

Homework #1

Pranav Jangir
Deep Learning Spring' 23
NYU

12th February 2023

Question 1.2

(a):

1. Forward Propagation
2. Loss computation
3. Clear Gradients (Reset to zero)
4. Backward propagation (Computing loss gradients for every layer wrt the weights)
5. Update weights

(b):

$$\begin{aligned} \text{Linear}_1 \text{ input} &= x & \text{Linear}_1 \text{ output} &= W^{(1)}x + b^{(1)} \\ f \text{ input} &= W^{(1)}x + b^{(1)} & f \text{ output} &= 3(W^{(1)}x + b^{(1)})^+ \\ \text{Linear}_2 \text{ input} &= 3(W^{(1)}x + b^{(1)})^+ & \text{Linear}_2 \text{ output} &= 3W^{(2)}(W^{(1)}x + b^{(1)})^+ + b^{(2)} \\ g \text{ input} &= 3W^{(2)}(W^{(1)}x + b^{(1)})^+ + b^{(2)} & g \text{ output} &= 3W^{(2)}(W^{(1)}x + b^{(1)})^+ + b^{(2)} \end{aligned}$$

Loss input = $3W^{(2)}(W^{(1)}x + b^{(1)})^+ + b^{(2)}$ and y

Loss output = $\| 3W^{(2)}(W^{(1)}x + b^{(1)})^+ + b^{(2)} - y \|^2$

(c):

We know that $x \in \mathbb{R}^n$ and $\tilde{y} \in \mathbb{R}^K$. If we assume that the first linear layer has a neurons, then $W^{(1)} \in \mathbb{R}^{a \times n}$ and $b^{(1)} \in \mathbb{R}^a$

also $W^{(2)} \in \mathbb{R}^{K \times a}$ and $b^{(2)} \in \mathbb{R}^K$

We have the following gradients that need to be computed :

$$\begin{aligned} \frac{\partial C}{\partial \tilde{y}} &\in \mathbb{R}^{1 \times K} \\ \frac{\partial \tilde{y}}{\partial s_2} &\in \mathbb{R}^{K \times K} \\ \frac{\partial s_2}{\partial a_1} &\in \mathbb{R}^{K \times a} \\ \frac{\partial a_1}{\partial s_1} &\in \mathbb{R}^{a \times a} \end{aligned}$$

Also, partial derivatives of the linear outputs z_1 and z_3 wrt to their layer's weights and bias that need to be computed :

$$\begin{aligned}\frac{\partial s_2}{\partial W^{(2)}} &\in \mathbb{R}^{K \times a \times K} \\ \frac{\partial s_2}{\partial b^{(2)}} &\in \mathbb{R}^{K \times K} \\ \frac{\partial s_1}{\partial W^{(1)}} &\in \mathbb{R}^{a \times n \times a} \\ \frac{\partial s_1}{\partial b^{(1)}} &\in \mathbb{R}^{a \times a}\end{aligned}$$

Calculating the gradients via chain rule :

We are going to use the two terms repeatedly :

$$\begin{aligned}\frac{\partial C}{\partial s_2} &= \frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} \\ \frac{\partial C}{\partial s_1} &= \frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} \frac{\partial s_2}{\partial a_1} \frac{\partial a_1}{\partial s_1} \\ &= \frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} W^{(2)} \frac{\partial a_1}{\partial s_1}\end{aligned}$$

For biases :

$$\begin{aligned}\frac{\partial C}{\partial b^{(2)}} &= \frac{\partial C}{\partial s_2} \frac{\partial s_2}{\partial b^{(2)}} \\ &= \boxed{\frac{\partial C}{\partial s_2}} \\ \frac{\partial C}{\partial b^{(1)}} &= \frac{\partial C}{\partial s_1} \frac{\partial s_1}{\partial b^{(1)}} \\ &= \frac{\partial C}{\partial s_1} \\ &= \boxed{\frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} W^{(2)} \frac{\partial a_1}{\partial s_1}}\end{aligned}$$

For weights we have :

$$\begin{aligned}
\frac{\partial C}{\partial W^{(2)}} &= \frac{\partial C}{\partial s_2} \frac{\partial s_2}{\partial W^{(2)}} \\
&= \sum_i \frac{\partial C}{\partial (s_2)_i} \times \frac{\partial (s_2)_i}{\partial W^{(2)}} \\
&= a_1 \frac{\partial C}{\partial s_2} \\
&= \boxed{3(W^{(1)}x + b^{(1)})^+ \frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2}} \\
\frac{\partial C}{\partial W^{(1)}} &= \frac{\partial C}{\partial s_1} \frac{\partial s_1}{\partial W^{(1)}} \\
&= \sum_i \frac{\partial C}{\partial (s_1)_i} \times \frac{\partial (s_1)_i}{\partial W^{(1)}} \\
&= x \frac{\partial C}{\partial s_1} \\
&= \boxed{x \frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} W^{(2)} \frac{\partial a_1}{\partial s_1}}
\end{aligned}$$

