

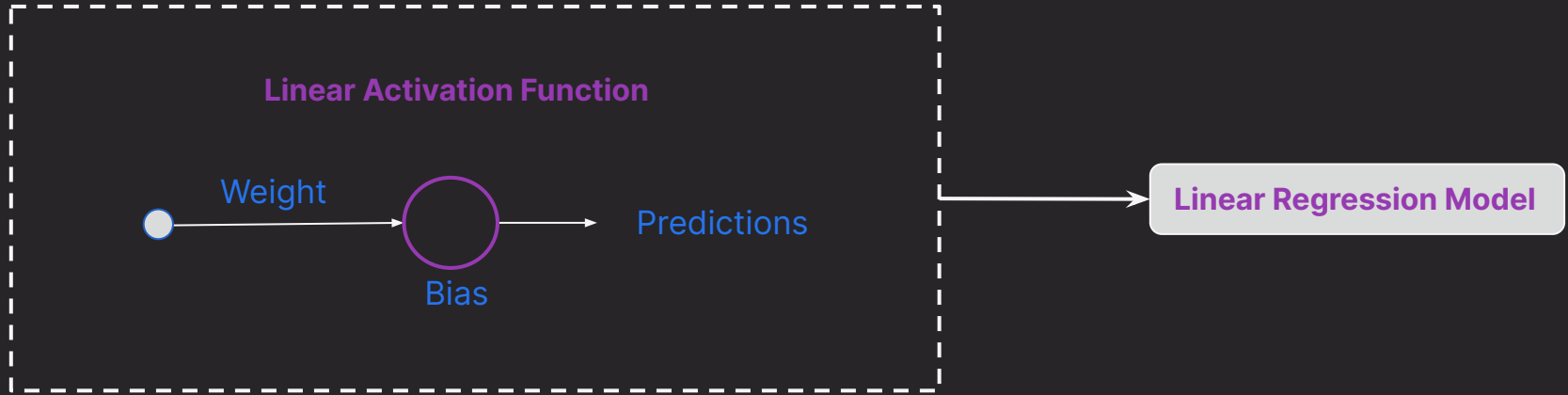


# Workings of Neural Networks

## Video 3: Optimization Techniques - Gradient Descent

In Air

# Minimizing Loss



# Minimizing Loss - LR

$$y = wx + b$$

Slope

Intercept

1

Closed form mathematical solution

2

Iterating  $w$  and  $b$  values for minimum loss

1

## Closed form mathematical solution



Direct method to calculate best-fit line



Approach complexity increases in multi-variable scenarios

## 2

Iterating  $w$  and  $b$  values for minimum loss

Infinite combinations of  $w$  and  $b$  make exhaustive search impractical.

