

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 \(https://github.com/adam-p/markdown-here/wiki/Markdown-Cheatsheet\)](https://github.com/adam-p/markdown-here/wiki/Markdown-Cheatsheet). 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 somewhat 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`.

Task 1

task 1a)

1)

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

$$C(w) = \frac{1}{N} \sum_{n=1}^N c^n(w)$$

$$\frac{\partial C^n(w)}{\partial w^i} = \frac{\partial C^n(w)}{\partial f(x^n)} \cdot \frac{\partial f(x^n)}{\partial w^i} = \frac{\partial C^n(w)}{\partial \hat{y}^n} \cdot \frac{\partial f(x^n)}{\partial w^i}$$

$$\frac{\partial C^n(w)}{\partial \hat{y}^n} = -\frac{y^n}{\hat{y}^n} + \frac{1 - y^n}{1 - \hat{y}^n}$$

$$\ln(-x) \Rightarrow -\frac{1}{x} \cdot -1$$

$$\ln(1-x) = \frac{1}{1-x} \cdot -1$$

$$\Rightarrow \frac{\partial C^n(w)}{\partial w^i} = -\left(\frac{y^n}{\hat{y}^n} - \frac{1 - y^n}{1 - \hat{y}^n}\right) \cdot \hat{y}^n \cdot (1 - \hat{y}^n) \cdot x_i^n$$

$$= (-y^n(1 - \hat{y}^n) - \hat{y}^n(1 - y^n)) x_i^n$$

$$= (y^n \hat{y}^n - y^n + \hat{y}^n - y^n \hat{y}^n) x_i^n = \underline{\underline{-(y^n - \hat{y}^n) x_i^n}}$$

task 1b)

$$b) C^n(w) = - \sum_{k=1}^K y_k^n \ln(\hat{y}_k^n) \quad \hat{y}_k^n = \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}}, \quad z_k = w_k^T X = \sum_i w_{k,i} \cdot x_i$$

$$\frac{\partial C^n(w)}{\partial w_{k,j}} = \frac{\partial C^n(w)}{\partial \hat{y}_k^n} \cdot \frac{\partial \hat{y}_k^n}{\partial w_{k,j}} = \frac{\partial C^n(w)}{\partial \hat{y}_k^n} \cdot \frac{\partial \hat{y}_k^n}{\partial z_k} \cdot \frac{\partial z_k}{\partial w_{k,j}}$$

$$\frac{\partial z_k}{\partial w_{k,j}} = x_j^n$$

$$\frac{\partial \hat{y}_k^n}{\partial z_k} = \frac{\partial}{\partial z_k} \left(\frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}} \right) = \frac{\partial}{\partial z_k} \left(\frac{e^{z_k}}{e^{z_k} + \sum_{k' \neq k} e^{z_{k'}}} \right) \quad \sum_{k'} e^{z_{k'}} = e^{z_k} + \sum_{k' \neq k} e^{z_{k'}}$$

$$= \frac{-e^{z_k} \cdot e^{z_k} + e^{z_k} \left(e^{z_k} + \sum_{k' \neq k} e^{z_{k'}} \right)}{\left(\sum_{k'} e^{z_{k'}} \right)^2} = \frac{e^{z_k} \sum_{k' \neq k} e^{z_{k'}}}{\left(\sum_{k'} e^{z_{k'}} \right)^2} = \hat{y}_k^n \cdot \frac{\sum_{k' \neq k} e^{z_{k'}}}{\sum_{k'} e^{z_{k'}}}$$

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

$$\frac{\partial C^n(w)}{\partial \hat{y}_k^n} = - \sum_{k'} y_{k'}^n \cdot \frac{1}{\hat{y}_{k'}^n}$$

