

Task 1:

a)

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

$$\hat{y}^n = f(z_k)$$

$$f(x) = \frac{1}{1 + e^{-w^T x}}$$

$$\hat{y} = \frac{1}{1 + e^{-z_k}}$$

$$z_k = w^T x = \sum w_i x_i$$

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

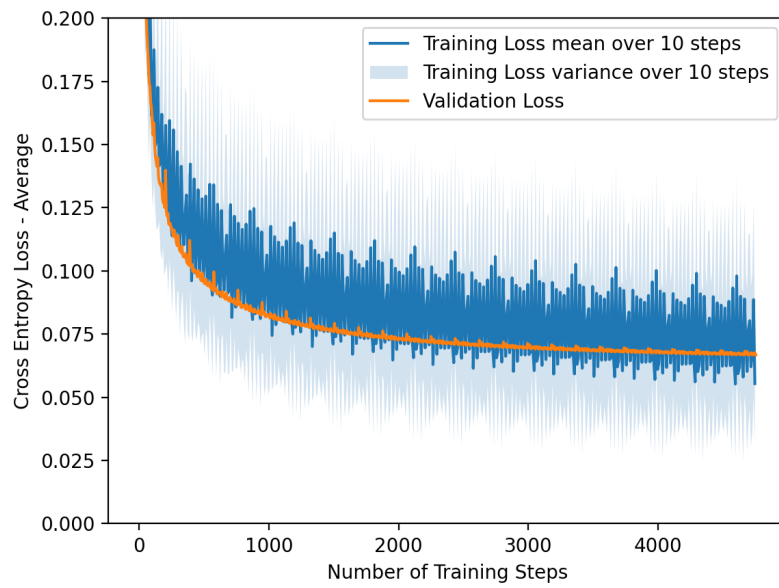
$$\begin{aligned} \frac{d\hat{y}^n}{dz_k} &= (1 + e^{-z_k})^{-1} = -(1 + e^{-z_k})^{-2} \cdot \frac{d}{dz_k}(1 + e^{-z_k}) = -(1 + e^{-z_k})^{-2} (-e^{-z_k}) = (1 + e^{-z_k})^{-2} e^{-z_k} = \frac{e^{-z_k}}{(1 + e^{-z_k})^2} = \frac{e^{-z_k}}{(1 + e^{-z_k})(1 + e^{-z_k})} \\ &= \frac{1}{(1 + e^{-z_k})} \cdot \frac{e^{-z_k}}{(1 + e^{-z_k})} = \frac{1}{(1 + e^{-z_k})} \cdot \frac{e^{-z_k} + 1 - 1}{(1 + e^{-z_k})} = \frac{1}{(1 + e^{-z_k})} \cdot \left(1 - \frac{1}{(1 + e^{-z_k})}\right) \end{aligned}$$

$$\frac{dz_k}{dw_i} = x_i$$

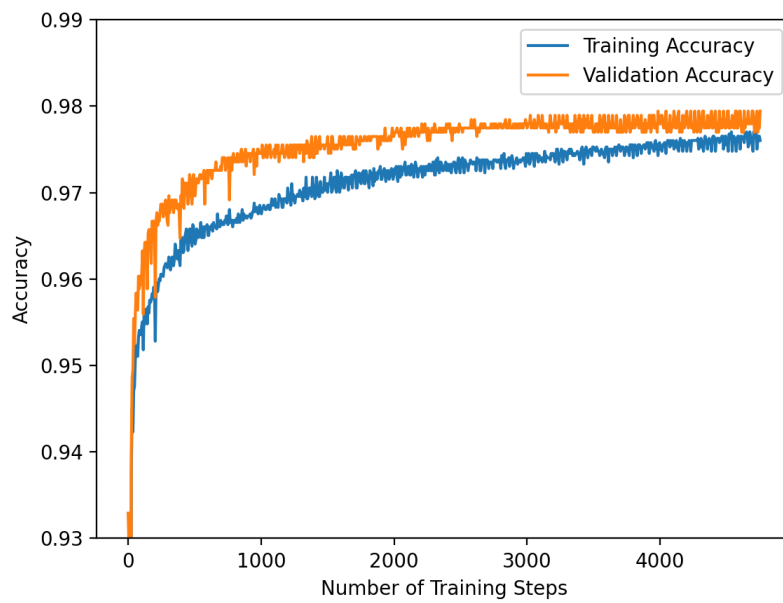
$$\begin{aligned} \text{Gradient: } \frac{dC(w)}{dw} &= \frac{dC(w)}{d\hat{y}^n} \cdot \frac{d\hat{y}^n}{dz_k} \cdot \frac{dz_k}{dw_i} = \left(-\frac{y^n}{\hat{y}^n} + \frac{(1 - y^n)}{(1 - \hat{y}^n)} \right) \cdot \left(\frac{1}{1 + e^{-z_k}} \cdot \left(1 - \frac{1}{(1 + e^{-z_k})}\right) \right) \cdot x_i \\ &= \left(-\frac{y^n}{\hat{y}^n} + \frac{(1 - y^n)}{(1 - \hat{y}^n)} \right) \cdot (\hat{y}^n \cdot (1 - \hat{y}^n)) \cdot x_i = \\ &= \left(-\frac{y^n \hat{y}^n (1 - \hat{y}^n)}{\hat{y}^n} + \frac{(1 - y^n) \hat{y}^n (1 - \hat{y}^n)}{(1 - \hat{y}^n)} \right) \cdot x_i = (-y^n (1 - \hat{y}^n) + (1 - y^n) \hat{y}^n) x_i \\ &= (-y^n + y^n \hat{y}^n + \hat{y}^n - y^n \hat{y}^n) x_i = \underline{(\hat{y}^n - y^n) x_i} \end{aligned}$$