(d):

$\frac{\partial a_1}{\partial s_1}$ as noted before is a $a \times a$ matrix with non diagonal entries equal to zero.

$$\boxed{\left(\frac{\partial a_1}{\partial s_1}\right)_{ii} = 3\mathbb{1}_{(s_1)_i > 0}}$$

for $i = 1 \cdots a$, where $\mathbb{1}$ is the indicator random variable.

Likewise, $\frac{\partial \tilde{y}}{\partial s_2}$ is a $K \times K$ matrix with non diagonal entries equal to zero.

$$\boxed{\left(\frac{\partial \tilde{y}}{\partial s_2}\right)_{ii} = 1}$$

$\frac{\partial C}{\partial \tilde{y}}$ is a row vector of dimensions $1 \times K$

$$\boxed{\left(\frac{\partial C}{\partial \tilde{y}}\right)_i = 2(\tilde{y}_i - y_i)}$$

Question 1.3

(a):

$$\begin{aligned}
 \text{Linear}_1 \text{ input} &= x & \text{Linear}_1 \text{ output} &= W^{(1)}x + b^{(1)} \\
 f \text{ input} &= W^{(1)}x + b^{(1)} & f \text{ output} &= \tanh(W^{(1)}x + b^{(1)}) \\
 \text{Linear}_2 \text{ input} &= \tanh(W^{(1)}x + b^{(1)}) \\
 \text{Linear}_2 \text{ output} &= W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)} \\
 g \text{ input} &= W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)} \\
 g \text{ output} &= \sigma(W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)})
 \end{aligned}$$

$$\begin{aligned}
 \text{Loss input} &= \sigma(W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)}) \text{ and } y \\
 \text{Loss output} &= \| \sigma(W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)}) - y \|^2
 \end{aligned}$$

Gradients can be computed just as they were computed for 1.2 using chain rule. The new gradients are :

For biases :

$$\begin{aligned}
 \frac{\partial C}{\partial b^{(2)}} &= \boxed{\frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2}} \\
 \frac{\partial C}{\partial b^{(1)}} &= \boxed{\frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} W^{(2)} \frac{\partial a_1}{\partial s_1}}
 \end{aligned}$$

For weights :

$$\begin{aligned}
 \frac{\partial C}{\partial W^{(2)}} &= \boxed{\tanh(W^{(1)}x + b^{(1)}) \frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2}} \\
 \frac{\partial C}{\partial W^{(1)}} &= \boxed{x \frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} W^{(2)} \frac{\partial a_1}{\partial s_1}}
 \end{aligned}$$

The values of the gradients are :

$\frac{\partial a_1}{\partial s_1}$ is a $a \times a$ matrix with non diagonal entries equal to zero.

$$\boxed{\left(\frac{\partial a_1}{\partial s_1}\right)_{ii} = (1 - \tanh(s_1)_i^2)}$$

for $i = 1 \cdots a$

Likewise, $\frac{\partial \tilde{y}}{\partial s_2}$ is a $K \times K$ matrix with non diagonal entries equal to zero.

$$\left(\frac{\partial \tilde{y}}{\partial s_2} \right)_{ii} = \sigma((s_2)_i)(1 - \sigma((s_2)_i))$$

$\frac{\partial C}{\partial \tilde{y}}$ is a row vector of dimensions $1 \times K$

$$\left(\frac{\partial C}{\partial \tilde{y}} \right)_i = 2(\tilde{y}_i - y_i)$$

(b):

Only the loss function has changed, therefore :

$$\begin{aligned} \text{Linear}_1 \text{ input} &= x & \text{Linear}_1 \text{ output} &= W^{(1)}x + b^{(1)} \\ f \text{ input} &= W^{(1)}x + b^{(1)} & f \text{ output} &= \tanh(W^{(1)}x + b^{(1)}) \\ \text{Linear}_2 \text{ input} &= \tanh(W^{(1)}x + b^{(1)}) \\ \text{Linear}_2 \text{ output} &= W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)} \\ g \text{ input} &= W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)} \\ g \text{ output} &= \sigma(W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)}) \end{aligned}$$

Loss input = $\sigma(W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)})$ and y

Loss output =

$$-\frac{1}{K} [y^T \log(\sigma(W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)})) + (1 - y)^T \log(1 - \sigma(W^{(2)} \tanh(W^{(1)}x + b^{(1)}) + b^{(2)}))]$$

