

Assignment 1 - Sigmund S. Mestad

February 4, 2022

1 Task 1

1.1 task 1a)

$$\begin{aligned}\frac{\partial C^n(w)}{\partial w_i} &= \frac{\partial}{\partial w_i} \left(- \left(y^n \ln(\hat{y}^n) + (1 - y^n) \ln(1 - \hat{y}^n) \right) \right) \\&= - \left(y^n \frac{\partial}{\partial w_i} \ln(\hat{y}^n) + (1 - y^n) \frac{\partial}{\partial w_i} \ln(1 - \hat{y}^n) \right) \\&= - \left(y^n \frac{1}{\hat{y}^n} \frac{\partial}{\partial w_i} \hat{y}^n + (1 - y^n) \frac{1}{1 - \hat{y}^n} \cdot \left(- \frac{\partial}{\partial w_i} \hat{y}^n \right) \right) \\&= - \left(y^n \frac{1}{\hat{y}^n} \cdot x_i^n \cdot \hat{y}^n (1 - \hat{y}^n) + (1 - y^n) \frac{1}{1 - \hat{y}^n} \cdot (-1) \cdot \hat{y}^n (1 - \hat{y}^n) \right) \\&= - \left(y^n x_i^n (1 - \hat{y}^n) - (1 - y^n) x_i^n \hat{y}^n \right) \\&= - \left(y^n x_i^n - y^n x_i^n \hat{y}^n - x_i^n \hat{y}^n + y^n x_i^n \hat{y}^n \right) \\&= - \left(y^n - \hat{y}^n \right) x_i^n\end{aligned}$$

Task 1a - The gradient for Logistic Regression

1.2 task 1b)

$$\begin{aligned}
 \textcircled{1} \frac{\partial C(\mathbf{w})}{\partial w_{kj}} &= \frac{\partial}{\partial w_{kj}} \left(- \sum_{k=1}^K y_k^n \ln(\hat{y}_k) \right) = - \sum_{k=1}^K y_k^n \frac{\partial}{\partial w_{kj}} \ln \left(\frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}} \right) \\
 &= - \sum_{k=1}^K y_k^n \frac{\partial}{\partial w_{kj}} \left[\ln(e^{z_k}) - \ln \left(\sum_{k'} e^{z_{k'}} \right) \right] \\
 &= - \sum_{k=1}^K \left[y_k^n \frac{\partial}{\partial w_{kj}} \ln(e^{z_k}) - y_k^n \frac{\partial}{\partial w_{kj}} \ln \left(\sum_{k'} e^{z_{k'}} \right) \right] \\
 &= - \left[\underbrace{\sum_{k=1}^K y_k^n \frac{\partial}{\partial w_{kj}} \ln(e^{z_k})}_{\textcircled{2}} - \underbrace{\sum_{k=1}^K y_k^n \frac{\partial}{\partial w_{kj}} \ln \left(\sum_{k'} e^{z_{k'}} \right)}_{\textcircled{3}} \right]
 \end{aligned}$$

But first: $\frac{\partial}{\partial w_{ki}} z_k = \sum_i \frac{\partial}{\partial w_{ki}} w_{ki} x_i = \begin{cases} x_i & \text{for } k=k^* \text{ and } i=j \\ 0 & \text{else} \end{cases}$

$\Rightarrow \frac{\partial}{\partial w_{kj}} z_k = \begin{cases} x_i & \text{for } k=k^* \\ 0 & \text{for } k \neq k^* \end{cases}$

$$\textcircled{2} \frac{\partial}{\partial w_{kj}} \ln(e^{z_k}) = \frac{1}{e^{z_k}} \cdot \frac{\partial}{\partial w_{kj}} e^{z_k} = \frac{1}{e^{z_k}} \cdot e^{z_k} \cdot \frac{\partial}{\partial w_{kj}} z_k = \begin{cases} x_i & \text{for } k=k^* \\ 0 & \text{for } k \neq k^* \end{cases}$$

$$\begin{aligned}
 \textcircled{3} \frac{\partial}{\partial w_{kj}} \ln \left(\sum_{k'} e^{z_{k'}} \right) &= \frac{1}{\sum_{k'} e^{z_{k'}}} \cdot \frac{\partial}{\partial w_{kj}} \sum_{k'} e^{z_{k'}} = \frac{1}{\sum_{k'} e^{z_{k'}}} \cdot \sum_{k'} e^{z_{k'}} \frac{\partial}{\partial w_{kj}} z_{k'} \\
 &= \frac{1}{\sum_{k'} e^{z_{k'}}} \cdot e^{z_{k^*}} \cdot x_i = \hat{y}_{k^*} x_i
 \end{aligned}$$