Task 2:

b)



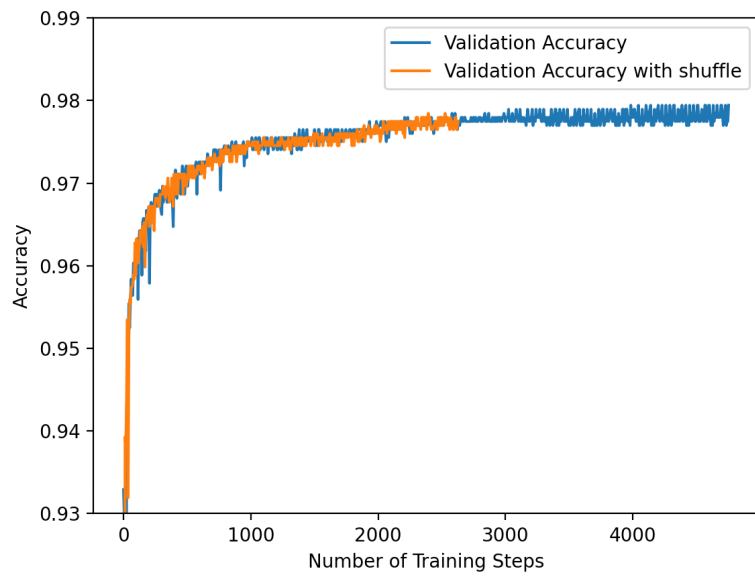
c)



d) The early stopping kicks in after 153 epochs

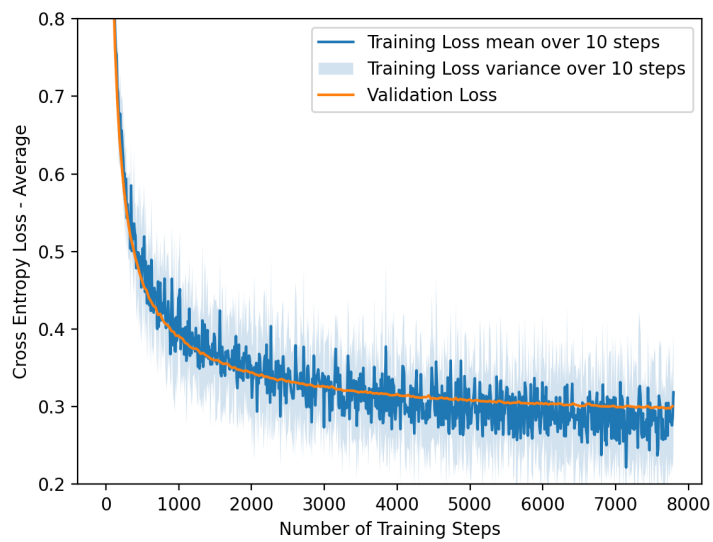
```
Train shape: X: (4005, 784), Y: (4005, 1)
Validation shape: X: (2042, 784), Y: (2042, 1)
Early stopping at epoch 153
Final Train Cross Entropy Loss: [0.07122371]
Final Validation Cross Entropy Loss: [0.06669121]
Train accuracy: 0.9760299625468165
Validation accuracy: 0.9794319294809011
```

e)

