qquad\qquad = (-y^n + y^n f(x^n) + f(x^n) - y^n f(x^n))\, x_i^n$

$\qquad\qquad = (-y^n + f(x^n))\, x_i^n$

$\qquad\qquad \overset{\hat{y}^n = f(x)}{=} -(y^n - \hat{y}^n)\, x_i^n$

1.b) $\dfrac{\partial C^n(w)}{\partial w_{kj}} = \dfrac{\partial C^n(w)}{\partial z_l} \dfrac{\partial z_l}{\partial w_{kj}}$

$z_l = \sum_i^I w_{ki} \cdot x_i^n$

$\dfrac{\partial z_l}{\partial w_{kj}} = \dfrac{\partial \sum_i^I w_{ki} x_i^n}{\partial w_{kj}} = \dfrac{\partial \left( \sum_{i \neq j}^I w_{ki} x_i^n + w_{kj} x_j^n \right)}{\partial w_{kj}} = x_j^n$

$C^n(w) = -\sum_{k=1}^K y_k^n \ln(\hat{y}_k^n)$

$\hat{y}_k^n = \dfrac{e^{z_k}}{\sum_{k'} e^{z_{k'}}}$

$\dfrac{\partial C^n(w)}{\partial z_l} = \dfrac{\partial \left( -\sum_{k=1}^K y_k^n \ln(\hat{y}_k^n) \right)}{\partial z_l}$

$= -\sum_{k=1}^K y_k^n \dfrac{\partial \ln(\hat{y}_k^n)}{\partial z_l} = -\sum_{k=1}^K y_k^n \dfrac{1}{\hat{y}_k^n} \dfrac{\partial \hat{y}_k^n}{\partial z_l}$

$\dfrac{\partial \hat{y}_k^n}{\partial z_l} = \dfrac{\partial \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}}}{\partial z_l}$

for $l = k$:

$\dfrac{\partial \frac{e^{z_k}}{\sum_{k'}^K e^{z_{k'}}}}{\partial z_k} = \dfrac{e^{z_k} \sum_{k'}^K e^{z_{k'}} - (e^{z_k})^2}{\left( \sum_{k'}^K e^{z_{k'}} \right)^2}$

$= \dfrac{e^{z_k}}{\sum_{k'}^K e^{z_{k'}}} \left( \dfrac{\sum_{k'} e^{z_{k'}} - e^{z_k}}{\sum_{k'}^K e^{z_{k'}}} \right) = \hat{y}_k^n(1 - \hat{y}_k^n)$

for $l \neq k$

$\dfrac{\partial \frac{e^{z_k}}{\sum_{k'}^K e^{z_{k'}}}}{\partial z_l} = \dfrac{0 \cdot \sum_{k'}^K e^{z_{k'}} - e^{z_k} e^{z_l}}{\left( \sum_{k'}^K e^{z_{k'}} \right)^2}$

$= -\dfrac{e^{z_k}}{\sum_{k'}^K e^{z_{k'}}} \cdot \dfrac{e^{z_l}}{\sum_{k'}^K e^{z_{k'}}} = -\hat{y}_k^n \hat{y}_l^n$

$\dfrac{\partial \hat{y}_k^n}{\partial z_l} = \begin{cases} \hat{y}_k^n(1 - \hat{y}_k^n) & \text{for } l = k \\ -\hat{y}_k^n \hat{y}_l^n & \text{for } l \neq k \end{cases}$

$\dfrac{\partial C^n(w)}{\partial z_l} = -\sum_{k=1}^K \dfrac{y_k^n}{\hat{y}_k^n} \dfrac{\partial \hat{y}_k^n}{\partial z_l}$

$= -\sum_{k \neq l}^K \dfrac{y_k^n}{\hat{y}_k^n} (-\hat{y}_k^n \hat{y}_l^n) - \dfrac{y_l^n}{\hat{y}_l^n} \hat{y}_l^n(1 - \hat{y}_l^n)$