**Example:** Weight of 3.3 with bias of 10.1, or weight of 3.31 with bias of 10.11, etc

# Gradient Descent

This technique is used to find the local minimum or optimize the loss function



A smarter approach to testing different  $w$  and  $b$  values

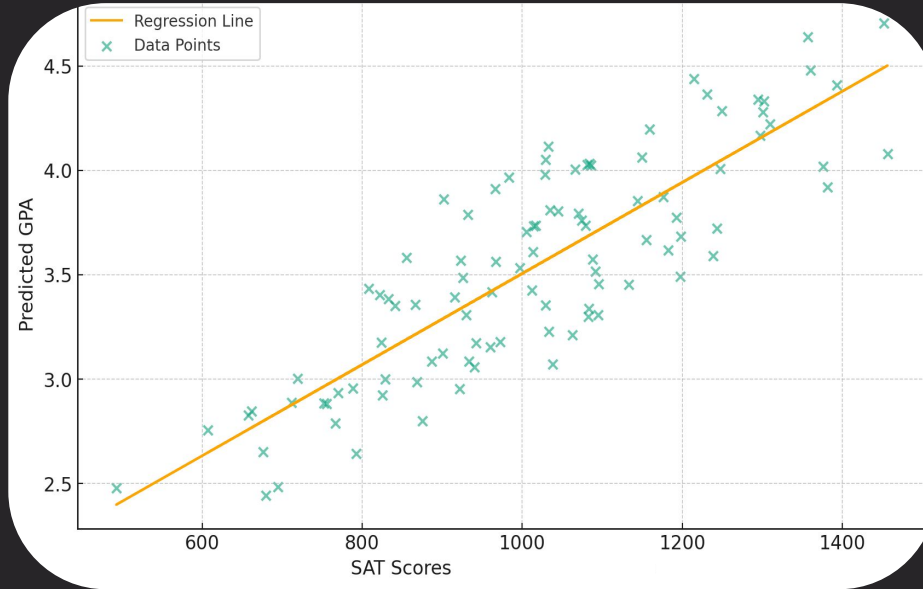


Backbone of neural networks



**Predicting first year GPA of students.**





$$y = w_i * x_i + b_i$$

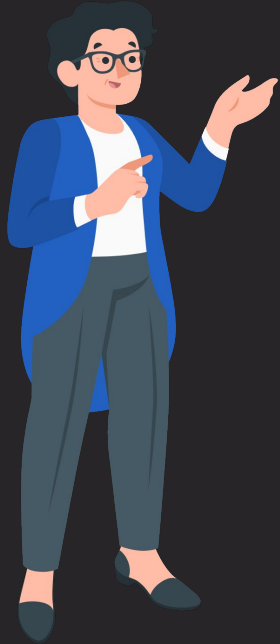
$w_i$  = weight of the slope

$x_i$  = input, sat\_score

$b_i$  = bias or intercept

Find the best fit line to minimize the loss

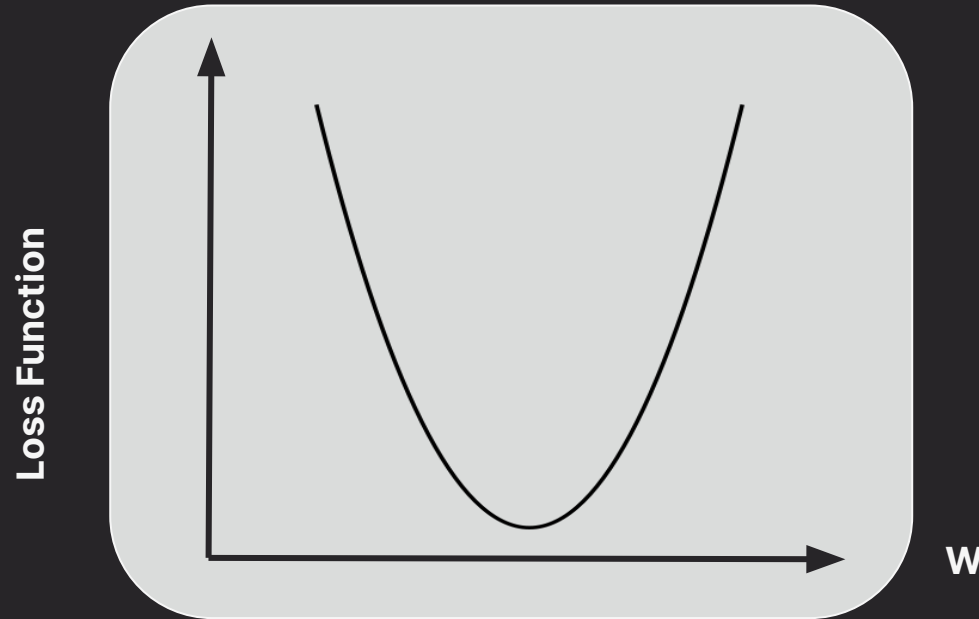
$$\text{MSE} = (1/n) \sum |\hat{y} - y|^2$$



