

# Assignment 1 Report TDT4265

## Group 16

Hanna Waage Hjelmeland, Siri Holde Hegsvold, MTTK

### ▼ Task 1

## ▼ task 1a)

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

$$1. \text{ Show } \frac{\partial C^n}{\partial \omega_i}(\omega) = -(y_i^n - \hat{y}_i^n) x_i^n$$

$$\text{Hint: } \frac{\partial f(x^n)}{\partial \omega_i} = x_i^n f(x^n) (1 - f(x^n)) = x_i^n \hat{y}_i^n (1 - \hat{y}_i^n)$$

$$\frac{\partial C^n}{\partial \omega_i}(\omega) = \frac{\partial}{\partial \omega_i} \left( \underbrace{- (y_i^n \ln(\hat{y}_i^n))}_I + \underbrace{(1 - y_i^n) \ln(1 - \hat{y}_i^n)}_{II} \right)$$

$$\text{Recall: } \hat{y}_i^n = f(x_i^n).$$

start with I:

$$\begin{aligned} \frac{\partial}{\partial \omega_i} (-y_i^n \ln(\hat{y}_i^n)) &= -y_i^n \frac{1}{\hat{y}_i^n} \cdot \frac{\partial}{\partial \omega_i} \hat{y}_i^n \\ &= -\frac{y_i^n}{\hat{y}_i^n} \frac{\partial f(x^n)}{\partial \omega_i} = -\frac{y_i^n}{\hat{y}_i^n} \cdot x_i^n \hat{y}_i^n (1 - \hat{y}_i^n) \\ &= -y_i^n x_i^n (1 - \hat{y}_i^n) \end{aligned}$$

$$\begin{aligned}
 \text{II: } \frac{\partial}{\partial \omega_i} (1 - y^n) \ln(1 - \hat{y}^n) &= -(1 - y^n) \frac{1}{(1 - \hat{y}^n)} \left( - \frac{\partial \hat{y}^n}{\partial \omega_i} \right) \\
 &= - \frac{1 - y^n}{1 - \hat{y}^n} \left( - x_i^n \hat{y}^n (1 - \hat{y}^n) \right) \\
 &= -(1 - y^n) (-x_i^n \hat{y}^n)
 \end{aligned}$$

$$\begin{aligned}
 \Rightarrow \frac{\partial C^n}{\partial \omega_i}(\omega) &= \text{I} + \text{II} = -y^n x_i^n (1 - \hat{y}^n) - (1 - y^n) (-x_i^n \hat{y}^n) \\
 &= -x_i^n (y^n (1 - \hat{y}^n) - \hat{y}^n (1 - y^n)) \\
 &= -x_i^n (y^n - \hat{y}^n y^n - \hat{y}^n + \hat{y}^n y^n) \\
 &= - (y^n - \hat{y}^n) x_i^n \quad \square
 \end{aligned}$$

### ▼ Task 1b)

1. b)

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

$$\hat{y}_k = \frac{e^{\varepsilon_k}}{\sum_{i=1}^K e^{\varepsilon_i}} \quad \varepsilon_k = \omega_k^T x^n$$

$$\hookrightarrow c$$

$$\begin{aligned} \frac{\partial C^n}{\partial \omega_{kj}}(\omega) &= \frac{\partial}{\partial \omega_{kj}} \left( - \sum_{k=1}^K y_k^n \ln(\hat{y}_k^n) \right) \\ &= - \sum_{k=1}^K \frac{y_k^n}{\hat{y}_k^n} \cdot \frac{\partial}{\partial \omega_{kj}} \hat{y}_k^n \end{aligned}$$

$$\frac{\partial}{\partial \omega_{kj}} \hat{y}_k^n = \frac{\partial}{\partial \omega_{kj}} \frac{e^{z_k}}{\sum_c e^{z_c}}$$

Using the quotient rule

$$\frac{d}{d\omega_j} \frac{f(\omega)}{g(\omega)} = \frac{\frac{\partial}{\partial \omega_j} [f(\omega)] g(\omega) - f(\omega) \frac{\partial}{\partial \omega_j} [g(\omega)]}{g^2(\omega)}$$

$$\text{with } f(\omega) = e^{z_k} = e^{\omega_k^T x^n}, g(\omega) = \sum_c e^{\omega_c^T x^n}$$

$$\begin{aligned} c = k : \quad \frac{\partial}{\partial \omega_{kj}} e^{z_k} &= \frac{\partial}{\partial \omega_{kj}} \left( e^{\omega_{k0}^T x^n} \cdot e^{\omega_{k1}^T x^n} \cdot e^{\omega_{kj}^T x^n} \right) \\ &= x_j^n \left( e^{\omega_{k0}^T x^n} \cdot e^{\omega_{k1}^T x^n} \cdot e^{\omega_{kj}^T x^n} \right) = x_j^n e^{\omega_k^T x^n} \end{aligned}$$

$$C \neq k \cdot \frac{\partial}{\partial \omega_{kj}} e^{z_{k'}} = 0$$

$$\text{Let } f(w) = e^{z_k} = e^{\omega_k^T x^n}, g(w) = \sum_{k'=1}^K e^{\omega_{k'}^T x^n}$$