$\begin{matrix} 0 & \text{for } k' \neq k^* \\ e^{z_{k^*}} x_i & \text{for } k' = k^* \end{matrix}$

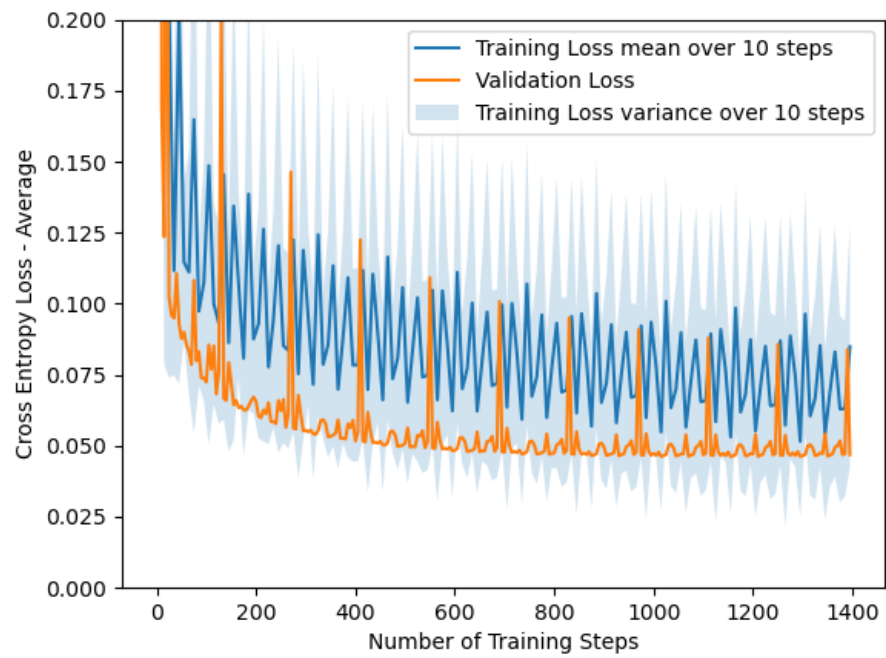
Substituting $\textcircled{2}$ and $\textcircled{3}$ in $\textcircled{1}$:

$$\begin{aligned}
 \textcircled{1} \frac{\partial C(\mathbf{w})}{\partial w_{kj}} &= - \left[\sum_{k=1}^K y_k^n \underbrace{\begin{cases} x_i & \text{for } k=k^* \\ 0 & \text{for } k \neq k^* \end{cases}}_{0 \text{ for } k \neq k^*} - \sum_{k=1}^K y_k^n \underbrace{\hat{y}_{k^*} x_i}_{\substack{\text{Independent of } k \\ = 1}} \right] \\
 &= - (y_{k^*}^n x_i - \hat{y}_{k^*}^n x_i) = \underline{\underline{-x_i (y_{k^*}^n - \hat{y}_{k^*}^n)}}
 \end{aligned}$$

