# Improving Deep Neural Networks

**Video 1:** Problems in Deep Neural Networks

Analytics
Vidhya

In Air

# Challenges in Neural Networks

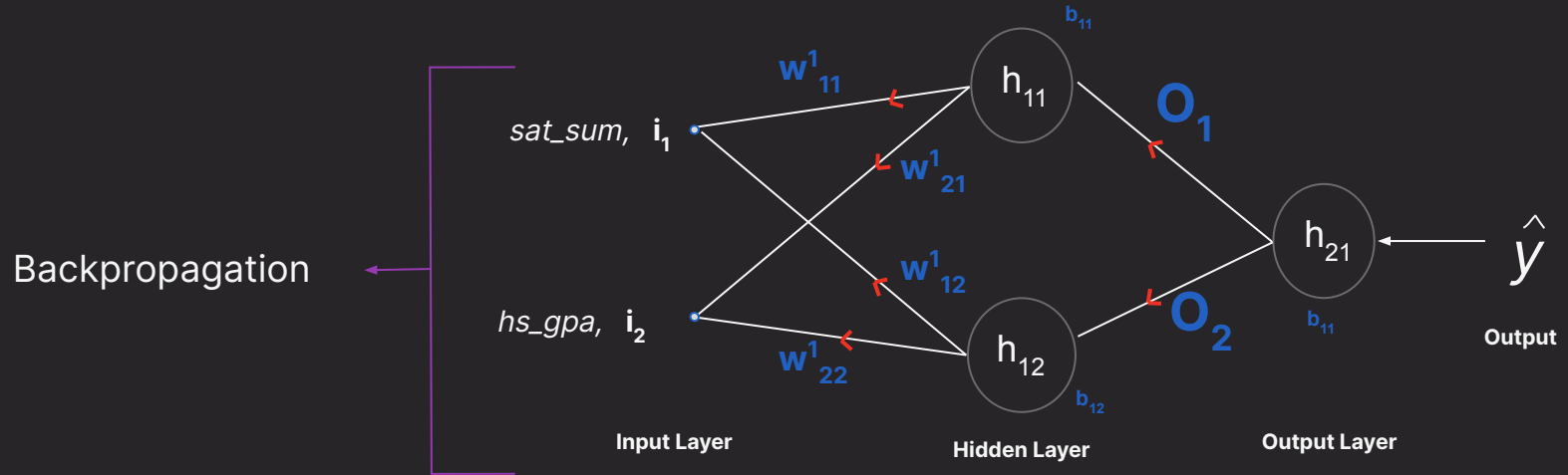## Problems with Large Neural Networks

▷ **Vanishing and Exploding Gradients**

▷ Overfitting

# Challenges in Neural Networks



Backpropagation

**How does too small or too large weight affect neural network learning ?**

# Vanishing Gradients



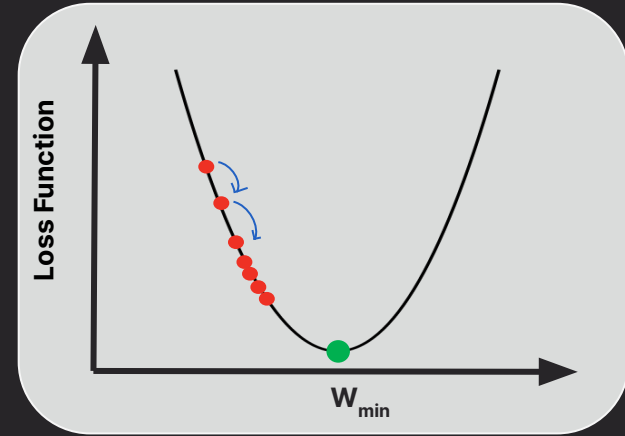▷ Update to the weights is small

$$w_{new} = w_{old} - eeta * dL/dw$$
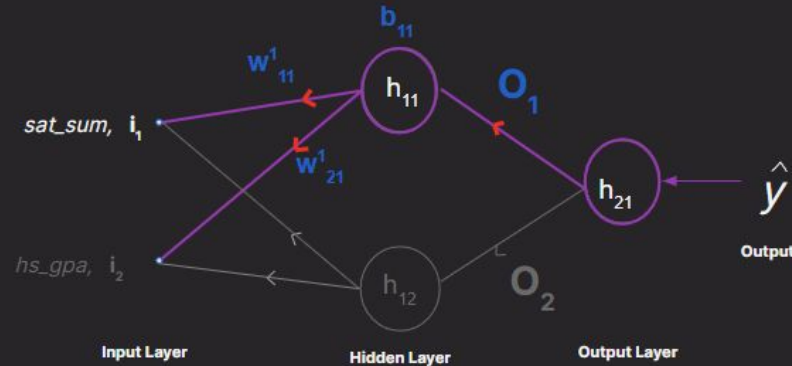
Becomes closer to 0

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

**Small Learning Rate**

Analytics Vidhya

Loss Function

$W_{min}$
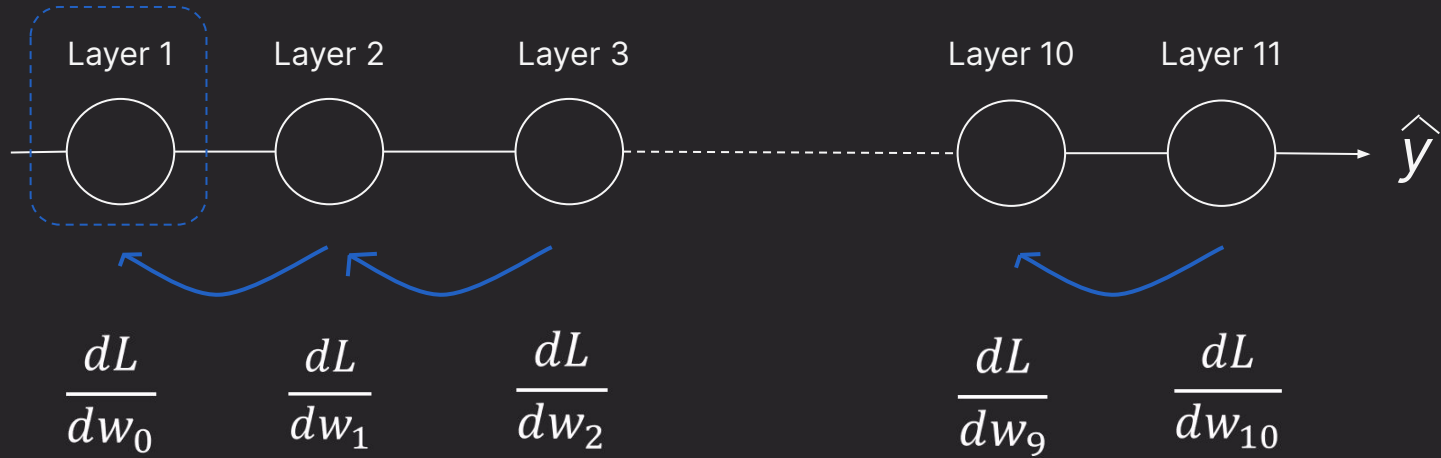
# Chain Rule of Differentiation



$$\frac{dL}{dw_{11}^1} = \frac{dL}{dO_1} * \frac{dO_1}{dw_{11}^1}$$

$$\frac{dL}{dw_{21}^1} = \frac{dL}{dO_1} * \frac{dO_1}{dw_{21}^1}$$

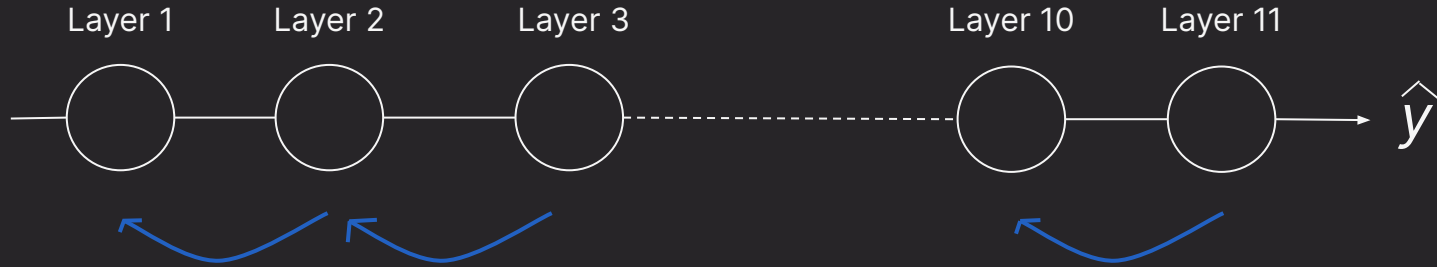$$\frac{dL}{db_{11}} = \frac{dL}{dO_1} * \frac{dO_1}{db_{11}}$$

# Chain Rule of Differentiation
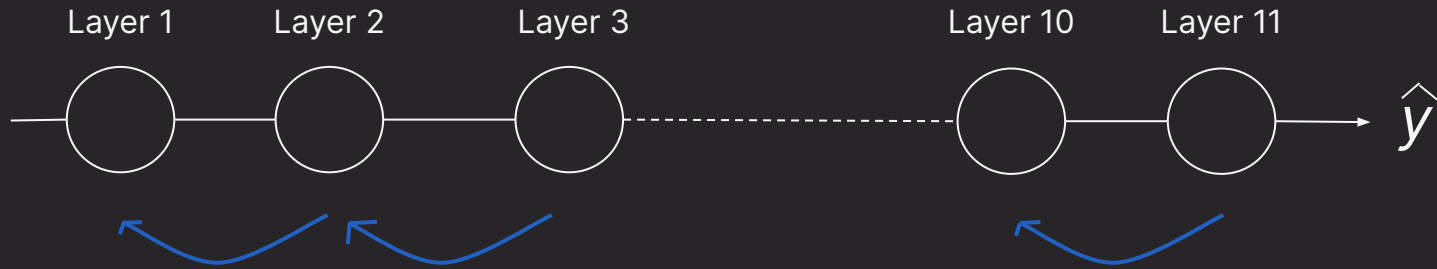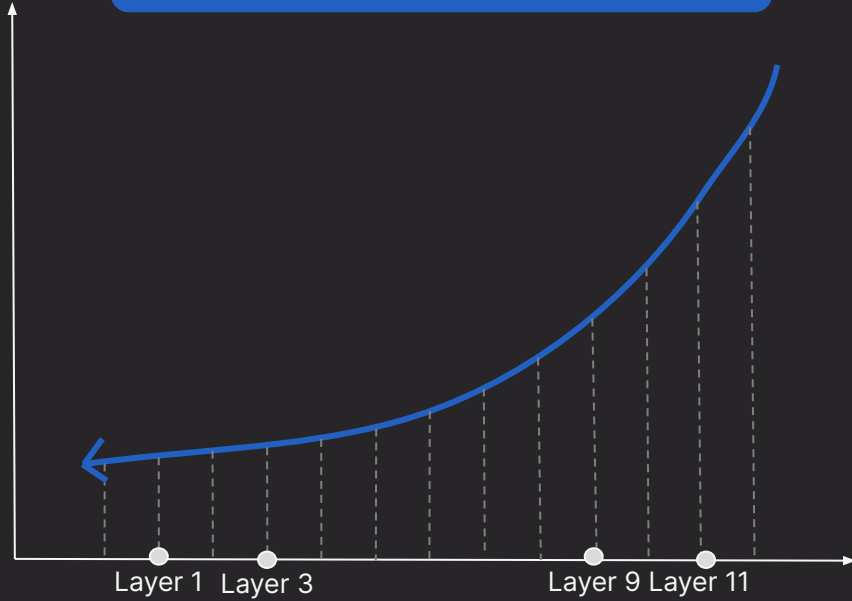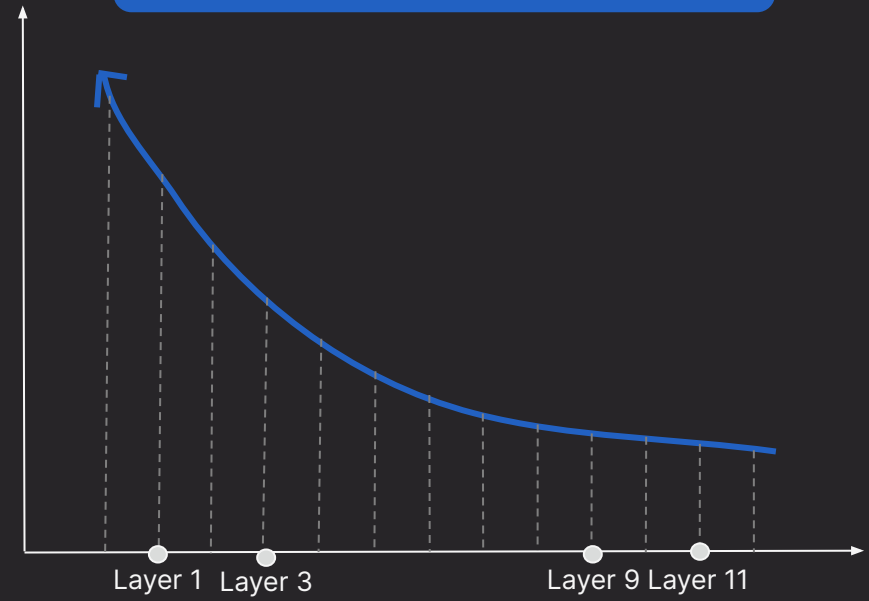


Could be represented as multiplication of all previous Gradients $\left(\dfrac{dL}{dw_n}\right)^{10}$

| Exploding Gradient | Vanishing Gradient |
|---|---|
| There is an exponential growth in the model parameters of the lower or initial layers. | The parameters of the higher layers change significantly whereas the parameters of lower or initial layers would not change much or not at all. |
| The model weights may become NaN during training | The model weights may become close to 0 during training. |
| During training the model may abruptly increase loss values by a large amount. | During training the model learns very slowly and the training stagnates at a very early stage. |