$$\Rightarrow \frac{\partial}{\partial \omega_{kj}} f(w) = x_j^n e^{\omega_k^T x^n}, \frac{\partial}{\partial \omega_{kj}} g(w) = x_j^n e^{\omega_k^T x^n}$$

$$\Rightarrow \frac{\partial}{\partial \omega_{kj}} \frac{f(w)}{g(w)} = \frac{\frac{\partial}{\partial \omega_{kj}} [f(w)] g(w) - f(w) \frac{\partial}{\partial \omega_{kj}} [g(w)]}{g^2(w)}$$

$$= \frac{x_j^n e^{\omega_k^T x^n} \sum_{k'=1}^K e^{\omega_{k'}^T x^n} - e^{\omega_k^T x^n} x_j^n e^{\omega_k^T x^n}}{\left( \sum_{k'=1}^K e^{\omega_{k'}^T x^n} \right)^2}$$

$$= x_j^n \left( \frac{e^{\omega_k^T x^n}}{\sum_{k'=1}^K e^{\omega_{k'}^T x^n}} \cdot \frac{\left( \sum_{k'=1}^K e^{\omega_{k'}^T x^n} - e^{\omega_k^T x^n} \right)}{\sum_{k'=1}^K e^{\omega_{k'}^T x^n}} \right)$$

$$= x_j^n \left( \hat{y}_k^n \cdot (1 - \hat{y}_k^n) \right)$$

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

$$\text{need to find } \frac{\partial}{\partial \omega_{kj}} \hat{y}_{k'}^n = \frac{\partial}{\partial \omega_{kj}} \frac{e^{z_{k'}}}{\sum_{c=1}^K e^{z_c}},$$

i.e. ...  $k' \neq k$ . Using  $f, g$  as

1, c. when  
previously, we obtain

$$\frac{\partial}{\partial \omega_{kj}} f(\omega) = \frac{\partial}{\partial \omega_{kj}} e^{z_{k'}} = 0. \quad e^{z_k} \text{ is contained here}$$

$$\frac{\partial}{\partial \omega_{kj}} g(\omega) = \frac{\partial}{\partial \omega_{kj}} \left( \sum_c^k e^{z_c'} \right) = x_j e^{\omega_k^T x^n}$$

$$\Rightarrow \frac{\partial}{\partial \omega_{kj}} \hat{y}_{k'}^n = \frac{0 - e^{\omega_{k'}^T x^n} x_j e^{\omega_k^T x^n}}{\left( \sum_c^k e^{\omega_{k'}^T x^n} \right)^2}$$

$$= - \frac{x_j \hat{y}_{k'}^n \hat{y}_k^n}{\hat{y}_k^n}$$

$$\Rightarrow \frac{\partial C^n}{\partial \omega_{kj}} = - \sum_{k=1}^K \frac{y_{k'}^n}{\hat{y}_{k'}^n} \frac{\partial \hat{y}_{k'}^n}{\partial \omega_{kj}}$$

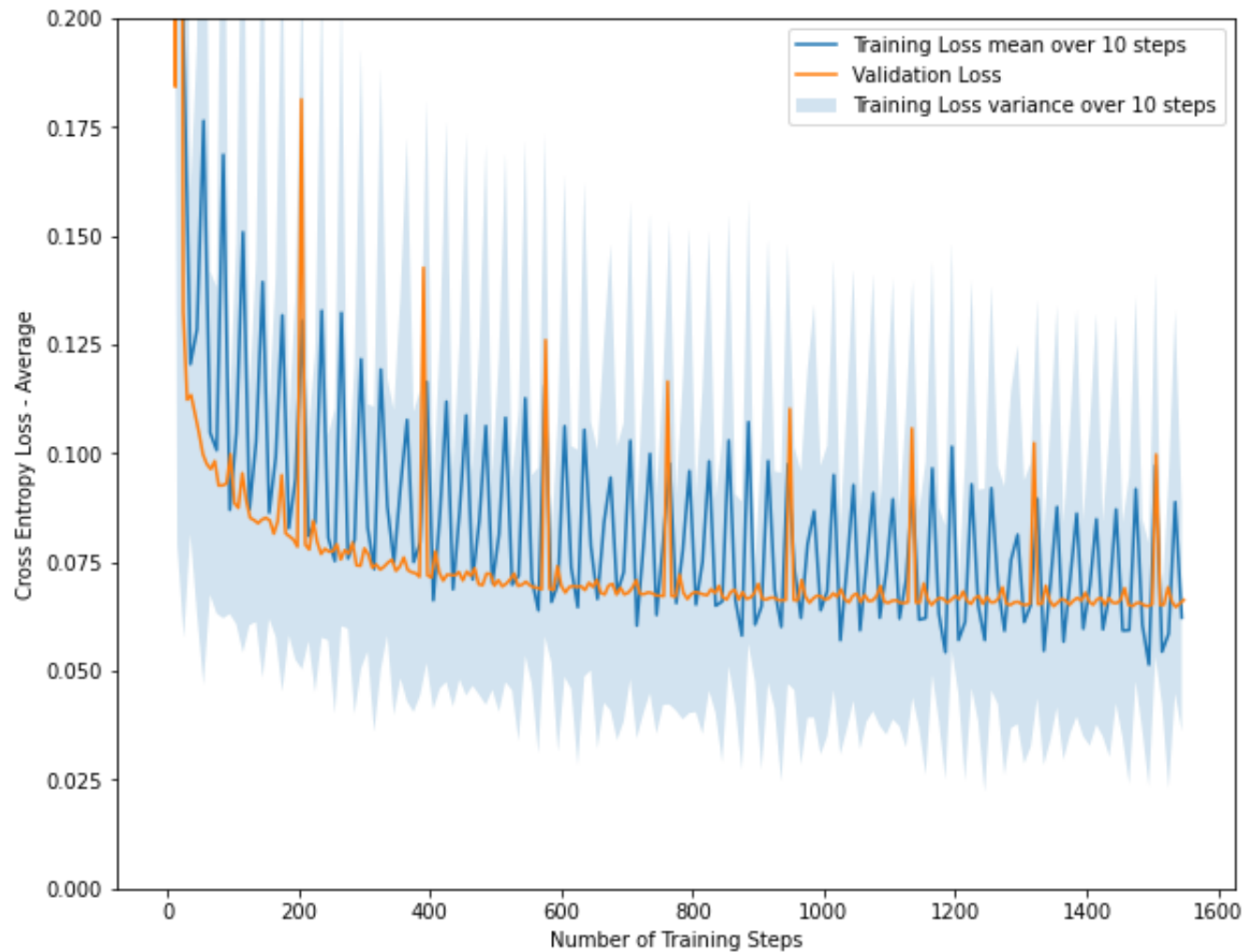
$$= - \left( \sum_{k=k}^n \frac{y_k^n}{\hat{y}_k^n} \frac{\partial \hat{y}_k^n}{\partial \omega_{kj}} + \sum_{k' \neq k}^n \frac{y_{k'}^n}{\hat{y}_{k'}^n} \frac{\partial \hat{y}_{k'}^n}{\partial \omega_{kj}} \right)$$

