

Assignment 1 Report

Authors:

- Alberto Robazza
- Ivan Garcia Alcaraz

Task 1

task 1a)

Equation 3

$$C(w) = \frac{1}{N} \sum_{n=1}^N C^n, \text{ where } C^n(w) = - \underbrace{(\gamma^n)}_{P1} \underbrace{\ln(\hat{y}^n)}_{P2} + \underbrace{(1-\gamma^n)}_{P3} \underbrace{\ln(1-\hat{y}^n)}_{P4}$$

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

For $f(x)$ and $\frac{\partial}{\partial w_i}$

$$\frac{\partial}{\partial w_i} \ln(f(x)) = \frac{\partial}{\partial w_i} \ln\left(\frac{1}{1+e^{-w^T x}}\right) = \frac{\partial}{\partial w_i} -\ln(1+e^{-w^T x})$$

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

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

$$\frac{\partial}{\partial w_i} \ln(1-f(x)) = x_i + x_i(1-f(x)) = \boxed{-f(x)x_i}$$

$$\frac{\partial}{\partial w_i} C^n(w) = - \left(\frac{\partial}{\partial w_i} \gamma^n \cdot \ln(\hat{y}^n) + \gamma^n \cdot \frac{\partial}{\partial w_i} \ln(\hat{y}^n) + \frac{\partial}{\partial w_i} (1-\gamma^n) \cdot \ln(1-\hat{y}^n) + (1-\gamma^n) \cdot \frac{\partial}{\partial w_i} \ln(1-\hat{y}^n) \right)$$

$$= -(\gamma^n x_i^n (1-f(x^n)) - (1-\gamma^n) x_i^n f(x^n)) =$$

$$= -(\gamma^n x_i^n - \gamma^n x_i^n f(x^n) - x_i^n f(x^n) + \gamma^n x_i^n f(x^n)) =$$

$$= -(\gamma^n x_i^n - x_i^n f(x^n)) = \boxed{-(\gamma^n - f(x^n)) x_i^n}$$

task 1b)

Calculus to makes the derivative more easily.

Equation 5

$$\hat{y}_k = \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}}, \text{ where } z_k = \mathbf{w}_k^T \cdot \mathbf{x} = \sum_i \mathbf{w}_{k,i} \cdot x_i$$

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

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

$$= - \sum_{k=1}^K y_k^n z_k^n + \sum_{k=1}^K y_k^n \ln(\sum_{k'} e^{z_{k'}^n}) = \underbrace{- \sum_{k=1}^K y_k^n z_k^n}_{-g^n(\mathbf{w})} + \underbrace{\ln(\sum_{k'} e^{z_{k'}^n})}_{+h^n(\mathbf{w})}$$

The derivative of a function $g(\mathbf{w})$ is the following:

$$\frac{\partial}{\partial w_{kj}} g(\mathbf{w}) = \frac{\partial}{\partial w_{kj}} \sum_{k=1}^K y_k^n z_k^n = \frac{\partial}{\partial w_{kj}} \sum_{k=1}^K y_k^n \sum_{i=1}^I w_{ki} \cdot x_i$$

The derivative of the previous formula takes the fact that only when $j=i$ the derivate of w_{jk} is non-zero so:

$$\frac{\partial}{\partial w_{kj}} g(\mathbf{w}) = \boxed{y_k^n \cdot x_j}$$

The derivative of a function $h(\mathbf{w})$ is the following:

$$\frac{\partial}{\partial w_{kj}} h(\mathbf{w}) = \frac{\partial}{\partial w_{kj}} \ln \left(\sum_{k'} e^{z_{k'}} \right) = \frac{\partial}{\partial w_{kj}} \ln \left(\sum_{k'} e^{\sum_i w_{ki} \cdot x_i} \right) = \frac{1}{\sum_{k'} e^{\sum_i w_{ki} \cdot x_i}} \frac{\partial}{\partial w_{kj}} \sum_{k'} e^{\sum_i w_{ki} \cdot x_i}$$

