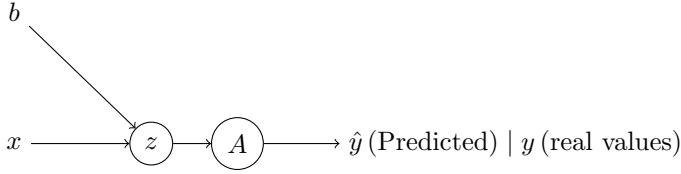


Summary of

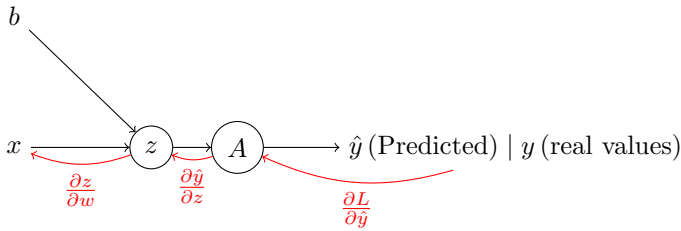
1 Neural Network Theory

Neural Network Structure:



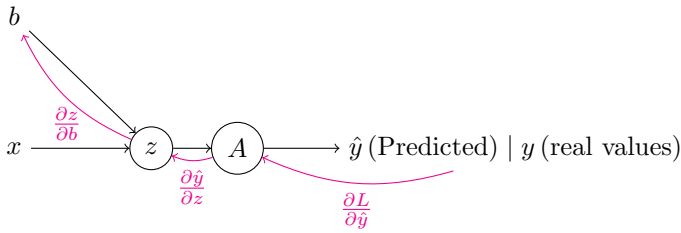
$$L(\hat{y}, y)$$

backward pass for w:



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

backward pass for b:



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

backward pass for w mathematically:

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

$$\underbrace{\frac{\partial \hat{y}}{\partial z}}_{\frac{\partial L}{\partial \hat{y}}}$$

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

as a Summary:

$$L = L(\hat{y}) = L(\hat{y}(z)) = L(\hat{y}(z(XW + b)))$$

$\frac{\partial L}{\partial \hat{y}} \quad \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \quad \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial W}$

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** ($\frac{\partial L}{\partial \hat{y}}$): This is computed in:

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

- `self.prediction = \hat{y}` (output of softmax)
- This gives the gradient from the loss with respect to the softmax output \hat{y} .
- This matches step 1 of the professor's chain rule.

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 $\left(\frac{\partial z}{\partial w}\right)$: 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) #  $\frac{\partial L}{\partial w}$ 
3     self.db = np.sum(up_grad, axis=0, keepdims=True) #  $\frac{\partial L}{\partial b}$ 
4     down_grad = np.dot(up_grad, self.w.T) #  $\frac{\partial L}{\partial \text{input}}$ 
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:

$$y \quad (7)$$

```
1 def greet(name):
2     print(f"Hello, {name}!")
3
4 greet("World")
```

Squared Error(SE): Most Common error function in linear regression is:

$$SE : (y^{(i)} - h(x^{(i)}, w))^2 \quad (8)$$

Sum of Squared Errors (SSE): Cost function should measure all predictions. Thus a choice could be Sum of Squared Error(SE)

$$SSE : \sum_{i=1}^N (y^{(i)} - h(x^{(i)}, w)) \quad (9)$$

Solve it analytically for one dimension: Predicted:

$$\hat{y} = w_0 + w_1 x \quad (10)$$

SSE or Cost Function:

$$J(w_0, w_1) := \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2 = \sum_{i=1}^N (y^{(i)} - w_0 + w_1 x^{(i)})^2 \quad (11)$$

Assumptions:

$$\frac{\partial J}{\partial w_0} = 0, \frac{\partial J}{\partial w_1} = 0 \quad (12)$$

$\frac{\partial J}{\partial w_0} = 0$: thus:

$$\frac{\partial}{\partial w_0} \left(\sum_{i=1}^N (y^{(i)} - (w_0 + w_1 x^{(i)}))^2 \right) = 0^1 \quad (13)$$

$$-2 \sum_{i=1}^N (y^{(i)} - w_0 - w_1 x^{(i)}) = 0 \quad (14)$$

For this equation to equal zero, the following condition must be met:

- $\sum_{i=1}^N y^{(i)} = 0 := Y$
- $\sum_{i=1}^N -w_0 = 0 := nw_0$
- $\sum_{i=1}^N -w_1 x^{(i)} = 0 := X$

Thus:

$$0 = Y - nw_0 - w_1 x^{(i)} \longrightarrow w_0 = \frac{(Y - w_1 x^{(i)})}{n} \quad (15)$$

¹ $f \circ g(x)' = f(g(x))' g(x)'$