

Homework 1

Alexandre Rodrigues - student number: 89404

Luís António Lima da Silva - student number: 102406

Rui Jorge Breda Perdigoto - student number: 102381

Question 1	2
Exercise 1	2
Exercise 2	3
Exercise 3	4
Exercise 4	6
Exercise 5	8
Question 2	9
Exercise 1	9
Exercise 2.a	11
Exercise 2.b	13
Question 3	14
Exercise 1.a	14
Exercise 2.a	15
Exercise 2.b	16
Question 4	16
Exercise 1	16
Exercise 2	20
Exercise 3	30

Question 1

Exercise 1

Question 1

$$1) \sigma(z) = \frac{1}{(1+e^{-z})}$$

$$\sigma'(z) = \frac{d}{dz} \sigma(z) = \frac{d}{dz} \frac{1}{1+e^{-z}} = \frac{d}{dz} (1+e^{-z})^{-1} = -(1+e^{-z})^{-2} \frac{d}{dz} (1+e^{-z})$$

$$= -(1+e^{-z})^2 \left(\frac{d}{dz} 1 + \frac{d}{dz} e^{-z} \right) = -(1+e^{-z})^2 \frac{d}{dz} e^{-z}$$

$$= -(1+e^{-z})^2 e^{-z} \cdot \frac{d}{dz} (-1) = -(1+e^{-z})^{-2} (e^{-z} - \frac{d}{dz} (z))$$

$$= -(1+e^{-z})^{-2} (e^{-z} - 1) = (1+e^{-z})^{-2} \cdot e^{-z}$$

$$= \frac{e^{-z}}{(1+e^{-z}) \cdot (1+e^{-z})} = \frac{e^{-z} + 1 - 1}{(1+e^{-z}) \cdot (1+e^{-z})} = \frac{1}{(1+e^{-z})} \cdot \frac{1+e^{-z}}{1+e^{-z}} - \frac{1}{1+e^{-z}}$$

$$= \underbrace{\frac{1}{1+e^{-z}}}_{\sigma(z)} \cdot 1 - \underbrace{\frac{1}{1+e^{-z}}}_{\sigma(z)} = \sigma(z) \cdot (1 - \sigma(z))$$

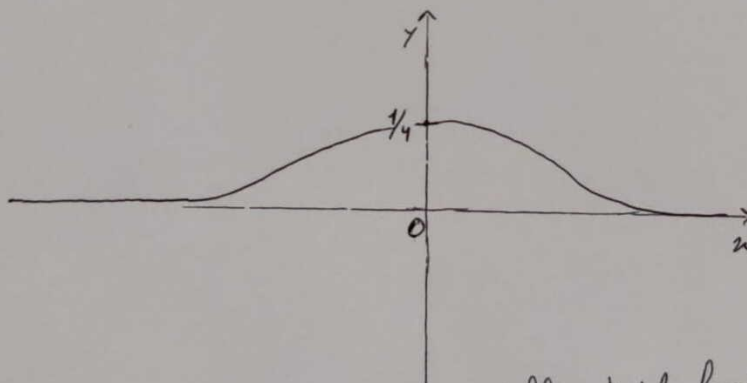
Exercise 2

$$\begin{aligned} 2) l'(z, y=+1) &= (-\log \sigma(z))' = \left(-\log \left(\frac{1}{1+e^{-z}} \right) \right)' = (-\log(1) - \log(1+e^{-z}))' \\ &= (\log(1+e^{-z}))' = \frac{(1+e^{-z})'}{1+e^{-z}} = \frac{-e^{-z}}{1+e^{-z}} \end{aligned}$$

$$l''(z, y=+1) = \left(-\frac{e^{-z}}{1+e^{-z}} \right)' = -\frac{(1+e^{-z})' \cdot (1+e^{-z}) - (1+e^{-z})' \cdot (1+e^{-z})'}{(1+e^{-z})^2}$$

$$= -\frac{(-e^{-z}) \cdot (1+e^{-z}) - (1+e^{-z}) \cdot (-e^{-z})}{(1+e^{-z})^2} = -\frac{-e^{-z} \cdot (e^{-z})^2 + (e^{-z})^2}{(1+e^{-z})^2}$$

$$= -\frac{-e^{-z}}{(1+e^{-z})^2} = \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{e^{+z}(\sigma(z))^2}$$



This function is convex, because a twice-differentiable function of a single variable is convex if and only its second derivative is nonnegative on its entire domain.

Exercise 3

Question 1.3

$$\text{softmax}: \mathbb{R}^K \rightarrow \mathbb{R}^K$$

$$[\text{softmax}(z)]_j = \frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)}$$

Jacobian matrix

$$J_{j,k} = \frac{\partial [\text{softmax}(z)]_j}{\partial z_k}$$

Two cases: $j = k$ and $j \neq k$

$$j = k: \quad \frac{\partial [\text{softmax}(z)]_j}{\partial z_k} = \frac{\partial}{\partial z_k} \frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)}$$

$$= \frac{\partial}{\partial z_j} \frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)} = \frac{\left(\frac{\partial}{\partial z_j} \exp(z_j) \right) \sum_{k=1}^K \exp(z_k) - \exp(z_j) \frac{\partial}{\partial z_j} \sum_{k=1}^K \exp(z_k)}{\left(\sum_{k=1}^K \exp(z_k) \right)^2}$$

$$= \frac{\exp(z_j) \cdot \sum_{k=1}^K \exp(z_k) - \exp(z_j) \cdot \exp(z_j)}{\sum_{k=1}^K \exp(z_k)^2} \quad \leftarrow (\text{since } j=1, \dots, K)$$

$$= \frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)} \cdot \frac{\sum_{k=1}^K \exp(z_k) - \exp(z_j)}{\sum_{k=1}^K \exp(z_k)}$$

$$= [\text{softmax}(z)]_j (1 - [\text{softmax}(z)]_j)$$

When $j \neq k$:

$$\frac{\partial}{\partial z_k} [\text{softmax}(z)]_j = \frac{\partial}{\partial z_k} \frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)}$$

$$= \frac{\left(\frac{\partial}{\partial z_k} \exp(z_j) \right) \sum_{k=1}^K \exp(z_k) - \exp(z_j) \cdot \frac{\partial}{\partial z_k} \sum_{k=1}^K \exp(z_k)}{\left(\sum_{k=1}^K \exp(z_k) \right)^2}$$

$$= \frac{-\exp(z_j) \cdot \exp(z_k)}{\left(\sum_{k=1}^K \exp(z_k) \right)^2}$$

$$= -\frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)} \cdot \frac{\exp(z_k)}{\sum_{k=1}^K \exp(z_k)}$$

$$= -[\text{softmax}(z)]_j \cdot [\text{softmax}(z)]_k$$

Note that when " $\exp(z_k)$ " is within a sum $\left(\sum_{k=1}^K \exp(z_k) \right)$, the index k refers to the index of the sum, whereas when " $\exp(z_k)$ " appears outside the sum (as a result of $\frac{\partial}{\partial z_k}$), the k refers to the index (j, k) in the jacobian matrix.

So, the jacobian matrix of the softmax function is defined as follows:

$$y_{(j,k)} = \begin{cases} [\text{softmax}(z)]_j \cdot (1 - [\text{softmax}(z)]_j) & \text{if } j = k \\ -[\text{softmax}(z)]_j \cdot [\text{softmax}(z)]_k & \text{if } j \neq k \end{cases}$$

Exercise 4

$$4. \quad L(z; y=j) = -\log[\text{softmax}(z)]_j$$

$$\frac{\partial}{\partial z_i} -\log[\text{softmax}(z)]_j$$

$$= -\frac{\partial}{\partial z_i} \log[\text{softmax}(z)]_j$$

$$= -\frac{1}{[\text{softmax}(z)]_j} \frac{\partial}{\partial z_i} [\text{softmax}(z)]_j$$

$$= -\frac{1}{[\text{softmax}(z)]_j} \cdot [\text{softmax}(z)]_j \cdot (1(i=j) - [\text{softmax}(z)]_i)$$

$$= -1(i=j) + [\text{softmax}(z)]_i$$

$$= [\text{softmax}(z)]_i - 1(i=j)$$

where the operation $1(i=j)$ is defined as: $1(i=j) = \begin{cases} 1, & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$

The Hessian matrix is the matrix containing all the possible second derivatives.

$$H_{(i,j)} = \begin{cases} \frac{\partial}{\partial z_j} [\text{softmax}(z)]_i - 1 & \text{if } i=j \\ \frac{\partial}{\partial z_j} [\text{softmax}(z)]_i & \text{if } i \neq j \end{cases}$$

$$= \frac{\partial}{\partial z_j} [\text{softmax}(z)]_i \quad (\text{since } \frac{\partial}{\partial z_j} (-1) = 0)$$

Notice that this matrix is the same as the jacobian of the softmax transformation calculated in the previous exercise.

$$\text{Hessian (loss function)} = \text{jacobian (softmax)}$$

So, the Hessian is

$$H_{(i,j)} = \begin{cases} [\text{softmax}(z)]_i \cdot (1 - \text{softmax}(z)_i) & , i=j \\ - [\text{softmax}(z)]_i [\text{softmax}(z)]_j & , i \neq j \end{cases}$$

~~A twice differentiable function of several variables is convex iff its Hessian is positive semidefinite, i.e.~~ $\forall v \in \mathbb{R}^n, v^T H v \geq 0$

For convenience, let's represent H as $H = \text{diag}(\text{softmax}(z)) - \text{softmax}(z) \cdot \text{softmax}(z)^T$

$$v^T H v = v^T \text{diag}(\text{softmax}(z)) v - v^T \text{softmax}(z) \text{softmax}(z)^T v$$

$$= \sum_{i=1}^K [\text{softmax}(z)]_i v_i^2 - \left(\sum_{i=1}^K v_i [\text{softmax}(z)]_i \right)^2$$

$$= \left(\sum_{i=1}^K [\text{softmax}(z)]_i \right) \left(\sum_{i=1}^K [\text{softmax}(z)]_i v_i^2 \right) - \left(\sum_{i=1}^K v_i [\text{softmax}(z)]_i \right)^2$$

due to the fact that this must sum to 1.

Based on the Cauchy-Schwartz inequality, we have that

$$\left(\sum_{i=1}^K [\text{softmax}(z)]_i \right) \left(\sum_{i=1}^K [\text{softmax}(z)]_i v_i^2 \right) - \left(\sum_{i=1}^K v_i \cdot [\text{softmax}(z)]_i \right)^2 \geq 0$$

so, the Hessian of the multinomial logistic loss function is positive semidefinite, implying that the loss function is convex with respect to z .

Exercise 5

Question 1.5

As it was proved in the previous exercise, the multinomial logistic loss function is convex with respect to z . Since convexity of a function is invariant under the composition with an affine map, then $L(z; y_j) \circ Z(w)$ is also convex w.r.t. the model parameters (w, b) .

The convexity of a composition of functions is preserved under the following circumstances:

- if f and g are convex and g is non-decreasing over an univariate domain, then $h(x) = g(f(x))$ is convex.
- if f is concave and g is convex and non-decreasing over a ~~univariate~~ univariate domain, then $h(x) = g(f(x))$ is convex.
- convexity is invariant under affine maps.

So, the logistic loss function can still be convex under specific circumstances when z is not a linear function of the model parameters.

