

1.2.

a)

1. Generate a prediction

$\tilde{y} = \text{model}(x)$

2. Compute the loss

$L(w, x, y) = f(x, y) = C(\tilde{y}, y) = ||y - \tilde{y}||^2$

3. Set the gradient parameters to zero

`optimizer.zero_grad()`

4. Compute and accumulate the gradient parameters

`L.backward()`

5. Step in towards the negative of the gradient parameters

`optimizer.step()`

b)

$x \rightarrow \text{Linear}_1 \rightarrow s_1 \rightarrow 3\text{ReLU}(\cdot) \rightarrow a_1 \rightarrow \text{Linear}_2 \rightarrow s_2 \rightarrow g \rightarrow \tilde{y}$

where,

$s_1 = W_1^T x + b_1,$

$a_1 = (3\text{ReLU})(W_1^T x + b_1)$

$s_2 = (W_2^T) \{3\text{ReLU}(W_1^T x + b_1)\} + b_2,$  and

$\tilde{y} = (W_2^T) \{3\text{ReLU}(W_1^T x + b_1)\} + b_2$

c)

$\partial C / \partial W_1 = (x) (\partial C / \partial \tilde{y}) (\partial \tilde{y} / \partial s_2) (W_2^T) (\partial a_1 / \partial s_1)$

$\partial C / \partial b_1 = (\partial C / \partial \tilde{y}) (\partial \tilde{y} / \partial s_2) (W_2^T) (\partial a_1 / \partial s_1)$

$\partial C / \partial W_2 = (3\text{ReLU}(W_1^T x + b_1)) (\partial C / \partial \tilde{y}) (\partial \tilde{y} / \partial s_2)$

$\partial C / \partial b_2 = (\partial C / \partial \tilde{y}) (\partial \tilde{y} / \partial s_2)$

Moreover,

$\partial s_2 / \partial W_2 = a_1$  (Used above)

$$\partial \mathbf{s}_1 / \partial \mathbf{W}_1 = \mathbf{x}$$

$$\partial \mathbf{s}_1 / \partial \mathbf{b}_1 = \mathbf{1}$$

$$\partial \mathbf{s}_2 / \partial \mathbf{b}_2 = \mathbf{1}$$

$$\text{And } (\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2) = \mathbf{1}$$

d)

L1 is the output size of the first linear layer

$$\partial \mathbf{a}_1 / \partial \mathbf{s}_1 = \mathbf{AS}_{ij} \text{ where } \mathbf{AS}_{ij} = 3 \text{ when } i = j \text{ and } \mathbf{S}_{ij} > 0$$

0 otherwise

$$\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2 = \mathbf{I}_{K \times K} \text{ where } \mathbf{I} \text{ is the identity matrix}$$

$$\partial C / \partial \tilde{\mathbf{y}} = 2[\tilde{\mathbf{y}} - \mathbf{y}]_{1 \times K}$$

1.3.

a.

$$\mathbf{x} \rightarrow \text{Linear}_1 \rightarrow \mathbf{s}_1 \rightarrow \tanh(\cdot) \rightarrow \mathbf{a}_1 \rightarrow \text{Linear}_2 \rightarrow \mathbf{s}_2 \rightarrow 1/(1 + \exp(-\mathbf{x})) \rightarrow \tilde{\mathbf{y}}$$

where,

$$\mathbf{s}_1 = \mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1,$$

$$\mathbf{a}_1 = \tanh(\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1) \text{ (Changed)}$$

$$\mathbf{s}_2 = \mathbf{W}_2^T \{\tanh(\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1)\} + \mathbf{b}_2, \text{ and (Changed)}$$

$$\tilde{\mathbf{y}} = 1/(1 + \exp(-\mathbf{W}_2^T \{\tanh(\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1)\} - \mathbf{b}_2)) \text{ (Changed)}$$

ii.

The value of  $(\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2)$  has changed to  $\text{sigmoid}(\mathbf{s}_2) * (1 - \text{sigmoid}(\mathbf{s}_2))$

And  $\mathbf{a}_1$  is  $\tanh(\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1)$

$$\partial C / \partial \mathbf{W}_1 = (x) (\partial C / \partial \tilde{\mathbf{y}}) (\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2) (\mathbf{W}_2^T) (\partial \mathbf{a}_1 / \partial \mathbf{s}_1)$$

$$\partial C / \partial \mathbf{b}_1 = (\partial C / \partial \tilde{\mathbf{y}}) (\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2) (\mathbf{W}_2^T) (\partial \mathbf{a}_1 / \partial \mathbf{s}_1)$$

$$\partial C / \partial \mathbf{W}_2 = (\tanh(\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1)) (\partial C / \partial \tilde{\mathbf{y}}) (\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2)$$

$$\partial C / \partial \mathbf{b}_2 = (\partial C / \partial \tilde{\mathbf{y}}) (\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2)$$

Moreover,

$$\partial \mathbf{s}_2 / \partial \mathbf{W}_2 = \mathbf{a}_1 \text{ (Used above)}$$

$$\partial \mathbf{s}_1 / \partial \mathbf{W}_1 = \mathbf{x}$$

$$\partial \mathbf{s}_1 / \partial \mathbf{b}_1 = 1$$

$$\partial \mathbf{s}_2 / \partial \mathbf{b}_2 = 1$$

The value of  $(\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2)$  has changed to  $\text{sigmoid}(\mathbf{s}_2) * (1 - \text{sigmoid}(\mathbf{s}_2))$

iii.