$$\text{MSE} = (1/n) \sum |\hat{y} - y|^2$$

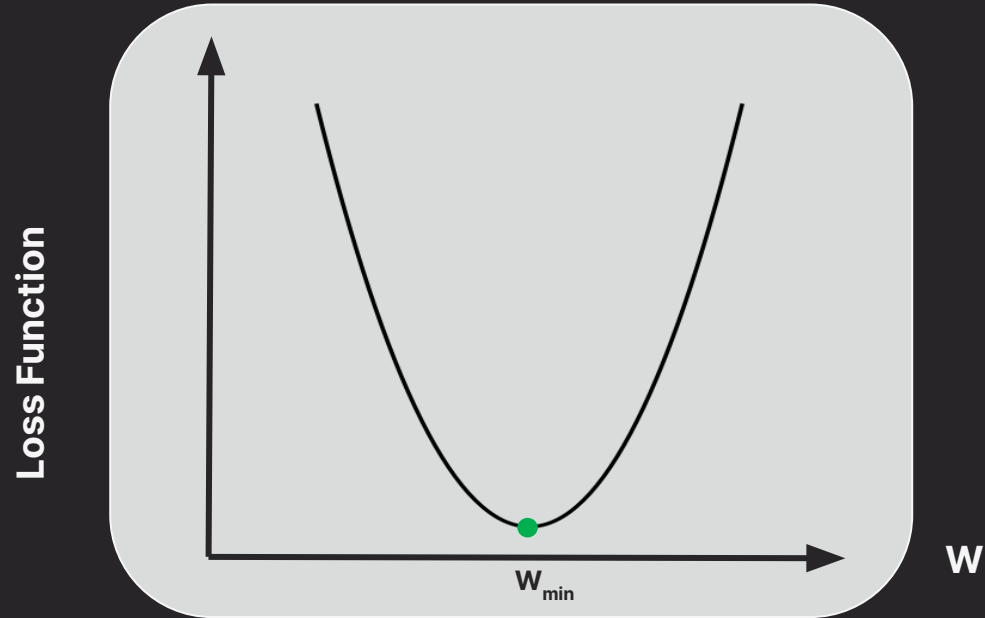
$$= (1/n) \sum |(wx + b) - y|^2$$

**Objective:** Find a value for weight where the loss is minimum

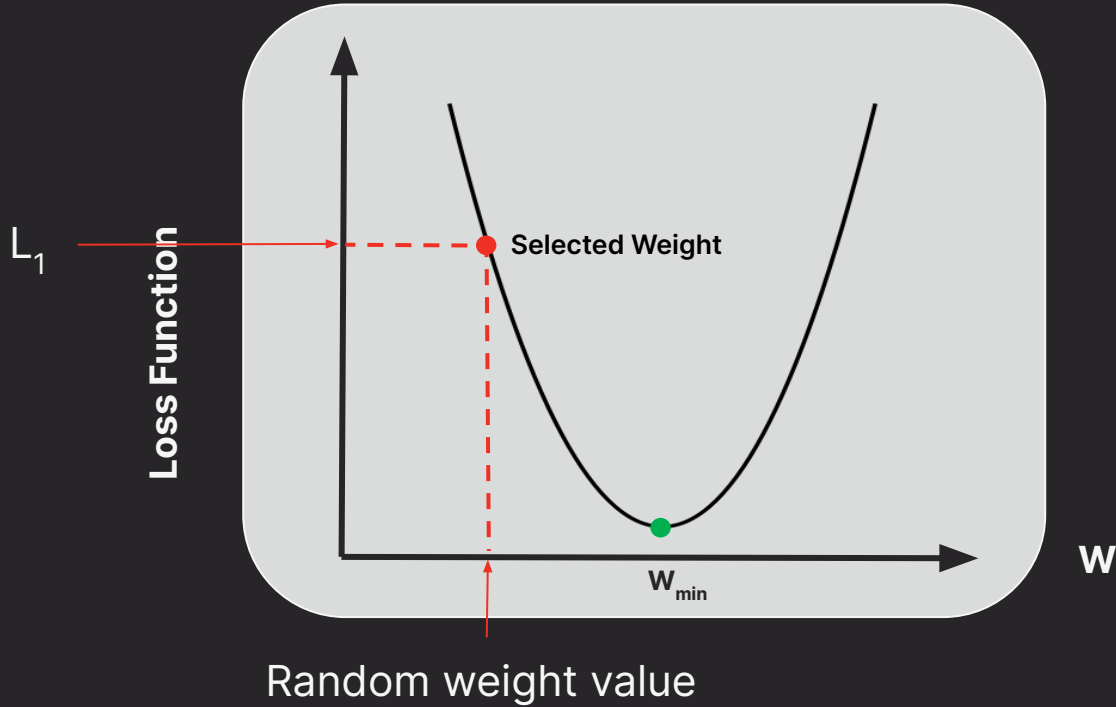


**Note:** Assume bias as constant to focus only on the impact of weight ( $w$ ) on loss.

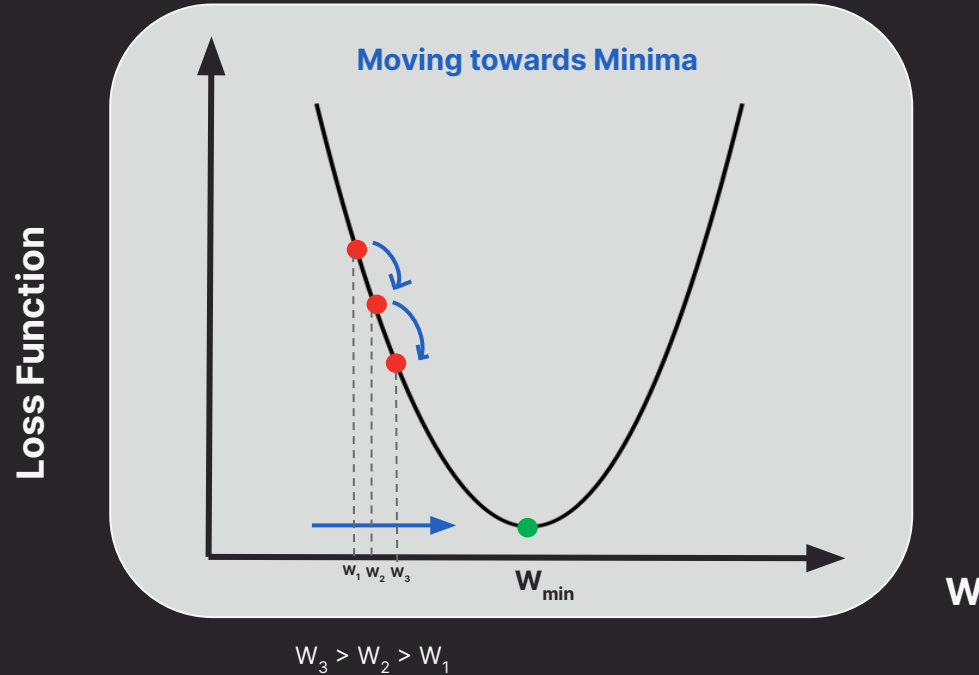
What is  $W_{min}$  ?