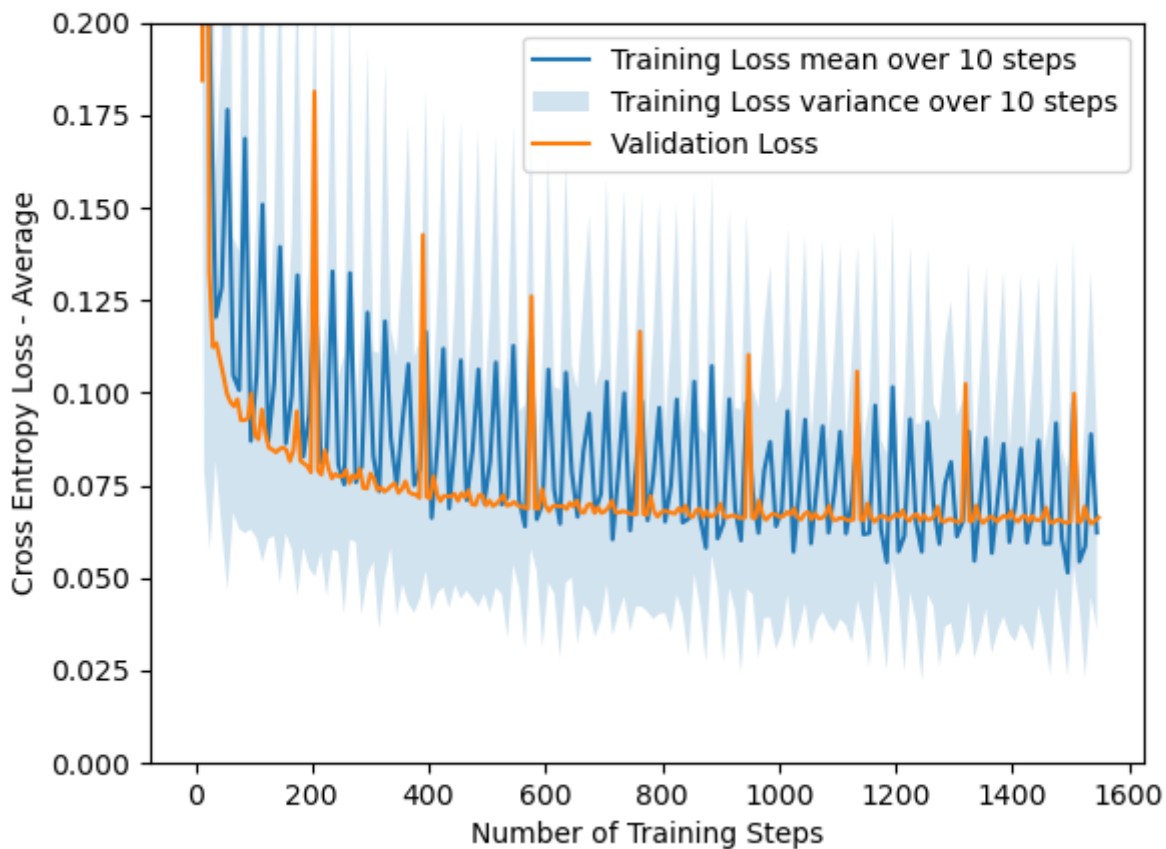
$$\Rightarrow \frac{\partial C^n(w)}{\partial w_{k,j}} = - \sum_{k'} \frac{y_{k'}^n}{\hat{y}_{k'}^n} \cdot \hat{y}_k^n (1 - \hat{y}_k^n) \cdot x_j^n = - \sum_{k'} y_{k'}^n (1 - \hat{y}_k^n) \cdot x_j^n$$

This was the closest I could get to the answer, I got the same answer when calculating dC/dz^k

