# 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 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`.

# Task 1

## task 1a)



## task 1b)

Softmax:

1b) $\quad \hat{y}_k = \dfrac{e^{z_k}}{\sum_{k'}^{K} e^{z_{k'}}} \quad , \quad z_k = w_k^T \cdot x = \sum_i^I w_{k,i} \cdot x_i \quad \Rightarrow \quad \dfrac{\partial z_k^n}{w_{kj}} = x_j^n$

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

$$\dfrac{\partial C^n(w)}{\partial w_{kj}} = \dfrac{\partial C^n}{\partial z_k^n} \dfrac{\partial z_k^n}{\partial w_{kj}}$$

$$= x_j^n \cdot \left( -\sum_{k=1}^{K} y_{kn} \dfrac{\partial}{\partial z_k^n} \left( \ln\left(e^{z^n}\right) \right) \right) = \quad \ln\left(\sum_{k'}^{K} e^{z_{k'}}\right)$$

$$\dfrac{\partial C^n}{\partial w_{kj}} = \dfrac{\partial}{\partial w_{kj}} \left( -\sum_{k''=1}^{K} y_{k''}^n \ln\left(\hat{y}_{k''}^n\right) \right)$$

$$= -\dfrac{\partial}{\partial w_{kj}} \sum_{k''=1}^{K} y_{k''}^n \ln\left( \dfrac{e^{z_{k''}^n}}{\sum_{k'}^{K} e^{z_{k'}^n}} \right)$$

$$\Rightarrow \dfrac{\partial C^n}{\partial w_{kj}} = +\dfrac{\partial}{\partial w_{kj}} \sum_{k''=1}^{K} y_{k''}^n \left( z_{k''}^n - \ln\left( \sum_{k'}^{K} e^{z_{k'}^n} \right) \right)$$

For $k'' = k$ we have:

$$\frac{\partial}{\partial w_{kj}} y_k^n \left( z_k^n - \ln\left( \sum_{h'}^{K} e^{z_{k'}^n} \right) \right)$$

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

For $k'' \neq k$ we have:

$$\frac{\partial}{\partial w_{kj}} \sum_{k'' \neq k}^{K} y_{k''}^n \left( z_{k''}^n - \ln\left( \sum_{k'}^{K} e^{z_{k'}^n} \right) \right)$$

$$= \sum_{h'' \neq h}^{K} 0 - \frac{y_{h''}^n}{\sum_{k'}^{K} e^{z_{k'}^n}} e^{z_k^n} x_j^n = -x_j^n \sum_{h'' \neq k}^{K} y_{h''}^n \hat{y}_{h''}^n$$

$$\Rightarrow + \frac{\partial c^n}{\partial w_{kj}} = -x_j^n \left( y_k^n - \sum_{h''}^{K} y_{h''}^n \hat{y}_{h''}^n \right)$$