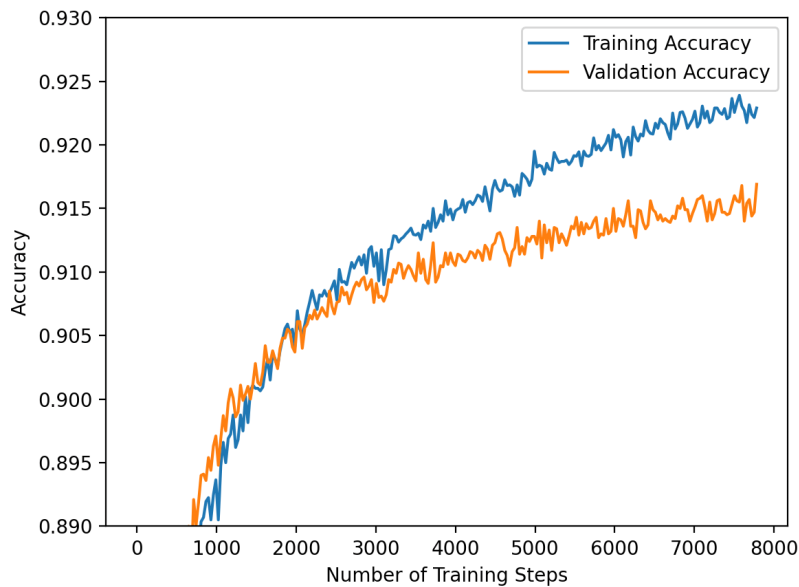


Task 3:

b)



c)



d)

We can see that the model starts to overfit in 3c because the validation accuracy stops increasing at the same rate as the training accuracy. The model likely overfits on the training data, so that when it must predict on unseen data (validation data) it expects the distribution of the data to look like the training data distribution (which it doesn't with some deviation).

Task 4

a)

$$a) \quad C(w) = \frac{1}{N} \sum_{k=1}^N C^k(w), \quad C^k(w) = - \sum_j y_k^j \ln(\hat{y}_k^j)$$

$$\hat{y}_k^j = \frac{e^{z_k^j}}{\sum_i e^{z_k^i}}, \quad z_k = w_k^T \cdot x$$

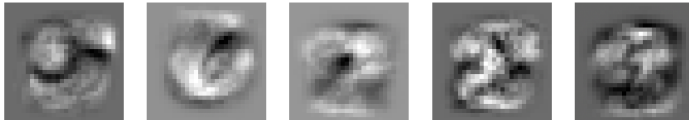
$$J(w) = C(w) + \lambda R(w), \quad R(w) = \|w\|^2$$

$$\frac{dJ}{dw} = \underbrace{\frac{dC}{d\hat{y}_k^j} \cdot \frac{d\hat{y}_k^j}{dz_k^j} \cdot \frac{dz_k^j}{dw}}_{\frac{dC}{dw}} + \frac{d\lambda R(w)}{dw} = \underbrace{-x_j^j (y_k^j - \hat{y}_k^j)}_{\frac{dC}{dw}} + \lambda 2w$$

$$\frac{dC}{dw} = -x_j^j (y_k^j - \hat{y}_k^j)$$

b)

Weights with $\lambda=0.0$:

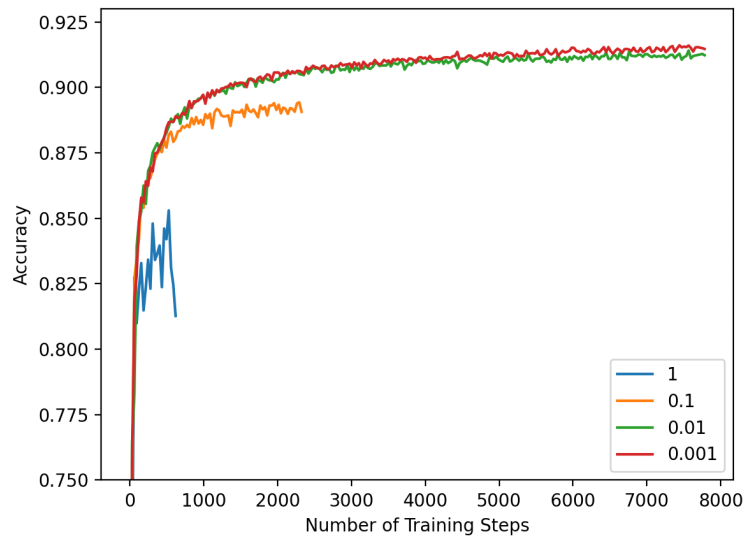


Weights with $\lambda=1.0$:



The weights for the model with $\lambda=1.0$ are less noisy because the regularization penalizes the model for learning the patterns for the training data too precisely.

c)



d)

