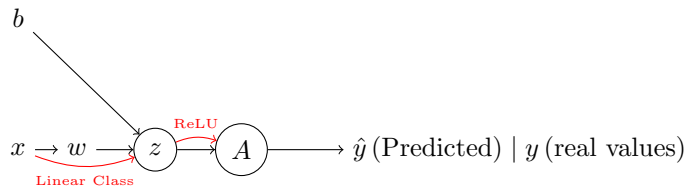


Summary of

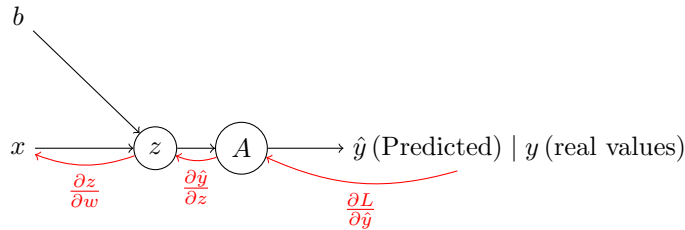
1 Neural Network Theory

Neural Network Structure:



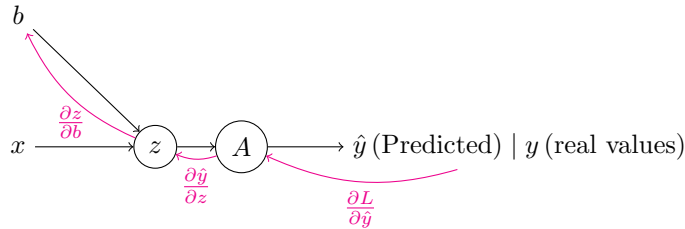
$$L(\hat{y}, y)$$

backward pass for w:



$$L(\hat{y}, y)$$

backward pass for b:



$$L(\hat{y}, y)$$

backward pass for w mathematically:

$$L = L(\underbrace{\hat{y}(z(Xw + b))}_{\underbrace{\frac{\partial z}{\partial w}}_{\underbrace{\frac{\partial \hat{y}}{\partial z}}_{\frac{\partial L}{\partial \hat{y}}}}}) \quad (1)$$

$$\frac{\partial L}{\partial w} = \underbrace{\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z}}_{\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z}} \cdot \frac{\partial z}{\partial w} \quad (2)$$

as a Summary:

$$L = \underbrace{L(\hat{y})}_{\frac{\partial L}{\partial \hat{y}}} = \underbrace{L(\hat{y}(z))}_{\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z}} = \underbrace{L(\hat{y}(z(XW + b)))}_{\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial W}}$$

Or more explicitly for the chain rule:

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z} \times \frac{\partial z}{\partial W}$$

- $z = XW + b$
- $\hat{y} = f(z)$  (activation function, e.g. softmax)
- $L = \text{loss}(\hat{y}, y)$

The derivatives  $\frac{\partial L}{\partial \hat{y}}$ ,  $\frac{\partial \hat{y}}{\partial z}$ , and  $\frac{\partial z}{\partial W}$  represent the gradient chain from output all the way down to weights.

## 2 Coding the Neural Network

Neural Network Structure:

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w}$$

where:

- **Loss derivative:**  $\frac{\partial L}{\partial \hat{y}}$
- **Activation derivative:**  $\frac{\partial \hat{y}}{\partial z}$
- **Linear model derivative:**  $\frac{\partial z}{\partial w}$

1. **Loss derivative**  $\left(\frac{\partial L}{\partial \hat{y}}\right)$ :

This is computed in:

```
1 def backward(self) -> np.ndarray: # CrossEntropy
2     grad = -self.target / self.prediction / self.target.shape[0]
3     return grad
```

- `self.prediction` =  $\hat{y}$  (output of softmax)
- This gives the gradient from the loss with respect to the softmax output  $\hat{y}$ .

**Activation derivative**  $\left(\frac{\partial \hat{y}}{\partial z}\right)$ :

If you have a Softmax layer:

```
1 def backward(self, up_grad: np.ndarray) -> np.ndarray: # Softmax
2     ...
3     down_grad[i] = np.dot(jacobian, up_grad[i])
```

- `up_grad` is  $\frac{\partial L}{\partial \hat{y}}$  from the loss.
- The softmax Jacobian gives  $\frac{\partial \hat{y}}{\partial z}$ .
- Output `down_grad` is  $\frac{\partial L}{\partial z}$ .
- This matches step 2 of the chain rule.

**Linear model derivative** ( $\frac{\partial z}{\partial w}$ ):

In your Linear layer:

```

1 def backward(self, up_grad: np.ndarray) -> np.ndarray: # Linear
2     self.dw = np.dot(self.inp.T, up_grad) #  $\partial L / \partial w$ 
3     self.db = np.sum(up_grad, axis=0, keepdims=True) #  $\partial L / \partial b$ 
4     down_grad = np.dot(up_grad, self.w.T) #  $\partial L / \partial \text{input}$ 
5     return down_grad

```

- up\_grad is  $\frac{\partial L}{\partial z}$  from the activation.
- Multiplying with inp.T applies  $\frac{\partial z}{\partial w}$  to get  $\frac{\partial L}{\partial w}$ .
- This is step 3 of the professor's chain rule.