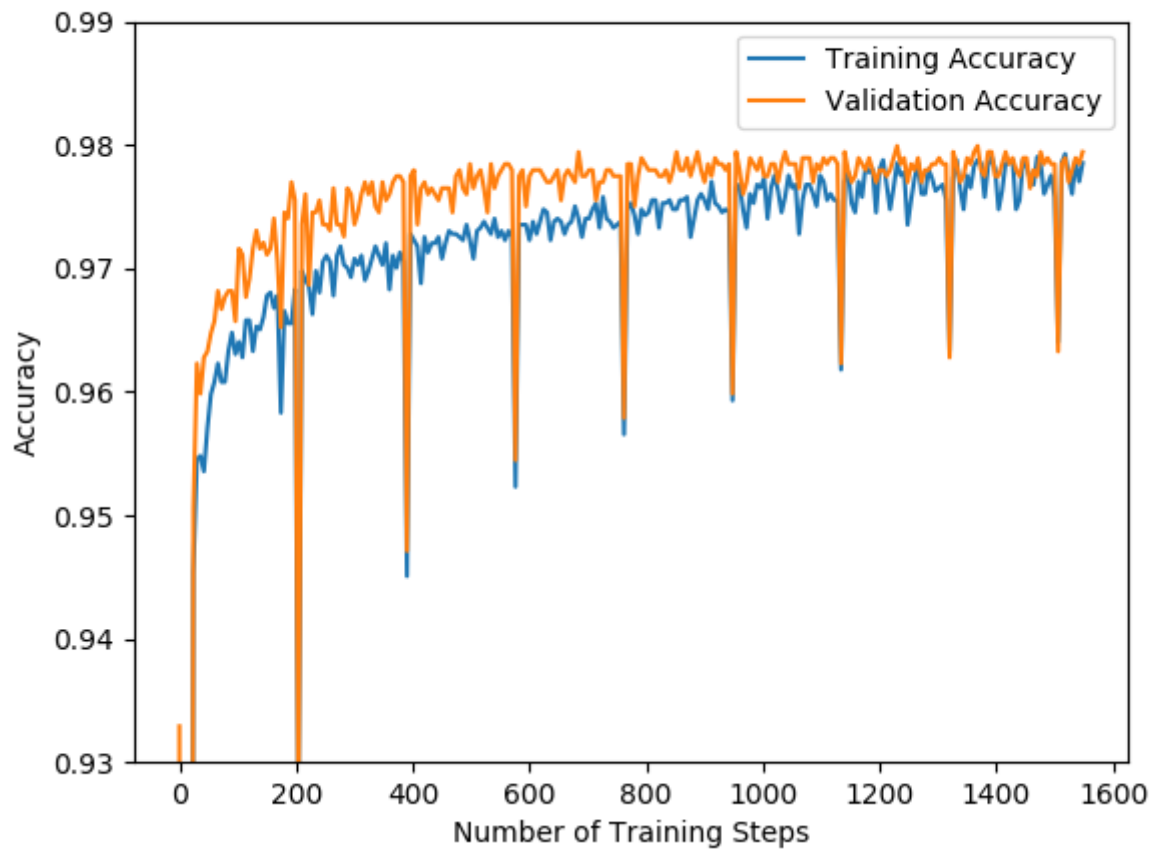$$= -x_i^n \left( y_k^n - \hat{y}_k^n \right)$$
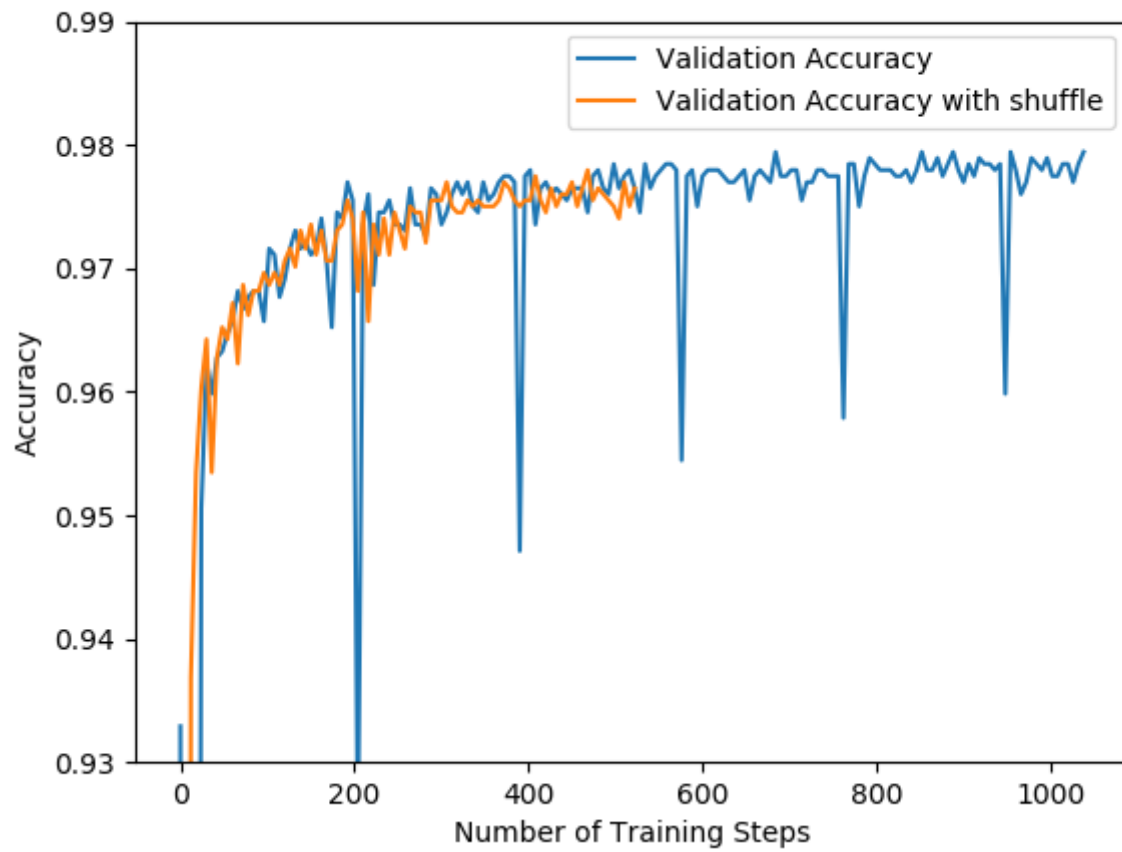
## Task 2