The gradients are the same as the gradients for 1.3 part (a) :

For biases :

$$\begin{aligned} \frac{\partial C}{\partial b^{(2)}} &= \left[\frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} \right] \\ \frac{\partial C}{\partial b^{(1)}} &= \left[\frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} W^{(2)} \frac{\partial a_1}{\partial s_1} \right] \end{aligned}$$

For weights :

$$\begin{aligned} \frac{\partial C}{\partial W^{(2)}} &= \left[\tanh(W^{(1)}x + b^{(1)}) \frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} \right] \\ \frac{\partial C}{\partial W^{(1)}} &= \left[x \frac{\partial C}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial s_2} W^{(2)} \frac{\partial a_1}{\partial s_1} \right] \end{aligned}$$

The values of the gradients are :

$\frac{\partial a_1}{\partial s_1}$ is a $a \times a$ matrix with non diagonal entries equal to zero.

$$\left(\frac{\partial a_1}{\partial s_1} \right)_{ii} = (1 - \tanh(s_1)_i^2)$$

for $i = 1 \cdots a$

Likewise, $\frac{\partial \tilde{y}}{\partial s_2}$ is a $K \times K$ matrix with non diagonal entries equal to zero.

$$\left(\frac{\partial \tilde{y}}{\partial s_2} \right)_{ii} = \sigma((s_2)_i)(1 - \sigma((s_2)_i))$$

$\frac{\partial C}{\partial \tilde{y}}$ is a row vector of dimensions $1 \times K$

$$\left(\frac{\partial C}{\partial \tilde{y}} \right)_i = \frac{1}{K} \frac{y_i - \tilde{y}_i}{\tilde{y}_i - \tilde{y}_i^2}$$

(c):

$f = ReLU$ is a much better choice than $f = \tanh$ and that is because \tanh suffers from the problem of vanishing gradients. The range of \tanh is $[-1, 1]$ and therefore, its derivative $1 - \tanh^2(x)$ has values between $[0, 1]$. During the training phase, the gradients may reach close to zero and thus slow down the weight change during backward pass, therefore slowing down the learning. $ReLU$ does not have this problem as its derivative is either 0 or 1 and as a result learning is faster.

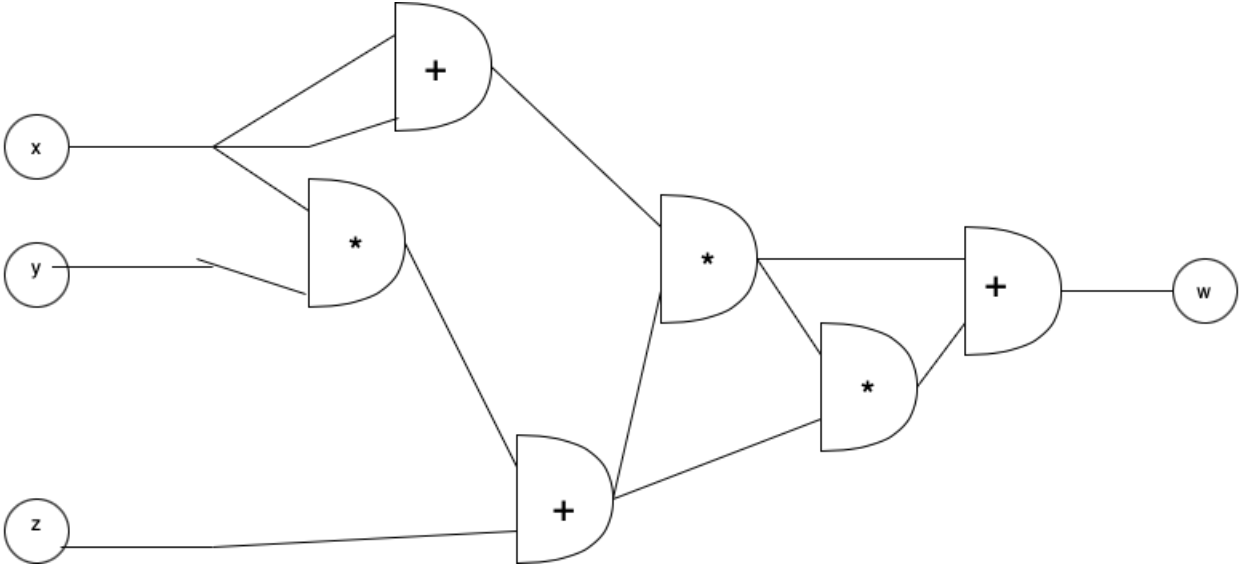
Question 1.4

(a):

Softmax is generally confused with Softargmax. What is colloquially called Softmax should be called Softargmax as it is a clearer name.

Softargmax gives a soft version of argmax. "Soft" here means differentiable everywhere. "Hard" functions like argmax are not differentiable everywhere. This makes it possible to give smooth gradients for the softargmax, which would be impossible with argmax.

(b):



(c):

(d):

Part (a) :

$$\frac{\partial f}{\partial \mathbf{x}} = J_f = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{x}_1} & \frac{\partial f_1}{\partial \mathbf{x}_2} & \dots & \frac{\partial f_1}{\partial \mathbf{x}_a} \\ \frac{\partial f_2}{\partial \mathbf{x}_1} & \frac{\partial f_2}{\partial \mathbf{x}_2} & \dots & \frac{\partial f_2}{\partial \mathbf{x}_a} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_b}{\partial \mathbf{x}_1} & \frac{\partial f_b}{\partial \mathbf{x}_2} & \dots & \frac{\partial f_b}{\partial \mathbf{x}_a} \end{bmatrix} = \mathbf{W}_1$$

$$\frac{\partial g}{\partial \mathbf{x}} = J_g = \begin{bmatrix} \frac{\partial g_1}{\partial \mathbf{x}_1} & \frac{\partial g_1}{\partial \mathbf{x}_2} & \dots & \frac{\partial g_1}{\partial \mathbf{x}_a} \\ \frac{\partial g_2}{\partial \mathbf{x}_1} & \frac{\partial g_2}{\partial \mathbf{x}_2} & \dots & \frac{\partial g_2}{\partial \mathbf{x}_a} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial g_b}{\partial \mathbf{x}_1} & \frac{\partial g_b}{\partial \mathbf{x}_2} & \dots & \frac{\partial g_b}{\partial \mathbf{x}_a} \end{bmatrix} = \mathbf{W}_2$$

Part (b):

$$J_h = \frac{\partial h}{\partial \mathbf{x}} = \frac{\partial f}{\partial \mathbf{x}} + \frac{\partial g}{\partial \mathbf{x}} = J_h = J_f + J_g = \mathbf{W}_1 + \mathbf{W}_2$$

Part (c):

If $\mathbf{W}_1 = \mathbf{W}_2$

$$J_h = \frac{\partial h}{\partial \mathbf{x}} = \frac{\partial f}{\partial \mathbf{x}} + \frac{\partial g}{\partial \mathbf{x}} = J_h = J_f + J_g = \mathbf{W}_1 + \mathbf{W}_2 = 2\mathbf{W}_1$$

Part (e):

(a)

The Jacobian of $f(\mathbf{x})$ is:

$$J_f = \frac{\partial f}{\partial \mathbf{x}} = J_f = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{x}_1} & \frac{\partial f_1}{\partial \mathbf{x}_2} & \dots & \frac{\partial f_1}{\partial \mathbf{x}_a} \\ \frac{\partial f_2}{\partial \mathbf{x}_1} & \ddots & \dots & \frac{\partial f_2}{\partial \mathbf{x}_a} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_b}{\partial \mathbf{x}_1} & \frac{\partial f_b}{\partial \mathbf{x}_2} & \dots & \frac{\partial f_b}{\partial \mathbf{x}_a} \end{bmatrix} = \mathbf{W}_1$$

The Jacobian of $g(\mathbf{x})$ is:

$$J_g = \frac{\partial g}{\partial \mathbf{x}} = J_g = \begin{bmatrix} \frac{\partial g_1}{\partial \mathbf{x}_1} & \frac{\partial g_1}{\partial \mathbf{x}_2} & \dots & \frac{\partial g_1}{\partial \mathbf{x}_b} \\ \frac{\partial g_2}{\partial \mathbf{x}_1} & \ddots & \dots & \frac{\partial g_2}{\partial \mathbf{x}_b} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial g_c}{\partial \mathbf{x}_1} & \frac{\partial g_c}{\partial \mathbf{x}_2} & \dots & \frac{\partial g_c}{\partial \mathbf{x}_b} \end{bmatrix} = \mathbf{W}_2$$

(b)

$$h(\mathbf{x}) = g(f(\mathbf{x})) = g(\mathbf{W}_1 \mathbf{x}) = \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$$

Hence, $J_h =$

$$\frac{\partial h}{\partial \mathbf{x}} = \mathbf{W}_2 \mathbf{W}_1$$

(c)

If $\mathbf{W}_1 = \mathbf{W}_2$,

$$J_h = \frac{\partial h}{\partial \mathbf{x}} = \mathbf{W}_2 \mathbf{W}_1 = \mathbf{W}_1^T \mathbf{W}_1$$