

Assignment 1 Report

Task 1

task 1a)

1a

y^n - label
 \hat{y} - result

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

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

$$\hat{y} = f(x) = \frac{1}{1 + e^{-w^T x}} = \frac{1}{1 + e^{-z}} = \sigma(z), \quad z = w^T x$$

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$\frac{\partial C}{\partial \hat{y}} = \frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} = \frac{(1 - y)\hat{y} - y(1 - \hat{y})}{(1 - \hat{y})\hat{y}} = \frac{\hat{y} - y}{(1 - \hat{y})\hat{y}}$$

$$\frac{\partial \hat{y}}{\partial z} = \frac{e^{-z}}{(e^{-z} + 1)^2} = \frac{1}{e^{-z} + 1} \cdot \frac{e^{-z}}{e^{-z} + 1} = \frac{1}{1 + e^{-z}} \left(\frac{e^{-z} + 1}{e^{-z} + 1} - \frac{1}{e^{-z} + 1} \right) = \hat{y}(1 - \hat{y})$$

$$\frac{\partial z}{\partial w} = x$$

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

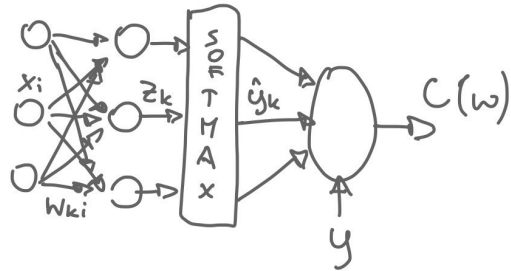
task 1b)

1b.

$$\hat{y}_k = \frac{e^{z_k}}{\sum_{k'}^K e^{z_{k'}}}, \quad z_k = w_k^T x = \sum_i^I w_{ki} \cdot x_i$$

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

$$\frac{\partial C^n}{\partial w_{ki}} = \frac{\partial C^n}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial z_k} \frac{\partial z_k}{\partial w_{ki}}$$



Derivative of \hat{y}_k has 2 cases when $k=k'$ and $k \neq k'$:

$$\text{case 1. } \frac{\partial \hat{y}_k}{\partial z_k} = \frac{e^{z_k} \sum_{k'}^K e^{z_{k'}} - e^{z_k} e^{z_k}}{(\sum_{k'}^K e^{z_{k'}})^2} = \hat{y}_k (1 - \hat{y}_k)$$

$$\text{case 2. } \frac{\partial \hat{y}_k}{\partial z_{k'}} = \frac{0 - e^{z_k} e^{z_{k'}}}{(\sum_{k'}^K e^{z_{k'}})^2} = -\hat{y}_k \hat{y}_{k'}$$

$$\frac{\partial C}{\partial \hat{y}_k} = - \sum_{k=1}^K \frac{y_k}{\hat{y}_k}$$

$$\frac{\partial C}{\partial z_k} = \frac{\partial C}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial z_k} = + \sum_{k' \neq k}^K \frac{y_k}{\hat{y}_k} \hat{y}_k \hat{y}_{k'} - \frac{y_{k'}}{\hat{y}_{k'}} \hat{y}_{k'} (1 - \hat{y}_k)$$