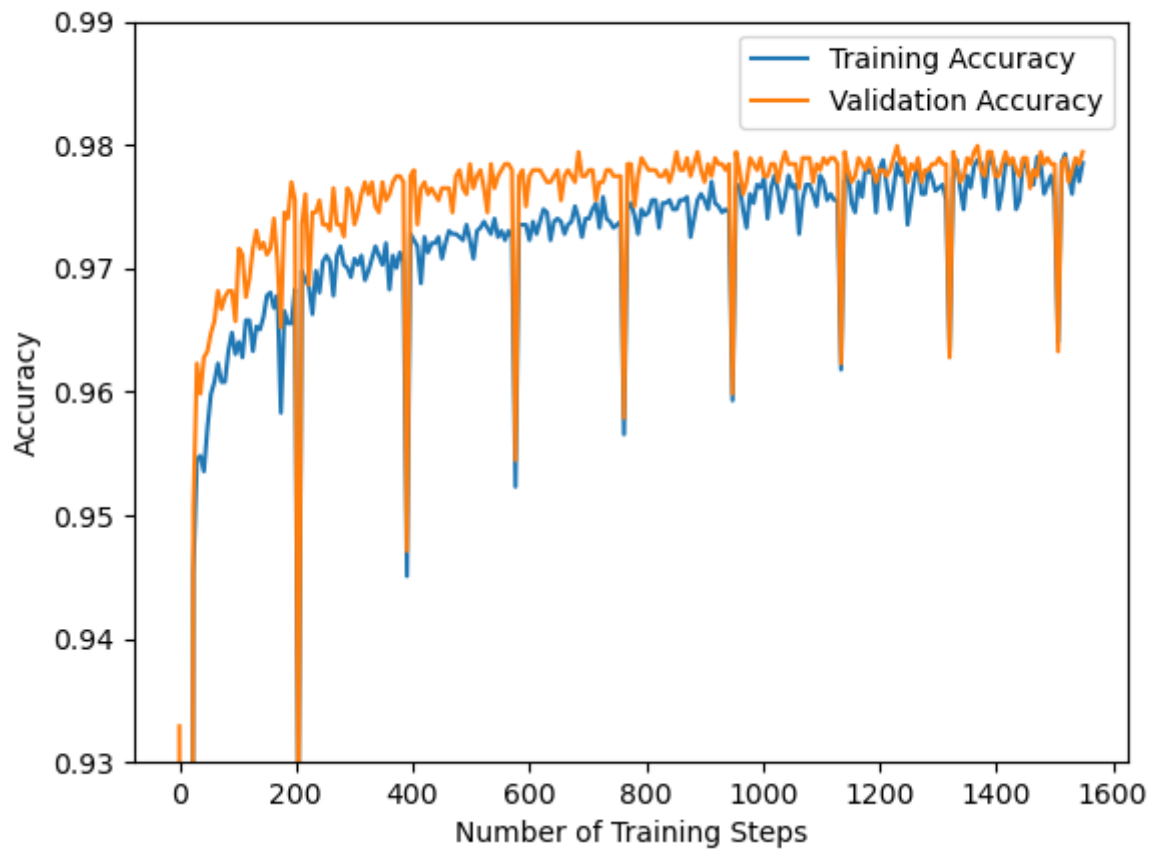
$$\begin{aligned}
 \frac{\partial \mathcal{L}(\omega)}{\partial z_n} &= \frac{\partial}{\partial z_n} \left(- \sum_{k=1}^K y_k^n \cdot \ln \left(\frac{e^{z_n}}{\sum_{k=1}^K e^{z_k}} \right) \right) \\
 &= - \frac{\partial}{\partial z_n} \left(\sum_{k=1}^K y_k^n \cdot \left(\underbrace{\ln(e^{z_n})}_{z_n} - \ln(\sum_{k=1}^K e^{z_k}) \right) \right), \\
 &= - \frac{\partial}{\partial z_n} \left(\sum_{k=1}^K y_k^n (z_n - \ln(\sum_{k=1}^K e^{z_k})) \right) = - \sum_{k=1}^K y_k^n \left(1 - \frac{e^{z_n}}{\sum_{k=1}^K e^{z_k}} \right) = \sum_{k=1}^K \left(y_k^n \frac{e^{z_n}}{\sum_{k=1}^K e^{z_k}} \right) - 1 \\
 &= \sum_{k=1}^K y_k^n \hat{y}_k^n - y_n^n
 \end{aligned}$$

Task 2

Task 2b)



Task 2c)