The derivative of the previous formula takes the fact that only when $j=i$ and $k=k'$ the derivate of w_{jk} is non-zero so:

$$\frac{\partial}{\partial w_{kj}} \sum_{k'} e^{\sum_i w_{ki} \cdot x_i} = \boxed{x_j e^{\sum_i w_{ki} \cdot x_i}}$$

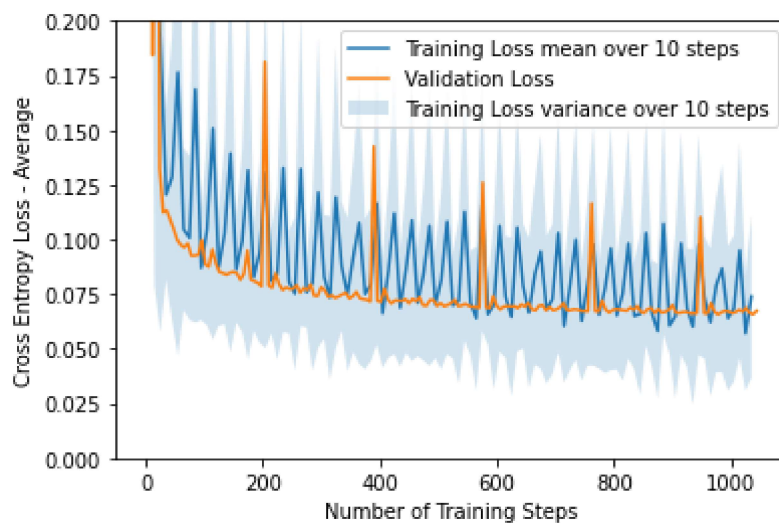
Now we merge the derivative of $g(\mathbf{w})$ with the derivative of $h(\mathbf{w})$ and we get the following:

$$\frac{1}{\sum_{k'} e^{\sum_i w_{ki} \cdot x_i}} \frac{\partial}{\partial w_{kj}} \sum_{k'} e^{\sum_i w_{ki} \cdot x_i} = \frac{e^{z_k}}{\sum_{k'} e^{\sum_i w_{ki} \cdot x_i}} x_j = \hat{y}_k^n x_j$$

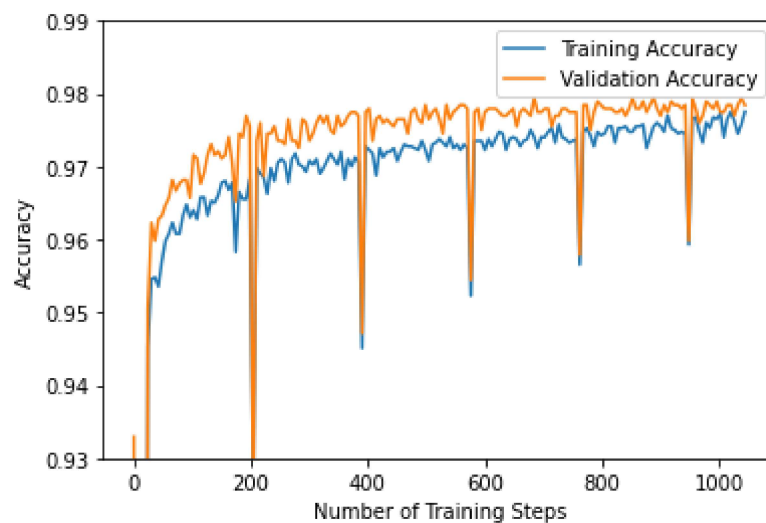
$$\frac{\partial}{\partial w_{kj}} C^n(\mathbf{w}) = \frac{\partial}{\partial w_{kj}} g^n(\mathbf{w}) + \frac{\partial}{\partial w_{kj}} h^n(\mathbf{w}) = -y_k^n x_j + \hat{y}_k^n x_j = \boxed{-x_j^n (y_k^n - \hat{y}_k^n)}$$

Task 2

Task 2b)



Task 2c)