$$= \sum_{k' \neq k}^K y_k \hat{y}_{k'} - y_{k'} + (y_{k'} \hat{y}_{k'})$$

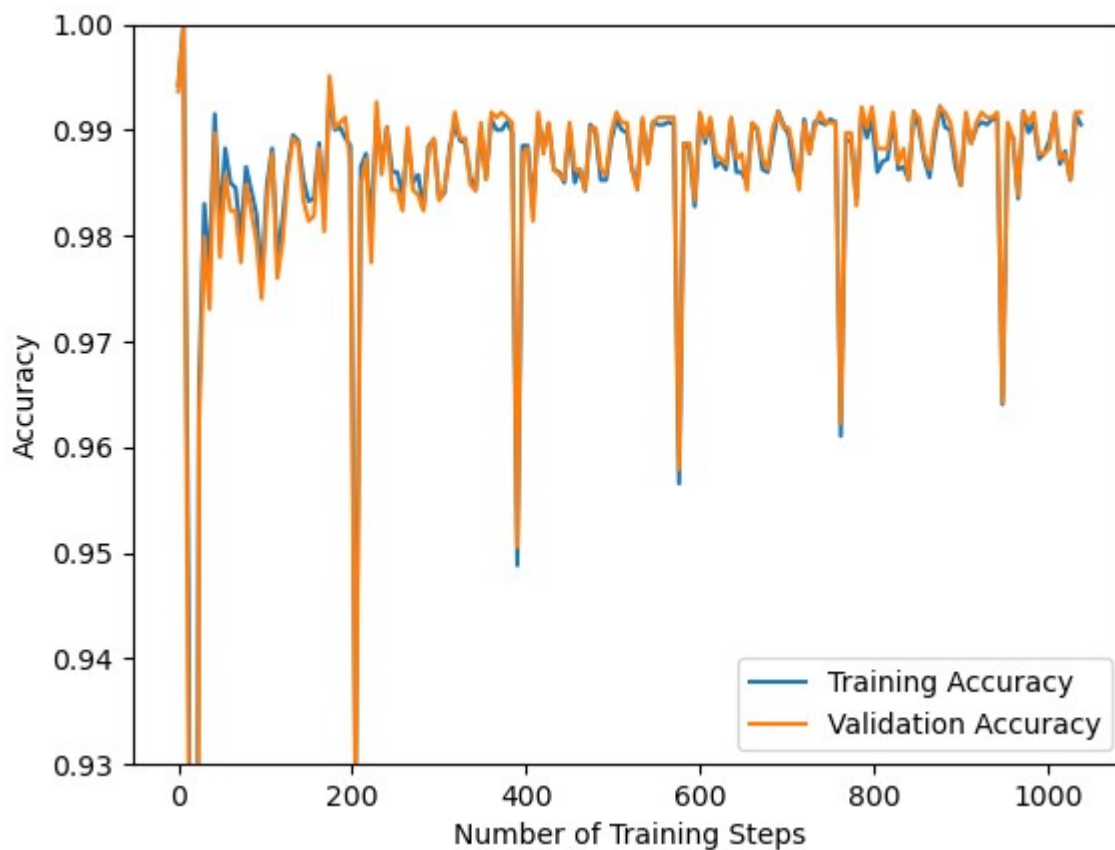
$$= \sum_k^K y_k \hat{y}_{k'} - y_{k'} = \hat{y}_{k'} \sum_{k=1}^K y_k - y_{k'} = \hat{y}_{k'} - y_{k'} = \hat{y}_k - y_k$$

* I think I may have indexed it other way around

$$z_k = \sum_i^I w_{ki} \cdot x_i \xrightarrow{i=j} \frac{\partial z_k}{\partial w_{kj}} = x_j$$

$$\Rightarrow \frac{\partial C^n(w)}{\partial w_{ij}} = \frac{\partial C^n}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial z_k} \frac{\partial z_k}{\partial w_{ij}} = x_i (\hat{y}_k - y_k) = -x_i (y_k - \hat{y}_k)$$