Task 1b - The gradient for Softmax Regression

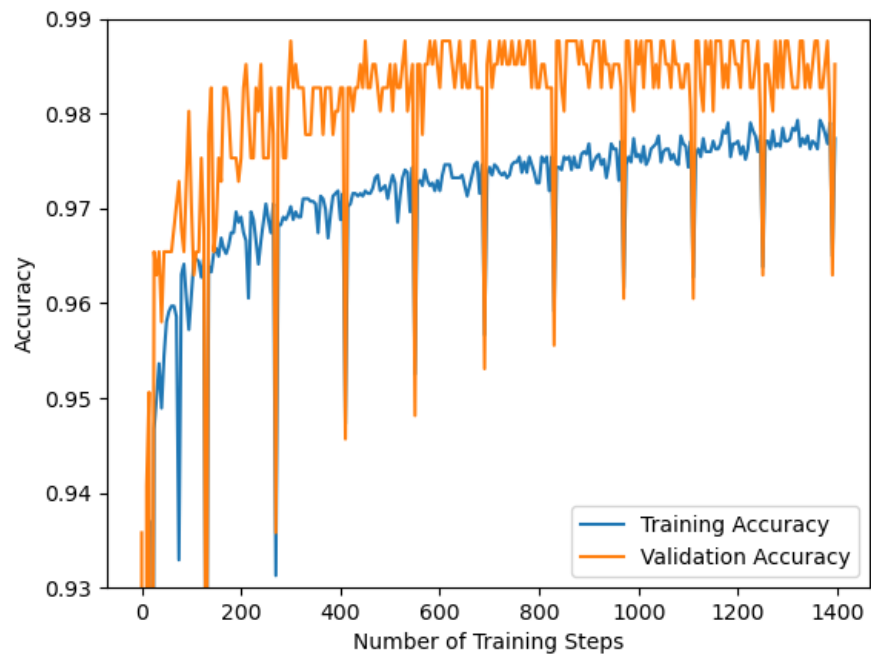
2 Task 2

2.1 Task 2b)



Task 2b - Loss

2.2 Task 2c)

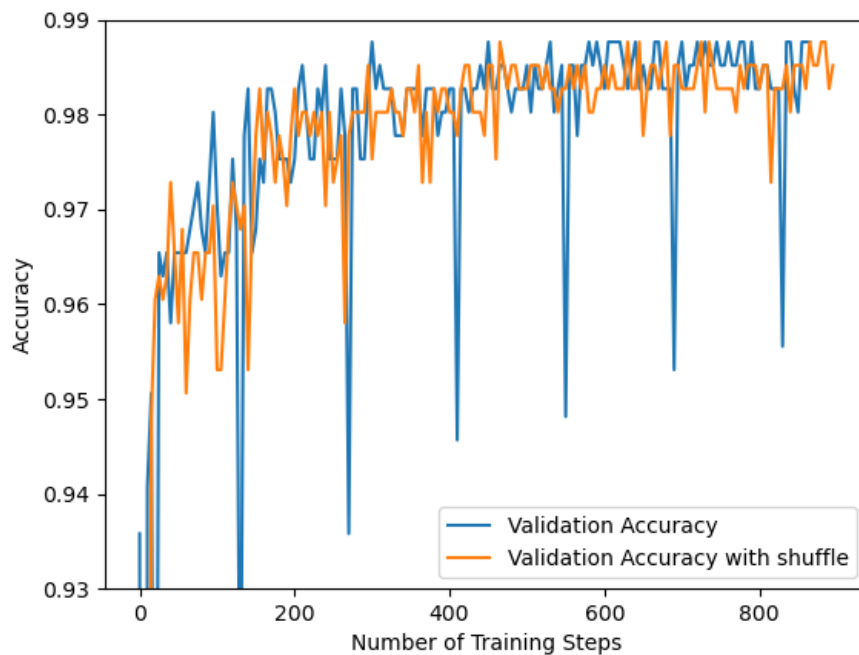
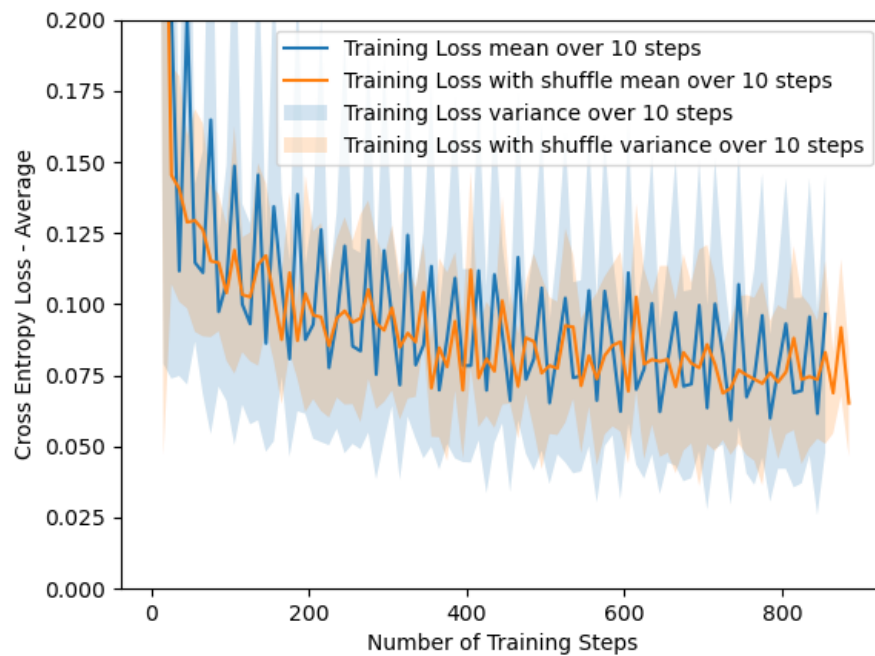