# Task 2b)



# Task 2c)

## Task 2d)

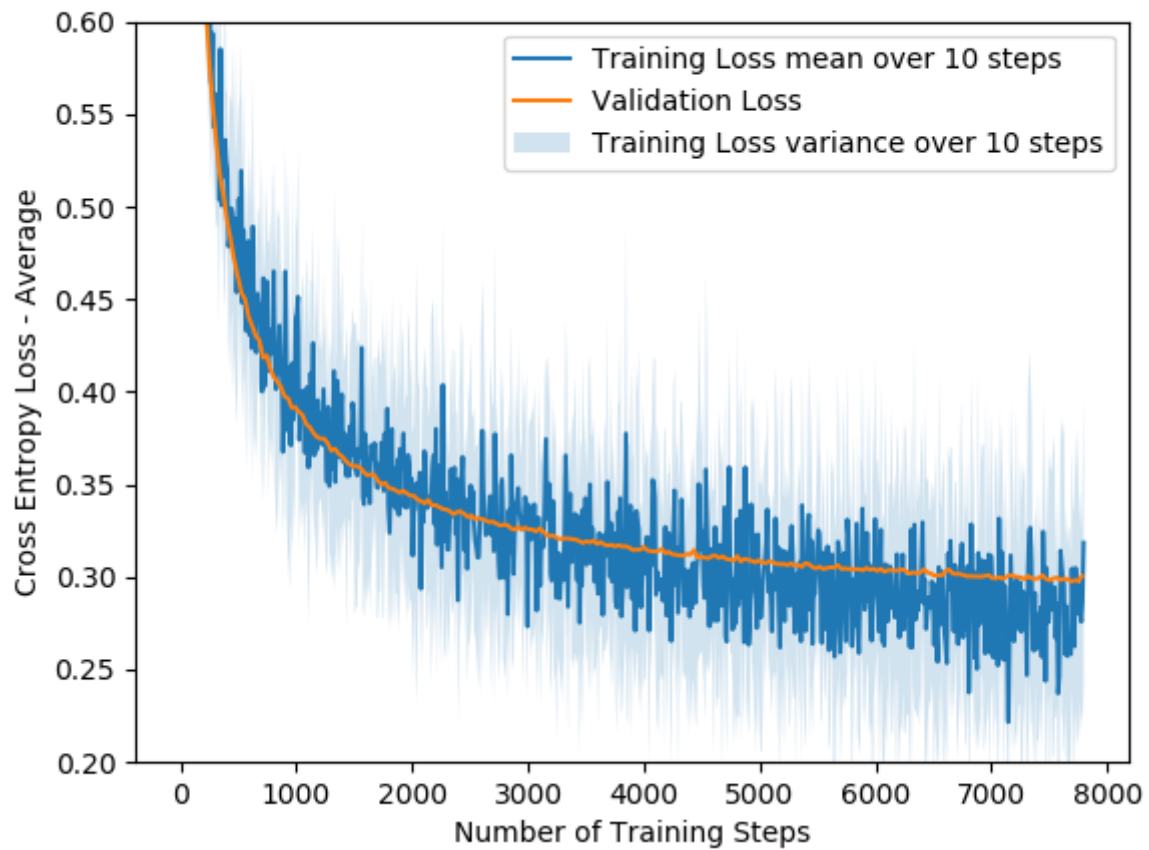The early stop kicked in after 33 epochs

## Task 2e)

The validation accuracy gets less spikes when shuffling because the