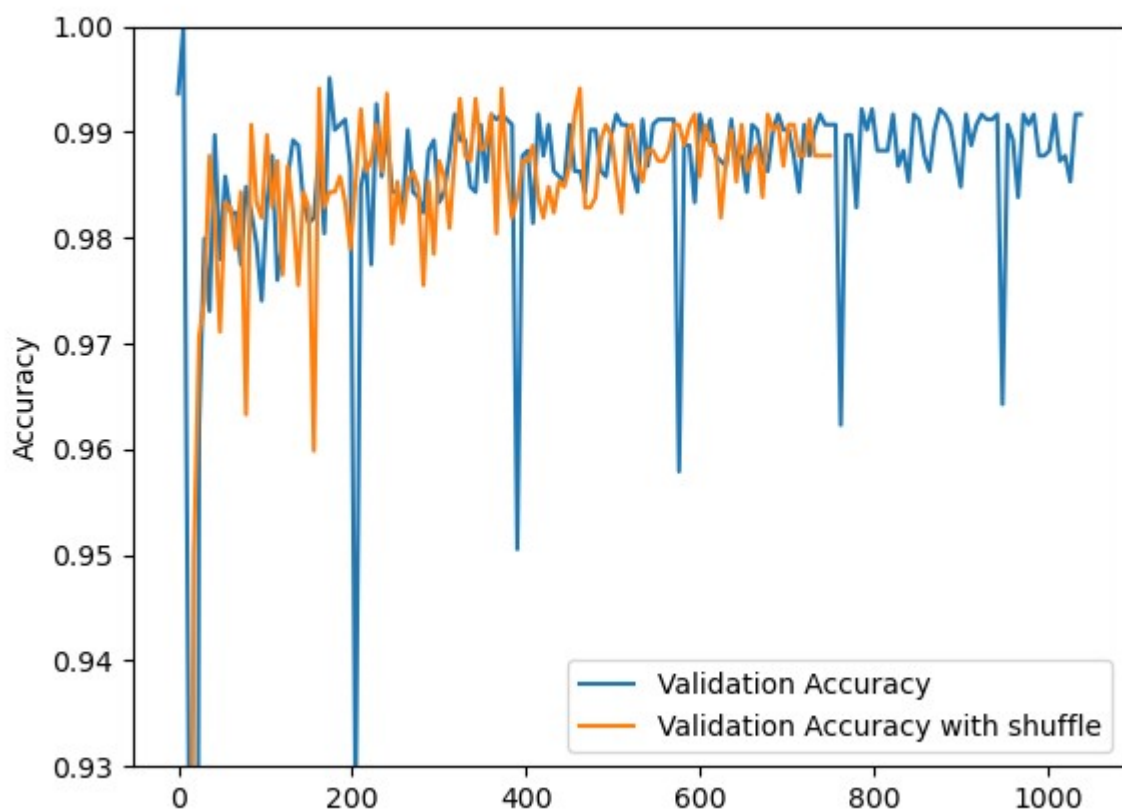
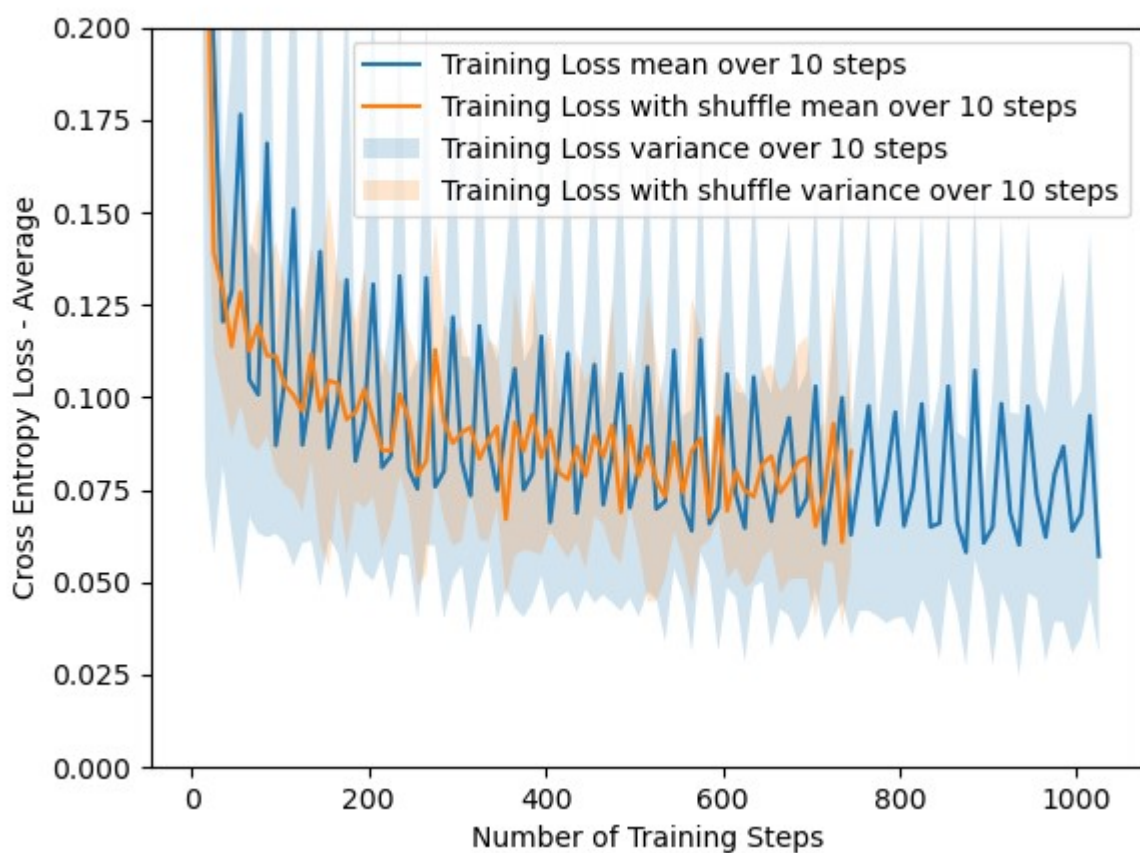
Task 2c)