$= \sum_{k \neq l}^K y_k^n \hat{y}_l^n - y_l^n + \hat{y}_l^n \hat{y}_l^n$
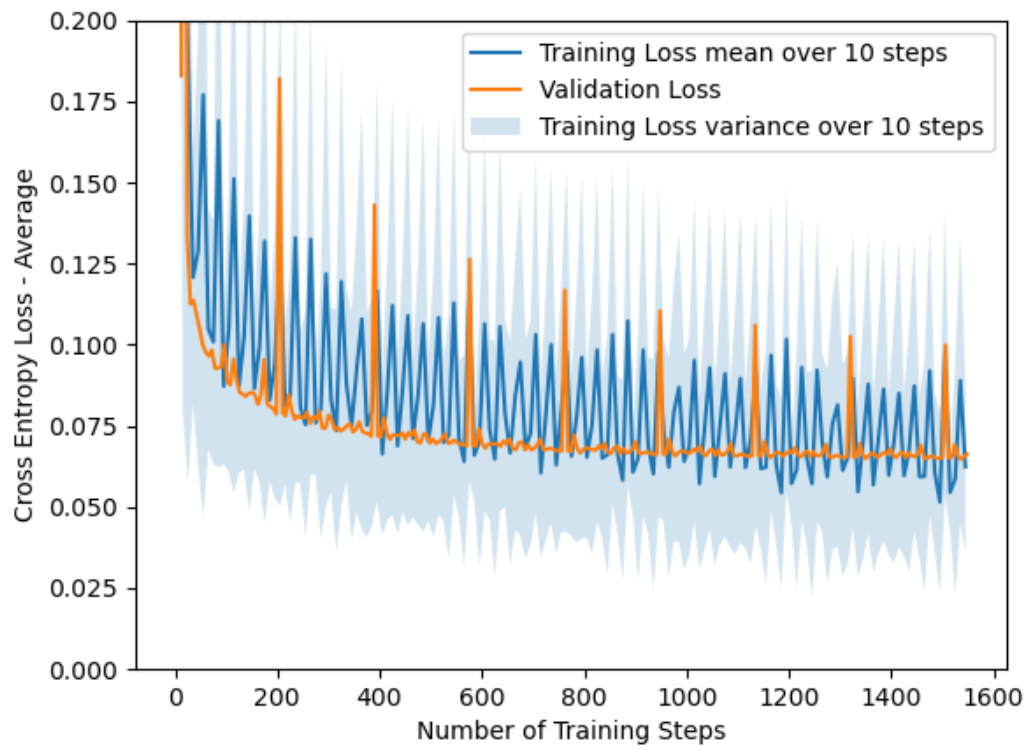
$= -y_l^n + \sum_{k=1}^K y_k^n \hat{y}_l^n$

$= -y_l^n + \hat{y}_l^n \sum_{k=1}^K y_k^n = -y_l^n + \hat{y}_l^n$

$\dfrac{\partial C^n(w)}{\partial w_{kj}} = \left( -y_l^n + \hat{y}_l^n \right) x_j^n = -x_j^n(y_l^n - \hat{y}_l^n)$
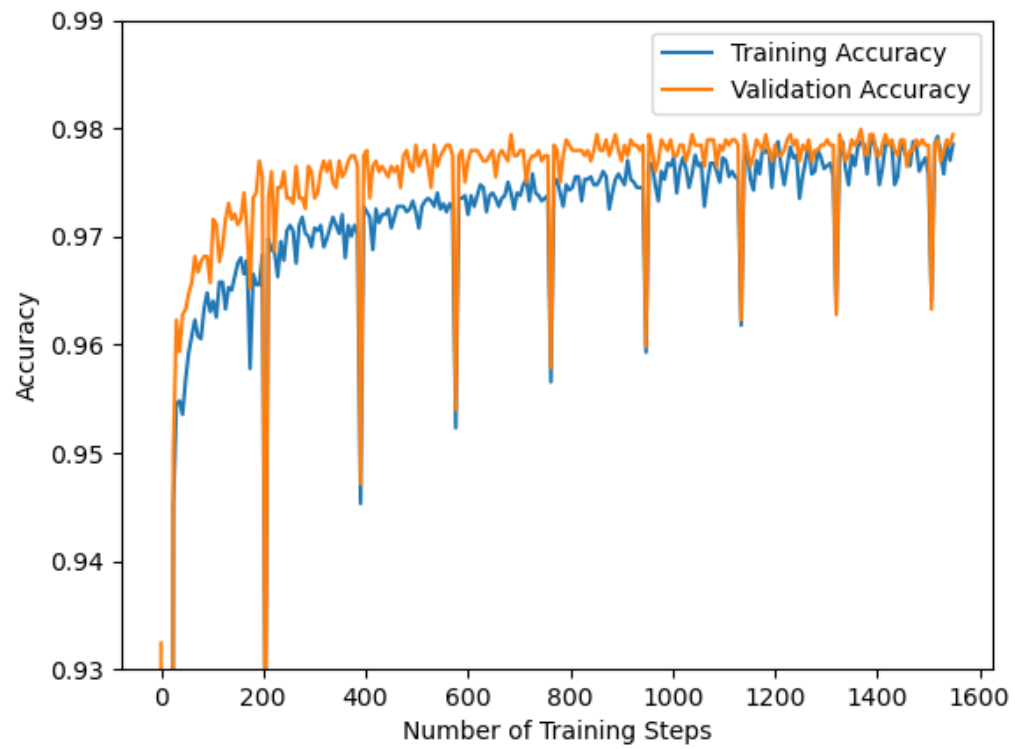
$= -x_j^n(y_k^n - \hat{y}_k^n)$