Question 2

Exercise 1

Question 2.1

$$L(z; \gamma) = \frac{1}{2} \|z - \gamma\|_2^2 = \frac{1}{2} (z - \gamma)^T (z - \gamma) \\ = \frac{1}{2} \sum_{i=1}^K (z_i - \gamma_i)^2, \quad |z| = K = |\gamma|$$

Consider the function $f: \mathbb{R}^K \rightarrow \mathbb{R}$ $f(z) = \frac{1}{2} \sum_{i=1}^K (z_i - \gamma_i)^2$

As an auxiliary result, we will first prove that f is convex w.r.t. z

$$\nabla f(z) = \begin{bmatrix} \frac{\partial f}{\partial z_1}(z) \\ \vdots \\ \frac{\partial f}{\partial z_K}(z) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial z_1} \left(\frac{1}{2} \sum_{i=1}^K (z_i - \gamma_i)^2 \right) \\ \vdots \\ \frac{\partial}{\partial z_K} \left(\frac{1}{2} \sum_{i=1}^K (z_i - \gamma_i)^2 \right) \end{bmatrix}$$

$\frac{\partial}{\partial z_i} (z_i^2 - 2z_i\gamma_i + \gamma_i^2) = 2z_i - 2\gamma_i$

$\frac{\partial}{\partial z_i} (z_j^2 - 2z_j\gamma_j + \gamma_j^2) = 0$

$$= \begin{bmatrix} z_1 - \gamma_1 \\ \vdots \\ z_K - \gamma_K \end{bmatrix}$$

In order to prove that f is convex, we must prove that its Hessian is positive semidefinite

$$H(z) = \begin{bmatrix} \frac{\partial^2 f}{\partial z_1 \partial z_1} & \frac{\partial^2 f}{\partial z_1 \partial z_2} & \dots & \frac{\partial^2 f}{\partial z_1 \partial z_K} \\ \frac{\partial^2 f}{\partial z_2 \partial z_1} & \frac{\partial^2 f}{\partial z_2 \partial z_2} & \dots & \frac{\partial^2 f}{\partial z_2 \partial z_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial z_K \partial z_1} & \dots & \dots & \frac{\partial^2 f}{\partial z_K \partial z_K} \end{bmatrix}$$

$$\frac{\partial^2 f}{\partial z_i \partial z_j} = \frac{\partial}{\partial z_i} \frac{\partial}{\partial z_j} \left(\frac{1}{2} \sum_{k=1}^K (z_k - \gamma_k)^2 \right)$$

if $i = j$,

$$= \frac{\partial}{\partial z_i} \left(\frac{1}{2} (2z_i - 2\gamma_i) \right) = \frac{\partial}{\partial z_i} (z_i - \gamma_i) = 1$$

$$\frac{\partial^2 f}{\partial z_i \partial z_j} = \frac{\partial}{\partial z_i} \left(\frac{1}{2} (2z_j - 2\gamma_j) \right)$$

if $i \neq j$,

$$\frac{\partial}{\partial z_i} (z_j - \gamma_j) = 0$$

$$H(z) = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$$

So, the Hessian matrix of f is the $K \times K$ identity matrix.

Definition The real symmetric $N \times N$ matrix V is said to be positive semidefinite

$$\text{if } a^T V a \geq 0$$

for any real $N \times 1$ vector a .

$$\begin{aligned} \text{So, } a^T H(u) a &= a^T I a \\ &= a^T a \\ &= \sum_{i=1}^K a_i^2 \geq 0 \end{aligned}$$

So, the hessian matrix of f is positive semidefinite, hence f is convex.

Since the composition of a affine map $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$ with a convex function $f: \mathbb{R}^m \rightarrow \mathbb{R}$ is convex, i.e., $(f \circ g)(u)$ is convex,

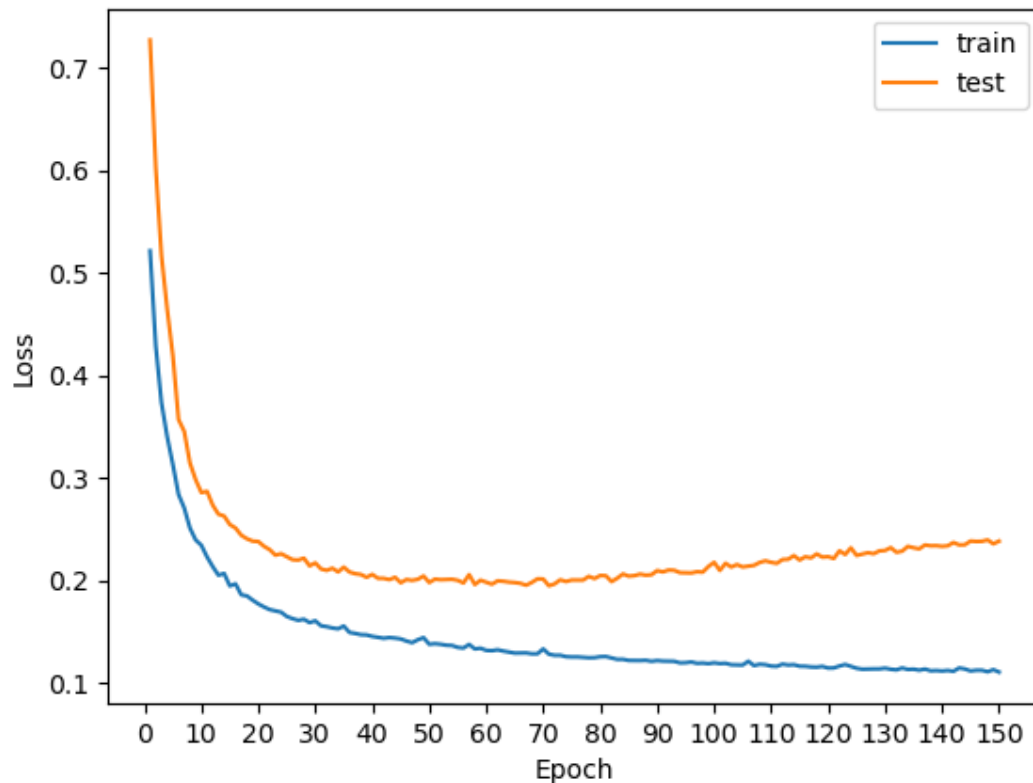
then the composition of $f = \frac{1}{2} \sum_{i=1}^K (\kappa_i - \tau_i)^2$ and $\kappa_i^{(w_i)} = w_i \phi(x) + b$

$$(f \circ g)(w_i) = \frac{1}{2} \sum_{i=1}^K (w_i \phi(x) + b - \tau_i)^2 \quad \text{is convex with}$$

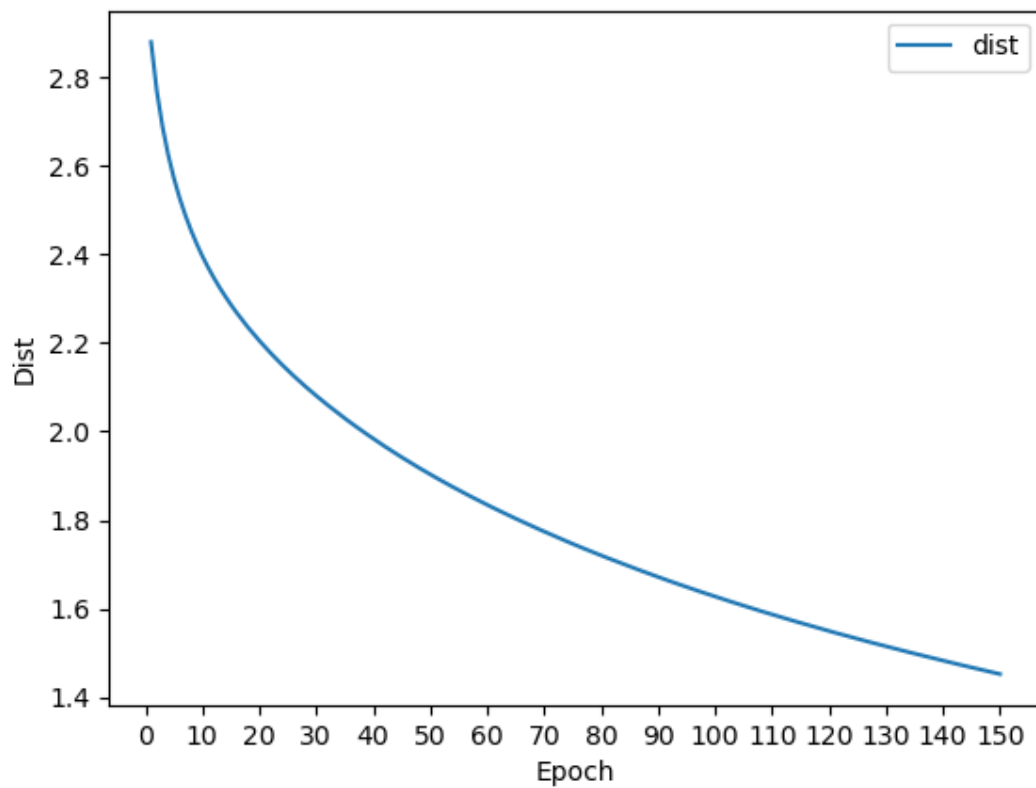
regards to w and b .

Exercise 2.a

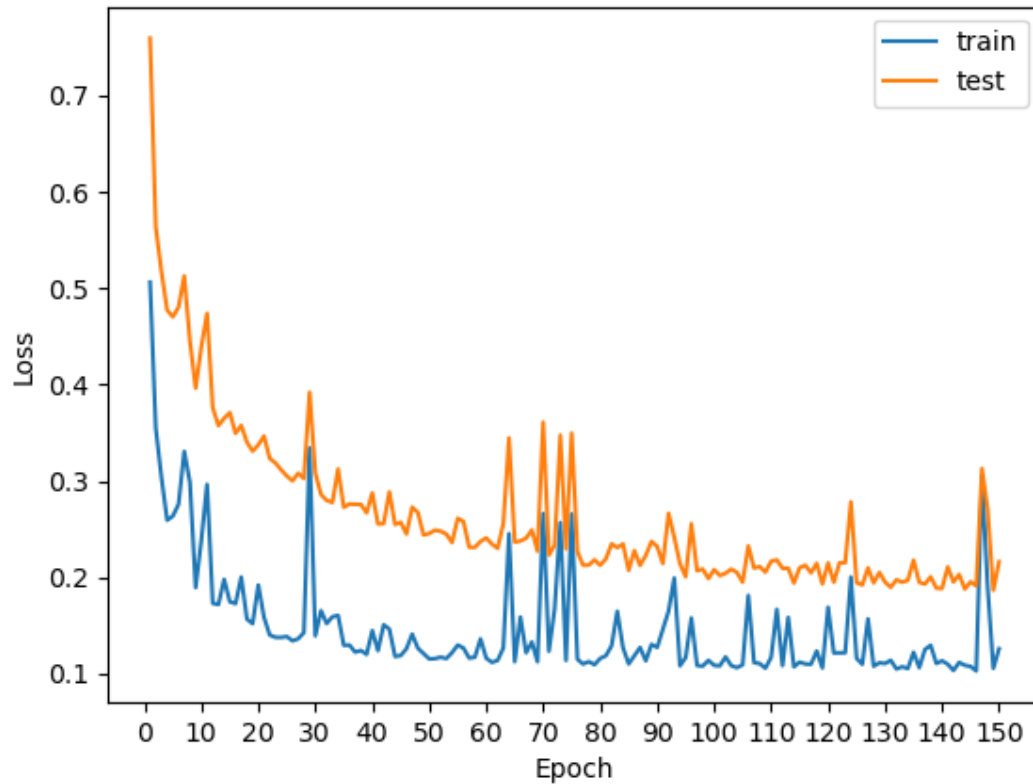
As the number of epochs increases the difference between the train and test sets' loss grows bigger. In the beginning both sets see a diminishing loss but at a certain point the train set's loss continues decreasing while the test set's loss increases. This is when our model starts to overfit. Overfitting is the reason why performance on the training and test sets differs as the number of epochs increases.



As we can see, the distance between our weight vector and the weight vector computed analytically decreases as the epochs increase.



Exercise 2.b

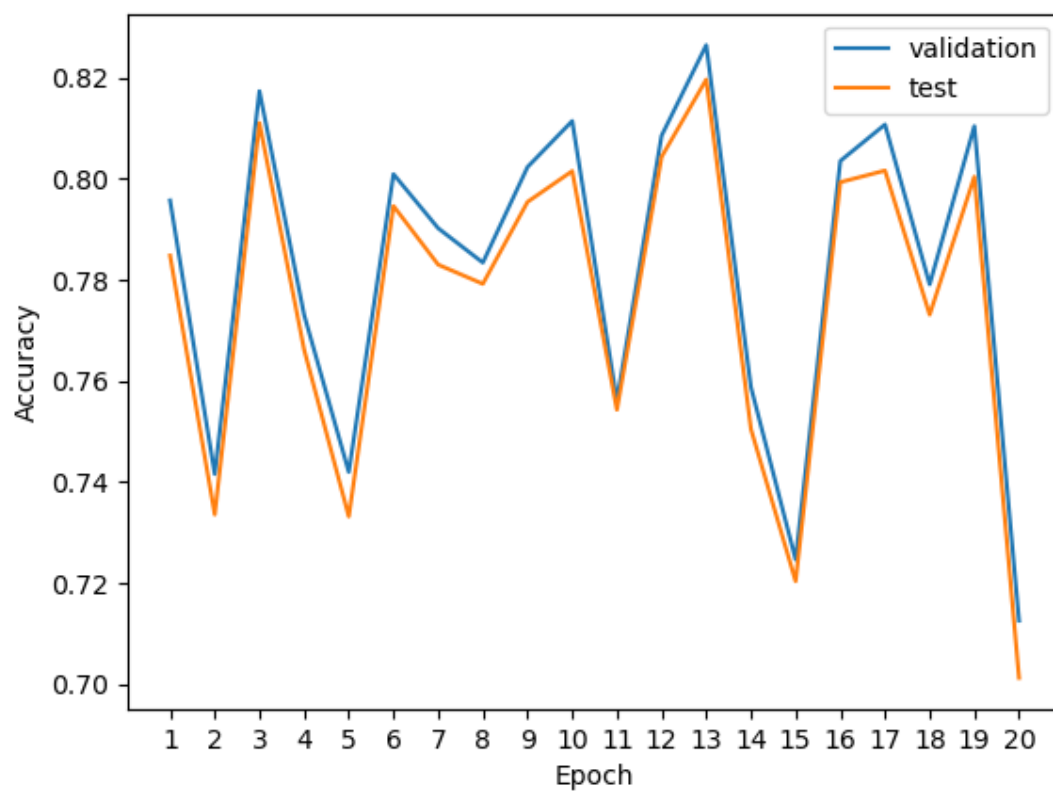


In comparison to the linear model, the observed difference between the performance in the training and test sets is that with a neural network we don't detect overfitting.