**Full Chain in Code:**

1. Loss backward  $\rightarrow$  CrossEntropy.backward()

$$\frac{\partial L}{\partial \hat{y}}$$

2. Activation backward  $\rightarrow$  Softmax.backward()

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z}$$

3. Linear backward  $\rightarrow$  Linear.backward()

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w}$$

**Linear:**

Abstract Python Code for Layer:

```

1 class Layer:
2     def __init__(self):
3         self.inp = None
4         self.out = None
5
6     def __call__(self, inp: np.ndarray) -> np.ndarray:
7         return self.forward(inp)
8
9     def forward(self, inp: np.ndarray) -> np.ndarray:
10        raise NotImplementedError
11
12    def backward(self, up_grad: np.ndarray) -> np.ndarray:
13        raise NotImplementedError
14
15    def step(self, lr: float) -> None:
16        pass

```

- Input features:

$$a \quad (3)$$

- Features weights:

$$a \quad (4)$$

- Bias term:

$$a \quad (5)$$

- Activation function :

$$a \quad (6)$$

- Output of the neuron:

$$y \quad (7)$$

**Loss Functions:**

types of loss functions:

- Mean Squared Error (MSE) for regression tasks:

$$L(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (8)$$

- Mean Absolute Error (MAE):

$$L(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i| \quad (9)$$

- Cross-Entropy Loss for classification tasks:

$$L(y, \hat{y}) = - \sum_{i=1}^C y_i \log(\hat{y}_i) \quad (10)$$

**explaining Forward in Linear and ReLU:**

in Linear and forward class we have:

```

1 class Linear(Layer):
2     ...
3     def forward(self, inp: np.ndarray) -> np.ndarray:
4         """Perform the linear transformation: output = inp * W + b"""
5         self.inp = inp
6         self.out = np.dot(inp, self.w) + self.b
7         return self.out
8     ...

```

and in ReLU we have:

```

1 class ReLU(Layer):
2     def forward(self, inp: np.ndarray) -> np.ndarray:
3         """ReLU Activation function: f(x) = max(0, x)"""
4         self.inp = inp
5         self.out = np.maximum(0, inp)
6         return self.out
7
8     def backward(self, up_grad: np.ndarray) -> np.ndarray:
9         """Backward pass for ReLU: derivative is 1 where input > 0, else 0."""
10        down_grad = up_grad * (self.inp > 0) # Efficient boolean indexing
11        return down_grad

```

then we make a list of layers like this:

```

1 layers = [
2     Linear(input_size, hidden_size),
3     ReLU(),
4     Linear(hidden_size, output_size)
5 ]
6 #passing layers to MLP Class:
7 model = MLP(layers, CrossEntropy(), learning_rate=0.001)

```

and this list of layers will be called in the MLP class as following:

```

1 class MLP:
2     def __init__(self, layers: list[Layer], loss_func: Loss, learning_rate: float) -> None:
3         self.layers = layers
4         self.loss_func = loss_func
5         self.learning_rate = learning_rate
6
7     def __call__(self, input: np.ndarray) -> np.ndarray:
8         return self.forward(input)
9
10    def forward(self, input: np.ndarray) -> np.ndarray:
11        """ pass input through each layer sequentially """
12        for layer in self.layers:
13            input = layer.forward(input)
14        return input

```

so the list of layers will be called one by one and the output of each layer will be feed to the next layer as input. the next simple example will show how it works:

## simple example of forward pass:

```

1 import numpy as np
2 # Input (batch of 2 samples, each with 3 features)
3 X = np.array([[1, 2, 3],
4               [4, 5, 6]])
5
6 # Linear layer
7 class Linear:
8     def __init__(self, in_features, out_features):
9         self.W = np.ones((in_features, out_features)) # just ones for simplicity
10        self.b = np.zeros((1, out_features))
11
12    def forward(self, x):
13        z = x @ self.W + self.b
14        print("Linear output z =", z)
15        return z
16
17 # ReLU layer
18 class ReLU:
19    def forward(self, x):
20        a = np.maximum(0, x)
21        print("ReLU output a =", a)
22        return a
23
24 # --- Build network ---
25 layers = [Linear(3, 2), ReLU()]
26
27 # --- Forward pass ---
28 out = X
29 for layer in layers:
30     print(out)
31     out = layer.forward(out)
32     print("After layer forward, out =", out)
33
34 print("Final output =", out)

```

First  $X$  is set to  $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$ ,  $W$  is set to ones, and  $b$  is set to zeros.

1. In `out = layer.forward(out)` the first layer is Linear, and `out` is equal to  $X$ . The matrix  $X$  will be fed into the Linear forward function.
2. In the Linear forward function: the output  $z$  will be calculated as:

$$z[0] = [1 + 2 + 3, 1 + 2 + 3] = [6, 6], \quad z[1] = [4 + 5 + 6, 4 + 5 + 6] = [15, 15]$$

3. In this step, `out` will be equal to the output of `layer.forward(out)`, which is  $z$ .
4. Now in the next step of the loop,  $z$  will be fed to the ReLU forward function, which is

$$a = \max(0, x)$$

and it will be equal to

$$\begin{bmatrix} 6 & 6 \\ 15 & 15 \end{bmatrix}$$

5. Then the result will be printed as:

Final output = `[[6. 6.] [15. 15.]]`