as both $k$ and $j$ are arbitrary integers between 1 and K.

# 3 Task 2

## 3.1 Task 2b)

## 3.2   Task 2c)



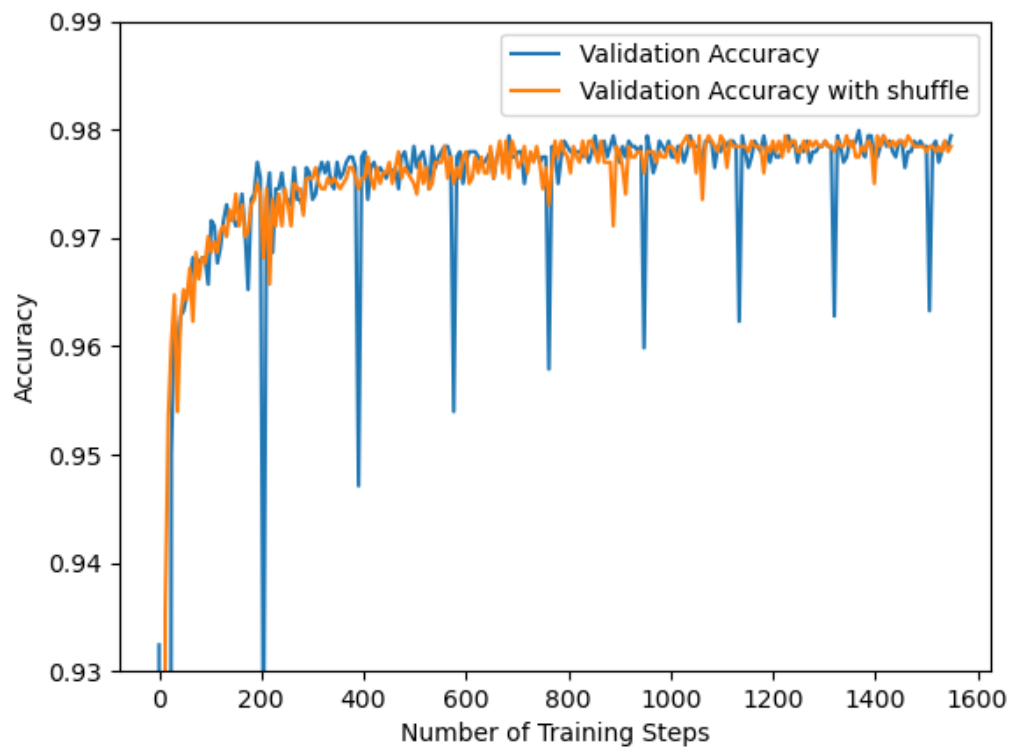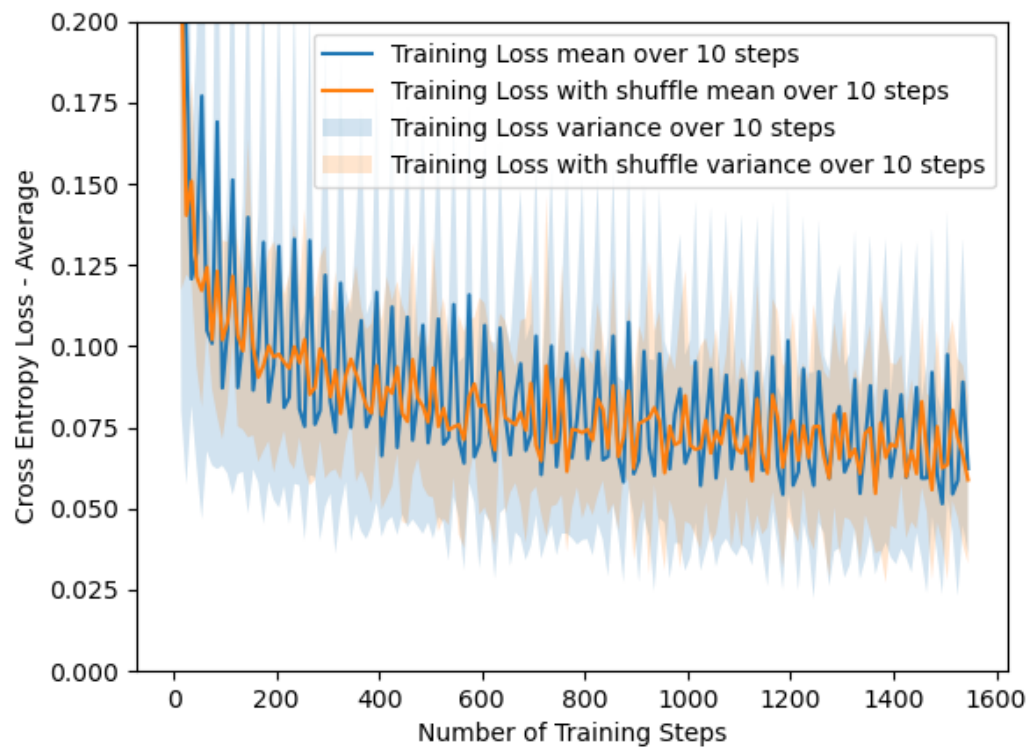## 3.3   Task 2d)