Task 2d)

The training stops after 33 epochs. The end accuracy results are also better than previously.

Task 2e)

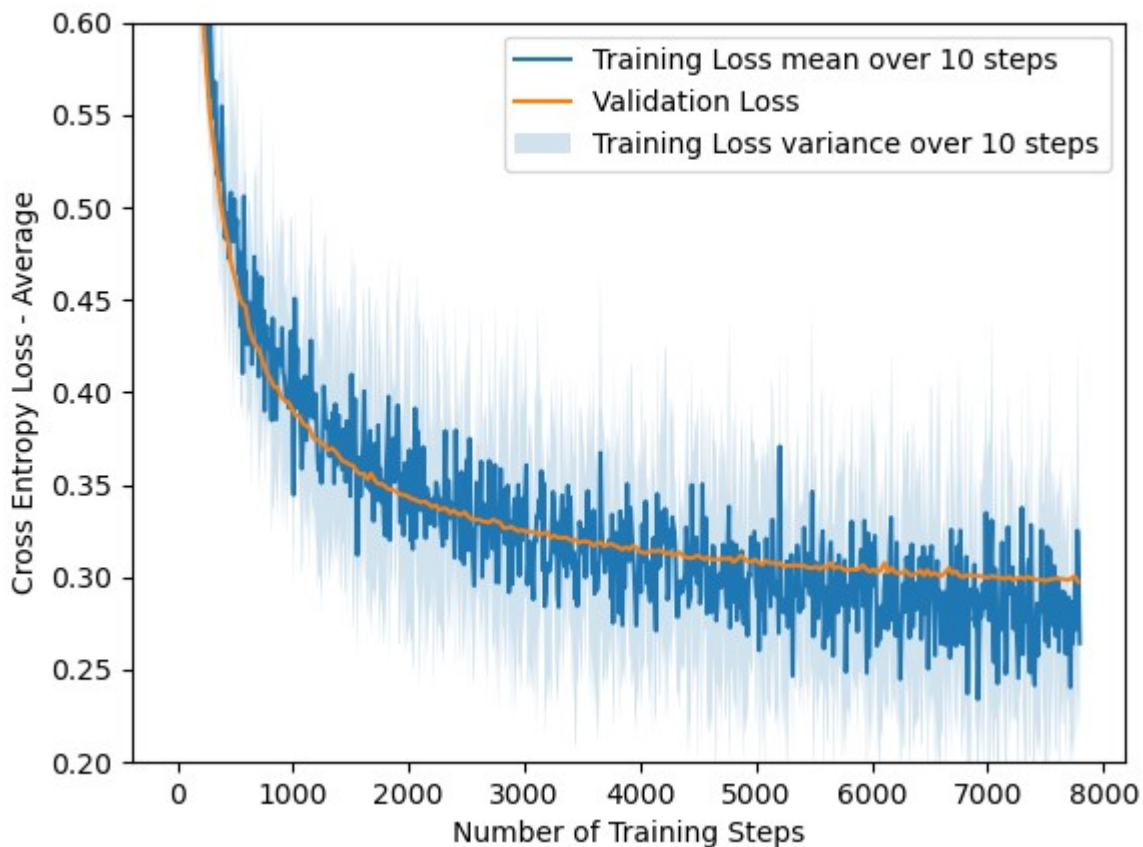


Number of Training Steps

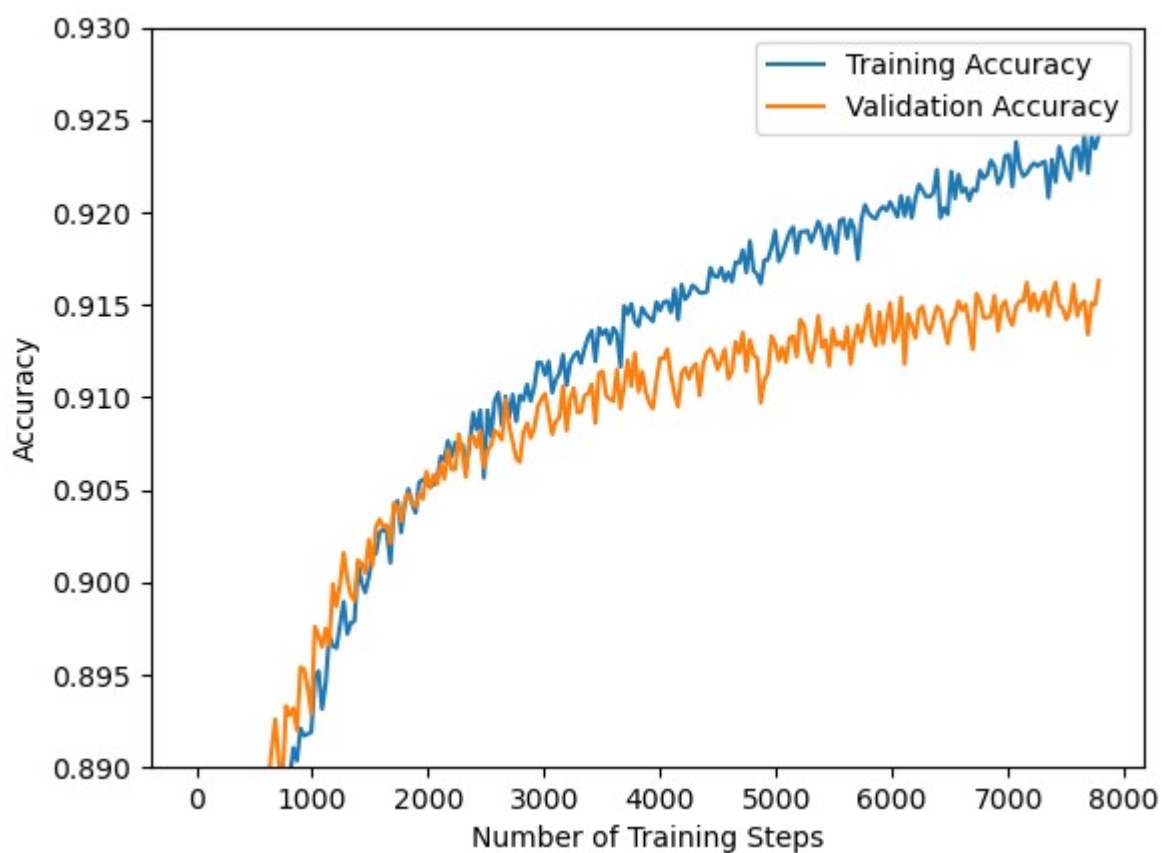
We can see that data shuffling removes many of the spikes that we could see previously. They were result of some difficult batch. When we shuffle the data the "difficult" images are spread around and the spikes disappear. We can also see that early stopping kicks in a bit earlier (after 24 batches). It will however vary from time to time since the data is shuffled randomly.