# Calculating the Slope

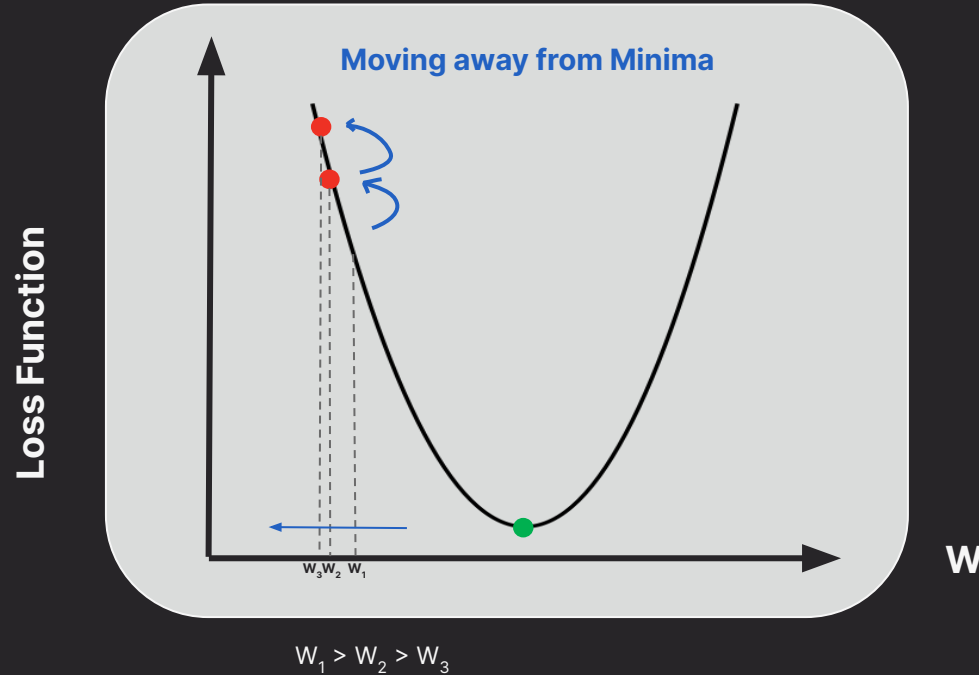


# Calculating the Slope

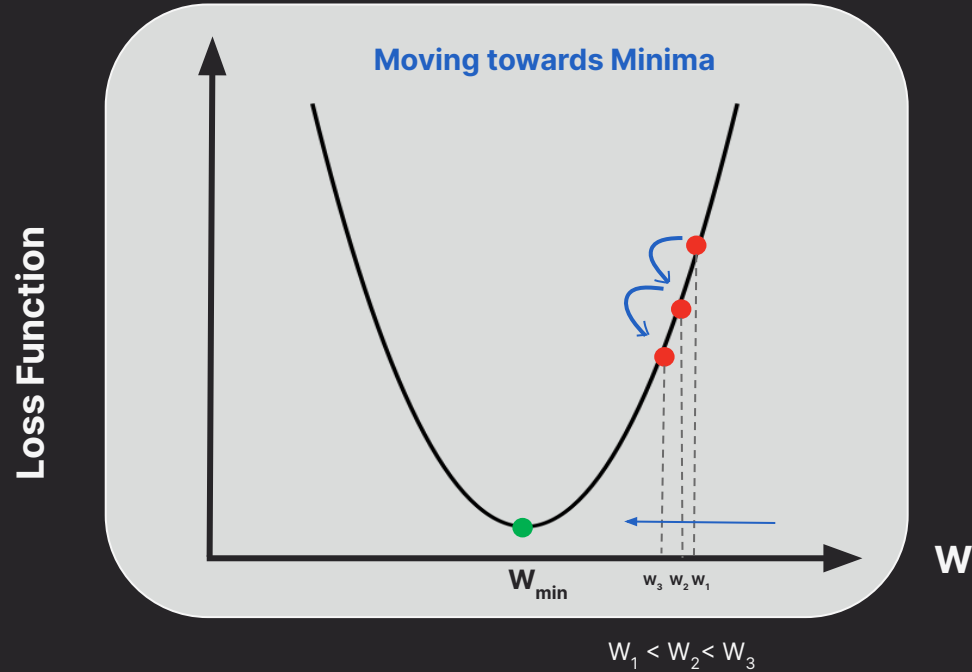