shuffling makes sure that we dont get the periodically minibatch combo, where one minibatch yields a poor validation accuracy for the other minibatch.
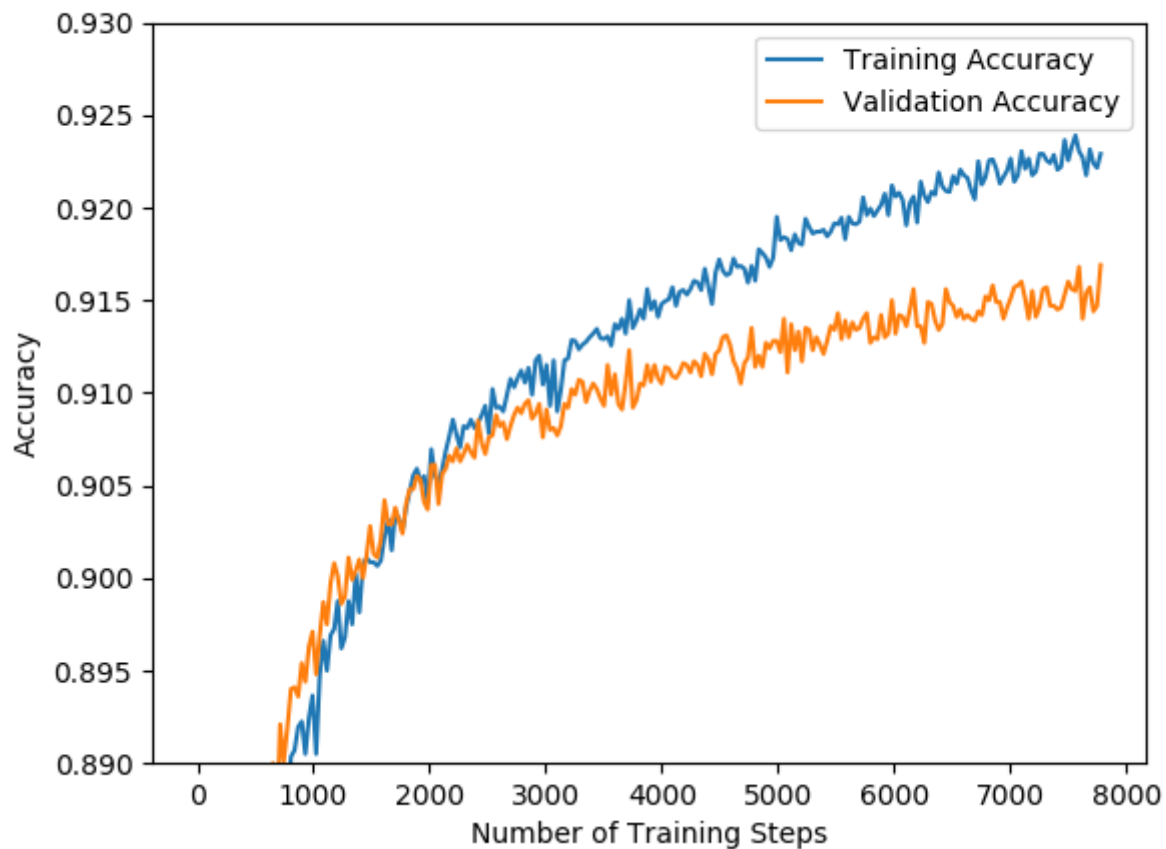
# Task 3