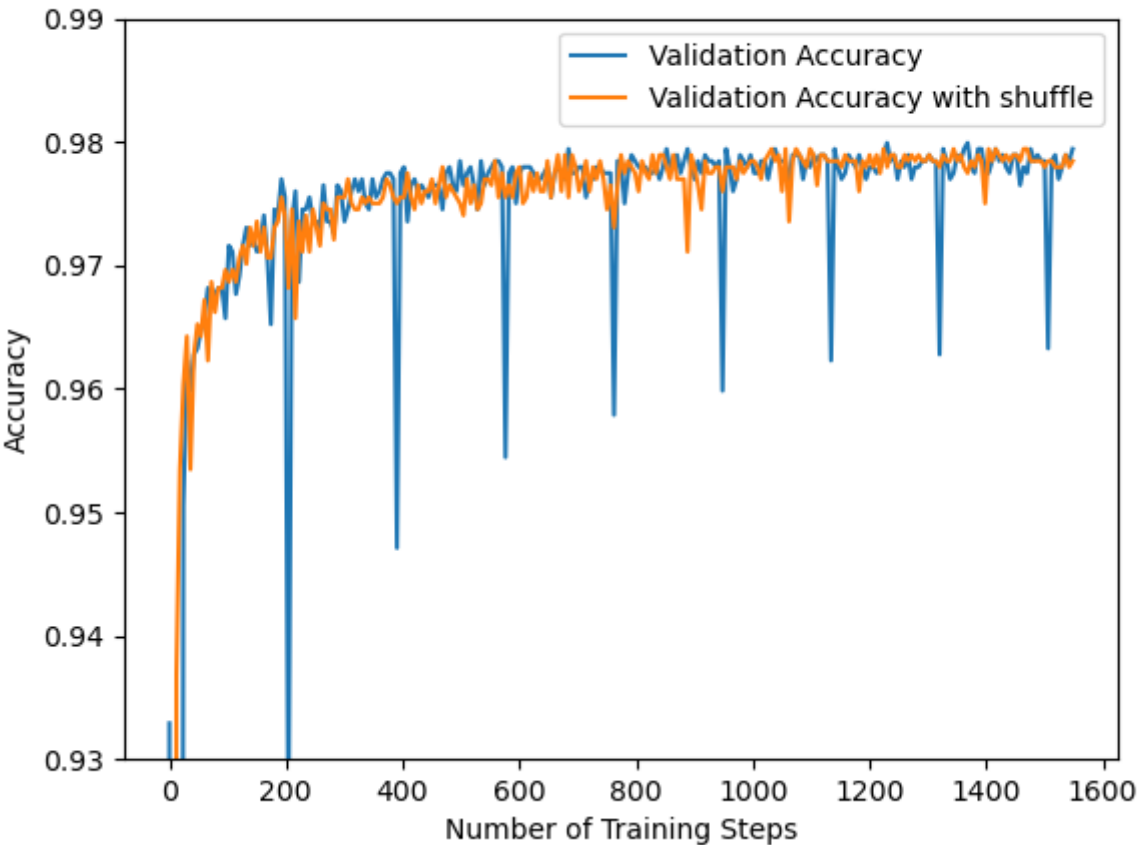
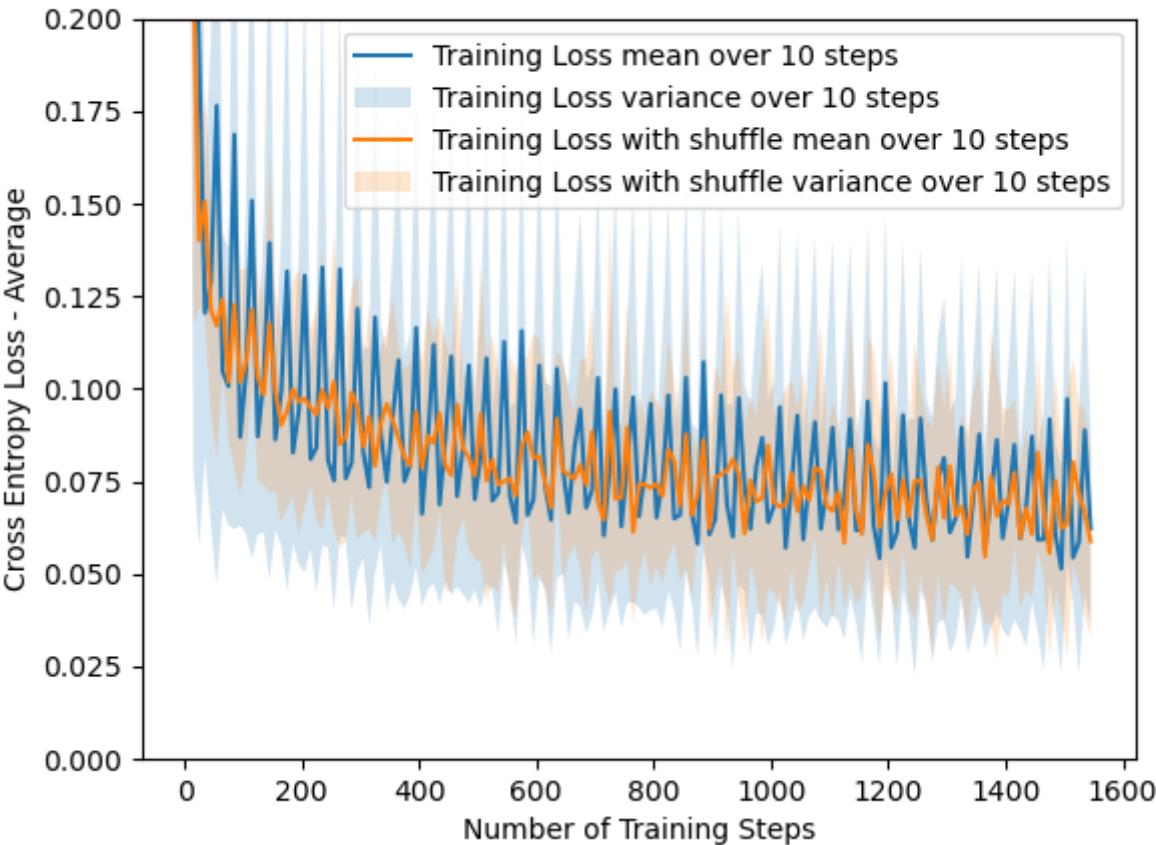