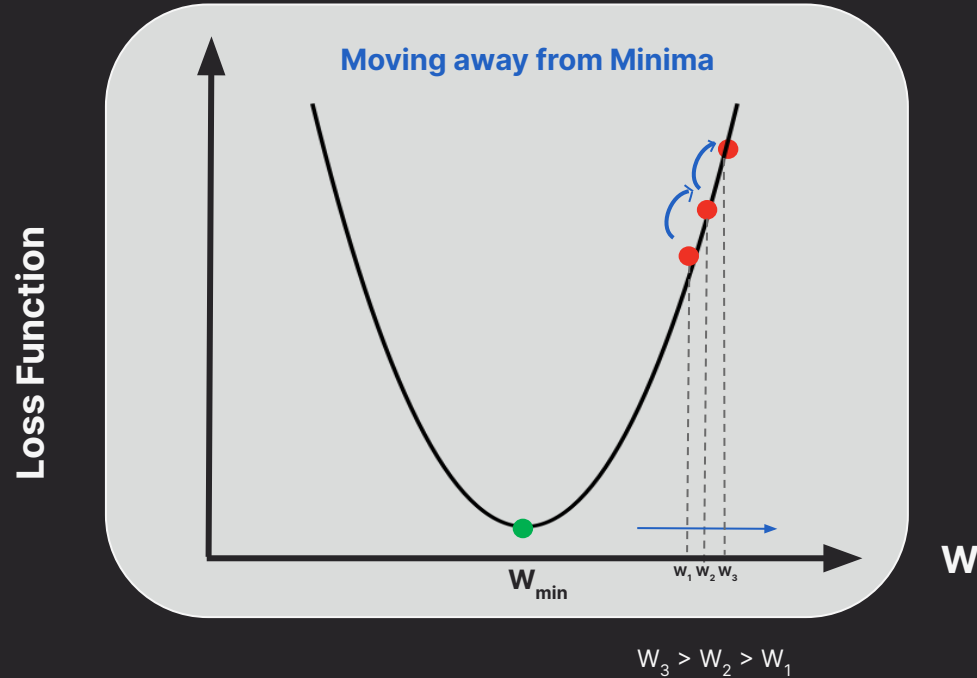
# Calculating the Slope



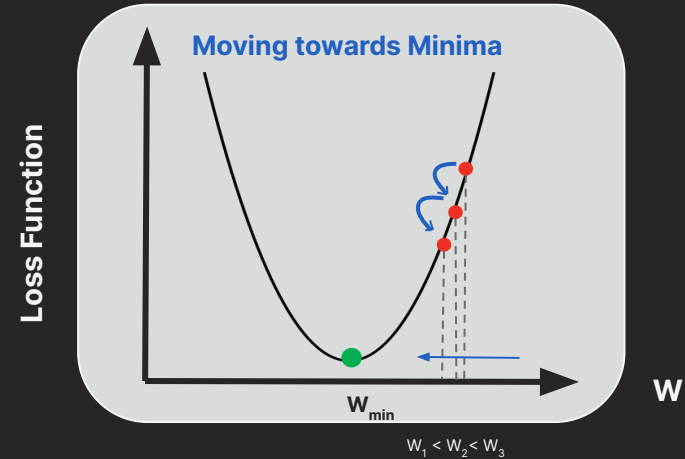
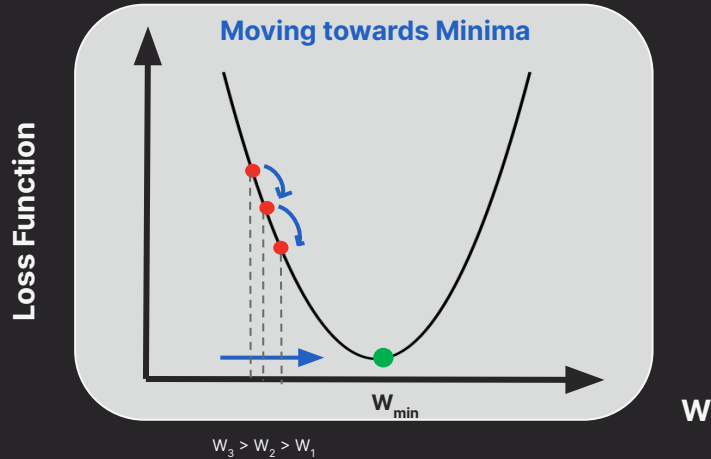
# Calculating the Slope