## Task 3b)

## Task 3c)

## Task 3d)
We can see signs of overfitting in the image above, where the validation accuracy starts flattening before the training accuracy.
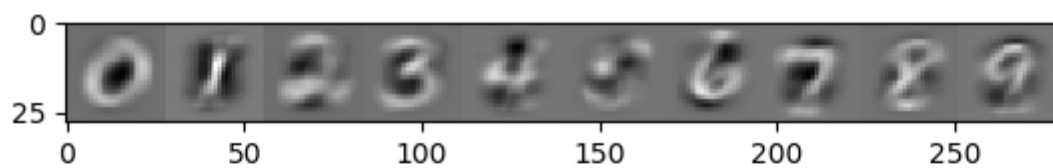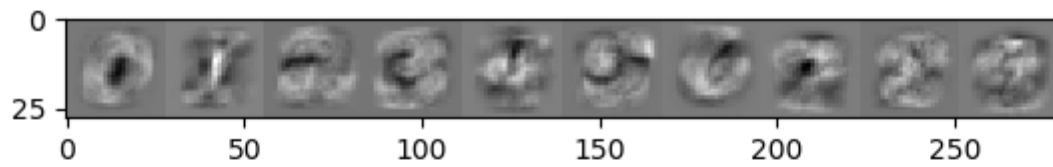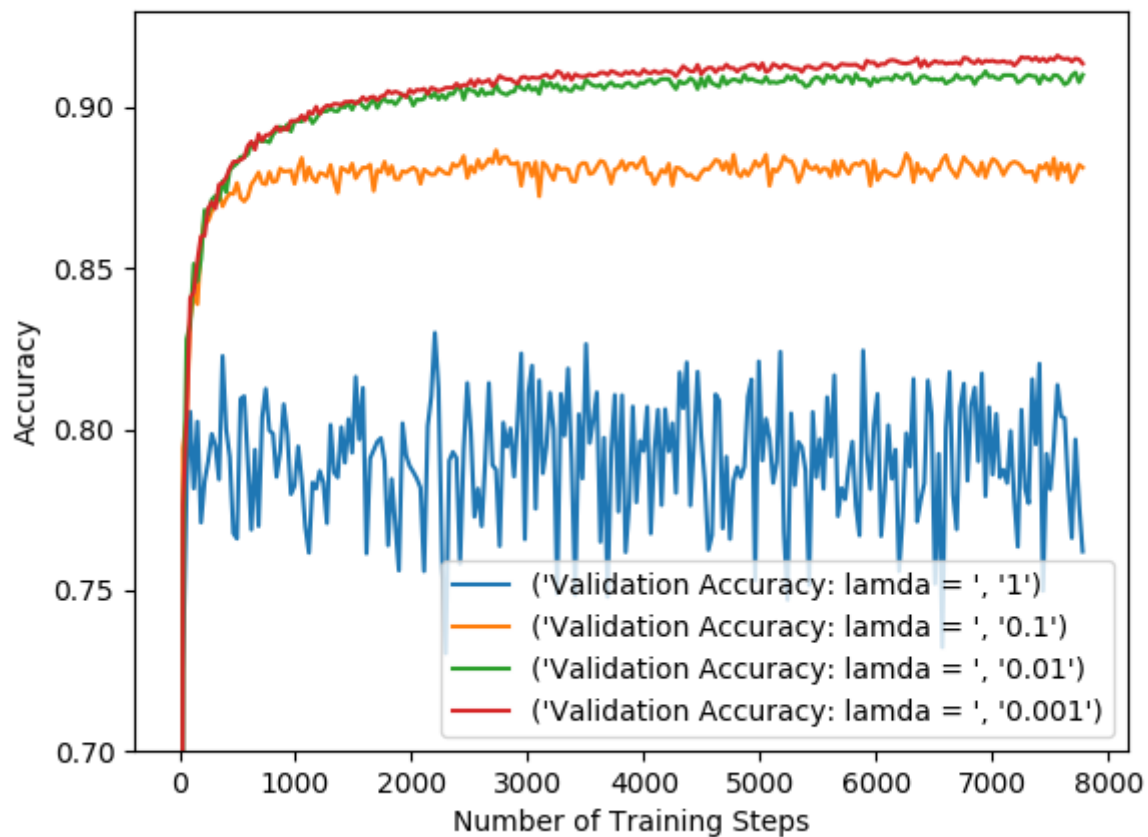
# Task 4

## Task 4a)

$$4_a) \quad J(w) = C(w) + \lambda R(w), \quad R(w) = ||w||^2 = \sum_{i,j} w_{i,j}^2$$

$$(= w^T w)$$

$$\frac{\partial R(w)}{\partial w} = 2w$$

$$\Rightarrow \frac{\partial J(w)}{\partial w} = \frac{\partial C}{\partial w} + \lambda \frac{\partial R}{\partial w} \quad \Rightarrow \quad = \frac{\partial C}{\partial w} + 2\lambda w$$

$$\text{where} \quad \frac{\partial C^n}{\partial w_{kj}} = x_j^n \left( y_k^n - \hat{y}_k^n \right)$$

## Task 4b)

The weights with lambda=1 is less noisy since the L2-regularisation reduces the complexity of the model, making it less smooth, hence giving a less noisy image.

## Task 4c)

## Task 4d)

The validation accuracy drops with lambda since it generalizes a low complexity problem which gives an underfitted result.

## Task 4e)

We see that higher lambda yields a lower L2-norm, which is what one could expect since the penalization is higher.