Task 2d)

It stops after 33 epochs, with shuffling after 16.

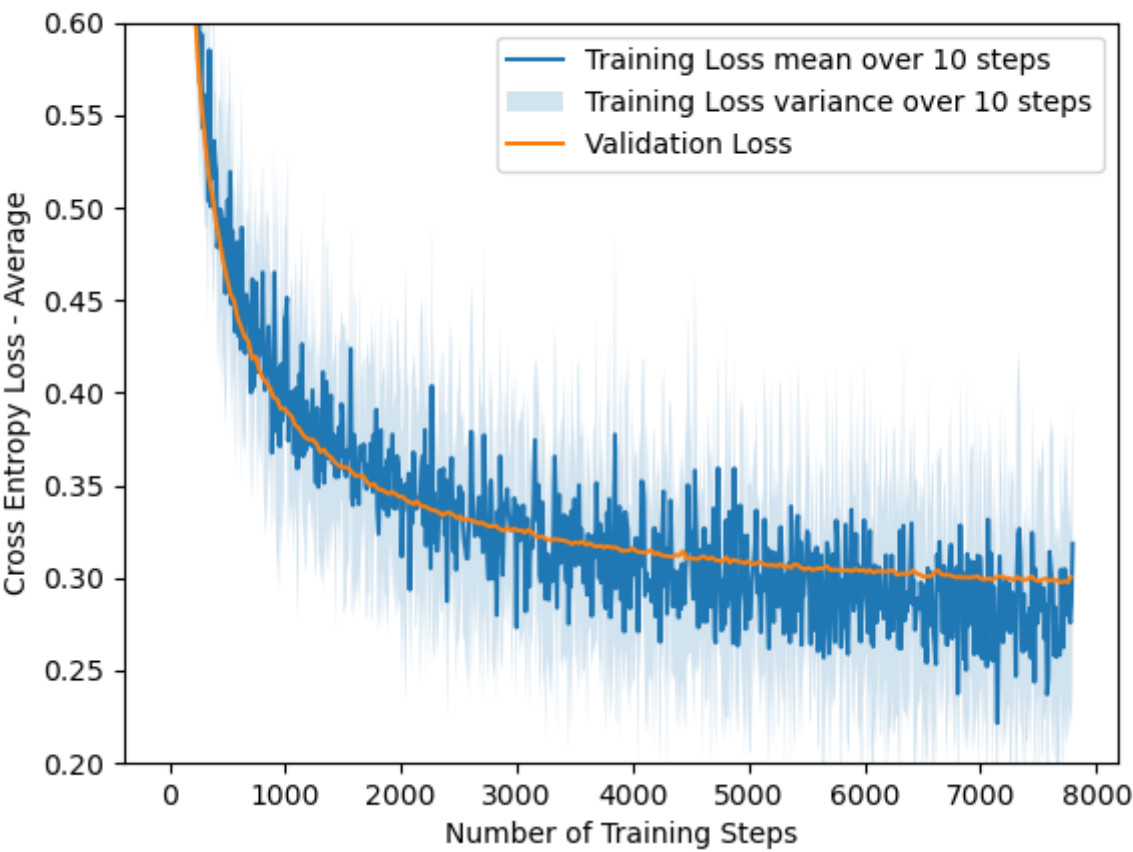
Task 2e)

My best guess would be that every 200-ish training steps there is a validation set that is particularly difficult for the neural net to classify, the shuffling of the data set will therefore spread out the difficult to classify data points. Generally the random shuffling will remove any deterministic bias of the data set.