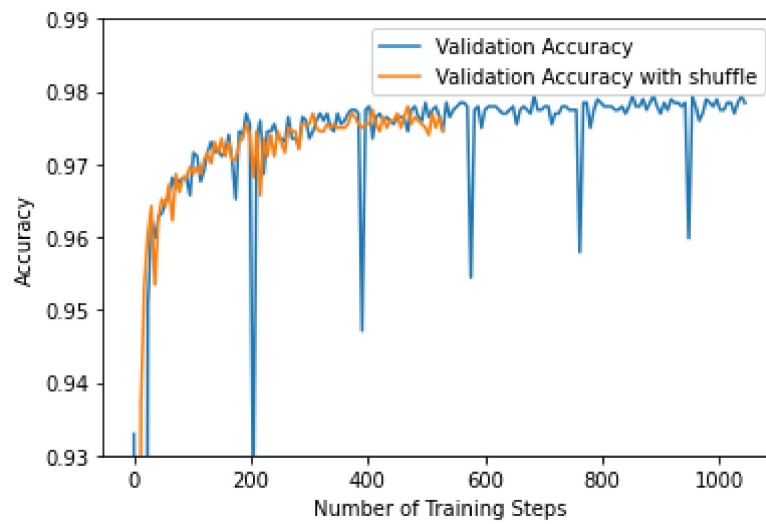


Task 2d)

In our test, we were bumped out after 17 epochs

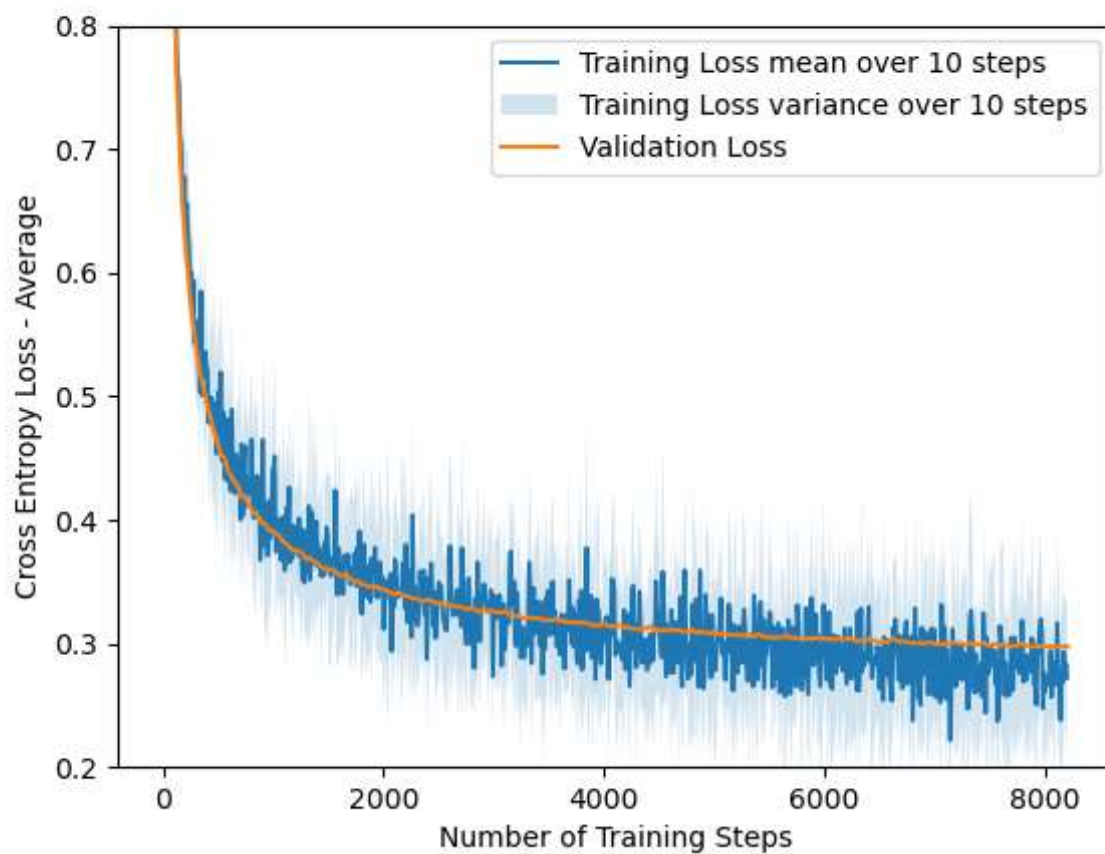
Task 2e)