Early stopping kicks in after 4 epochs.
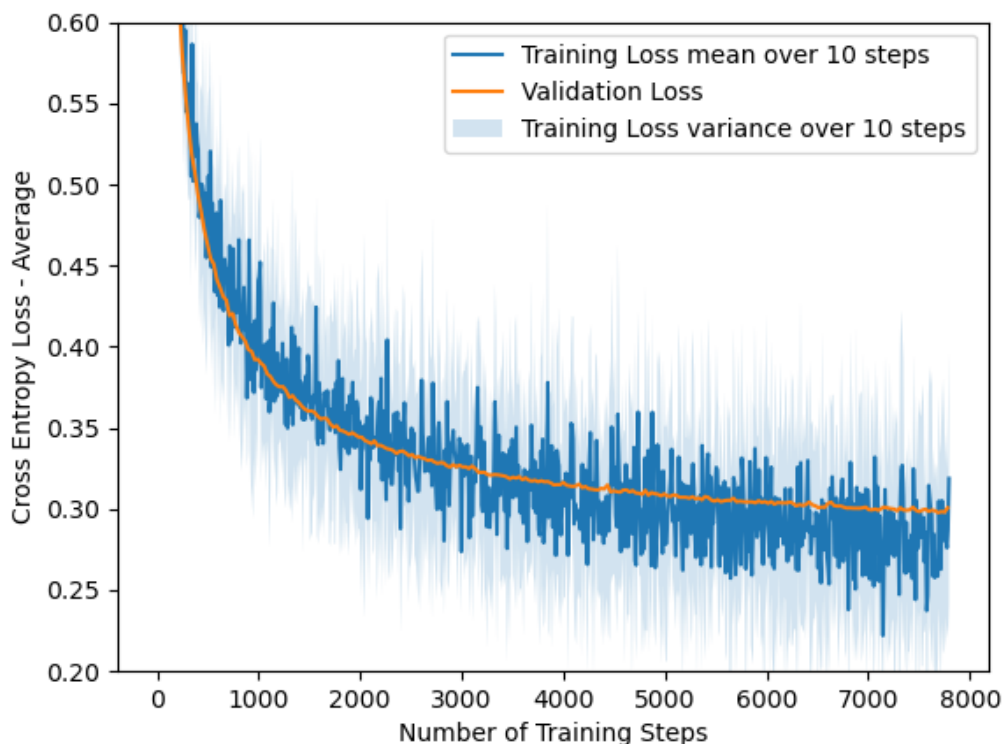
## 3.4   Task 2e)

By shuffling the dataset we avoid the model being trained for only 2s at the start of every epoch, which would result in an accuracy drop when we hit the 3s until the model learns how to correctly classify both. With a shuffled dataset we encounter 2s and 3s randomly and the model learns that there exists more than one number to classify.