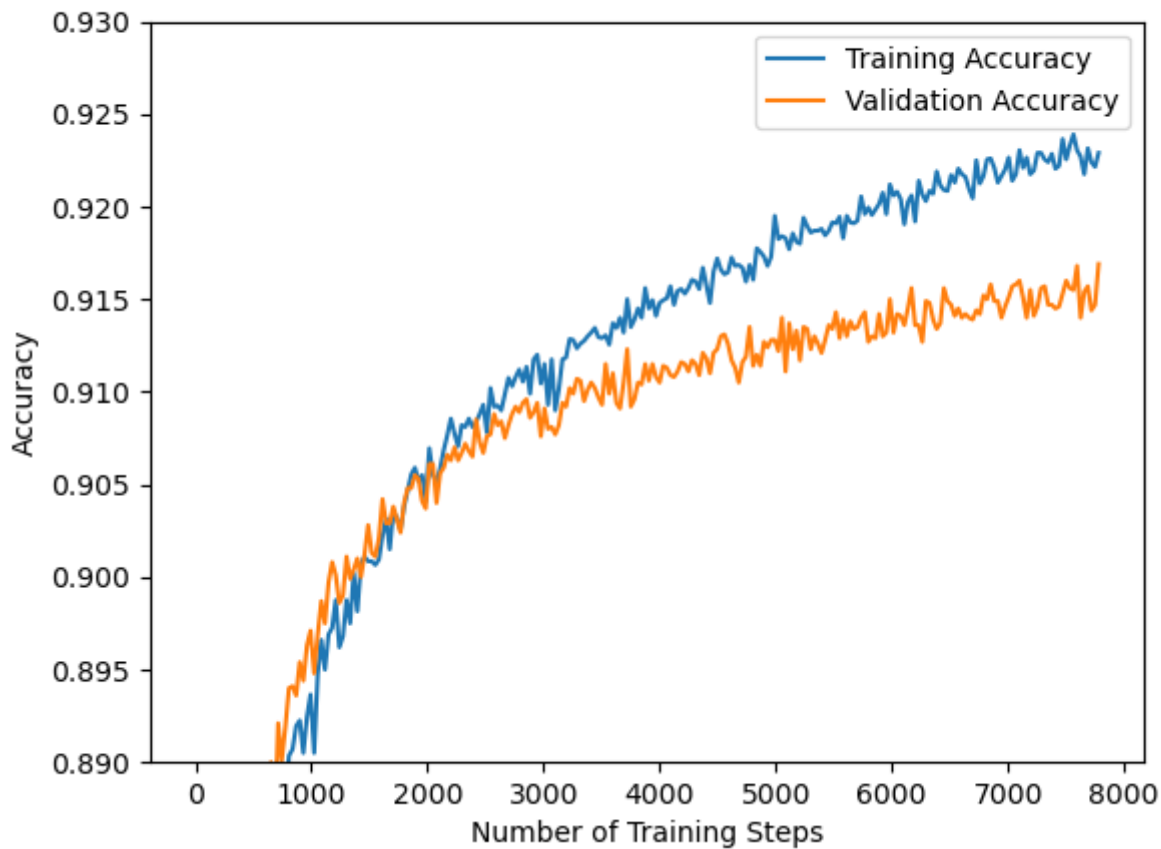


Task 3

Task 3b)



Task 3c)



Task 3d)

The training accuracy is growing at a much faster rate than the validation accuracy. Although the validation accuracy is still rising it may start to fall off as the training accuracy continues to grow. This is because the neural net is very accurately fitting itself to the training set while not generalizing well to the validation set.

Task 4

Task 4a)

4)

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

$$R(w) = \|w\|^2 = \sum_{ij} w_{ij}^2$$

$$\frac{\partial J(w)}{\partial w} = \frac{\partial C(w)}{\partial w} + \lambda \frac{\partial R(w)}{\partial w}$$

$$\frac{\partial C(w)}{\partial w} = \frac{\partial}{\partial w} \left(\frac{1}{N} \sum_{n=1}^N C^n(w) \right) = \frac{1}{N} \sum \frac{\partial C^n(w)}{\partial w}$$

$$\frac{\partial R(w)}{\partial w} = \frac{\partial}{\partial w} \sum w_{ij}^2 = \begin{bmatrix} 2w_{00} & \dots & 2w_{0j} \\ \vdots & \ddots & \vdots \\ 2w_{i0} & \dots & 2w_{ij} \end{bmatrix}$$

