

ECE 284 HW2 (Jing-Hua Chang A69040678)

Q2:

$$\text{MSE Loss: } L(w, b) = \frac{1}{m} \sum_{j=1}^m (y_j - \hat{y}_j)^2 = \frac{1}{2} [(Y_{3_1} - T_1)^2 + (Y_{3_2} - T_2)^2]$$

$$S_2 = [1, 2] \begin{bmatrix} 2 \\ 0 \end{bmatrix} = [2, 3]$$

$$Y_2 = \phi([2, 3]) = [2, 3]$$

$$S_3 = [2, 3] \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix} = [5, 10]$$

$$Y_3 = \phi([5, 10]) = [5, 10]$$

$$L = \frac{1}{2} [(5 - 0)^2 + (10 - 0)^2] = \frac{1}{2} (25 + 100) = 62.5$$

```
print(W1.grad)
print(W2.grad)
print(Y3)

tensor([[25., 50.], [50., 50.]])
tensor([[10., 20.], [15., 30.]])
tensor([5., 10.], grad_fn=<ReluBackward0>)
```

$$\frac{\partial L}{\partial S_3} = \frac{\partial L}{\partial Y_3} \odot \text{ReLU}'(S_3) = (Y_3 - T) \odot [1, 1] = [5, 10]$$

$$S_3 = Y_2 W_2 \Rightarrow \frac{1}{W_2} = Y_2 \frac{1}{S_3} \frac{\partial L}{\partial L} \Rightarrow \frac{\partial L}{\partial W_2} = Y_2^T \cdot \frac{\partial L}{\partial S_3}$$

$$= \begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot [5 \ 10] = \begin{bmatrix} 10 & 20 \\ 15 & 30 \end{bmatrix}$$

* gradient of W_2

$$S_3 = Y_2 W_2 \Rightarrow \frac{\partial L}{\partial S_3} = \frac{\partial L}{\partial Y_2} W_2 \Rightarrow \frac{\partial L}{\partial Y_2} = \frac{\partial L}{\partial S_3} W_2^T$$

$$= [5, 10] \begin{bmatrix} 1 & 1 \\ 2 & 2 \end{bmatrix} = [25, 25]$$

$$\frac{\partial L}{\partial S_2} = \frac{\partial L}{\partial Y_2} \odot \text{ReLU}'(S_2) = [25, 25] \odot [1, 1] = [25, 25]$$

$$S_2 = X W_1 \Rightarrow \frac{\partial L}{\partial W_1} = X^T \frac{\partial L}{\partial S_2} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \cdot [25 \ 25] = \begin{bmatrix} 25 & 25 \\ 50 & 50 \end{bmatrix}$$

* gradient of W_1

Q3:

$$loss = \frac{1}{2} [(Y_3_1 - T_1)^2 + (Y_3_2 - T_2)^2]$$

$$S2 = XW1 = \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 0 & -1 \end{bmatrix} = \begin{bmatrix} 2 & -1 \end{bmatrix}$$

$$Y2 = \begin{bmatrix} 2 & 0 \end{bmatrix}$$

$$S3 = \begin{bmatrix} 2 & 0 \end{bmatrix} \begin{bmatrix} -1 & 2 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} -2 & 4 \end{bmatrix}$$

$$Y3 = \begin{bmatrix} 0 & 4 \end{bmatrix}$$

$$S3 = Y2W2 \Rightarrow \frac{\partial L}{\partial W2} = Y2^T \frac{\partial L}{\partial S3} = \begin{bmatrix} 2 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 0 & 4 \end{bmatrix} = \begin{bmatrix} 0 & 8 \\ 0 & 0 \end{bmatrix}$$

$$\frac{\partial L}{\partial S3} = \frac{\partial L}{\partial Y3} \circ ReLU'(S3) = (Y3 - T) \cdot \begin{bmatrix} 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 4 \end{bmatrix} \quad *$$

$$S2 = XW1 \Rightarrow \frac{\partial L}{\partial W1} = X^T \frac{\partial L}{\partial S2} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \begin{bmatrix} 8 & 0 \end{bmatrix} = \begin{bmatrix} 8 & 0 \\ 16 & 0 \end{bmatrix} \quad *$$

$$S3 = Y2W2 \Rightarrow \frac{\partial L}{\partial Y2} = \frac{\partial L}{\partial S3} W2^T = \begin{bmatrix} 0 & 4 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} 8 & 8 \end{bmatrix}$$

$$\frac{\partial L}{\partial S2} = \frac{\partial L}{\partial Y2} \circ ReLU'(S2) = \begin{bmatrix} 8 & 8 \end{bmatrix} \circ \begin{bmatrix} 1 & 0 \end{bmatrix} = \begin{bmatrix} 8 & 0 \end{bmatrix}$$

Q 4:

Parameter

	Weights	bias
f _{c1}	$(28 \times 28) \times 512 = 401408$	512
f _{c2}	$512 \times 512 = 262144$	512
f _{c3}	$512 \times 10 = 5120$	10

Total parameter: $(401408 + 262144 + 5120) + (512 + 512 + 10) = 669706$ *

Operation

	includes input, weight, bias MAC (input \times output)	ReLU
f _{c1}	$(28 \times 28) \times 512 = 401408$	512
f _{c2}	$512 \times 512 = 262144$	512
f _{c3}	$512 \times 10 = 5120$	10

total operations 668672×2

$1034 = 1338378$

MAC has 2 operations (add & multiplication)

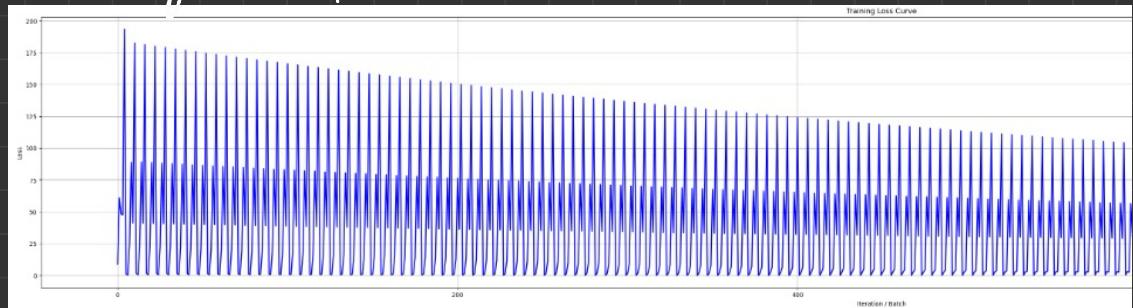
The default bias is set to True

ReLU is counted as 1 operation

There are 1338378 operations / image *

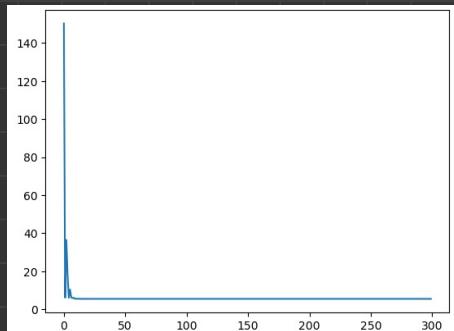
Q6:

SGD, single data point



Weights are updated using only one data at each iteration.
So the gradient goes up and down to update different points,
making the oscillation sharp and produce noisy gradient.
But the overall gradient is still smooth

Batch-based



Loss drops and stabilize fast, for it update all data
point a once. Allowing it to stabilize faster.
The gradient looks smooth