Question 3

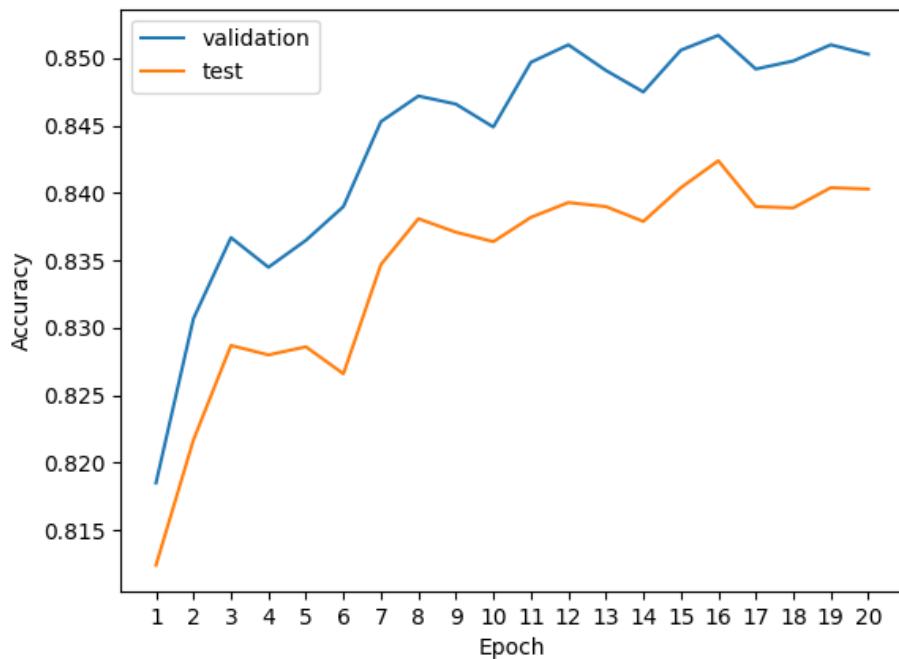
Exercise 1.a

Train with 20 epochs



Exercise 1.b

Train with 20 epochs

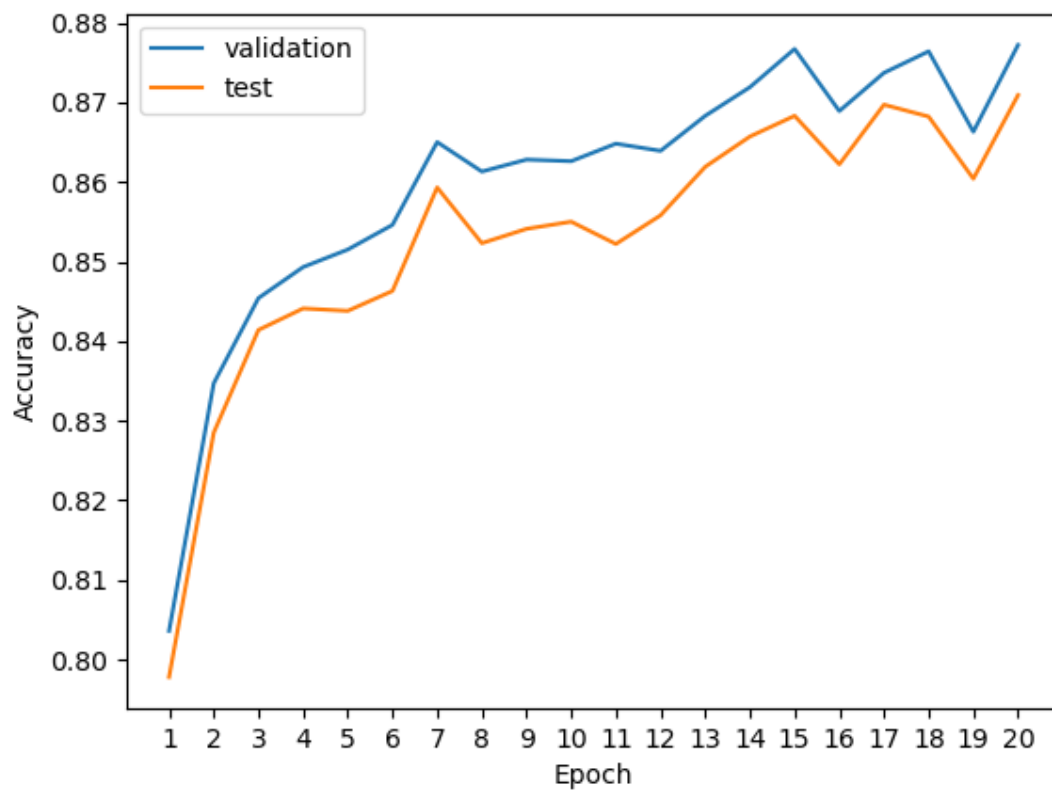


Exercise 2.a

The simple perceptron is a threshold-based feed-forward network that can only classify linearly separable cases with a binary target. Does not include hidden layers that allow the neural network to create a feature hierarchy.

Multi-layer perceptrons with non-linear activations are more expressive than the simple perceptron because they allow the model to create complex connections between the inputs of the Neural Network and outputs. They can solve the problems of linear activation functions because they allow backpropagation because they have a derived function that is related to the inputs and allow making a kind of "stack" of several layers of neurons constituted by hidden layers that are essential for learning complex datasets with high levels of accuracy. The limitation is that a Multi-layer perceptron with non-linear activations is different from any structure with one layer if the activation function is the same.

Exercise 2.b



Question 4

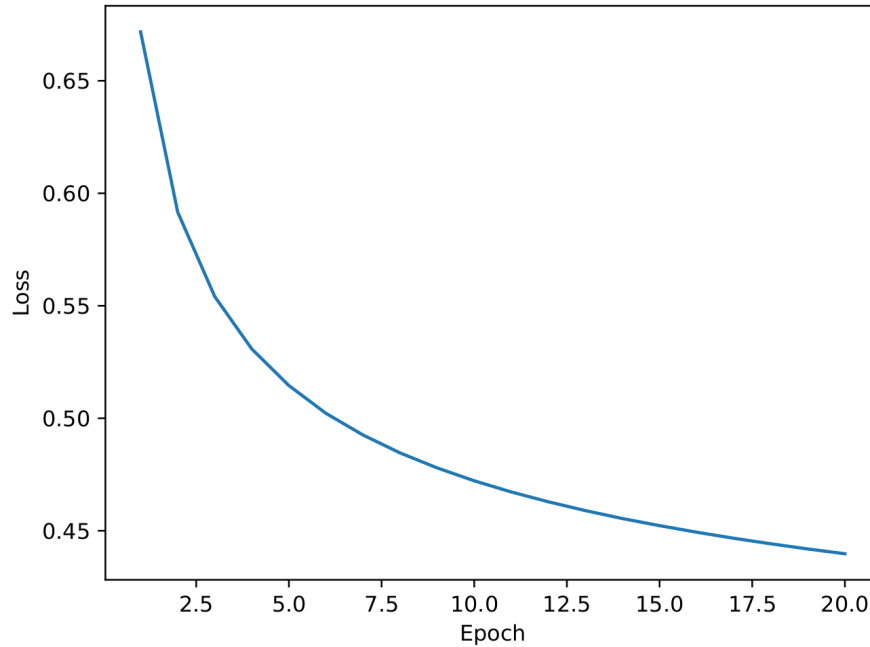
Exercise 1

This exercise consisted of building a logistic regression model using pytorch and training it using the MNIST dataset with 20 epochs. Three different learning rates were used {0.001,

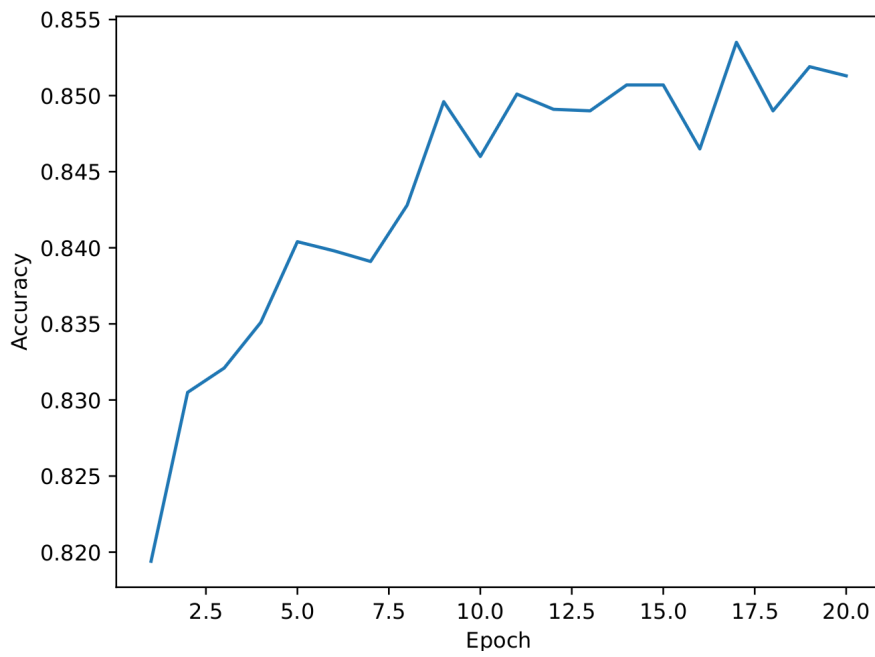
0.01 and 0.1} and their respective accuracies and losses were compared. As it can be observed, the configuration which had a better performance was the one using a learning rate of 0.001.

The final test accuracy with $lr = 0.001$: 0.8402

Training loss with $lr = 0.001$:

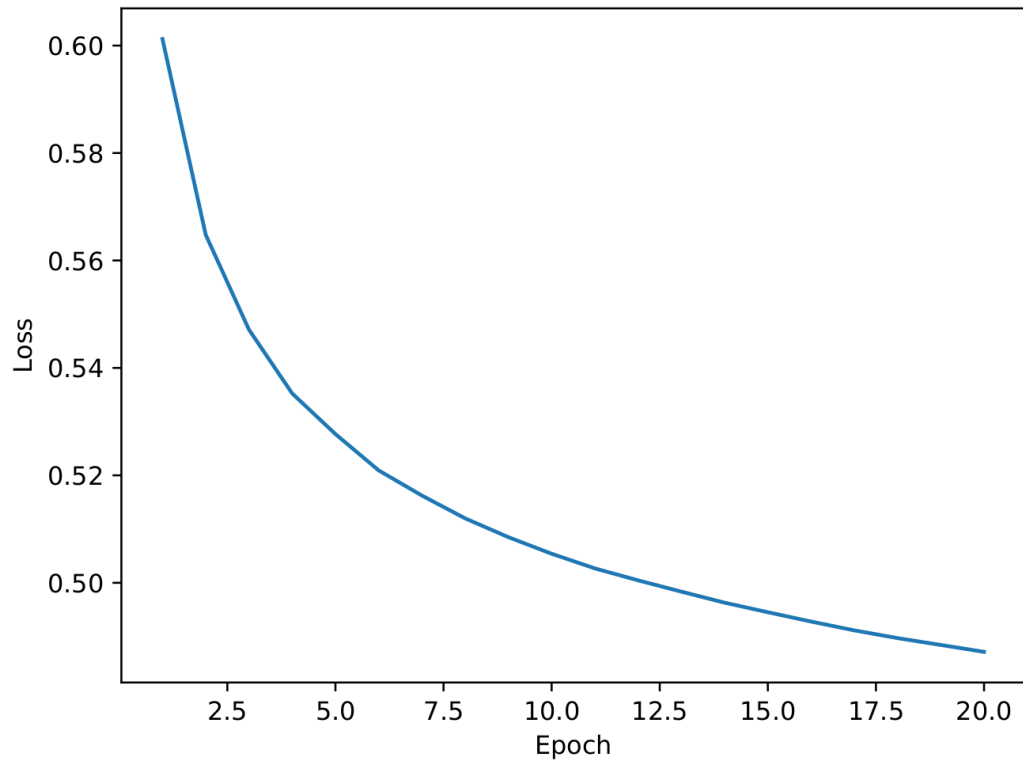


Validation accuracy with $lr = 0.001$:

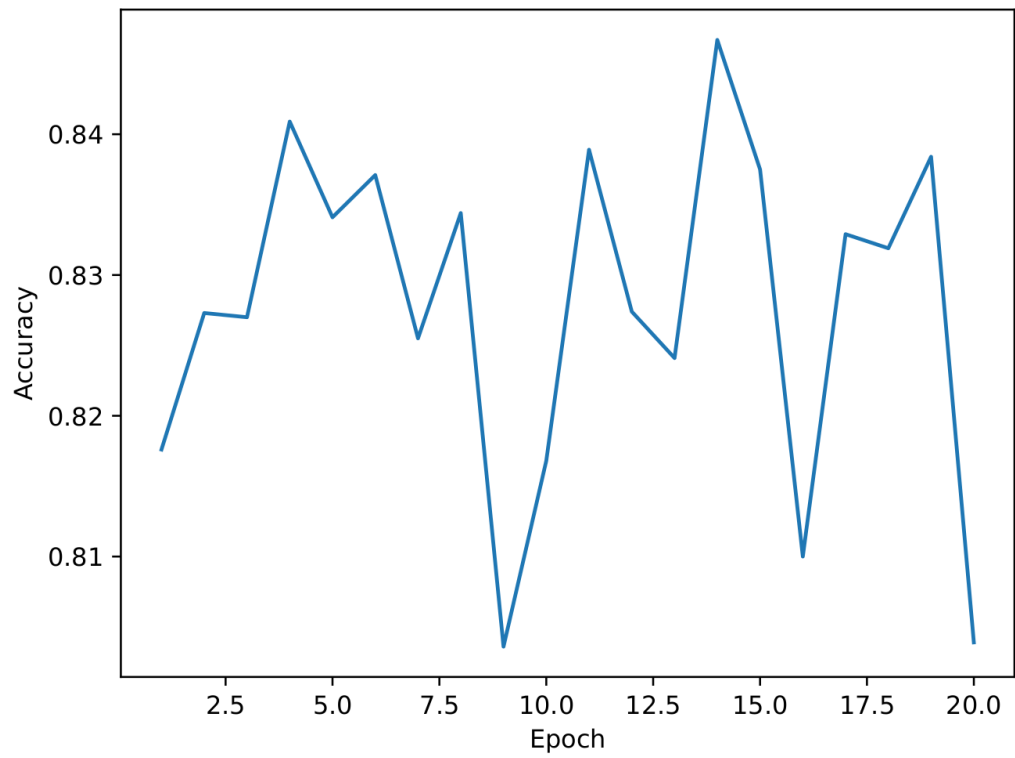


The final test accuracy with $lr = 0.01$: 0.7944

Training loss with $lr = 0.01$:

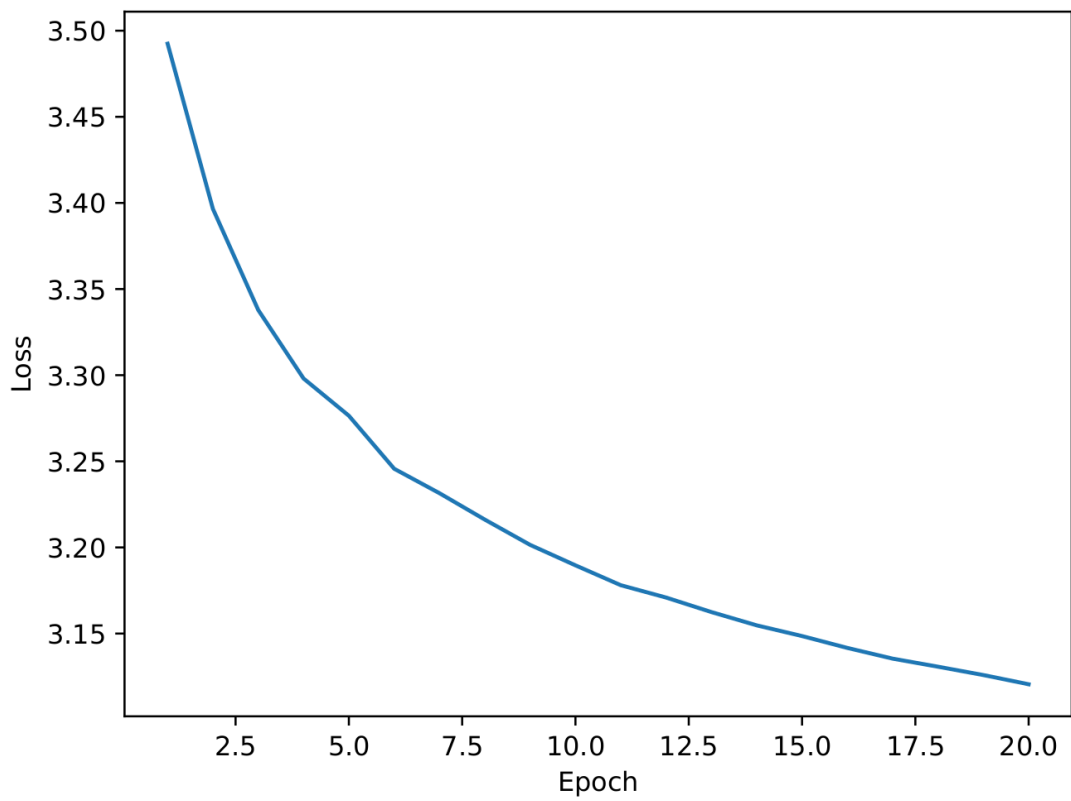


Validation accuracy with $lr = 0.01$:

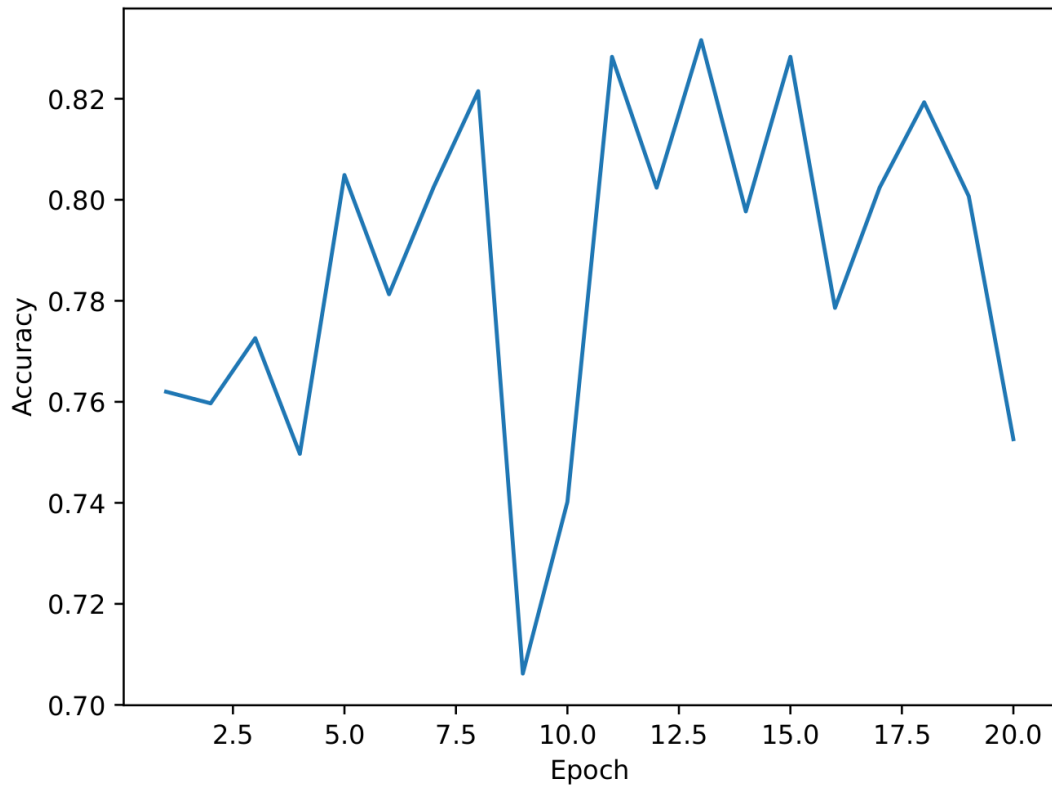


The final test accuracy with $\text{lr} = 0.1$: 0.7500

Training loss with $\text{lr} = 0.1$:



Validation accuracy with $\text{lr} = 0.1$:

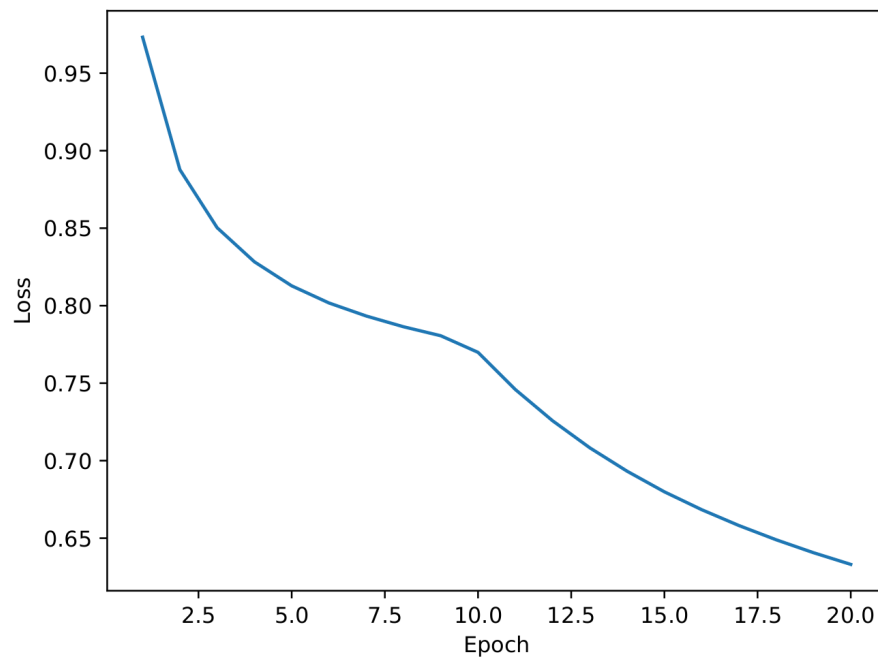


Exercise 2

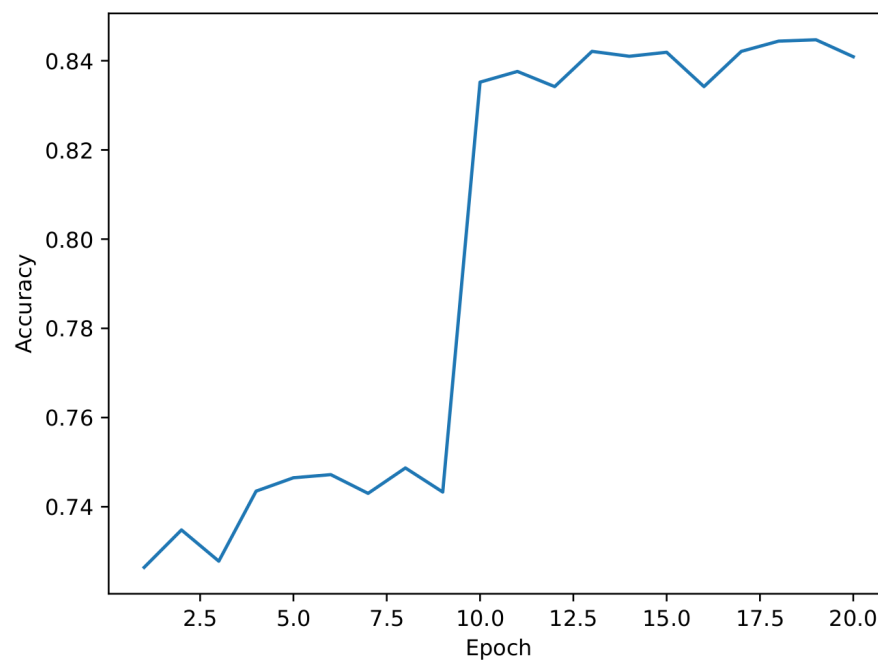
In this section we implemented a feed-forward neural network with one layer and compared the effects of different hyperparameters.