# Calculating the Slope

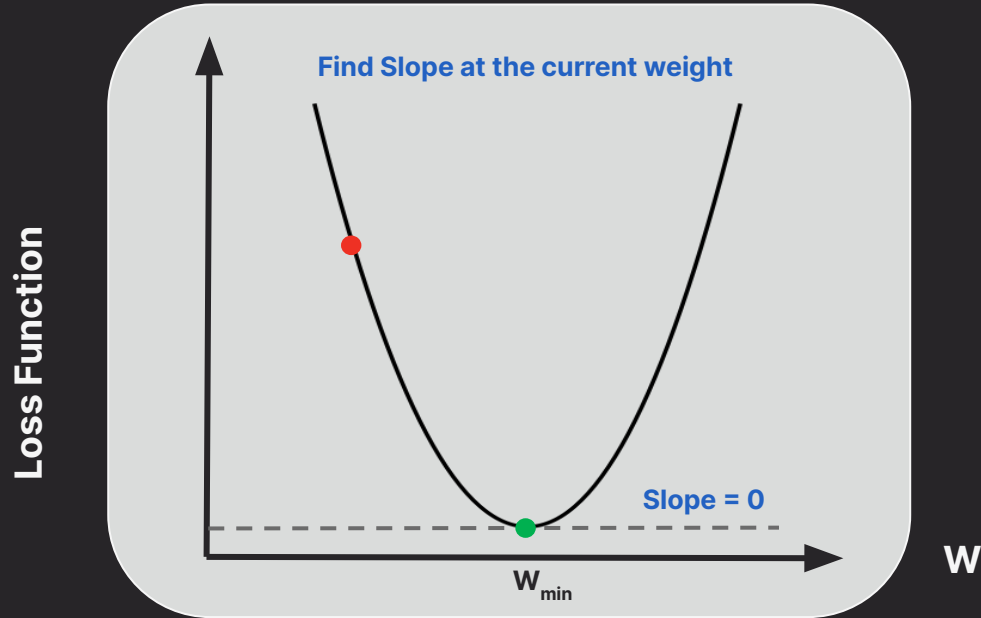


# Calculating the Slope

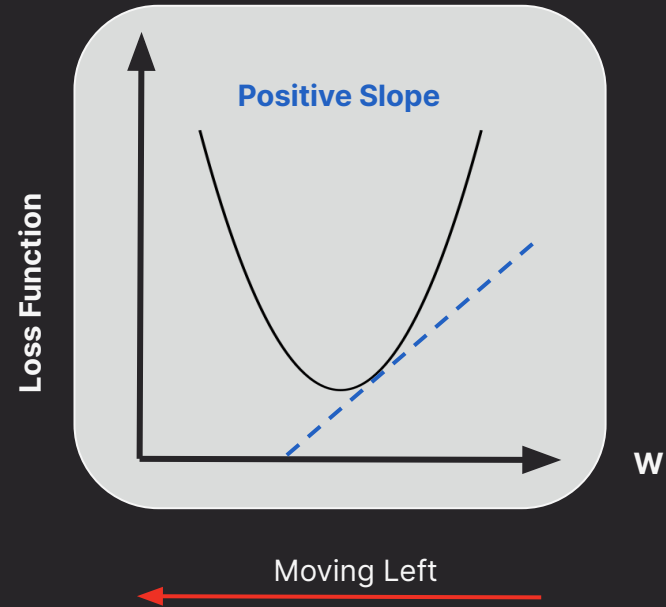
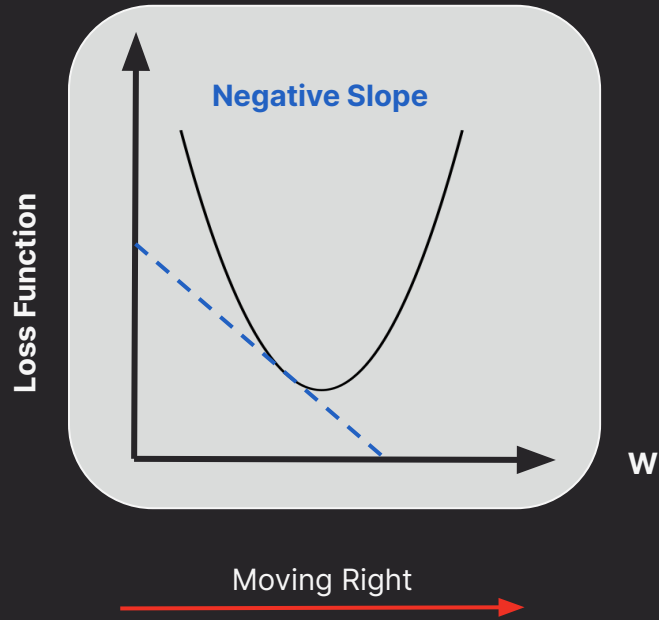


How would the algorithm know which side to move towards?

# Calculating the Slope

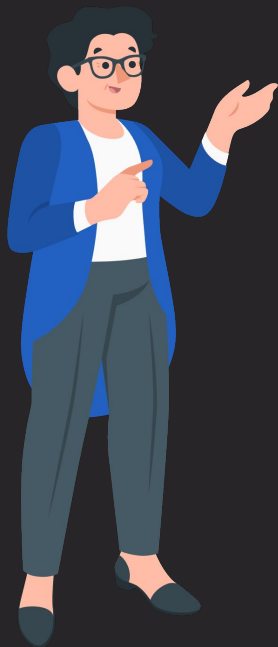


# Calculating the Slope



# Calculating the Slope

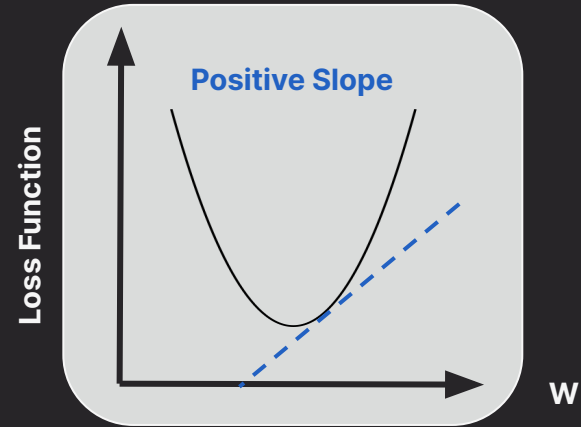
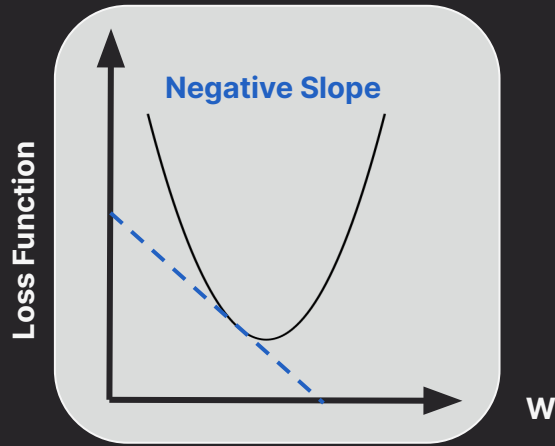
$$\frac{d}{dw}(MSE)$$



$$MSE = \frac{1}{N} \sum_{i=1}^n ((wx_i + b) - y_i)^2$$

# Updating Weights

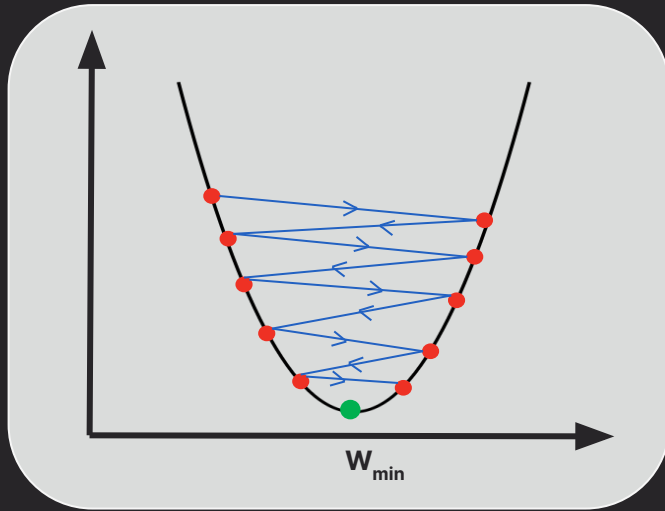
$$\text{New weight } (w_{\text{new}}) = \text{Old Weight } (w_{\text{old}}) - \text{Slope}$$





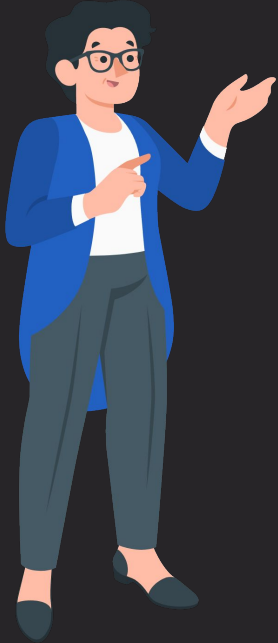
# Challenges in Weight Update

$$\text{New weight } (w_{new}) = \text{Old Weight } (w_{old}) - \text{Slope}$$



Subtracting large slope can jump over the minimum.

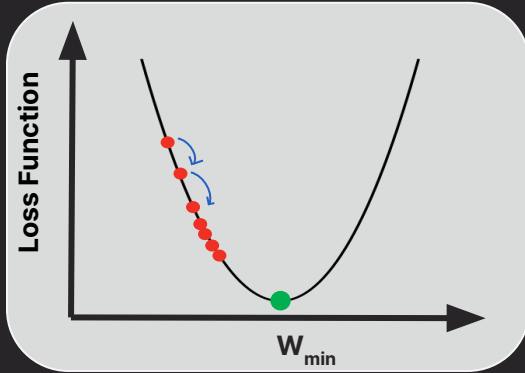
# Using a Learning Rate



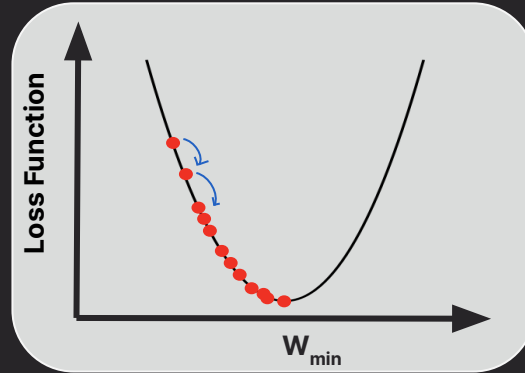
**Learning Rate** : A small positive number that scales the slope to refine weight updates.

$$w_{new} = w_{old} - \text{Learning Rate} * \text{slope}$$

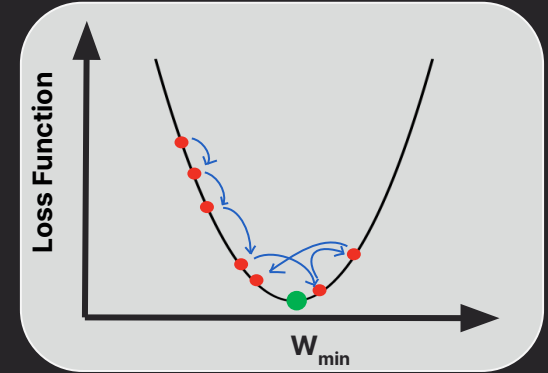
# Choosing a Learning Rate



Small Learning Rate



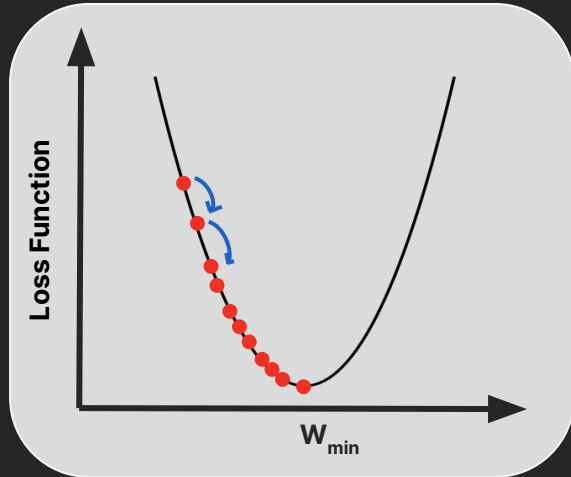
Ideal Learning Rate



Large Learning Rate

# Using a Learning Rate

$$w_{new} = w_{old} - \text{Learning Rate} * \text{slope}$$



**Learning rate = 0.001** (Good starting point to train models)



This value can be increased or decreased based on observations.

**When do we  
stop updating  
the weights ?**



# When to Stop Updating?

1. Stop when  $w_{new} = w_{old}$

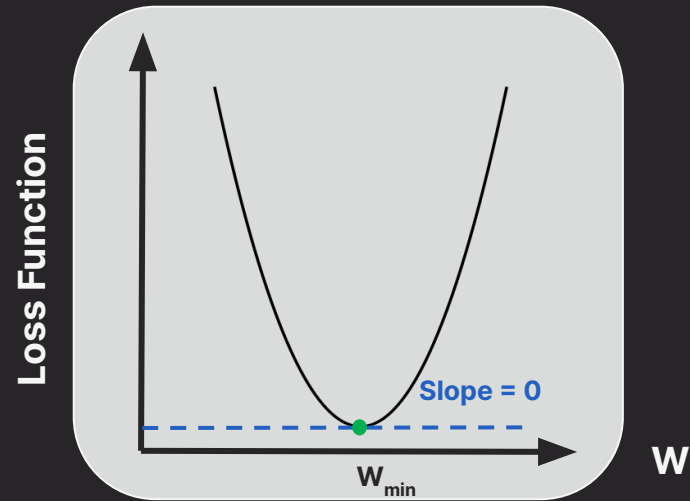
$$w_{new} = w_{old} - \text{learning rate} * \text{slope}$$