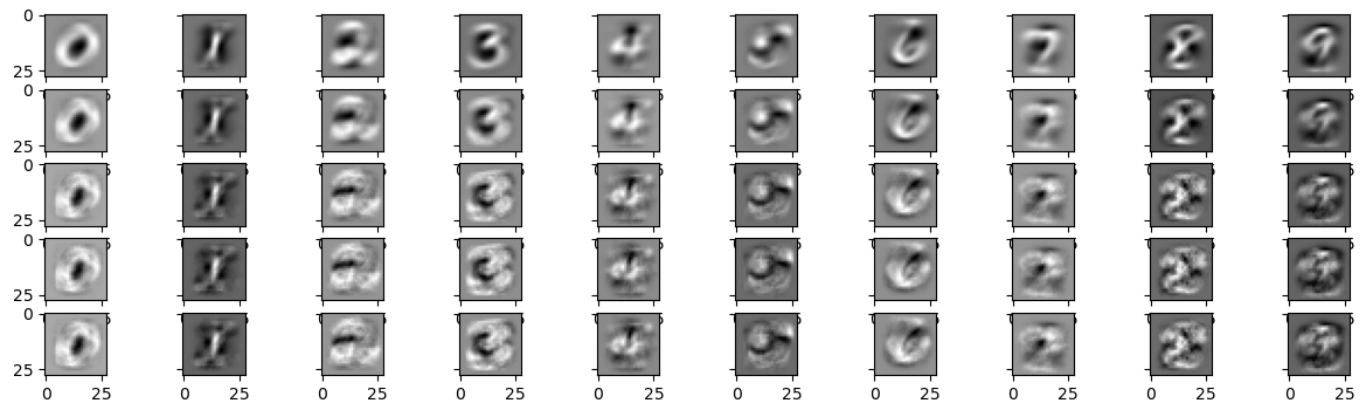
$$\frac{\partial J(w)}{\partial w} = \frac{\partial C(w)}{\partial w} + \lambda \cdot w$$

Same as
before

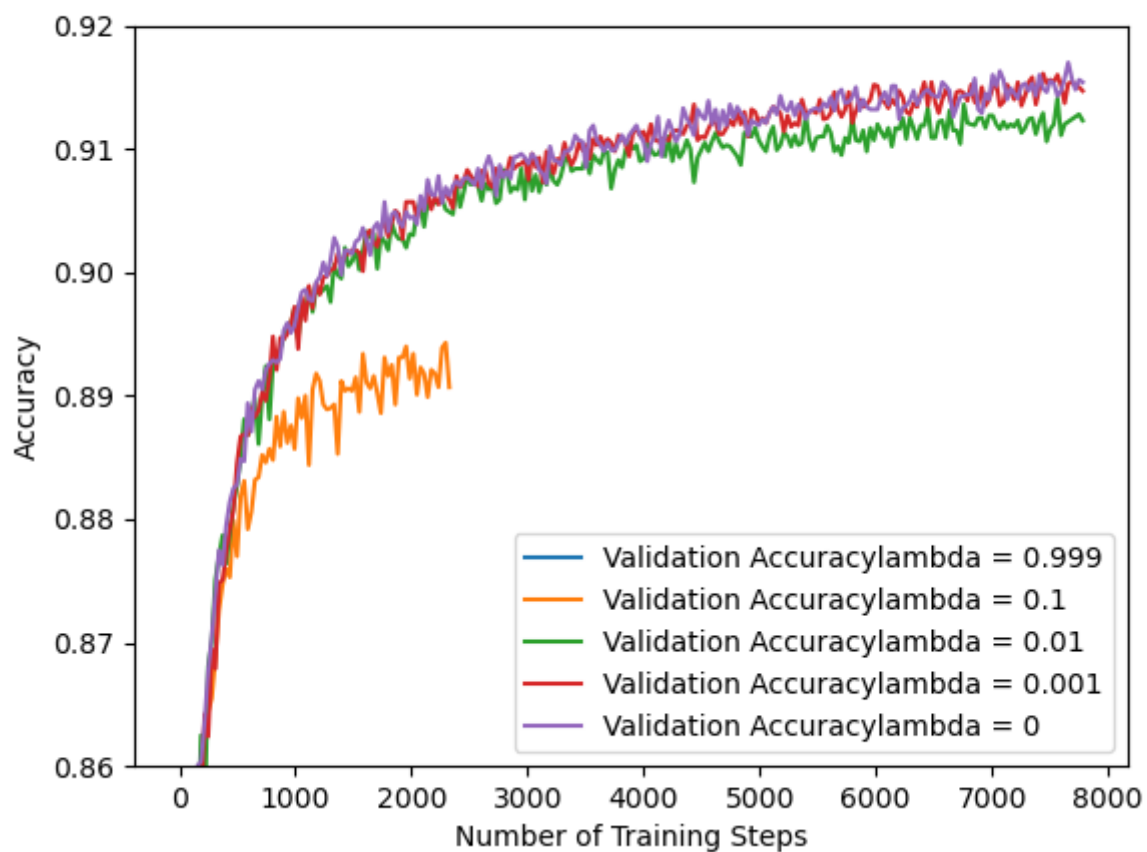
the 2 is absorbed into the λ

Task 4b)

Visualization of weights using $\lambda = [1, 0.1, 0.01, 0.001, 0]$. Zero at the bottom



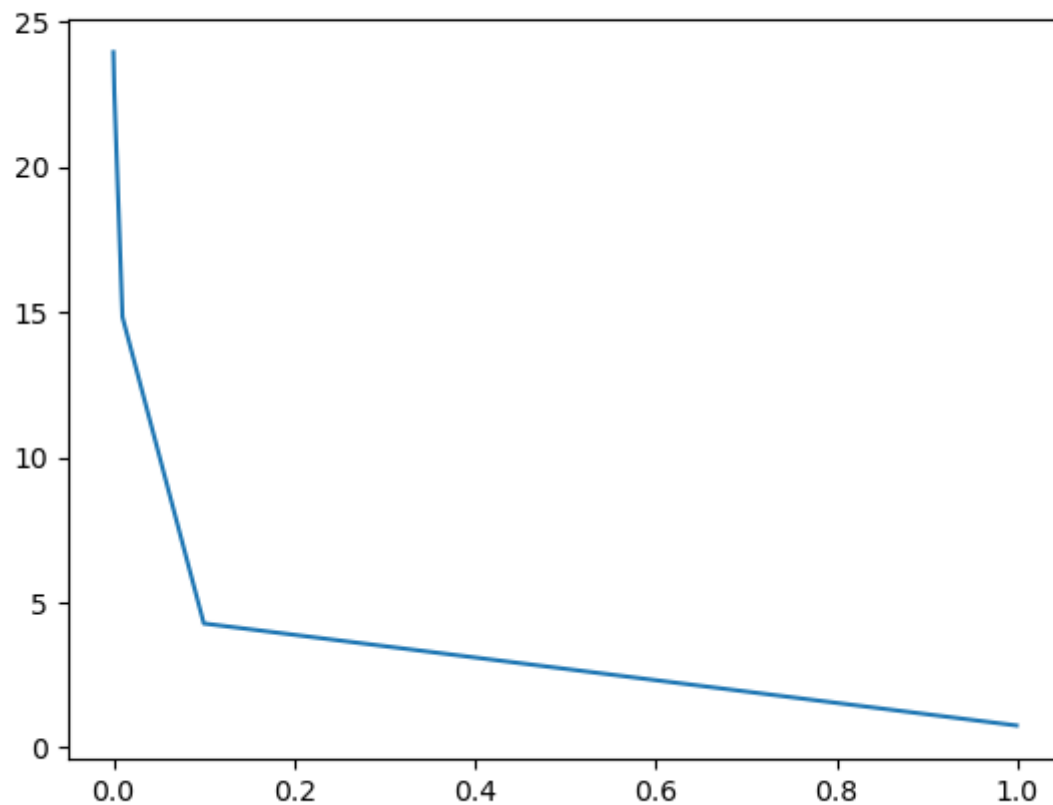
Task 4c)



Task 4d)

It makes sense that the validation accuracy drops off whenever we use l2 regularization as the weights that were trained using what is currently the validation set will diminish due to the weights having a cost.

Task 4e)



X axis is L2 and y axis is the norm of the weights squared.

As expected adding a cost for the size of the weights will reduce the size of the weights. We move in the opposite direction of the gradient which is the opposite direction of an increase in weights. Therefore high values of lambda will result in smaller weights.