$$\begin{aligned}
&= - \left( \frac{y_k}{\hat{y}_k^n} \left( x_j^n \left( \hat{y}_k^n \cdot (1 - \hat{y}_k^n) \right) - \sum_{k' \neq k} \frac{y_{k'}}{\hat{y}_{k'}^n} x_j^n \hat{y}_{k'}^n \hat{y}_k^n \right) \right) \\
&= \sum_{k' \neq k} y_{k'}^n x_j^n \hat{y}_k^n - x_j^n y_k^n (1 - \hat{y}_k^n) \\
&= -x_j^n \left( y_k^n - \left( y_k^n \hat{y}_k^n + \sum_{k' \neq k} y_{k'}^n \hat{y}_k^n \right) \right) \\
&= -x_j^n \left( y_k^n - \sum_{k'=1}^K y_{k'}^n \hat{y}_k^n \right) = -x_j^n \left( y_k^n - \hat{y}_k^n \sum_{k'=1}^K y_{k'}^n \right) \\
&= \underline{\underline{-x_j^n (y_k^n - \hat{y}_k^n)}} \quad \square
\end{aligned}$$

## ▼ Task 2

## ▼ Task 2b)

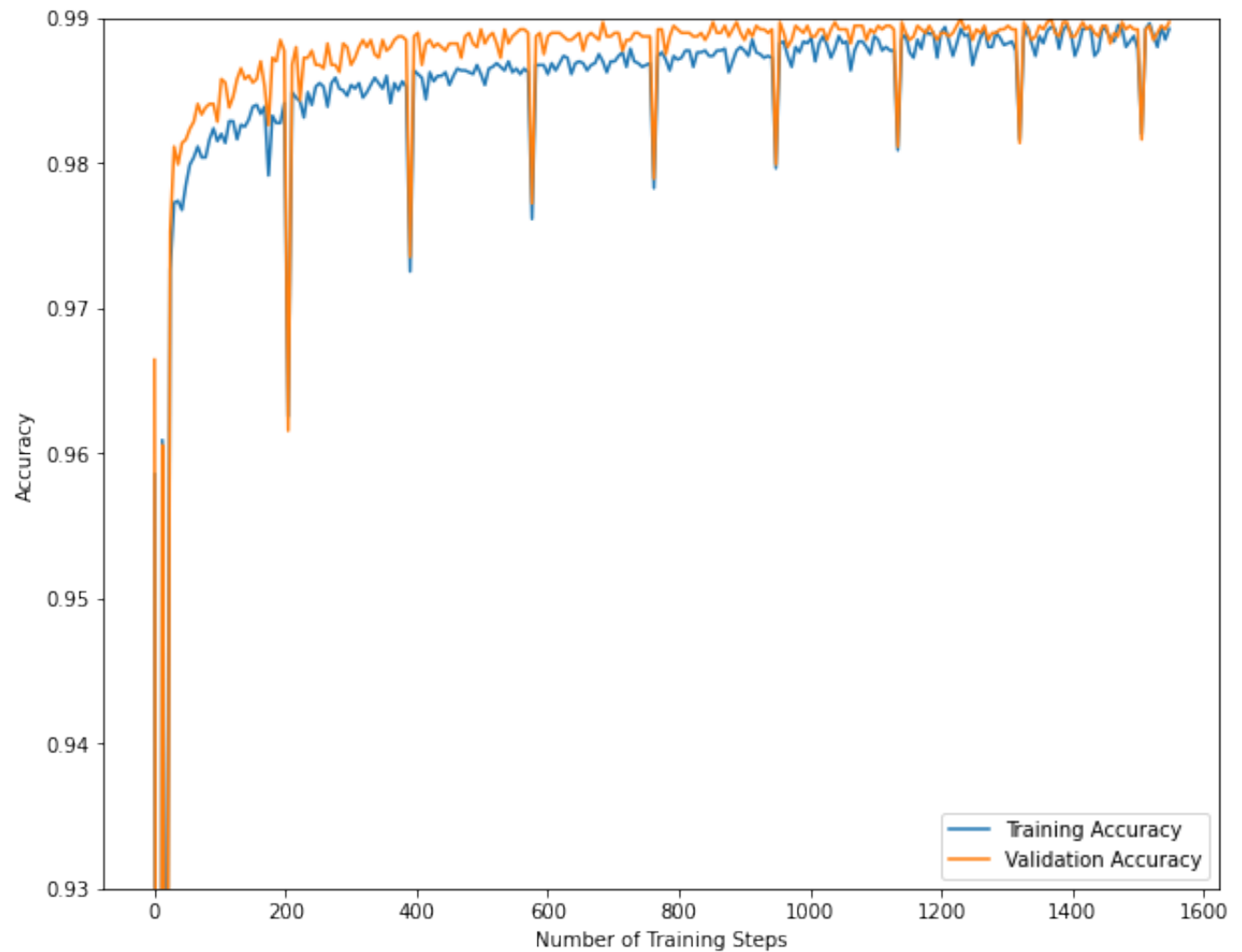
Train shape: X: (4005, 784), Y: (4005, 1)  
Validation shape: X: (2042, 784), Y: (2042, 1)  
Final Train Cross Entropy Loss: 0.06509106268432195  
Final Validation Cross Entropy Loss: 0.06522830574814967  
Train accuracy: 0.9792759051186017  
Validation accuracy: 0.9789422135161606



## Task 2c)

- ▼ Implemented `calculate_accuracy`. Resulting plot:





## Task 2d)

- With early stopping, we get:

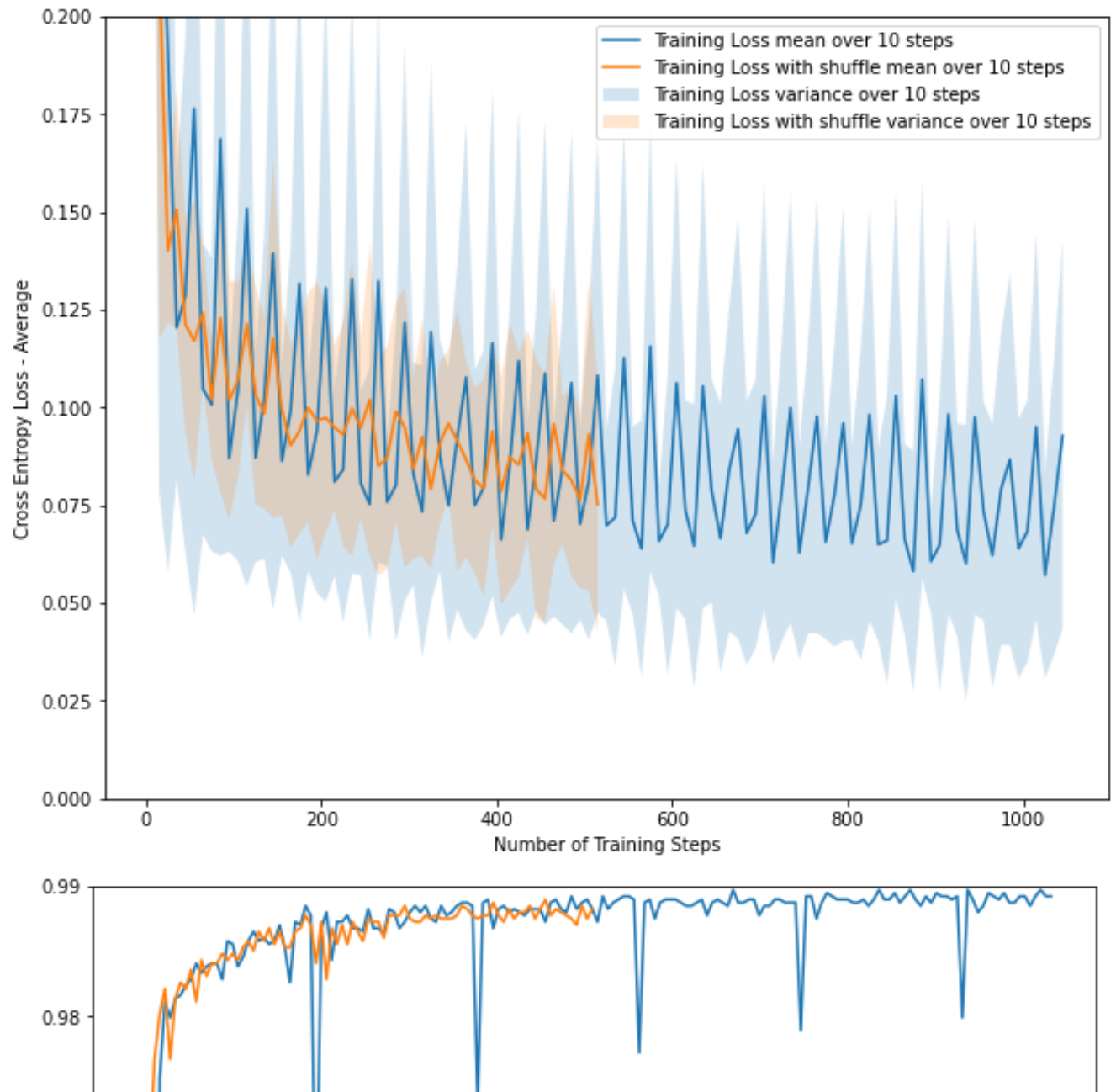
```
Train shape: X: (4005, 784), Y: (4005, 1)
Validation shape: X: (2042, 784), Y: (2042, 1)
Early stopping at epoch: 33
```

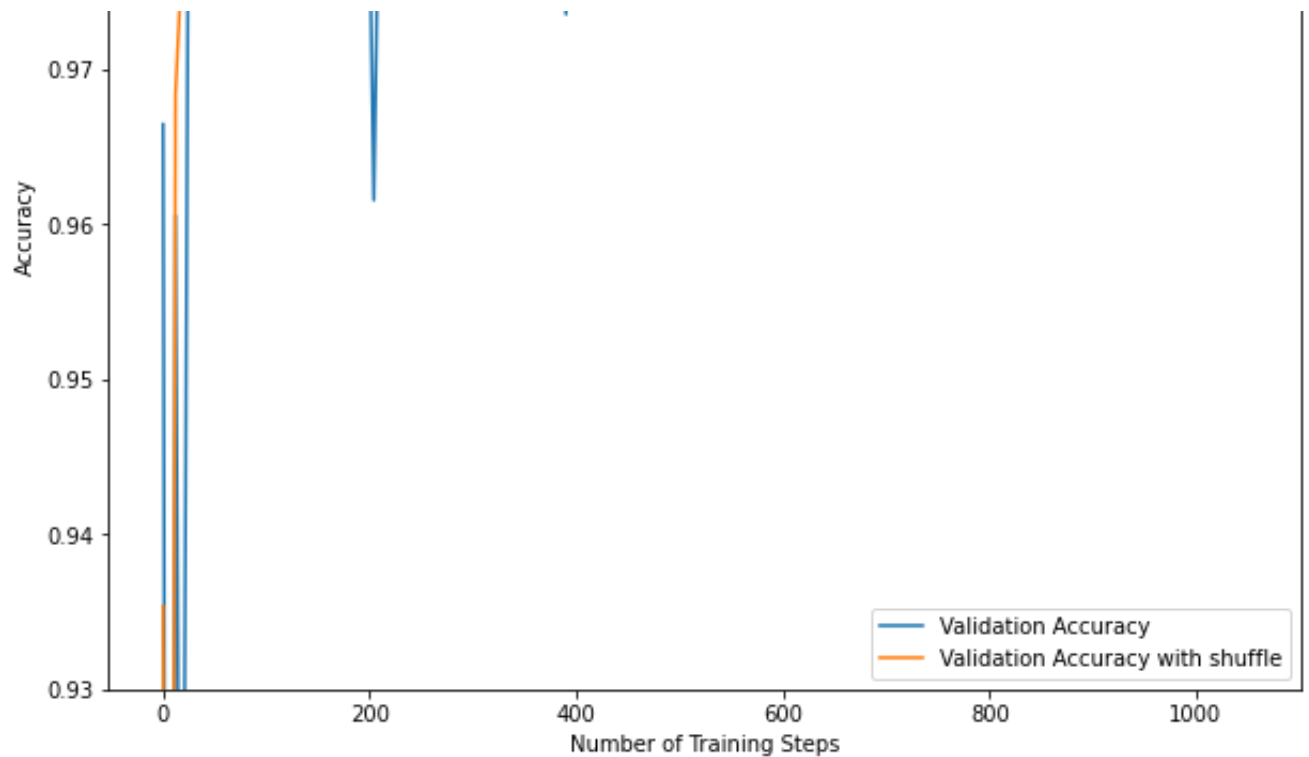
Result is, evidently: Early stopping at epoch 33.

## Task 2e)

### ▼ Result:

Train shape: X: (4005, 784), Y: (4005, 1)  
Validation shape: X: (2042, 784), Y: (2042, 1)  
Early stopping at epoch: 33  
Early stopping at epoch: 16





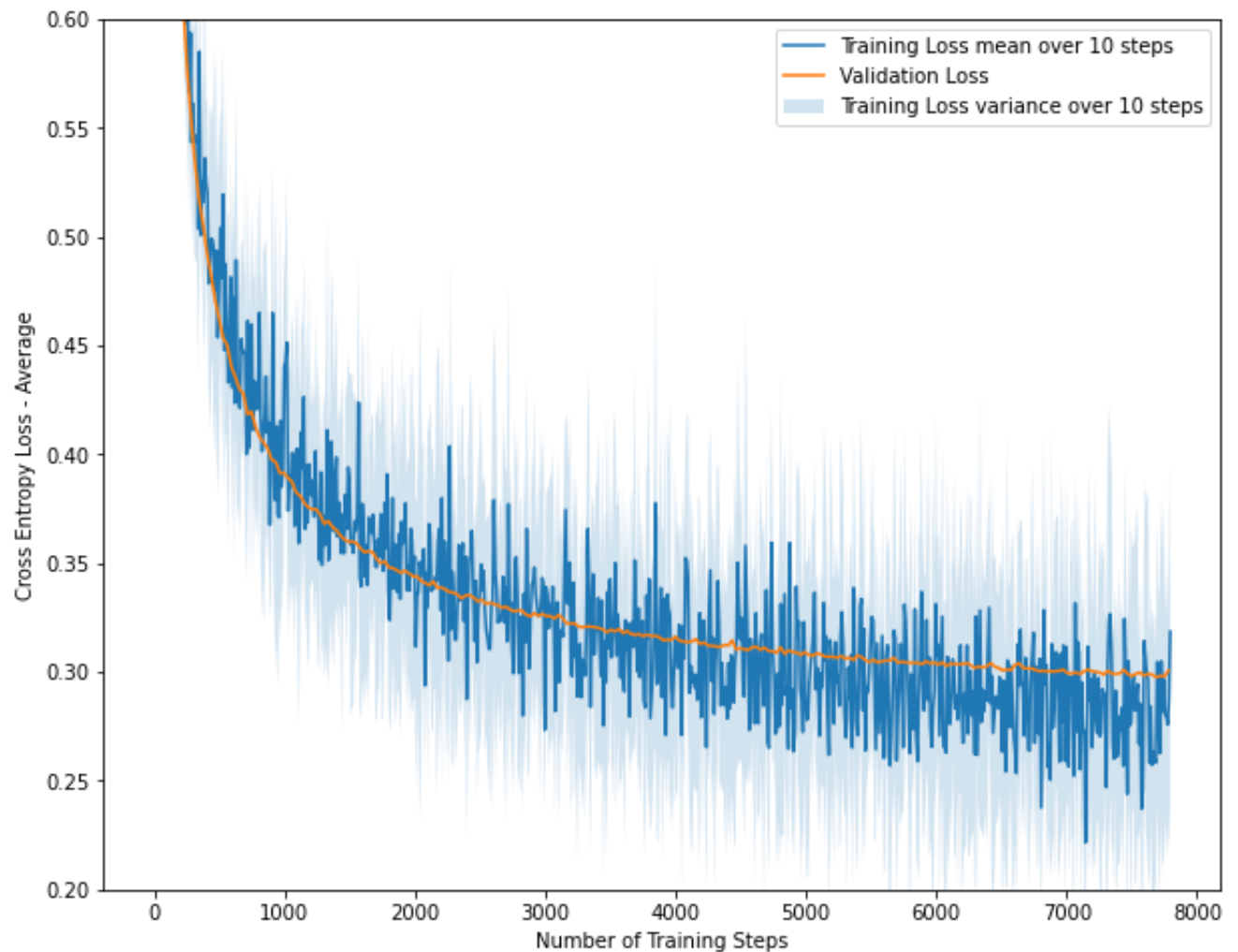
When using shuffle before every epoch, we get less "spikey behaviour". If we do not do this, at every epoch the network tries the updated weights on the same set as before, and these naturally fits badly. But, when using shuffling, it is a "new set" that the weights have no relation to, and this leads to the spikes disappearing as we get more generalization from the start. We also see that it leads to faster training, as the early stopping kicks in at 16 as opposed to 33.

### ▼ Task 3

#### Task 3b)

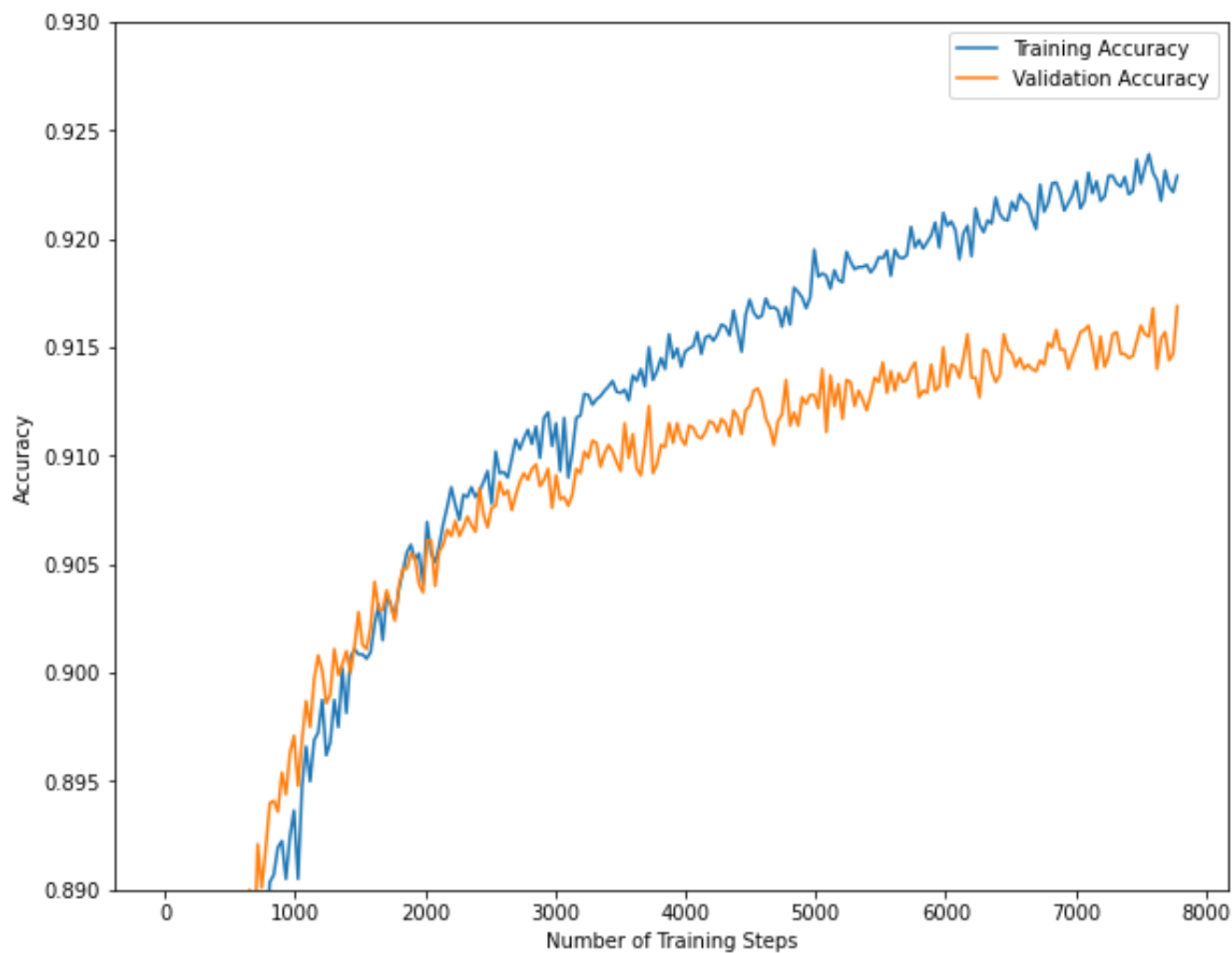
## ▼ Result, Softmax

Final Train Cross Entropy Loss: 0.2800335875350492  
Final Validation Cross Entropy Loss: 0.29902402393704086  
Final Train accuracy: 0.9225  
Final Validation accuracy: 0.9157



Task 3c)

## ▼ Plot accuracy



## Task 3d)

Yes, we see signs of overfitting from around training step 2000. From there, the validation accuracy *does increase*, hence it is not stopped in the early stopping-mechanism, but it increases a lot slower than the training accuracy which indicates that we are on our way to overfitting.

## ▼ Task 4

## Task 4a)

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

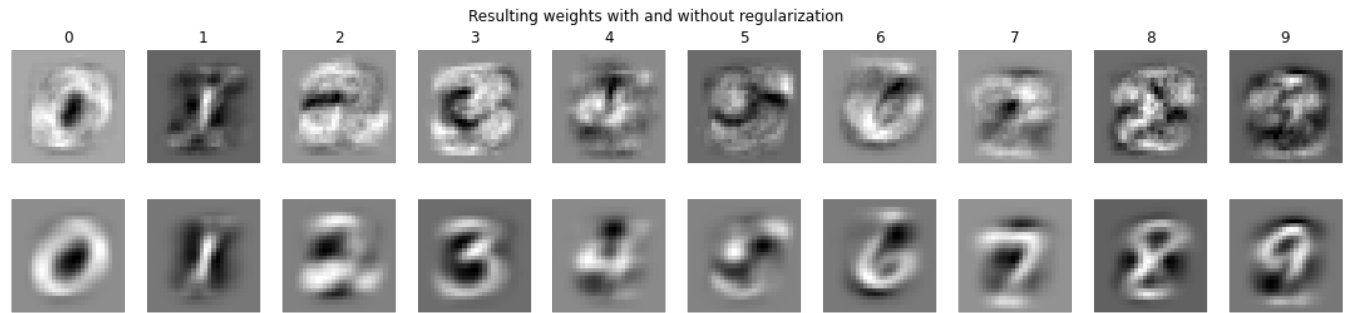
$$\frac{\partial R}{\partial \omega}(\omega) = \frac{\partial}{\partial \omega} \|\omega\|^2 = 2\lambda \omega$$

$$\Rightarrow \frac{\partial J}{\partial \omega}(\omega) = -x_j^{\wedge} (y_k^{\wedge} - \hat{y}_k^{\wedge}) + 2\lambda \omega$$

## Task 4b)

## ▼ Plotting of softmax weights (Task 4b)

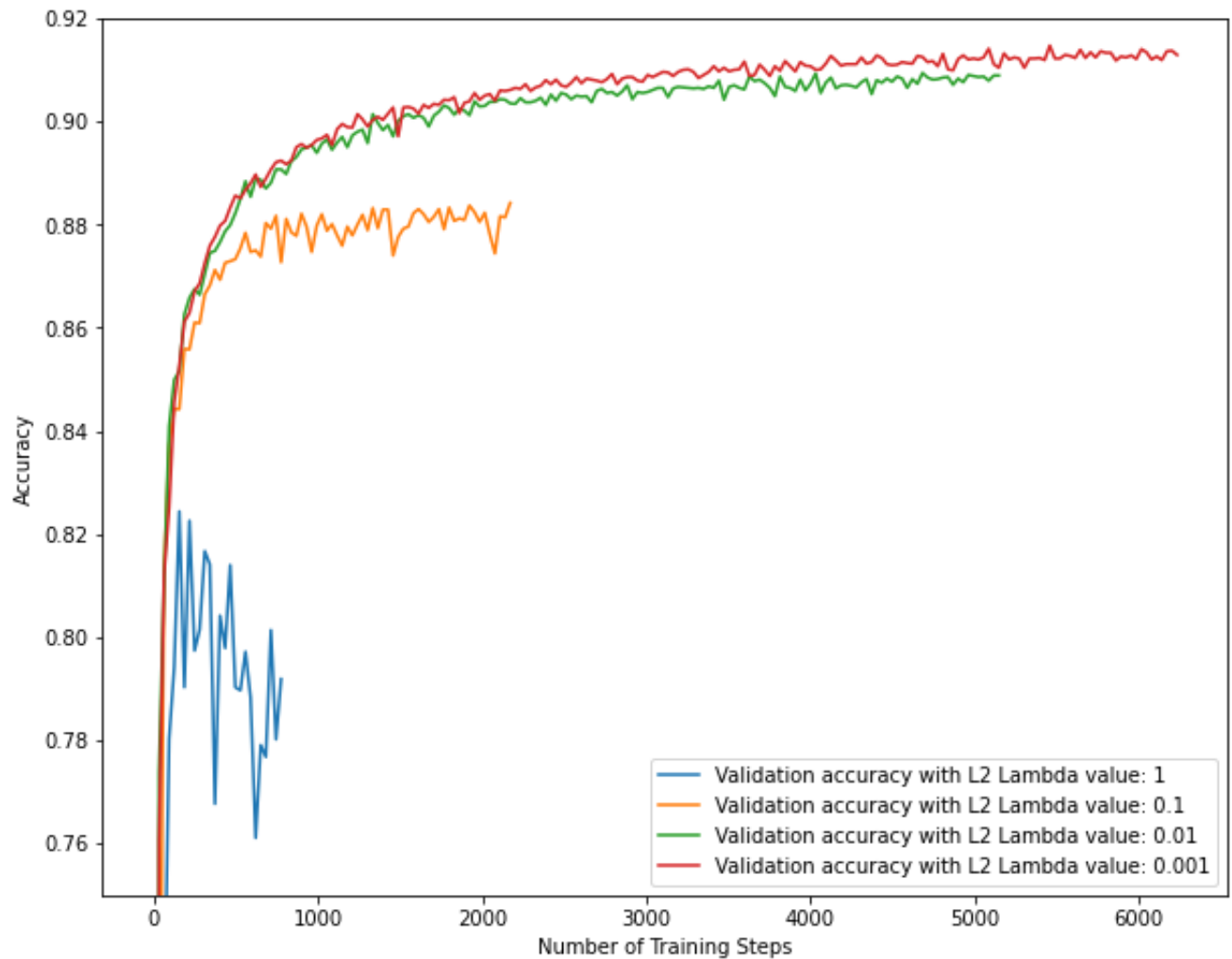
NB: the first row is the original without regularization. The second row is with regularization.



The weights for the model with  $\lambda = 1$  is more noisy. This is due to the fact that large values for the weights are penalized, and thus we do not see hard edges and very definite shapes.

## Task 4c)

- Validation accuracies for different values of  $\lambda$



## ▼ Task 4d)

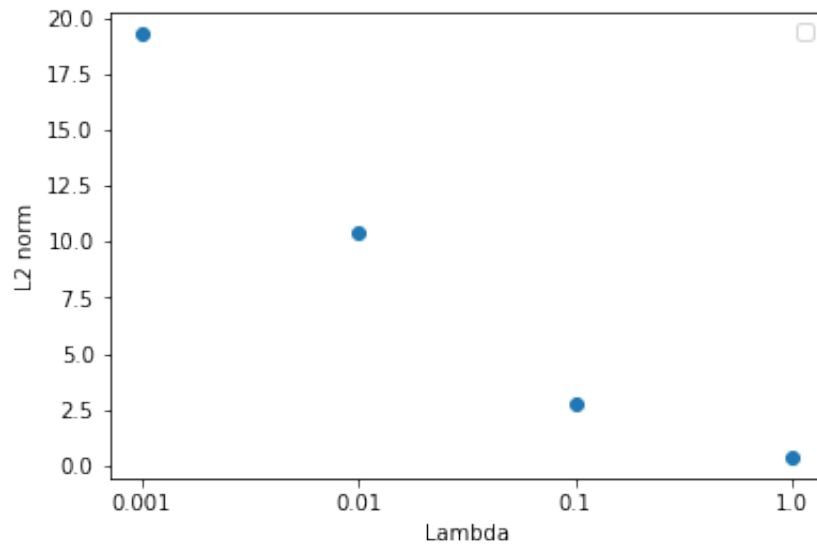
The accuracy degrades when we apply regularization. This can be connected to task 4b), where we see that the weights are more noisy with  $\lambda = 1$ . Thus, as we apply more regularization, the different predictions will get more blurry and, with a higher likelihood, blend in with each other. In other words, it becomes harder for the model to separate the digits.

## Task 4e)



## ▼ Plotting of the l2 norm for each weight

No handles with labels found to put in legend.



We observe, as expected, that the L2 norm of the weights decrease as we introduce more and more penalization of the size of the weights.