Task 3

Task 3b)



Task 3c)



Task 3d)

The first sign of overfitting is that the training loss is not improving very much after about 3000 samples. We can also see that the validation accuracy is getting quite lower than the training accuracy (also after about 3000 samples). It means that the network has learned the training set good, but is worse at recognising new data. Which is overfitting.

Task 4

Task 4a)

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

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

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

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

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

$$= -x_j^n (y_k^n - \hat{y}_k^n) + 2\lambda w$$

Task 4b)

lambda = 0.0

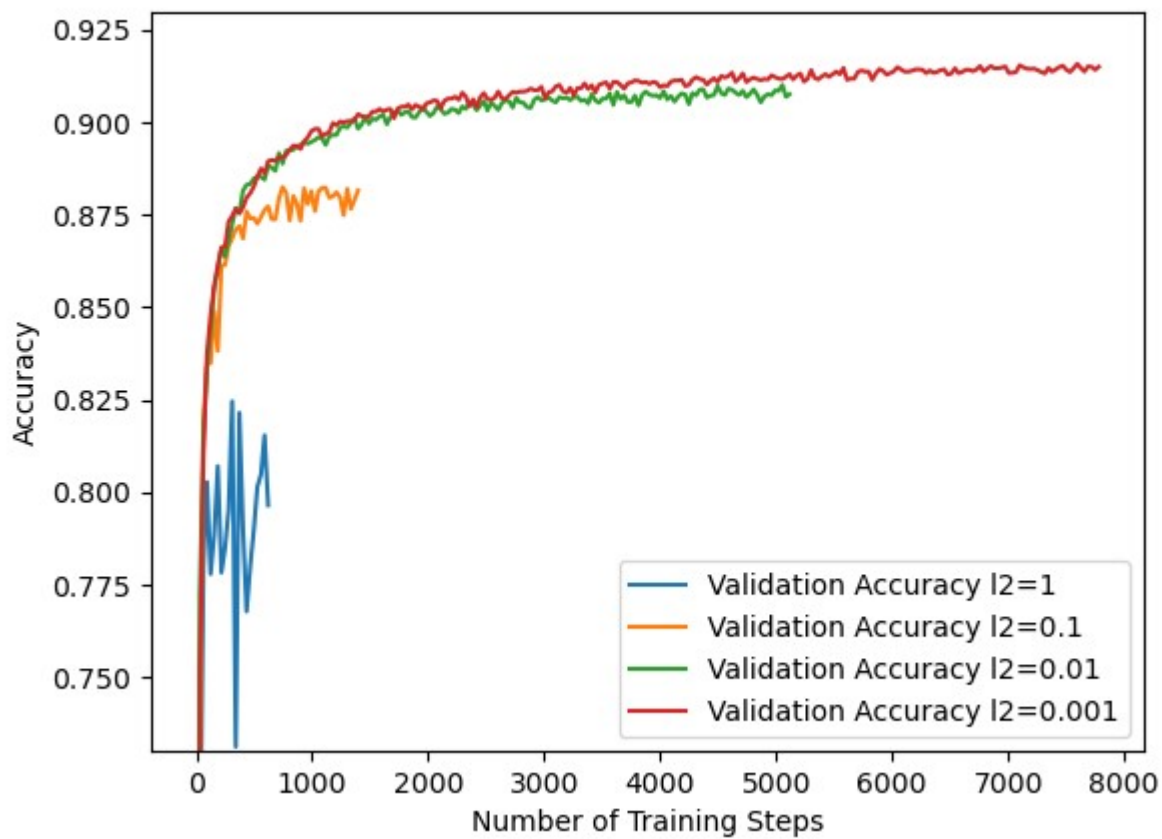


lambda = 1.0