L1 is the output size of the first linear layer

$$\partial \mathbf{a}_1 / \partial \mathbf{s}_1 = 1 - \tanh^2(s_{ij}) \text{ where } i \text{ and } j \text{ go to } L1 \times L1$$

$$\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2 = \tilde{\mathbf{y}} * (1 - \tilde{\mathbf{y}}) \text{ of dimension } K \times K$$

$$\partial C / \partial \tilde{\mathbf{y}} = 2[\tilde{\mathbf{y}} - \mathbf{y}]_{1 \times K}$$

b)

The value of  $\partial C / \partial \tilde{\mathbf{y}}$  is  $-1 * \tilde{\mathbf{y}} / \tilde{\mathbf{y}} + (1 - \tilde{\mathbf{y}}) / (1 - \tilde{\mathbf{y}})$  where

$$\tilde{\mathbf{y}} = 1 / (1 + \exp(-\mathbf{W}_2 * \tanh(\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1) - \mathbf{b}_2))$$

Hence, the expressions for  $\mathbf{s}_1$ ,  $\mathbf{a}_1$ ,  $\mathbf{s}_2$  and  $\tilde{\mathbf{y}}$  do not change over the previous part but the substitution of  $\partial C / \partial \tilde{\mathbf{y}}$  would change in the values

$$\partial C / \partial \mathbf{W}_1 = (x) (\partial C / \partial \tilde{\mathbf{y}}) (\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2) (\mathbf{W}_2^T) (\partial \mathbf{a}_1 / \partial \mathbf{s}_1)$$

$$\partial C / \partial \mathbf{b}_1 = (\partial C / \partial \tilde{\mathbf{y}}) (\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2) (\mathbf{W}_2^T) (\partial \mathbf{a}_1 / \partial \mathbf{s}_1)$$

$$\partial C / \partial \mathbf{W}_2 = (3\text{ReLU}(\mathbf{W}_1^T \mathbf{x} + \mathbf{b}_1)) (\partial C / \partial \tilde{\mathbf{y}}) (\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2)$$

$$\partial C / \partial \mathbf{b}_2 = (\partial C / \partial \tilde{\mathbf{y}}) (\partial \tilde{\mathbf{y}} / \partial \mathbf{s}_2)$$

c)

$$\partial C / \partial \tilde{\mathbf{y}} = [(1-y)/(1-\tilde{y}) - y/\tilde{y}]_{1 \times K}$$

Rest everything is same as the part b

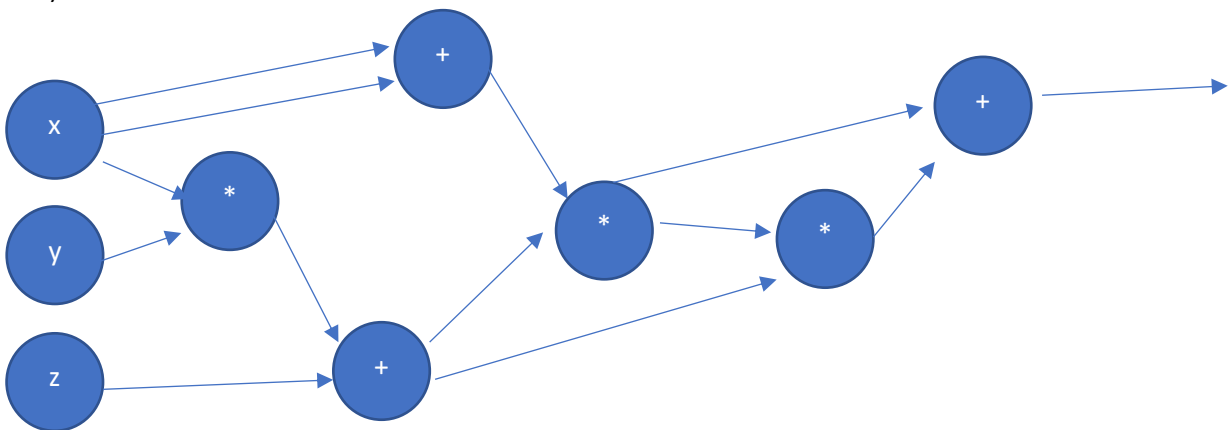
1.4.

Because softmax implemented in most machine learning libraries actually gives the one-hot encoded vector with 1 at the index with the maximum argument, while softmax mathematically gives

$\text{Softmax}_b(\mathbf{e}) = (1/b) * \log \sum \exp(b * \mathbf{e})$  ( $= \max(\mathbf{e})$  as  $b$  tends to infinity) and

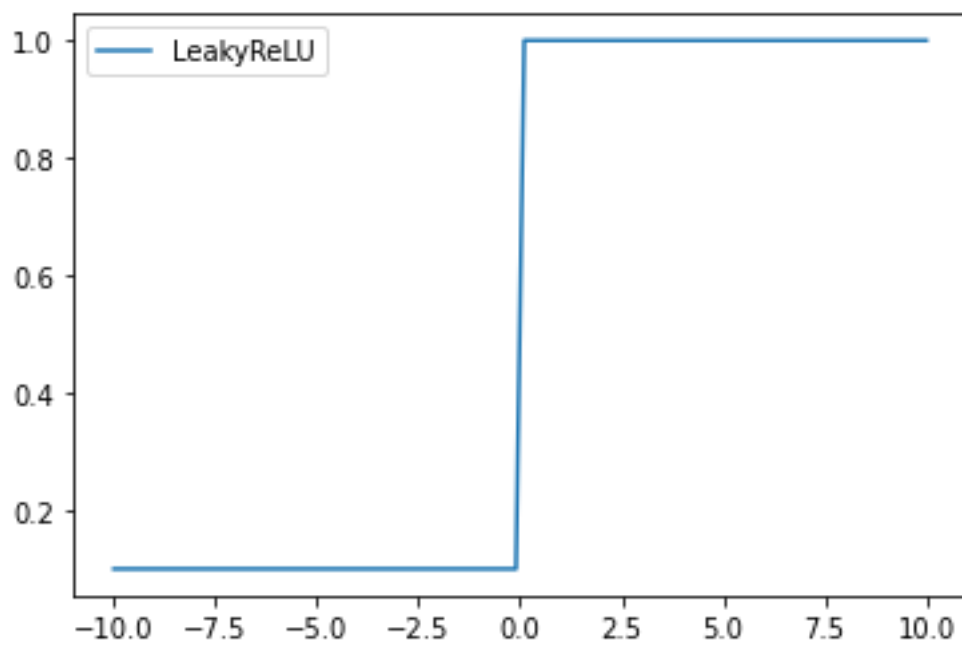
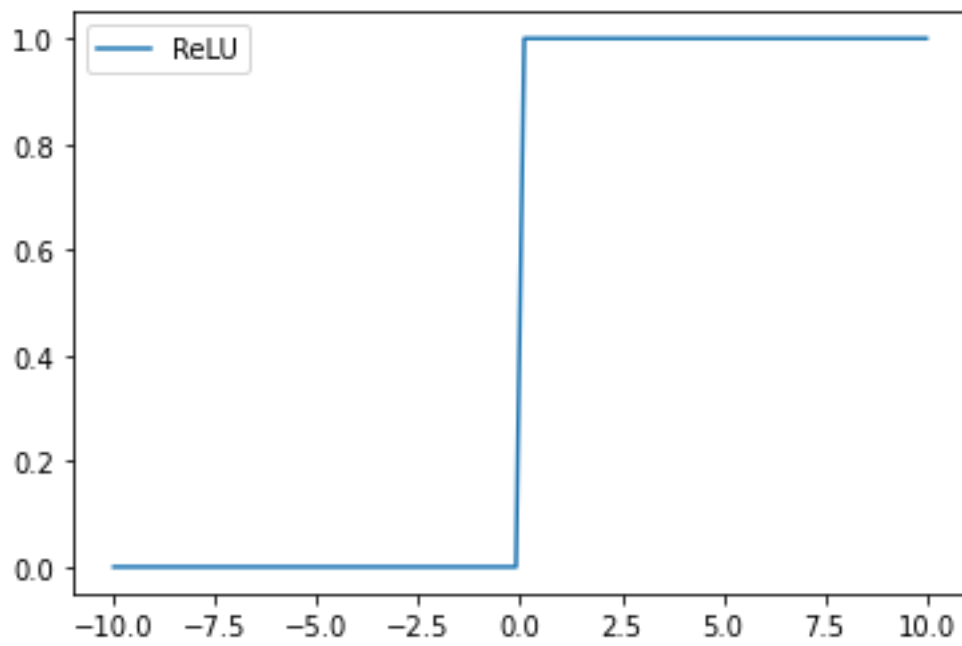
$\text{Softargmax}_b(\mathbf{e}) = \exp(b \mathbf{e}) / \sum \exp(b * \mathbf{e})$  actually mathematically corresponds to the maximum argument (as  $b$  tends to infinity)

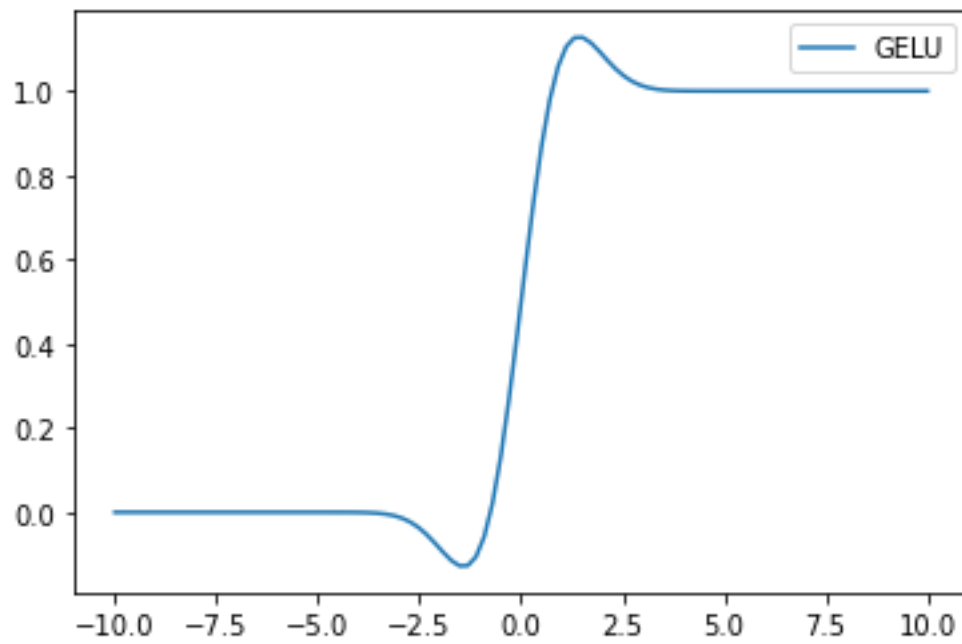
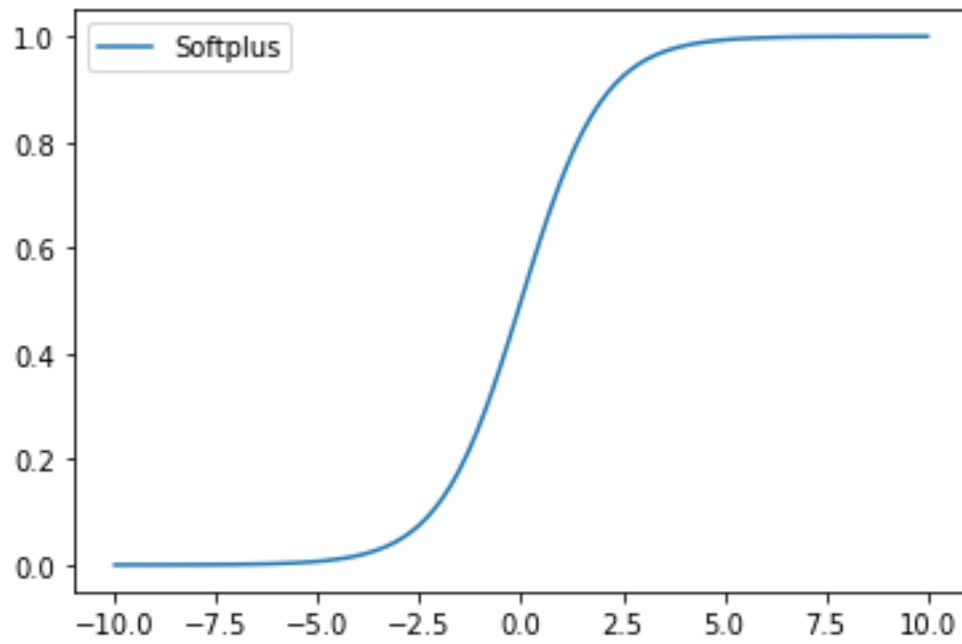
b)



1.4.

c)





d)

Jacobian of  $f = \partial f / \partial x = W1$

Jacobian of  $g = \partial g / \partial x = W2$

ii.

Jacobian of  $h = \partial h / \partial x = W_1 + W_2$

If  $W_1 = W_2$

$h = \partial h / \partial x = 2 (W_2)$

e)

Jacobian of  $f = \partial f / \partial x = W_1$

Jacobian of  $g = \partial g / \partial x = W_2$

Jacobian of  $h = \partial h / \partial x = (W_2) (W_1)$

If  $W_1 = W_2$

$h = \partial h / \partial x = (W_1)^2 = (W_2)^2$