The final test accuracy with $lr = 0.001$: 0.8324

Training loss with $lr = 0.001$:

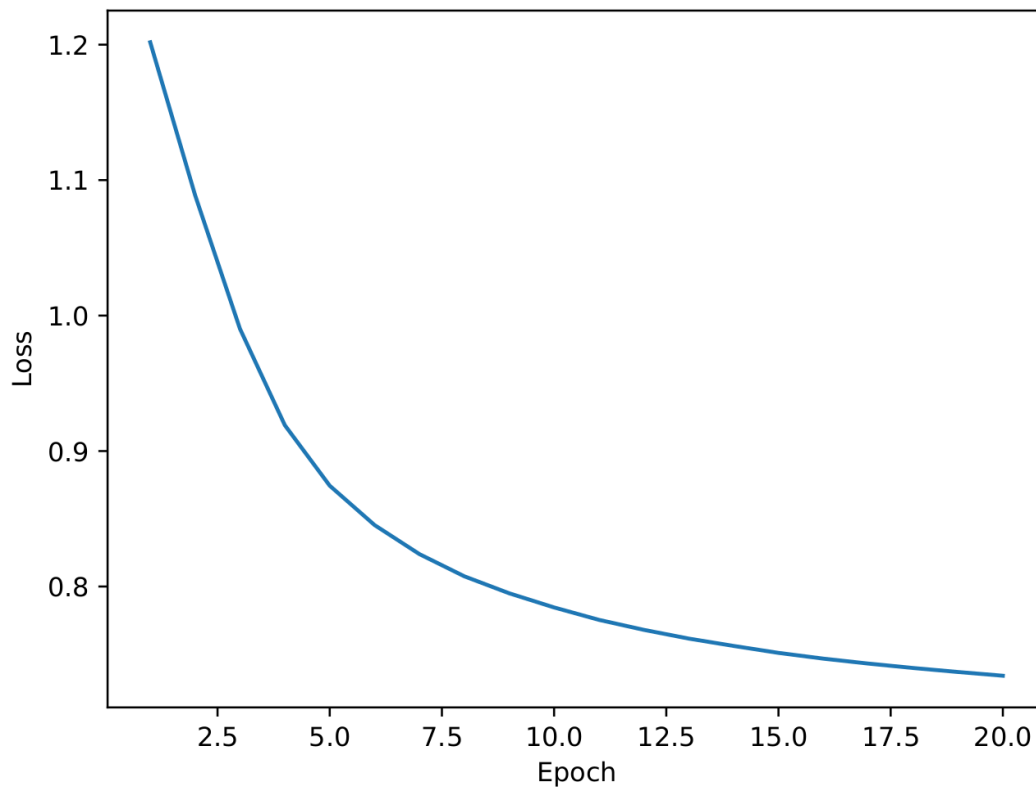


Validation accuracy with $\text{lr} = 0.001$:

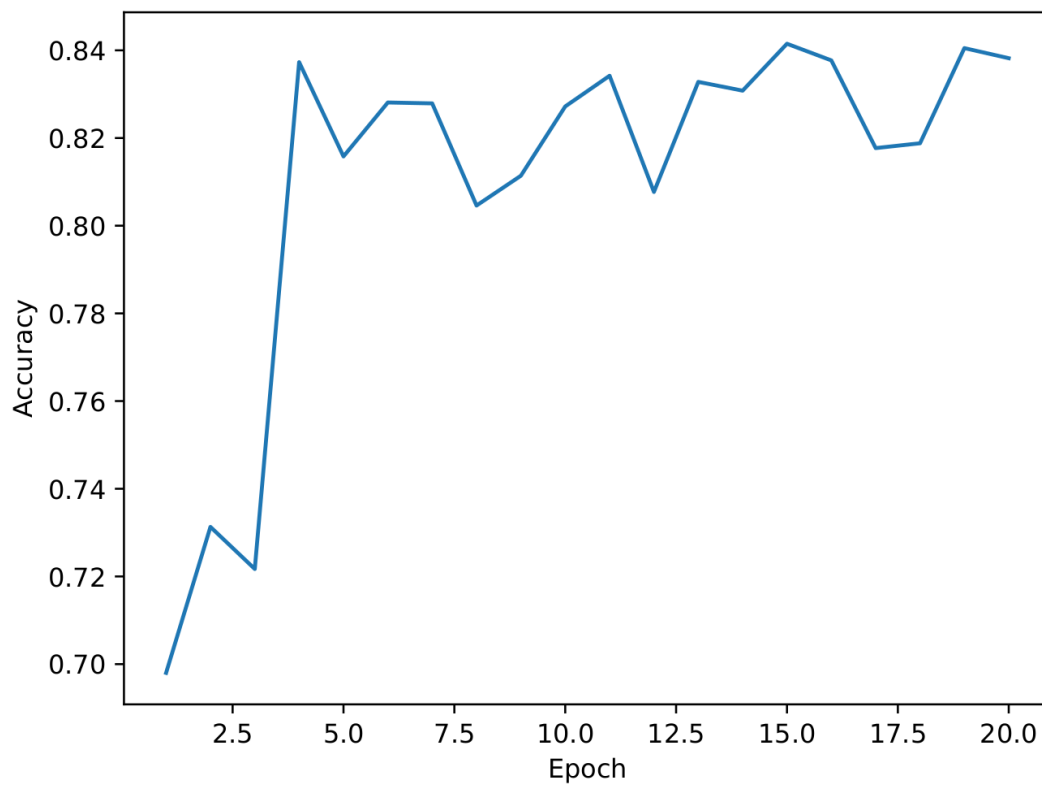


The final test accuracy with $\text{lr} = 0.01$: 0.8284

Training loss with $\text{lr} = 0.01$:

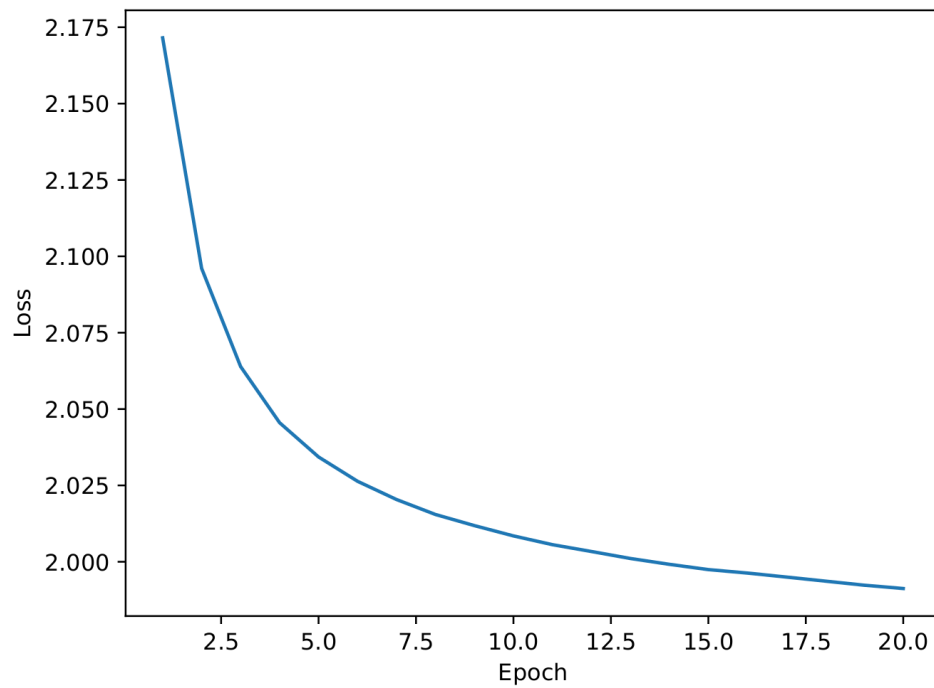


Validation accuracy with $\text{lr} = 0.01$:

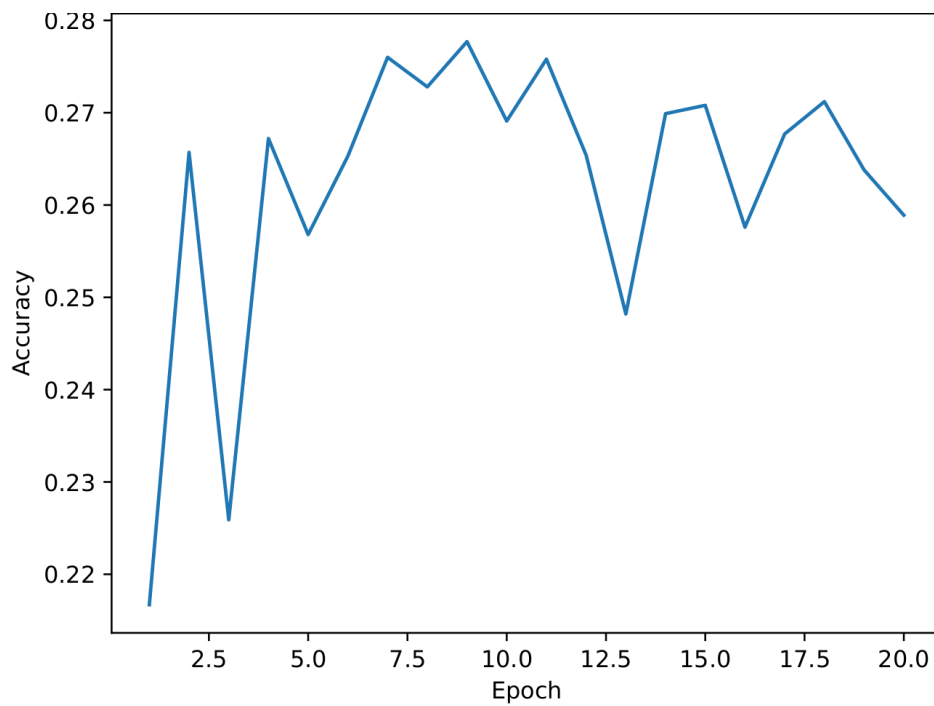


The final test accuracy with $\text{lr} = 0.1$: 0.2552

Training loss with $\text{lr} = 0.1$:

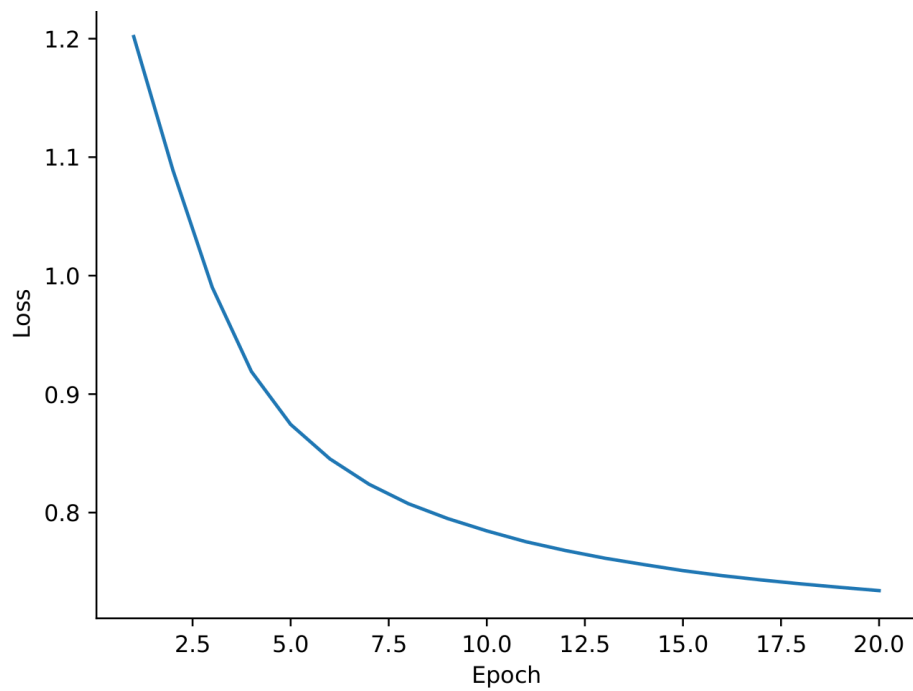


Validation accuracy with $\text{lr} = 0.1$:

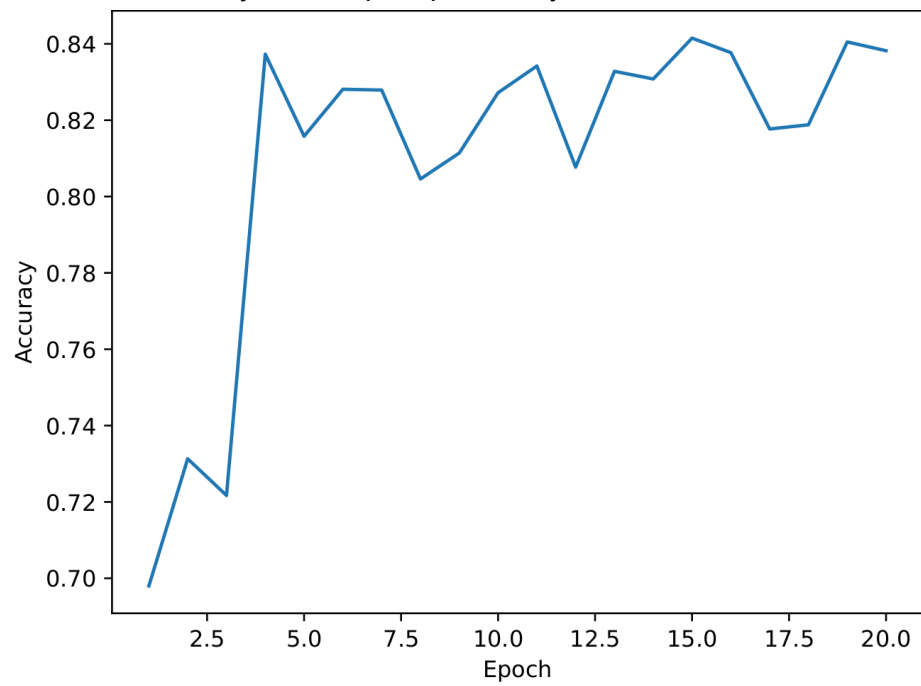


The final test accuracy with dropout probability = 0.3: 0.8284

Training loss with dropout probability = 0.3:

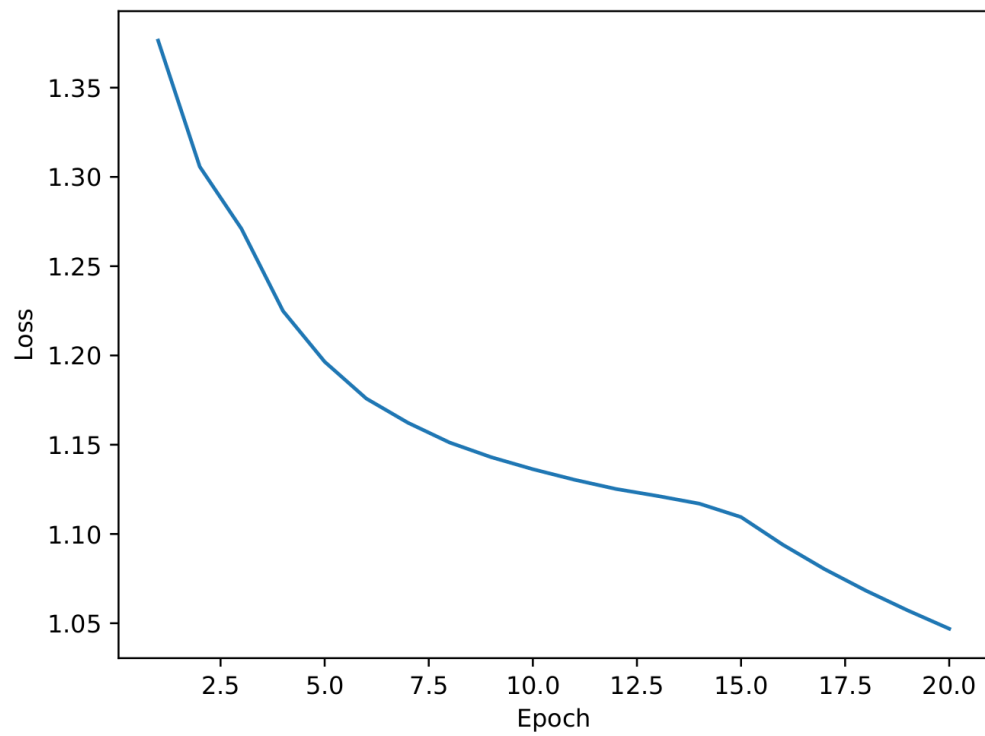


Validation accuracy with dropout probability = 0.3:

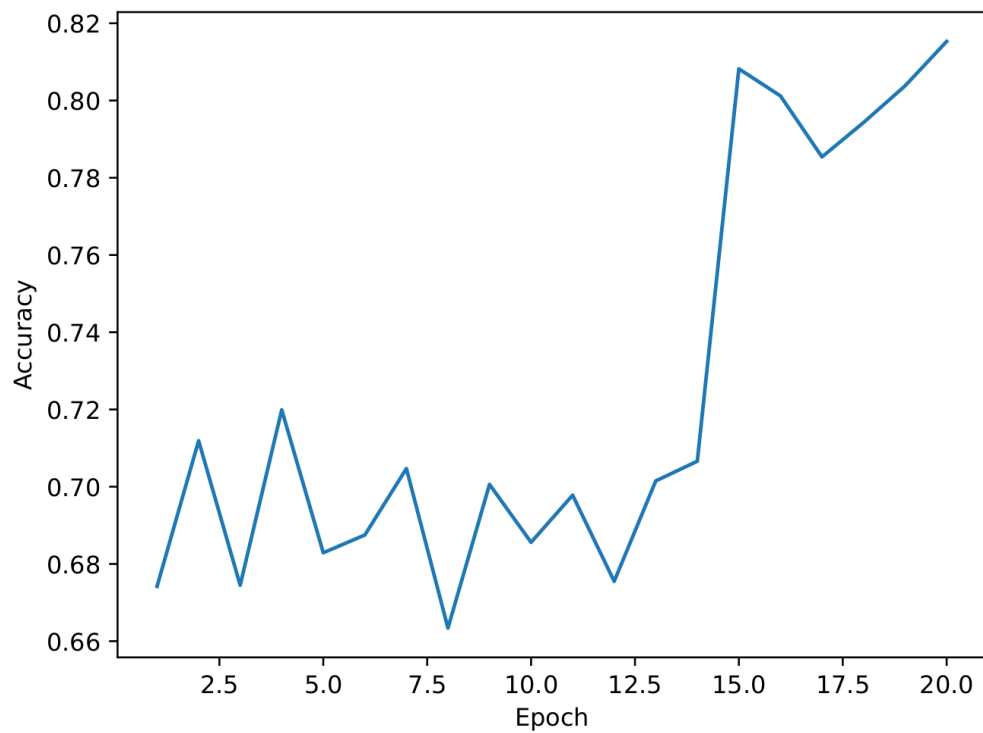


The final test accuracy with dropout probability = 0.5: 0.8114

Training loss with dropout probability = 0.5:

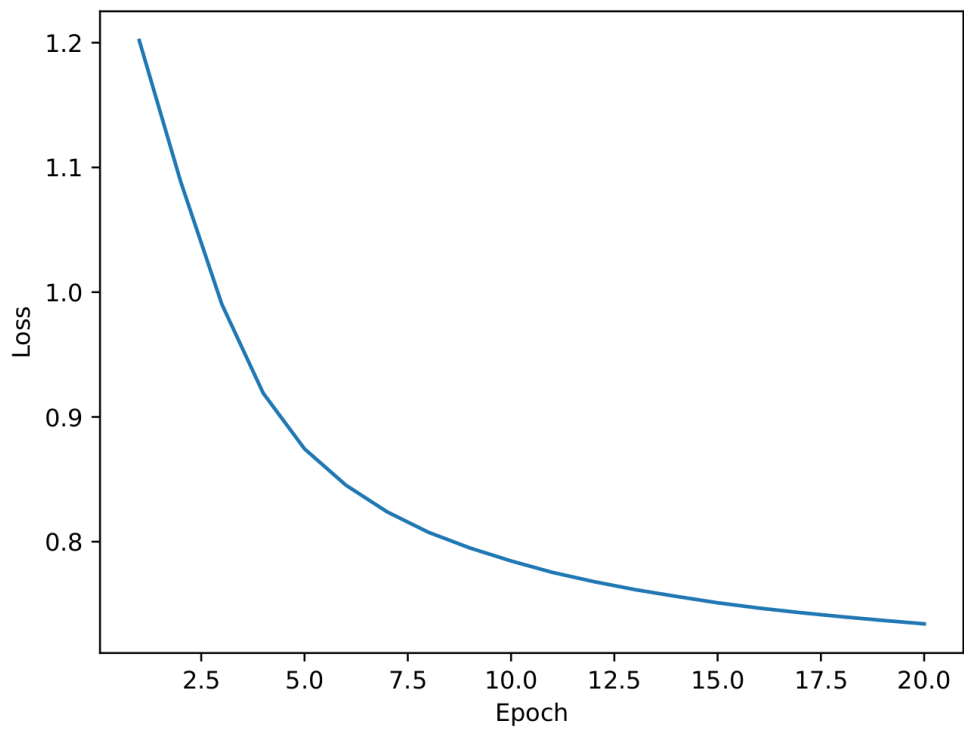


Validation accuracy with dropout probability = 0.5:

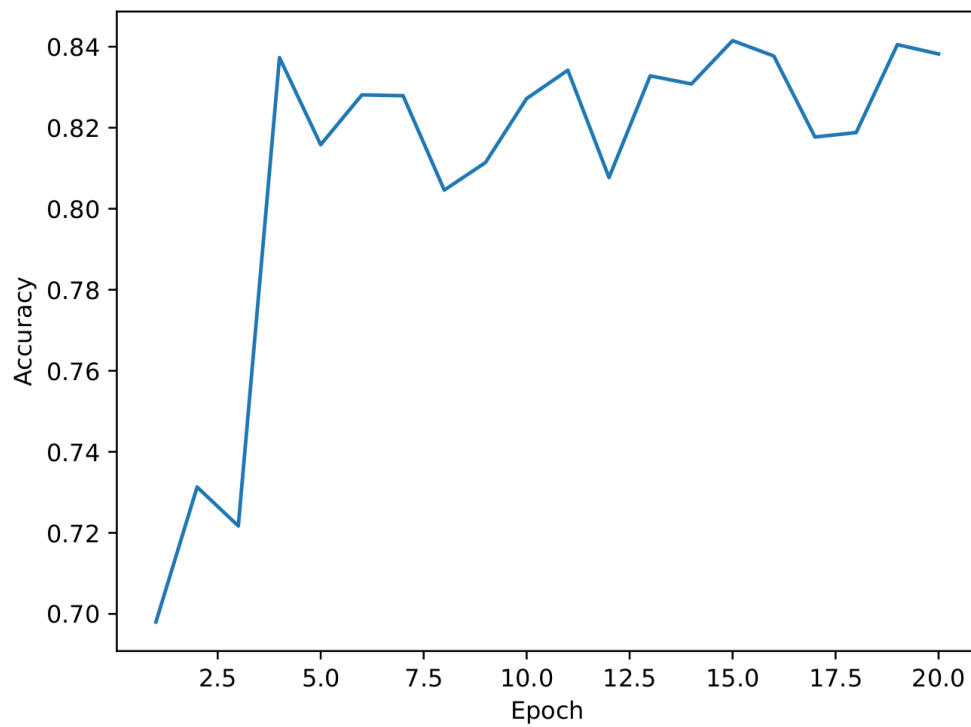


The final test accuracy with activation function ReLU: 0.8284

Training loss with activation function ReLU:

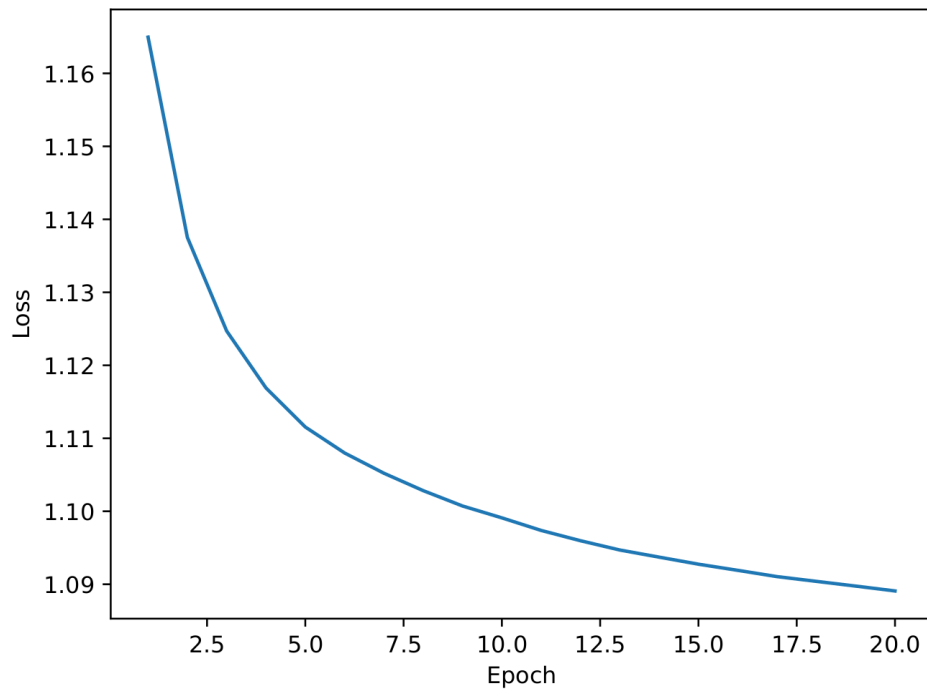


Validation accuracy with activation function ReLU:

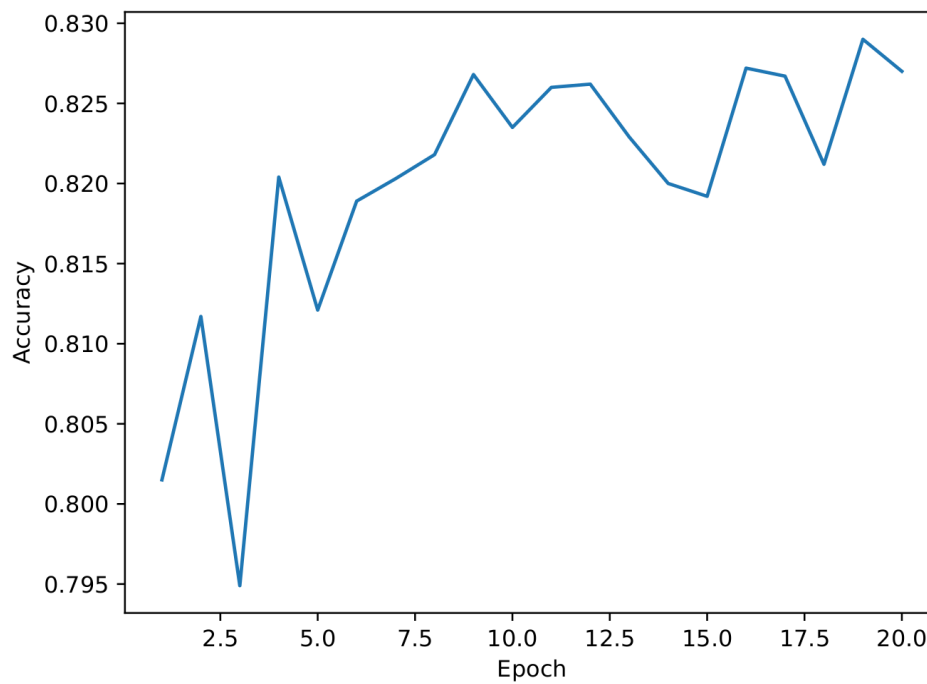


The final test accuracy with activation function Tanh: 0.8207

Training loss with activation function Tanh:

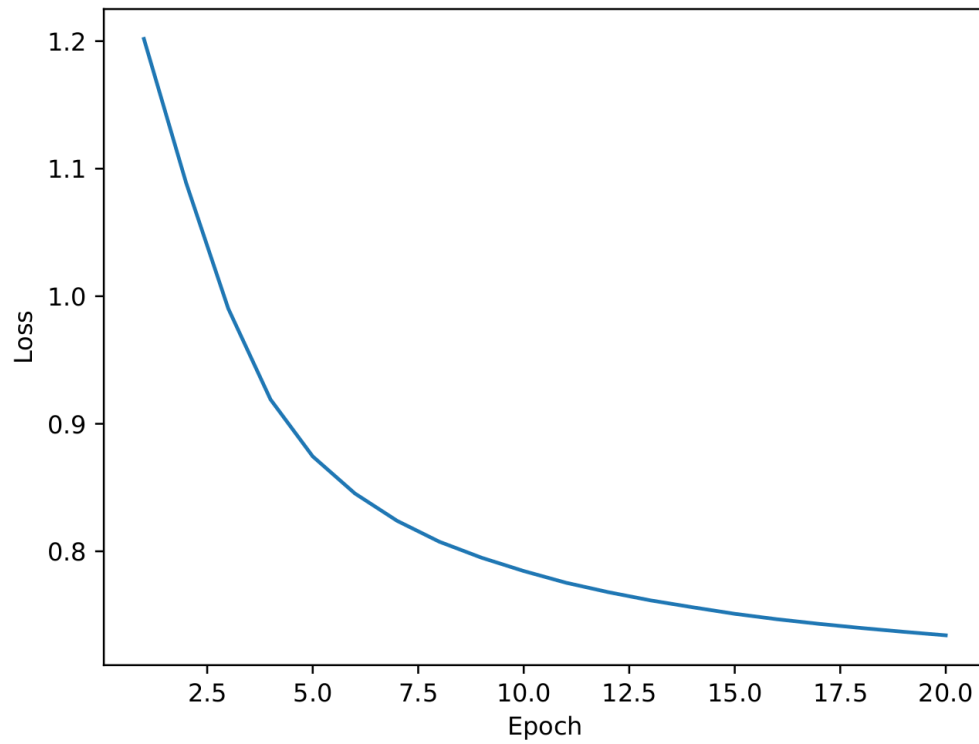


Validation accuracy with activation function Tanh:

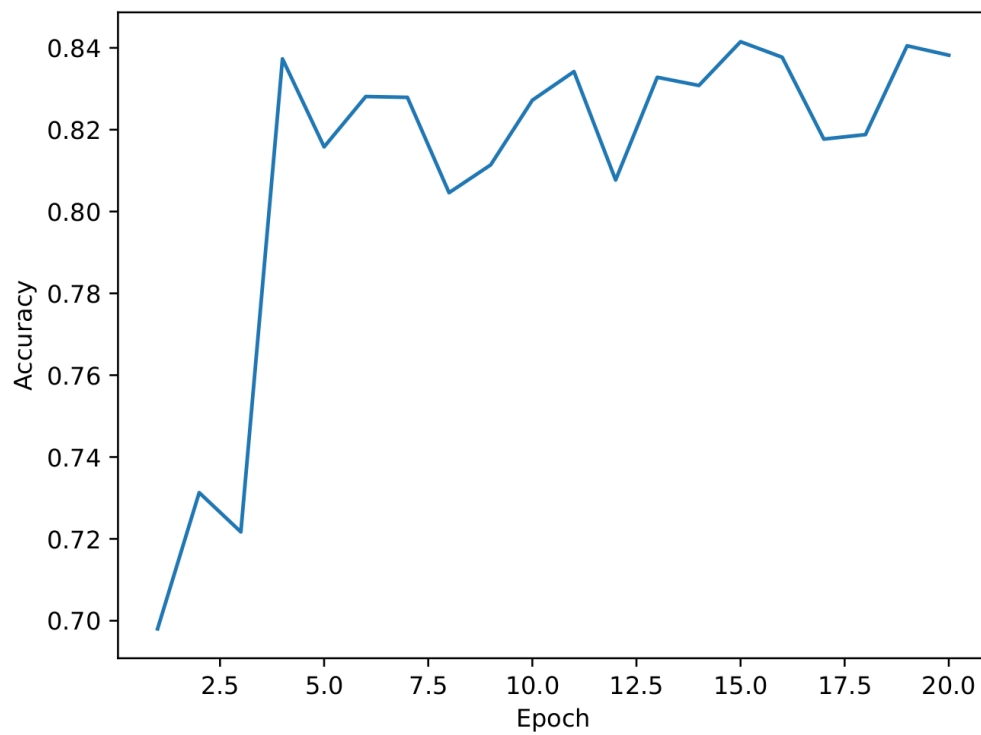


The final test accuracy with optimizer SGD: 0.8284

Training loss with optimizer SGD:

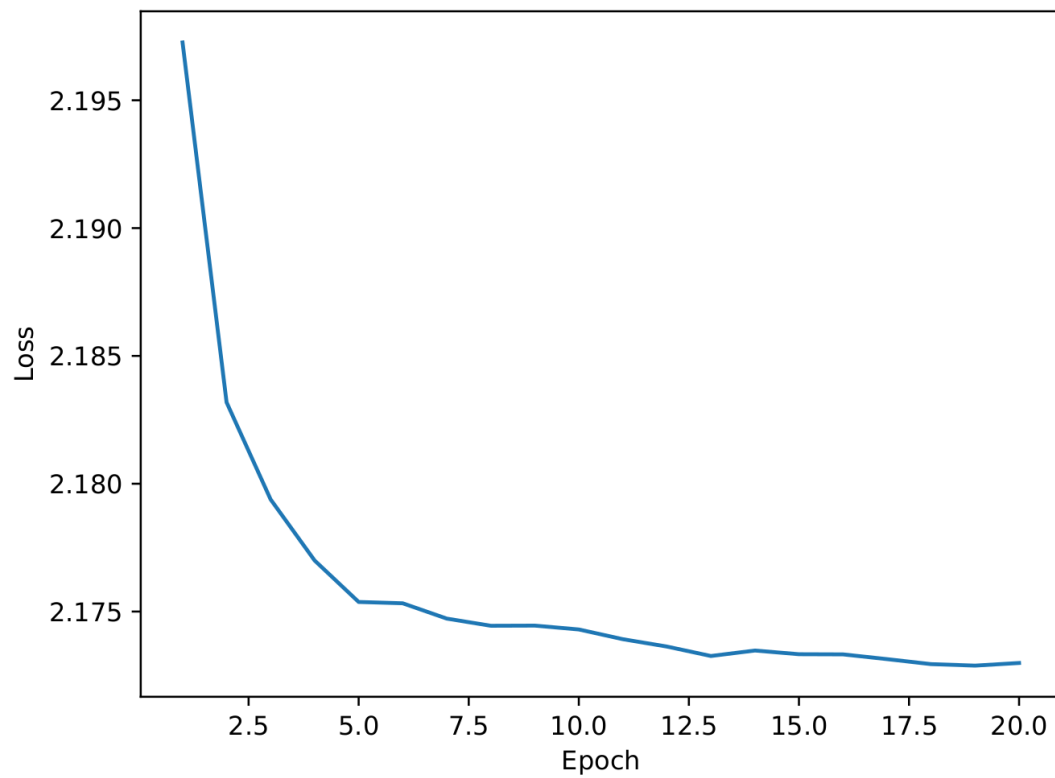


Validation accuracy with optimizer SGD:

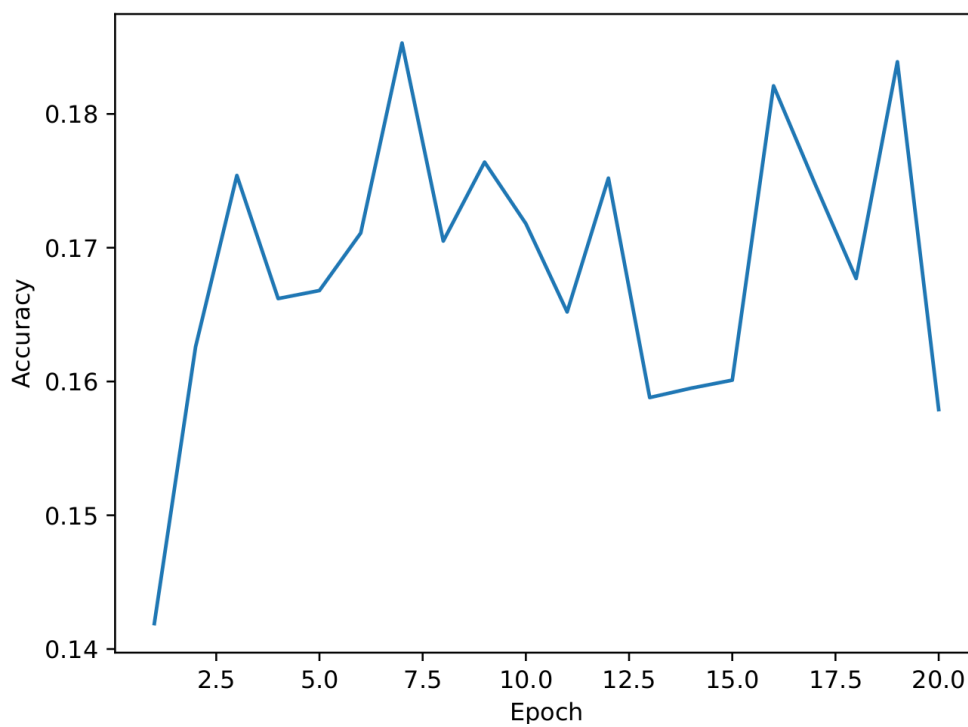


The final test accuracy with optimizer ADAM: 0.1579

Training loss with optimizer ADAM:



Validation accuracy with optimizer ADAM:



Taking into account all the test accuracies measured, it is clear that the highest one was observed when implementing the neural network with a learning rate of 0.001. Although