With lambda=1 the gradient will be smaller and weights will also change slower and be smaller. Smaller weights lead to less complex models that don't overfit that much. We can often see that overfitting model will take the noise into account and becomes less general. The L2 regression makes model more general and we can see that clearly from the images of the weights. The weight without regularization are much more noisy than those with.

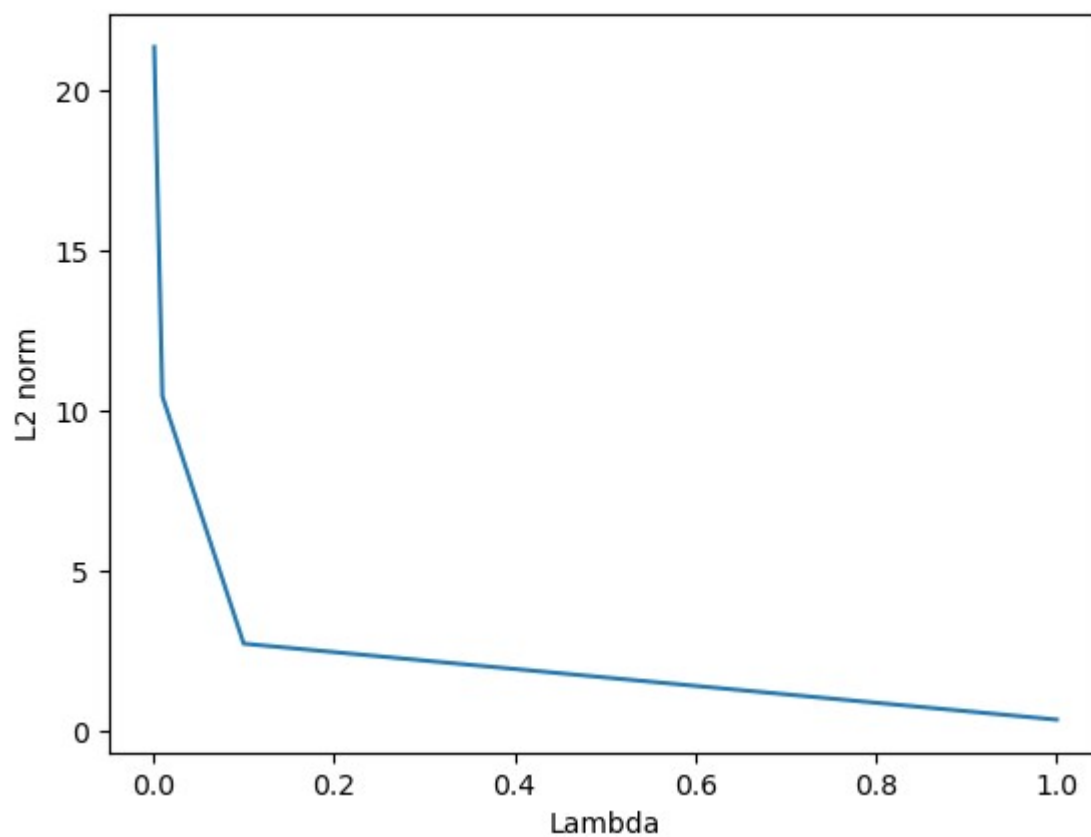
Task 4c)



Task 4d)

With regularization the model might become too "simple" for the task and won't be able to respond well for more varying inputs. This may not be a problem when working with more complex problems.

Task 4e)



We can see that the higher the λ is, the smaller the weights get. The smaller weights will lead to simpler model and less overfitting. This confirms what we see both in task b and c.