↓

$$w_{new} = w_{old} - \text{learning rate} * 0$$

↓

$$w_{new} = w_{old}$$



# When to Stop Updating?

## 2. Limit number of iterations

Ex: Limiting iterations to 1000





# The Math Behind Gradient Descent



Find the derivative of the slope with respect to the weight  $w$

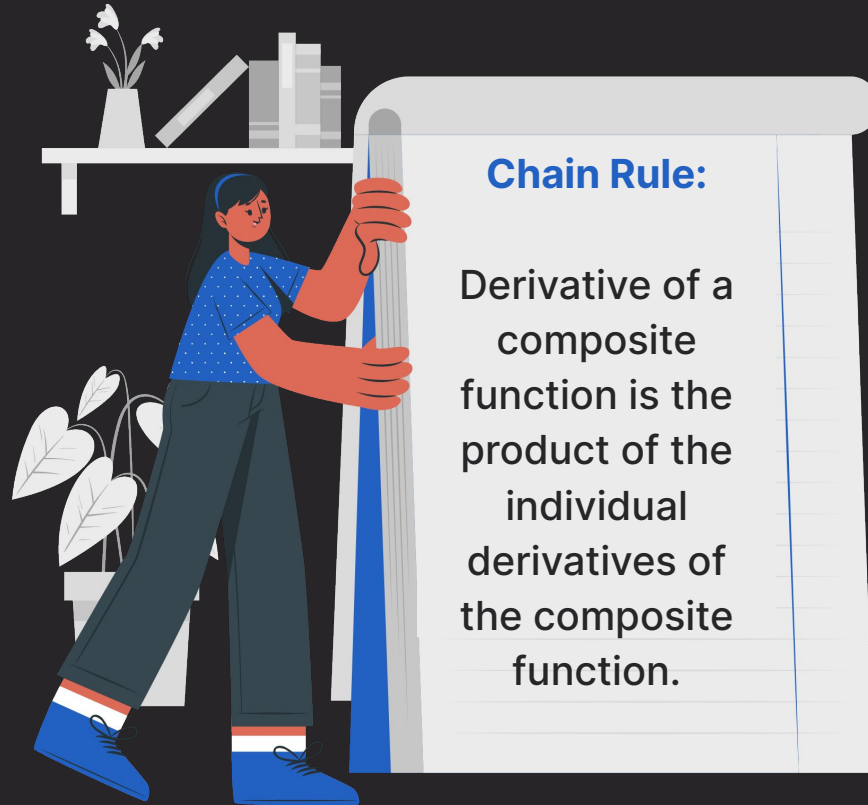
$$L = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$\hat{y} = \omega_i^* x_i + b$$

$x_i$  = features from dataset



# Chain Rule of Differentiation



## Chain Rule:

Derivative of a composite function is the product of the individual derivatives of the composite function.

$$L = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$\hat{y} - y = u$$

$$L = \frac{1}{N} \sum u^2$$

# Chain Rule of Differentiation

$$\frac{dL}{du} = \frac{1}{N} \frac{d(u^2)}{du}$$

$$\frac{dL}{du} = \frac{2u}{N}$$

# Chain Rule of Differentiation

$$\frac{dL}{d\hat{y}} = \frac{du}{d\hat{y}} * \frac{dL}{du}$$

How  $u$  impacts the **loss value** and then how  $\hat{y}$  impacts  $u$ .

# Chain Rule of Differentiation

$$\frac{du}{d\hat{y}} = \frac{d}{d\hat{y}} (\hat{y} - y) = 1 \text{ ————— } \textcircled{1}$$

$$\frac{dL}{du} = \frac{2u}{N} \text{ ————— } \textcircled{2}$$

$$\frac{dL}{d\hat{y}} = \frac{du}{d\hat{y}} * \frac{dL}{du} \longrightarrow \frac{dL}{d\hat{y}} = \frac{2u}{N} \text{ ————— Combining 1 \& 2}$$

# Chain Rule of Differentiation

$$y_i = w \cdot x_i + b$$

$$\frac{dL}{dw} = \frac{dL}{d\hat{y}} * \frac{d\hat{y}}{dw} \quad \text{—————} \textcircled{1}$$

$$\frac{dL}{db} = \frac{dL}{d\hat{y}} * \frac{d\hat{y}}{db} \quad \text{—————} \textcircled{2}$$

# Chain Rule of Differentiation

$$\frac{dL}{dw} = 2u \cdot x_i \quad \text{or} \quad 2(\hat{y} - y) \cdot x_i$$

$$2((w^*x + b) - y) \cdot x_i$$

# Chain Rule of Differentiation

$$w_{new} = w_{old} - \text{learning rate} * \text{slope}$$

$$b_{new} = b_{old} - \text{learning rate} * \text{slope}$$

$$\frac{dL}{dB} = 2u - y = 2(w_{old} * x + b)$$





# Backpropagation

Find how weights in earlier layers affect the loss function in a multi neural network.