In general shuffling helps reduce variance and overfitting by ensuring the data used is a representation of the whole dataset. In the MNIST database a sorted batch of data will never include all 10 digits and thereby make the model overfitted to the digits represented and completely oblivious to other digits. While a shuffled batch of data will almost always contain all digits such that the model trains on a more general batch.
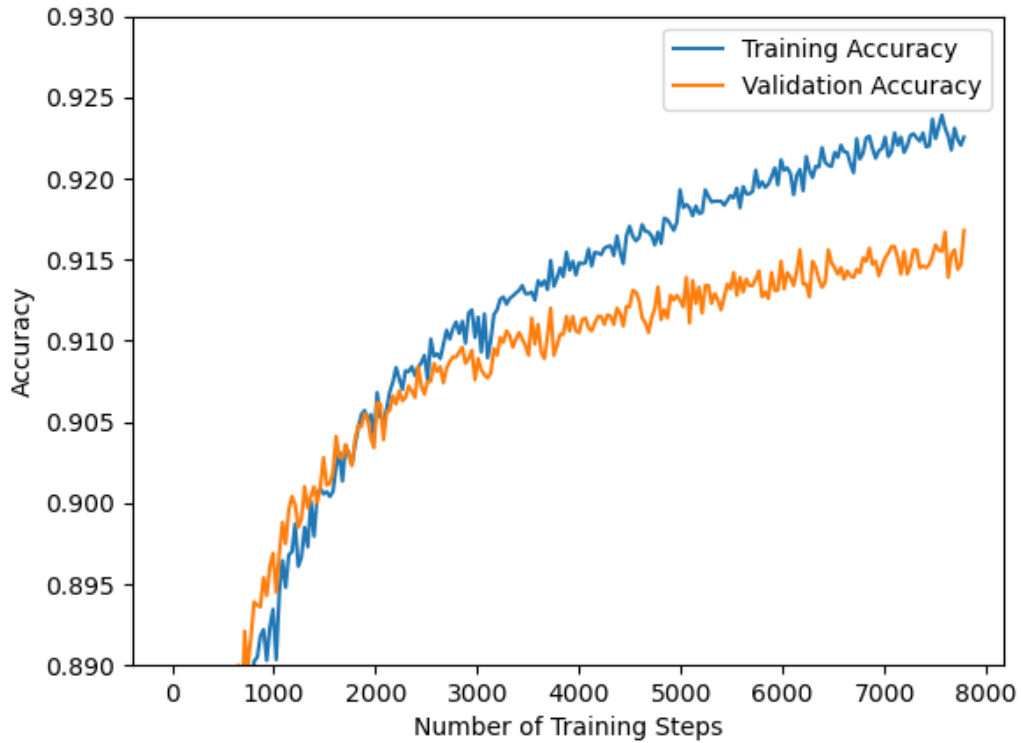
# 4 Task 3

## 4.1 Task 3b)

## 4.2 Task 3c)



## 4.3 Task 3d)

The training accuracy is rising faster than the validation accuracy, indicating some overfitting. The validation accuracy is still increasing which implies more training could give a higher validation accuracy, but this higher accuracy would cost a lot of computational time compared to earlier training where the training and validation accuracies were rising almost identically.