Input is more stable, the model is trained on different pictures compared to the immediatly validatd ones

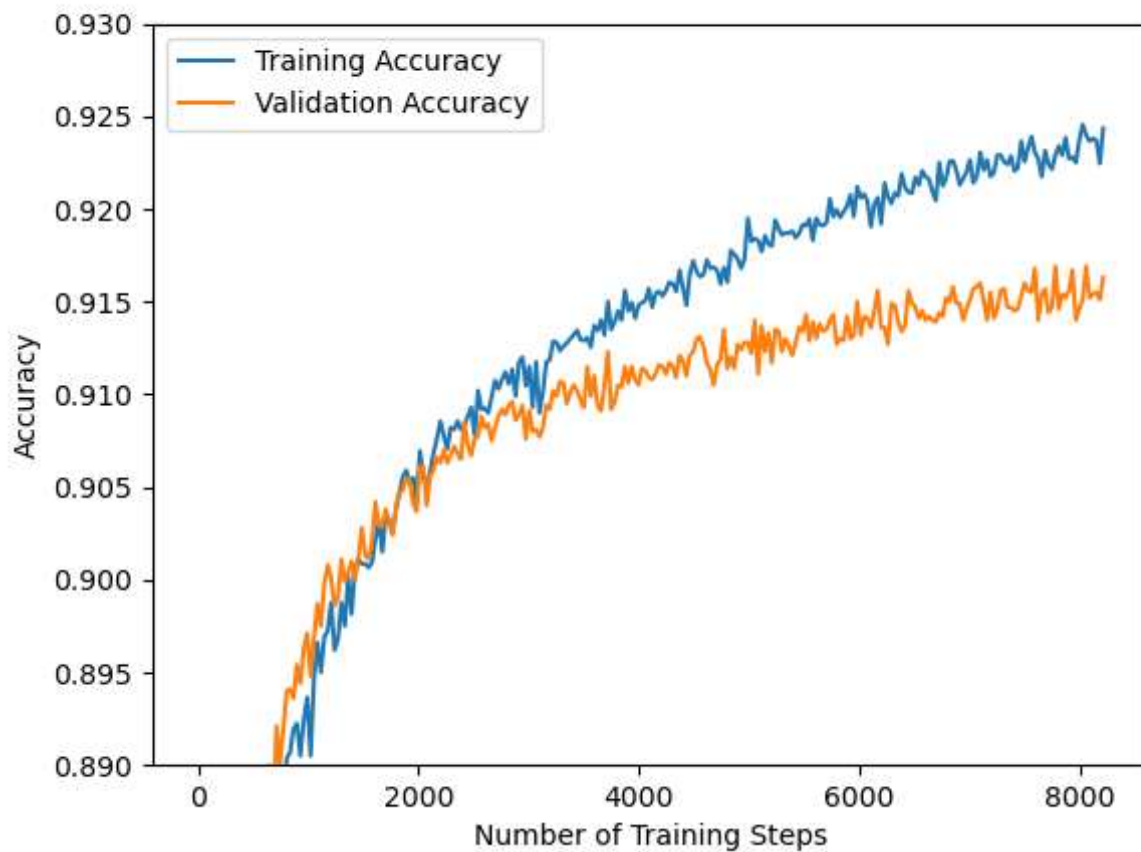


Task 3

Task 3b)



Task 3c)



Task 3d)

In Figure 6 the training accuracy is better than the validation accuracy after 3000 training steps. And the increase distance between the training accuracy and the validation accuracy could be a sign of overfitting, because the model performs better in one than the other.

