

Q 1.2.

a)

1. Import the required libraries, and define the model architecture.
2. Initialize the model and define the loss function.
3. Forward pass, and compute loss.
4. Perform backward pass, and calculate derivatives.
5. Update the parameters.

b)

Layer	Input	Output
Linear1	x	$W(1)x + b(1)$
f	$W(1)x + b(1)$	$\text{ReLU}(W(1)x + b(1))$
Linear2	$\text{ReLU}(W(1)x + b(1))$	$W(2) \text{ReLU}(W(1)x + b(1)) + b(2)$
g	$W(2) \text{ReLU}(W(1)x + b(1)) + b(2)$	$W(2) \text{ReLU}(W(1)x + b(1)) + b(2)$
MSE	$(W(2) \text{ReLU}(W(1)x + b(1)) + b(2)), y)$	

c)

Parameter	Gradient
$W(1)$	$\delta l / \delta W(1) = \delta l / \delta y * \delta y / \delta z3 * \delta z3 / \delta z2 * \delta z2 / \delta z1 * \delta z1 / \delta W(1)$
$b(1)$	$\delta l / \delta b(1) = \delta l / \delta y * \delta y / \delta z3 * \delta z3 / \delta z2 * \delta z2 / \delta z1 * \delta z1 / \delta b(1)$
$W(2)$	$\delta l / \delta W(2) = \delta l / \delta y * \delta y / \delta z3 * \delta z3 / \delta W(2)$
$b(2)$	$\delta l / \delta b(2) = \delta l / \delta y * \delta y / \delta z3 * \delta z3 / \delta b(2)$

d)

$$\delta l / \delta y = 2(y^{\wedge} - y)$$

$$\delta y^{\wedge} / \delta z3 = 1$$

$$\delta z2 / \delta z1 = 0 \text{ if } z1 \leq 0 \text{ and } 1 \text{ if } z1 > 0$$

Q 1.3.

a)

Layer	Input	Output
Linear1	x	$W(1)x + b(1)$
f	$W(1)x + b(1)$	$\text{sigmoid}(W(1)x + b(1))$
Linear2	$\text{sigmoid}(W(1)x + b(1))$	$W(2) \text{ sigmoid}(W(1)x + b(1)) + b(2)$
g	$Z2 = W(2) \text{ sigmoid}(W(1)x + b(1)) + b(2)$	$\text{Sigmoid}(Z2)$
MSE	$\text{Sigmoid}(Z2)$	

Parameter	Gradient
$W(1)$	$\delta l / \delta W(1) = \delta l / \delta y * \delta y^{\wedge} / \delta z3 * \delta z3 / \delta z2 * \delta z2 / \delta z1 * \delta z1 / \delta W(1)$
$b(1)$	$\delta l / \delta b(1) = \delta l / \delta y * \delta y^{\wedge} / \delta z3 * \delta z3 / \delta z2 * \delta z2 / \delta z1 * \delta z1 / \delta b(1)$
$W(2)$	$\delta l / \delta W(2) = \delta l / \delta y * \delta y^{\wedge} / \delta z3 * \delta z3 / \delta W(2)$
$b(2)$	$\delta l / \delta b(2) = \delta l / \delta y * \delta y^{\wedge} / \delta z3 * \delta z3 / \delta b(2)$

$$\delta l / \delta y = 2(y^{\wedge} - y)$$

$$\delta y^{\wedge} / \delta z3 = \sigma(z3)(1 - \sigma(z3))$$

$$\delta z2 / \delta z1 = \sigma(z1)(1 - \sigma(z1))$$

b)

for Q 1.2. b, we need to change the Loss from MSE to BCE.

Also we need to modify the gradient of the loss w.r.t. the predicted output  $y^{\wedge}$ :

$$\delta_{bce}/\delta y^{\wedge} = -1/K * (y/y^{\wedge} - (1-y)/(1-y^{\wedge}))$$

c)

Sigmoid activation function tends to saturate for large positive or negative inputs, resulting in very small gradients. When backpropagating the gradients through multiple layers, the gradients can diminish significantly, leading to the vanishing gradient problem. This can make it difficult for the earlier layers in the network to learn effectively. ReLU, on the other hand, does not suffer from this saturation issue, allowing gradients to flow more freely and mitigate the vanishing gradient problem.