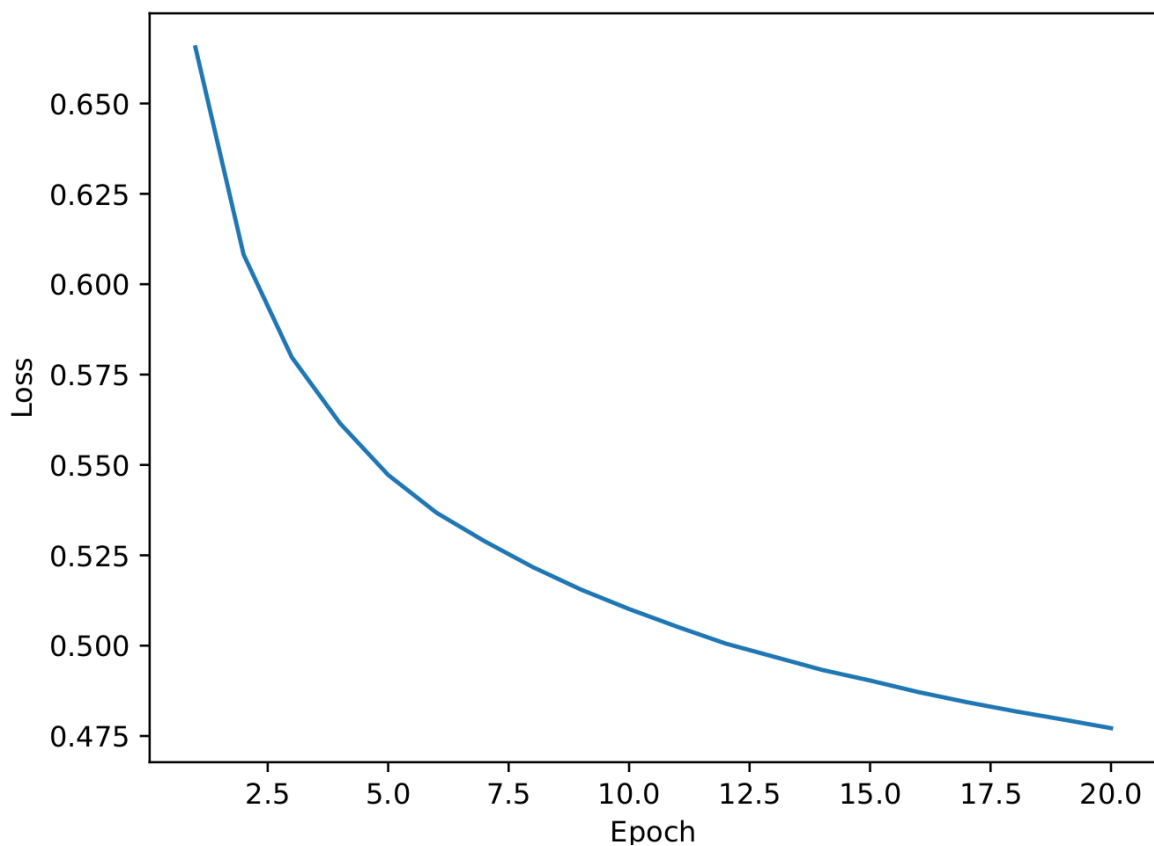
most of them have similar performances, one could speculate that smaller "steps" during the optimization process leads to function approximators that better generalize to the test data.

Exercise 3

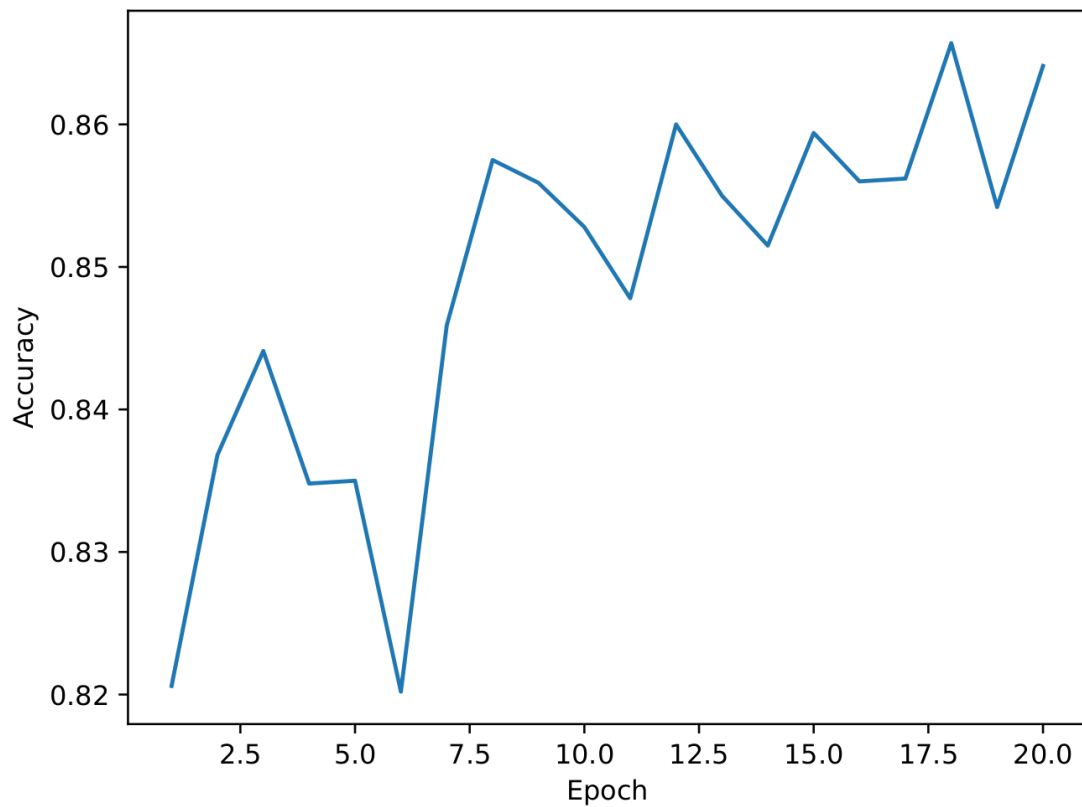
In this final exercise, we implemented a 2-layer and a 3-layer neural network using the same hyperparameters as the previous exercise. As before, performances were plotted and tested.

The final test accuracy with 2 layers: 0.8576

Training loss with 2 layers:

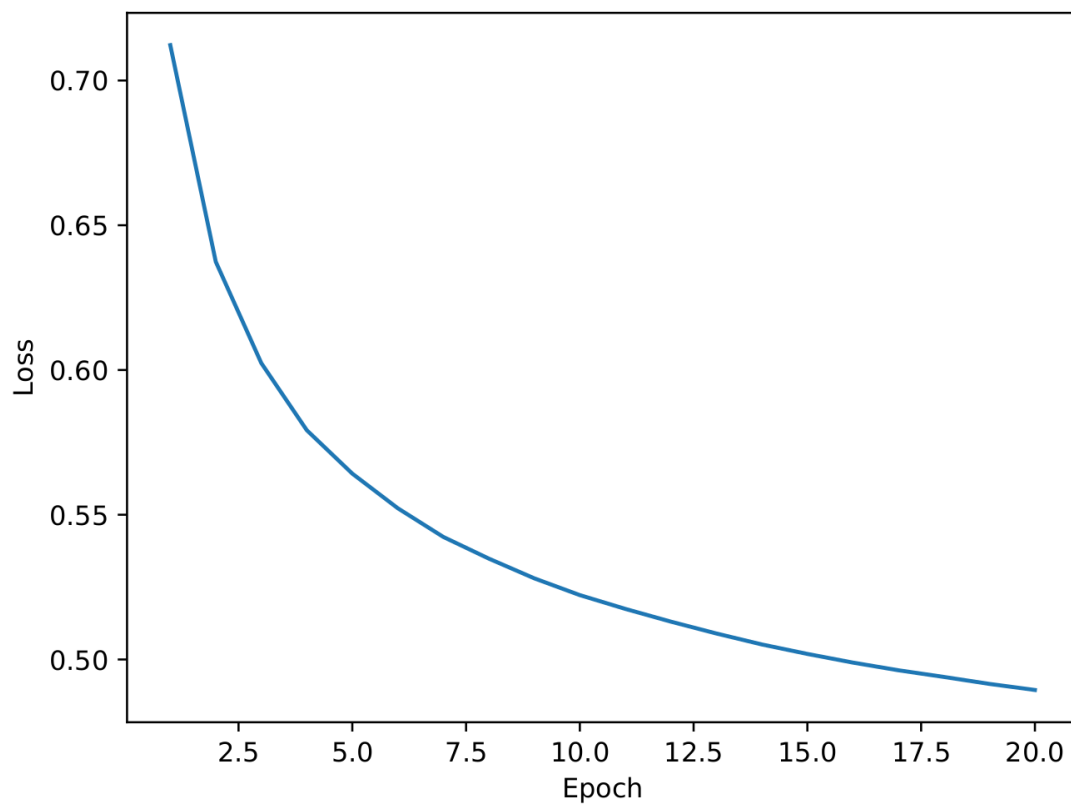


Validation accuracy with 2 layers:

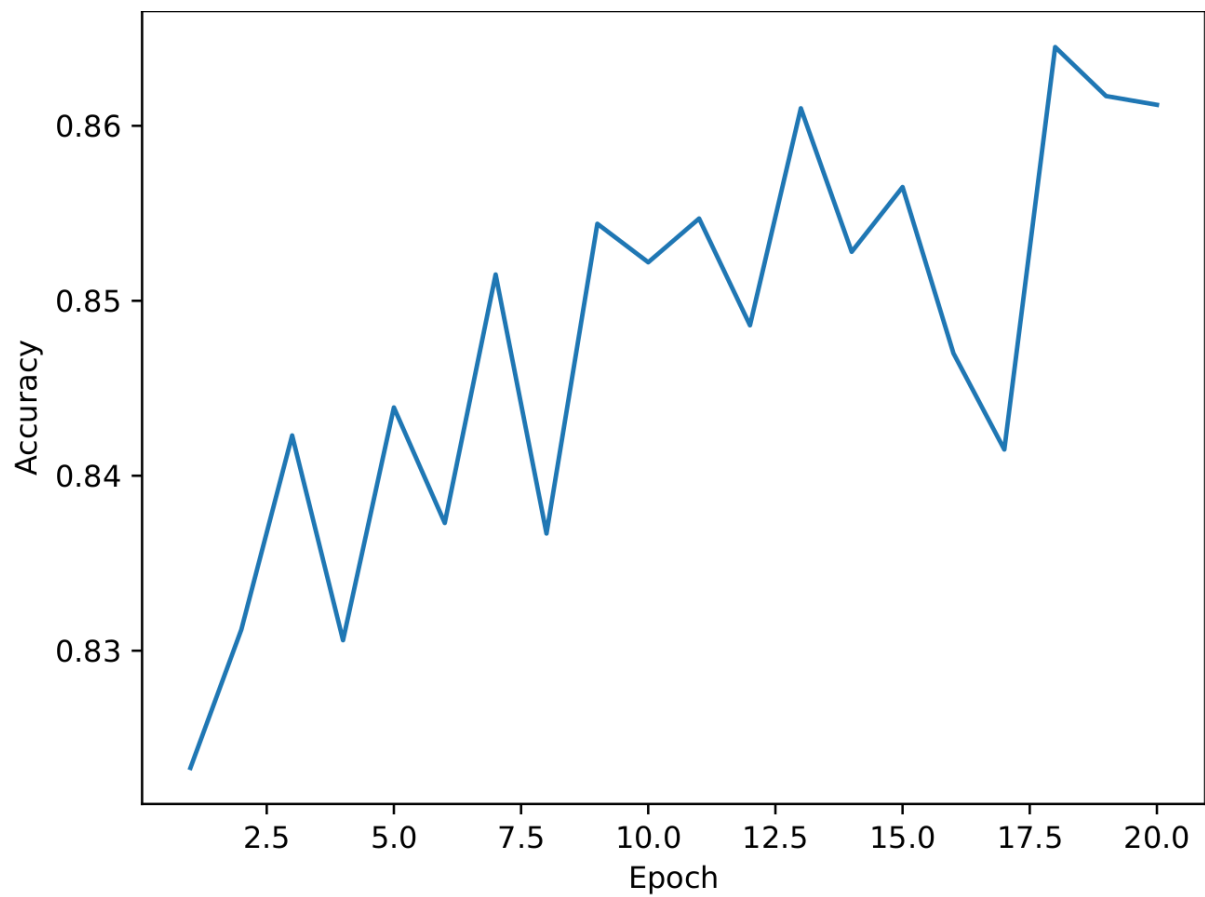


The final test accuracy with 3 layers: 0.8573

Training loss with 3 layers:



Validation accuracy with 3 layers:



As expected, the difference in the test accuracies of 1-layer and 2-layer perceptrons are noticeable, with the 3-layer perceptron having similar performance to the 2-layer, with slightly less test accuracy but generally higher valid accuracies during training.