We can also see that the validation loss is plateauing faster than the training loss, which signals the same as the plateauing in validation accuracy compared to training accuracy.
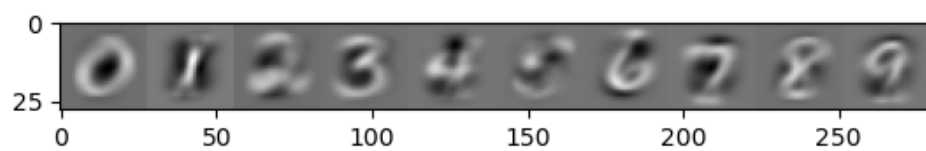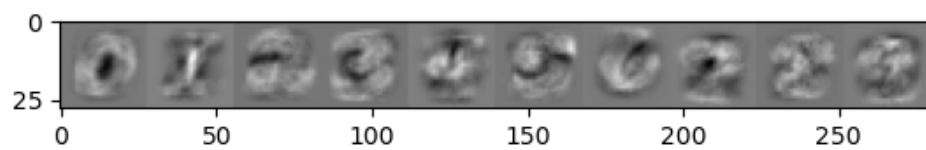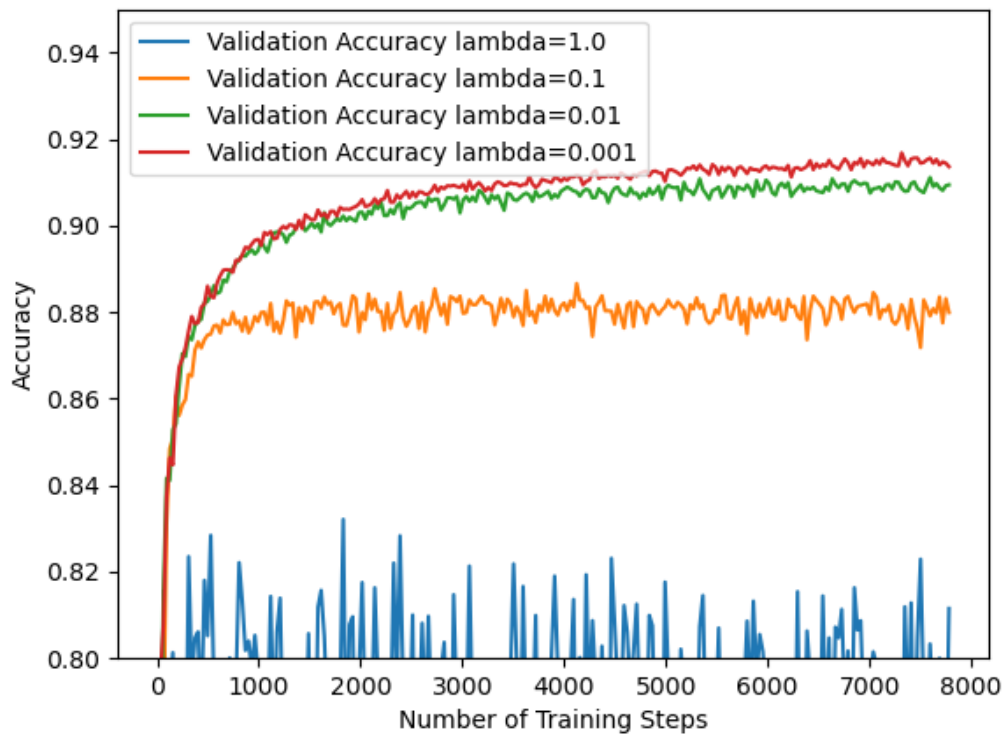
# 5 Task 4

## 5.1 Task 4a)

$$J(w) = C(w) + \lambda R(w)$$

$$\frac{\partial J(w)}{\partial w} = \frac{1}{N} \sum_{n=1}^{N} \left[ \frac{\partial C^n(w)}{\partial w} \right] + \lambda \frac{\partial \|w\|^2}{\partial w}$$

$$= \frac{1}{N} \sum_{n=1}^{N} \sum_{j=1}^{I} \sum_{k=1}^{K} \left[ \frac{\partial C^n(w)}{\partial w_{kj}} \right] + 2\lambda w$$

$$= \frac{1}{N} \sum_{n=1}^{N} \sum_{j=1}^{I} \sum_{k=1}^{K} \left[ x_j^n (y_k^n - \hat{y}_k^n) \right] + 2\lambda w$$

$$= \frac{1}{N} \sum_{n=1}^{N} \left[ -x(y^n - \hat{y}) \right] + 2\lambda w$$

$$= -\frac{1}{N} \sum_{n=1}^{N} \left[ x^n (y^n - \hat{y}^n) \right] + 2\lambda w$$

## 5.2 Task 4b)

The L2 regularization reduces the model complexity which results in less of the noise from the dataset leaking through the model. The reduction of model complexity increases bias but reduces noise in the bias/noise tradeoff.

## 5.3 Task 4c)



The zoom for this plot is a bit off, as the plot with edited range for ylim did not have time to finish running, same as task4c.

## 5.4 Task 4d)

More regularization reduces complexity which trades less noise for more bias, which might reduce accuracy, especially for high regularization constants where more bias is introduced, ie. lambda = 1.

## 5.5 Task 4e)

Could not finish running the code in time for the submission deadline, started running at about 23:00 but had not yet finished at 23:55 so that i could get the final plots into the report. Code for plotting the L2 norms is included in the code zip.

Also had to download Tex to deliver jupyter as pdf so my estimate of 5 min to be able to deliver was too short. Would have delivered 1 hr earlier if i knew code would not finish in the timeframe.

[ ]: