# Back Propagation

**Back Propagation makes Neural Networks Trainable**
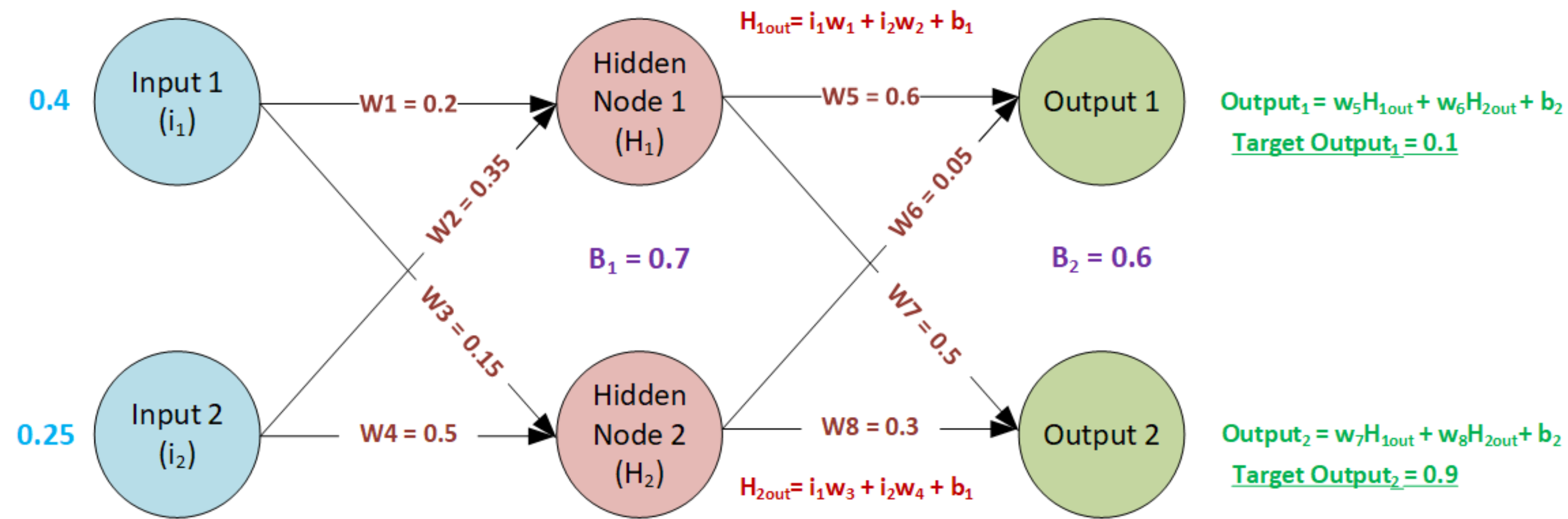
# Back Propagation
## This is what makes Neural Networks Trainable :)

- The importance of Back Propagation cannot be understated

- Using the loss, it tells us how much to change/update the gradients by so that we reduce the overall loss
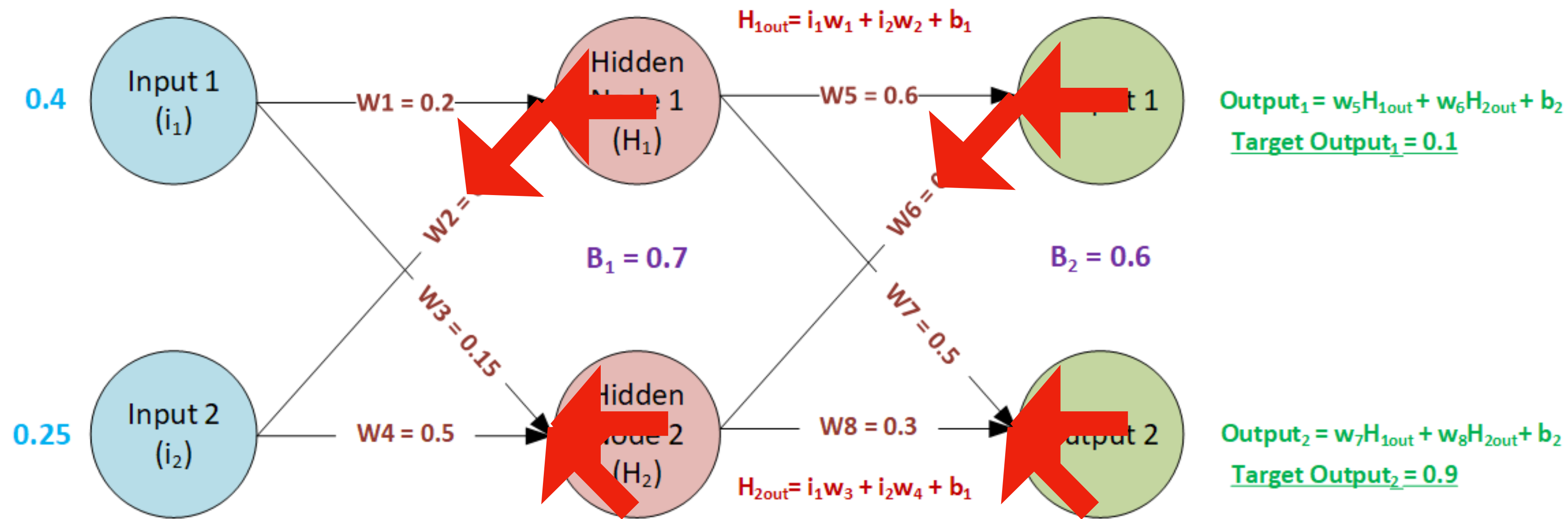
# Back Propagation Example
## Explained Using Neural Networks



- Let's look at a regular Neural Network above.

- Using the Loss value, Back Propagation can tells us whether a **small increase** of $W_5$ to 0.6001 or a **small decrease** 0.5999 will lead to a **reduction in the overall loss**

# Back Propagation Example



$H_{1out} = i_1 w_1 + i_2 w_2 + b_1$

$Output_1 = w_5 H_{1out} + w_6 H_{2out} + b_2$

Target Output$_1$ = 0.1

Input 1 ($i_1$) — 0.4

W1 = 0.2

Hidden Node 1 ($H_1$)

W5 = 0.6

Output 1

$B_1 = 0.7$

$B_2 = 0.6$

W2 =

W6 =

W3 = 0.15

W7 = 0.5

Input 2 ($i_2$) — 0.25

W4 = 0.5

Hidden Node 2 ($H_2$)

W8 = 0.3

Output 2

$H_{2out} = i_1 w_3 + i_2 w_4 + b_1$

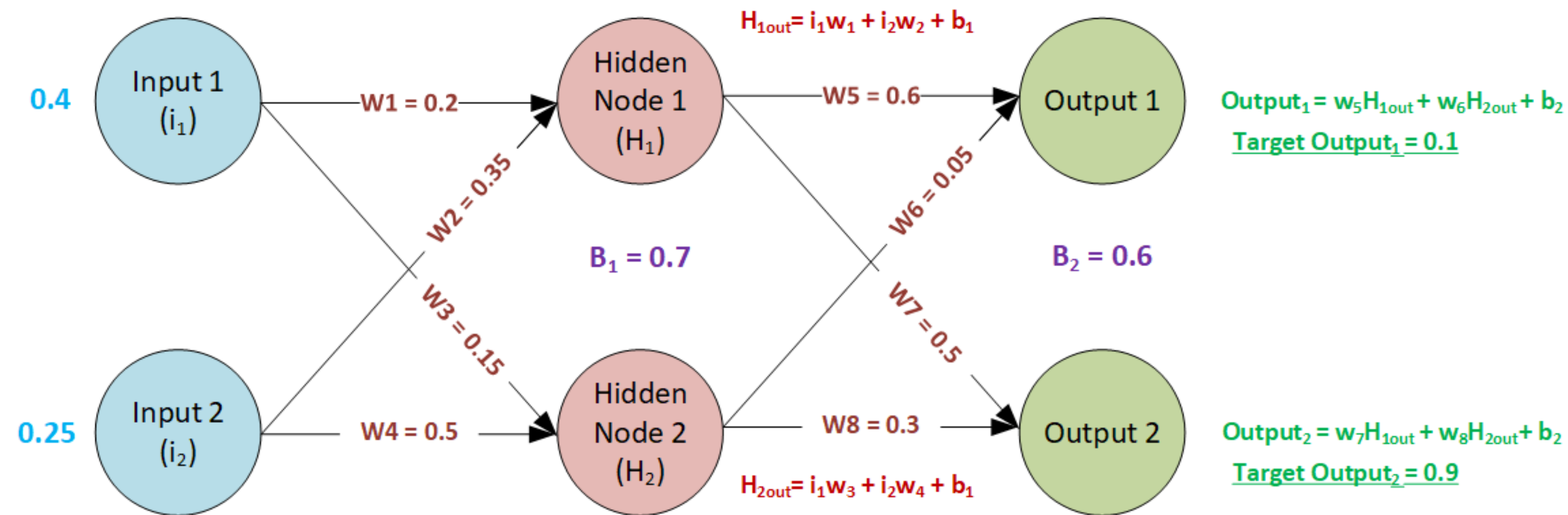$Output_2 = w_7 H_{1out} + w_8 H_{2out} + b_2$

Target Output$_2$ = 0.9

- Moving **right to left**

- Back Propagation gives us the new gradient or weight values for each node so that the overall loss is decreased

- This is done for all nodes

# Back Propagation Process

- By **forward** propagating input data we can use back propagation to **lower** the weights to **lower the loss**

- **But**, this simply tunes the weights for that particular input (or batch of inputs)

- We improve **Generalisation** (ability to make good predictions on unseen data) by using all data in our training dataset

- By continuously changing the weights for each data input (or batch of images) we are lowering the overall loss for our training data.

# What do our Weights or Gradients Look Like?

## Let's look at a Simple Neural Network



- The output from Hidden Node 1 is:

$$H_1 out = i_1 w_1 + i_2 w_2 + b_1$$

# For a Convolutional Neural Network

| 1 | 0 | 1 | 0 | 1 |
|---|---|---|---|---|
| 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 |

\*

| 0 | 1 | 0 |
|---|---|---|
| 1 | 0 | -1 |
| 0 | 1 | 0 |

=

| 2 | 1 | -1 |
|---|---|---|
| -1 | 1 | 3 |
| 2 | 1 | 1 |

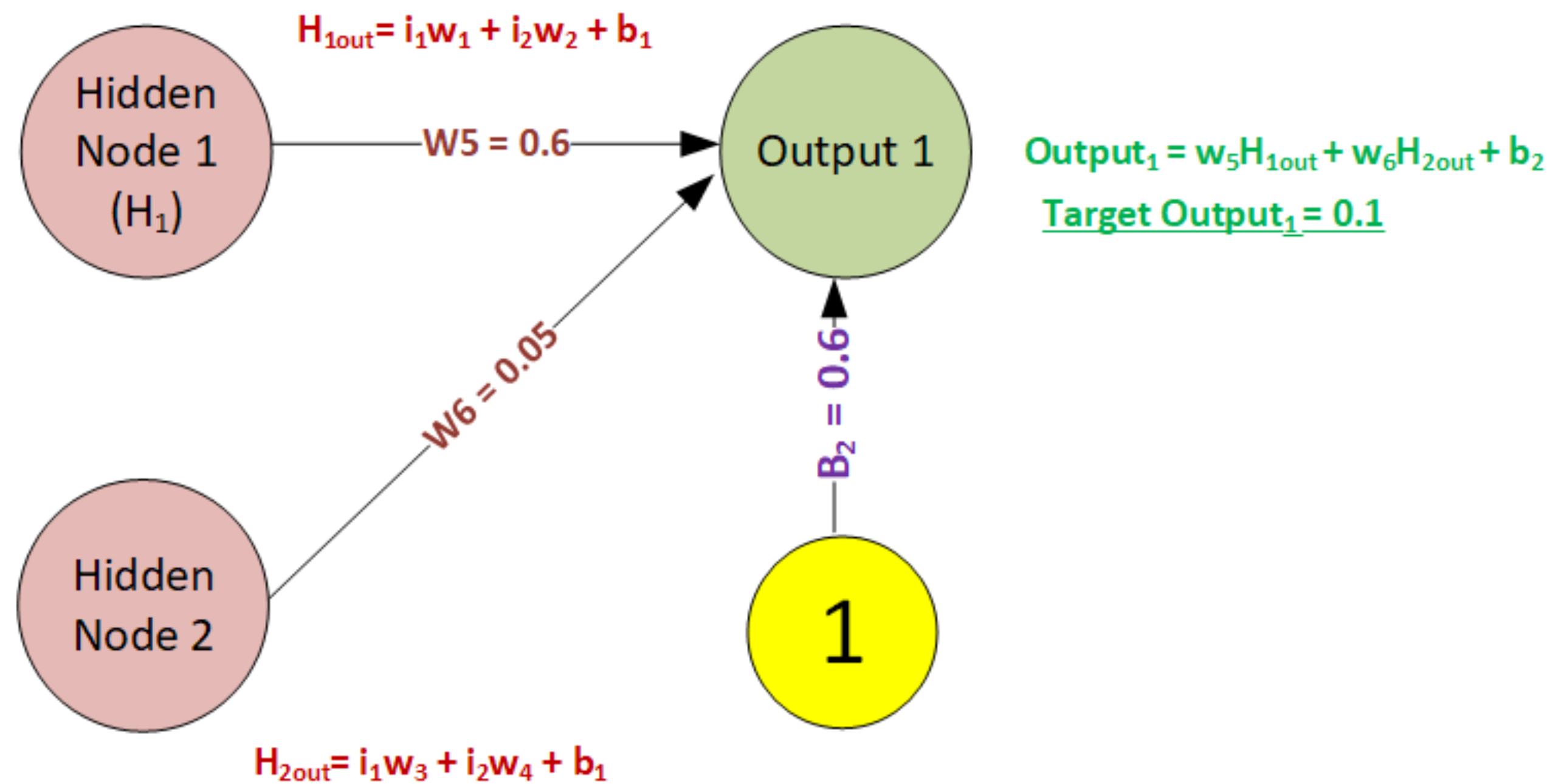Input Image          Filter or Kernel          Output or Feature Map

- The values of our Filter/Kernel are the weights!

# How does Back Propagation Work?

- **Chain Rule!**

- If we have two functions $y = f(u)$ and $u = g(x)$ then the derivative of $y$ is:

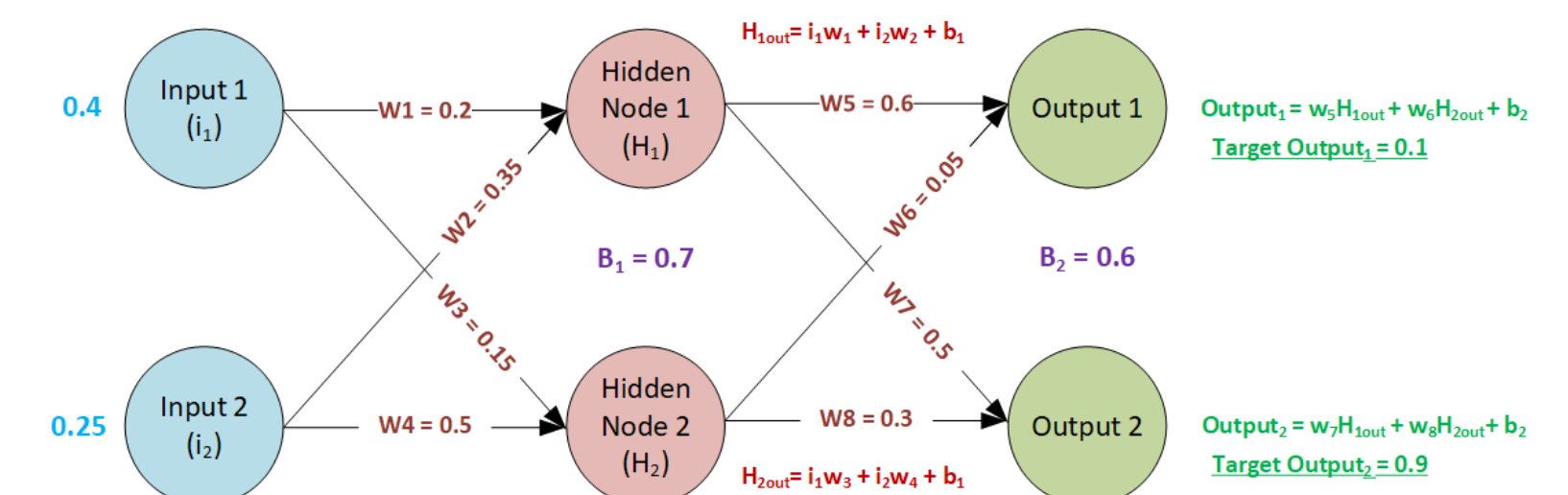- $$\frac{dy}{dx} = \frac{dy}{du} \times \frac{du}{dx}$$

# A Simple Back Propagation Example



$H_{1out} = i_1w_1 + i_2w_2 + b_1$

W5 = 0.6

W6 = 0.05

$B_2 = 0.6$

Hidden Node 1 ($H_1$)

Hidden Node 2

Output 1

1

$Output_1 = w_5H_{1out} + w_6H_{2out} + b_2$

Target Output$_1$ = 0.1

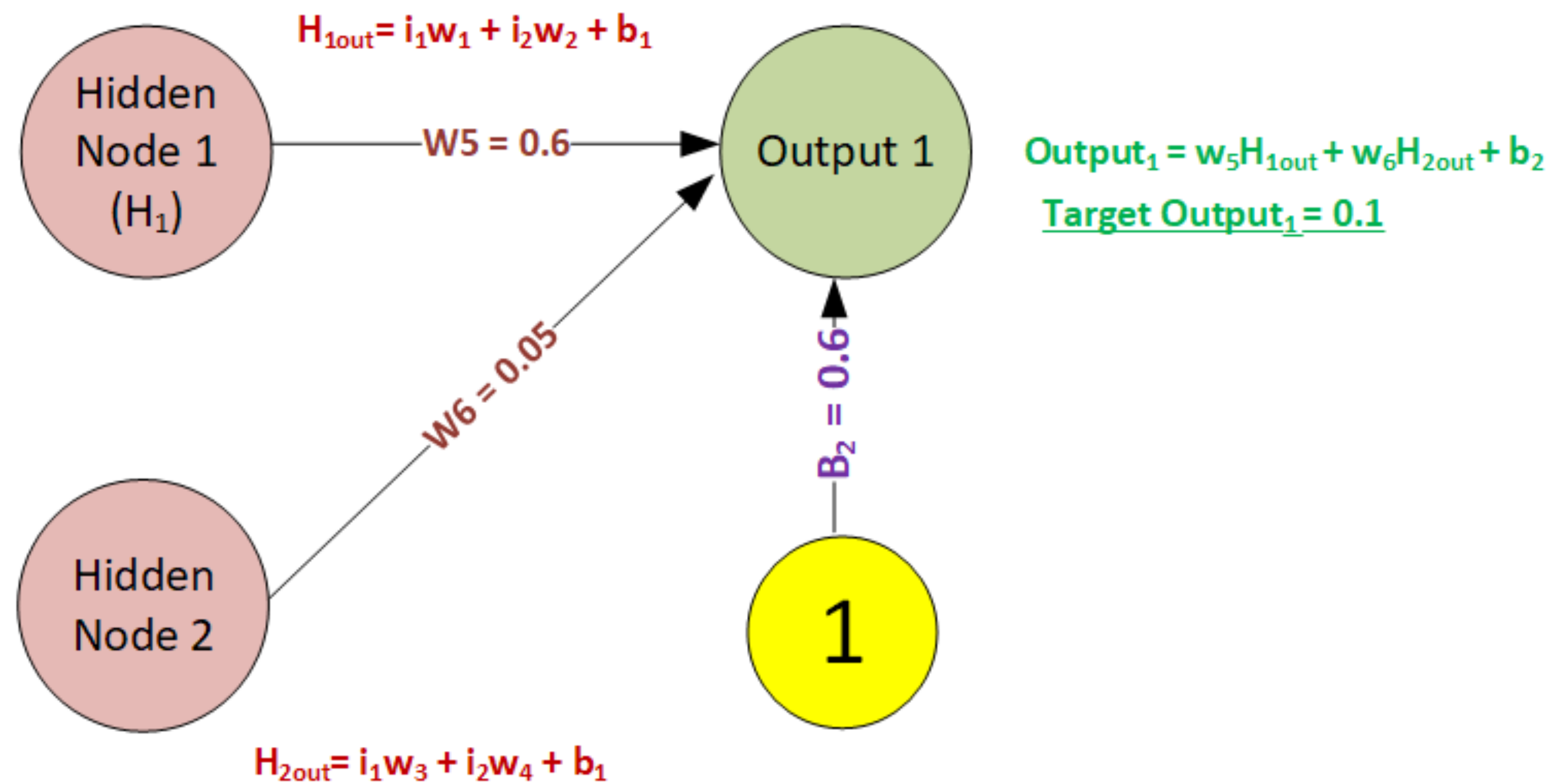$H_{2out} = i_1w_3 + i_2w_4 + b_1$

- We want to know how much changing $W_5$ changes the **Total Error**.

- That is given by:

$$\frac{dE_T}{dW_5}$$

- Where $E_T$ is the sum of the error from Outputs 1 and 2 (see below)

# A Simple Back Propagation Example



New $W_5 = -\lambda \times \dfrac{dE_T}{dW_5}$

- Note we introduced a new parameter $\lambda$

- $\lambda$ is our learning rate

- It controls how a big a jump (positive or negative) we take when updating $W_5$

- Large learning rates train faster, but can get stuck in a Global Minimum

- Small learning rates train more slowly

# Next...

**Gradient Descent**