Task 4

Task 4a)

Task 4a:

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

The $\frac{\partial C}{\partial w}$ was done in Task 1

$$\text{The } \frac{\partial \|w\|^2}{\partial w} = \frac{\partial}{\partial w} \sum_{i,j} w_{i,j}^2$$

In the matrix generated by $w_{i,j}$ only the terms where $i'=i$ or $j'=j$ are the ones where the derivative is non-zero
So

$$\frac{\partial}{\partial w} = \boxed{2w}$$

$$\frac{\partial J}{\partial w} = \boxed{- \frac{(\bar{y}_k^n - \hat{y}_k^n) x_j^n}{N} + 2\lambda w}$$

Task 4b)

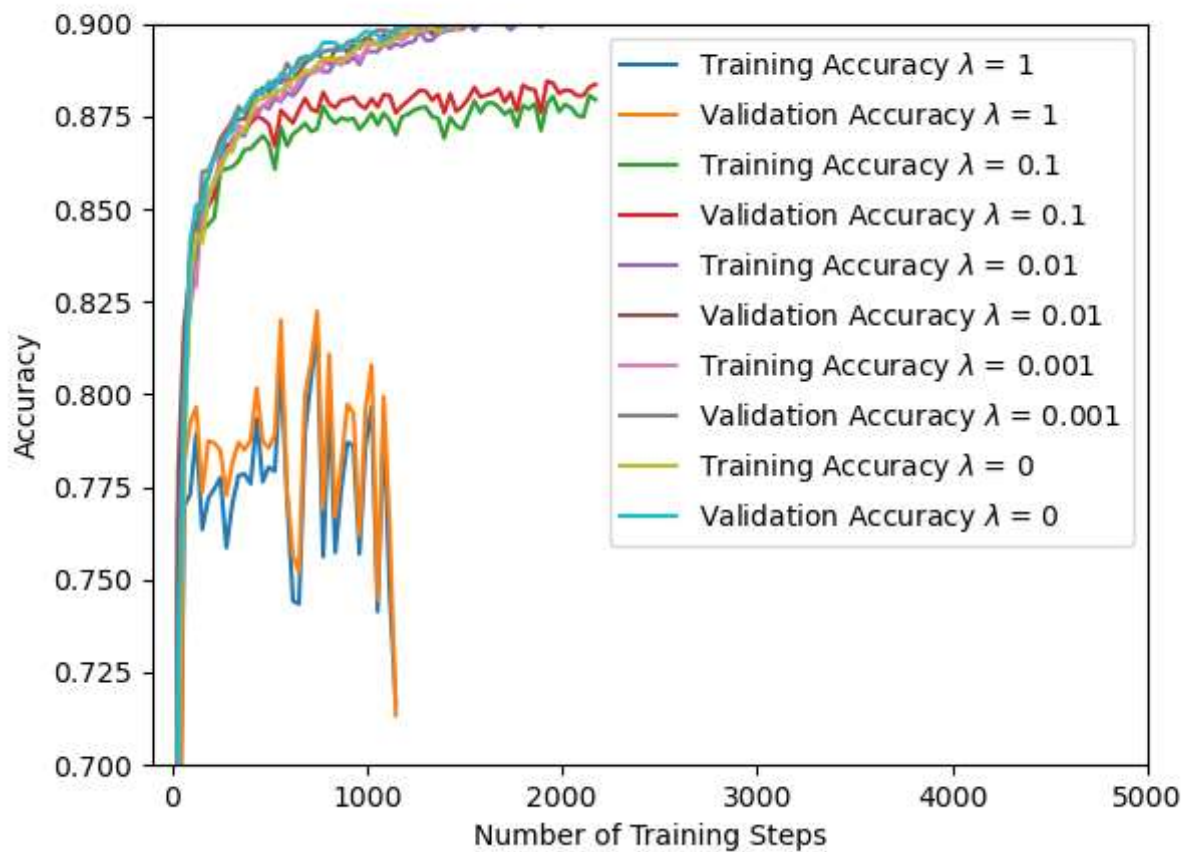
FILL IN ANSWER



With stronger L2 normalization (high values for lambda), the model tries to be simpler, the difference within different values would be lower, than the noise is reduced

Task 4c)

FILL IN ANSWER



Task 4d)

Since L2 normalization penalize high values, the model tries to reduce them. It brings to a worse fitting model since it wouldn't align to the points

Task 4e)

On increasing lambda, l2 value lowers down. This is probabilly caused by the penalization of L2 on higher lambda values

Plot of length of w vs λ