Task 2c - Accuracy

2.3 Task 2d)

Early stop kicks in after 32 epochs

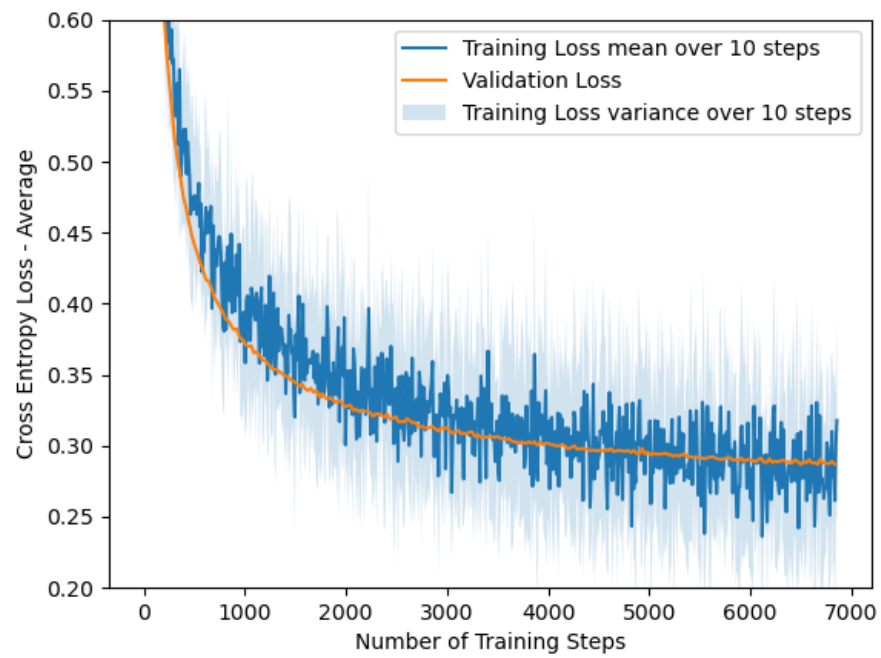
2.4 Task 2e)



The accuracy have less spikes when shuffling because the network is robust and generalized when presented with a new batch of images every time. This makes it less vulnerable for some “bad” batch confusing the network.

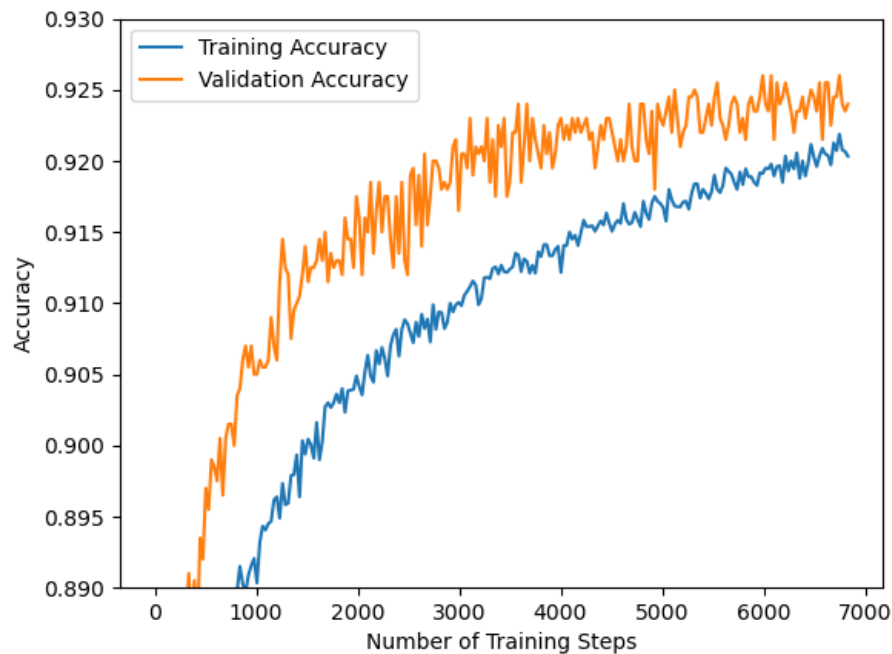
3 Task 3

3.1 Task 3b)



Task 3b - Training and validation loss

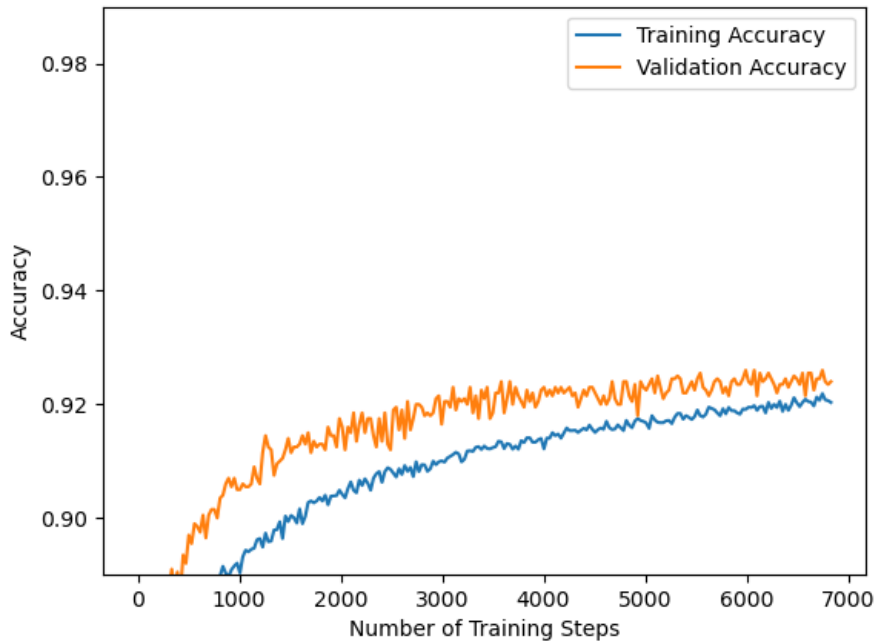
3.2 Task 3c)



Task 3c - Training and validation accuracy

3.3 Task 3d)

The model doesn't show serious signs of overfitting. The validation accuracy seems to reach a plateau at around 40-50 epochs (the early stop kicks in at 49). Doubling the number of epochs, we see that the training accuracy continues to grow, while the validation accuracy is completely flat (see plot below). A possible explanation of why the network doesn't seem to overfit could be because the model is too simple to learn all the details and noise of the training data which normally would lead to worse performance.



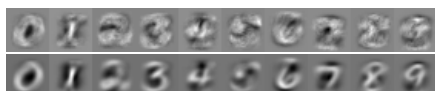
4 Task 4

4.1 Task 4a)

$$\begin{aligned}
 \textcircled{1} \quad \frac{\partial J(w)}{\partial w} &= \frac{\partial C(w)}{\partial w} + \lambda \frac{\partial R(w)}{\partial w} \\
 \textcircled{2} \quad \frac{\partial R(w)}{\partial w} &= \frac{\partial}{\partial w} \frac{1}{2} \sum_{i,j} w_{i,j}^2 = \frac{1}{2} \sum_{i,j} 2w_{i,j} = \sum_{i,j} w_{i,j} \\
 \textcircled{3} \quad \text{We know from task 1b that } \frac{\partial C(w)}{\partial w} &= -x_j^n (y_k^n - \hat{y}_k^n) \\
 \Rightarrow \frac{\partial C(w)}{\partial w} &= \frac{1}{N} \sum_{n=1}^N -x_j^n (y_k^n - \hat{y}_k^n) \\
 \textcircled{1} \quad \frac{\partial J(w)}{\partial w} &= \frac{1}{N} \sum_{n=1}^N -x_j^n (y_k^n - \hat{y}_k^n) + \lambda \sum_{i,j} w_{i,j}
 \end{aligned}$$

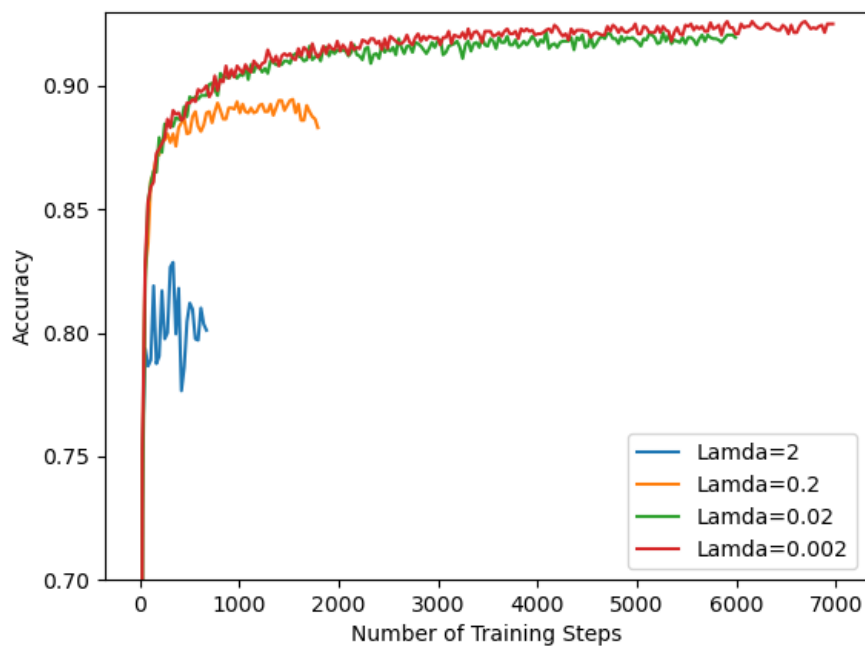
Task 4a - Update term for softmax regression with L2 regularization

4.2 Task 4b)



The regularization term reduces all weights, reducing the effect of less common patterns. Weights with low consensus will diminish.

4.3 Task 4c)

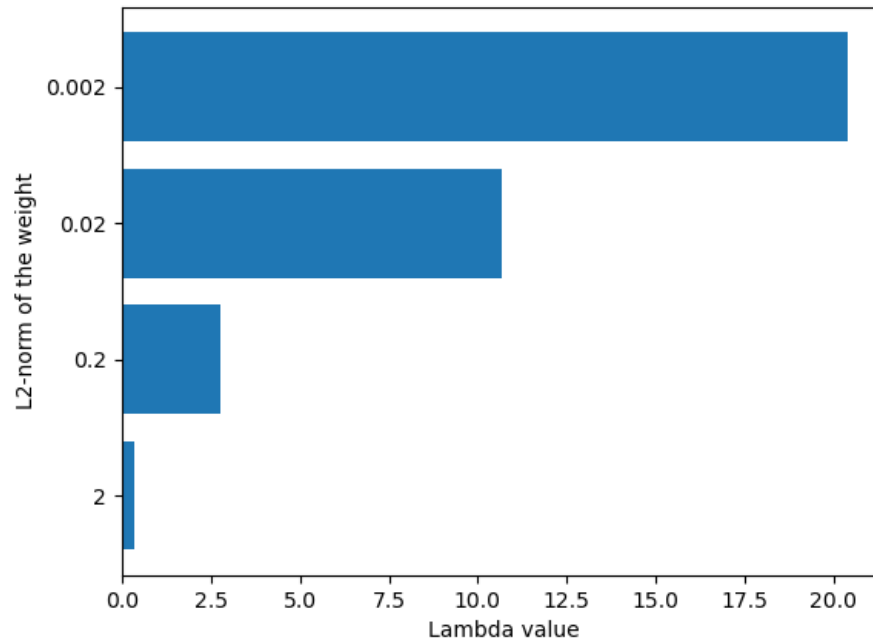


Task 4c - Accuracy with different lambda values

4.4 Task 4d)

The validation accuracy degrades when applying regularization because the model architecture is so simple. Making the model even more general, reduces its capacity to differentiate the numbers.

4.5 Task 4e)



Task 4e - L2-normalization of weights

It's clear that more regularization leads to smaller weights.