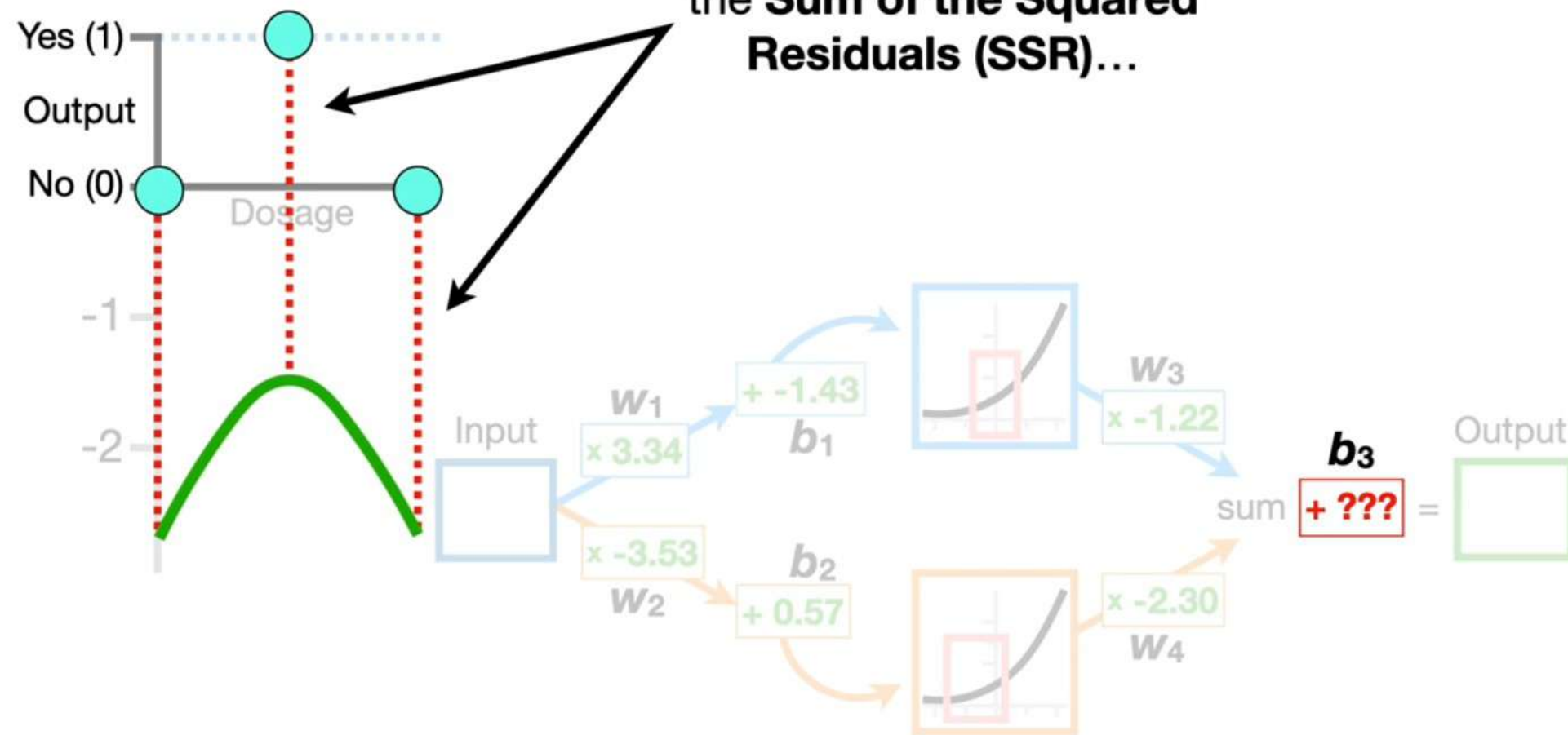| Name | Plot | Equation | Derivative |
|---|---|---|---|
| <span style="color:red">Activation functions with small ranges are usually used for solving classification problems.</span> | | | |
| Identity | | $f(x) = x$ | $f'(x) = 1$ |
| Binary step | | $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$ |
| Logistic (a.k.a Soft step) | | $f(x) = \dfrac{1}{1 + e^{-x}}$ | $f'(x) = f(x)(1 - f(x))$ |
| TanH | | $f(x) = \tanh(x) = \dfrac{2}{1 + e^{-2x}} - 1$ | $f'(x) = 1 - f(x)^2$ |
| ArcTan | | $f(x) = \tan^{-1}(x)$ | $f'(x) = \dfrac{1}{x^2 + 1}$ |
| Rectified Linear Unit (ReLU) | | $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ |
| Parameteric Rectified Linear Unit (PReLU) [2] | | $f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ |
| Exponential Linear Unit (ELU) [3] | | $f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ |
| SoftPlus | | $f(x) = \log_e(1 + e^x)$ | $f'(x) = \dfrac{1}{1 + e^{-x}}$ |

We first used **The Chain Rule** to calculate the derivative of the **Sum of the Squared Residuals (SSR)**…

$$\frac{d\ SSR}{d\ b_3} = \frac{d\ SSR}{d\ Predicted} \times \frac{d\ Predicted}{d\ b_3}$$

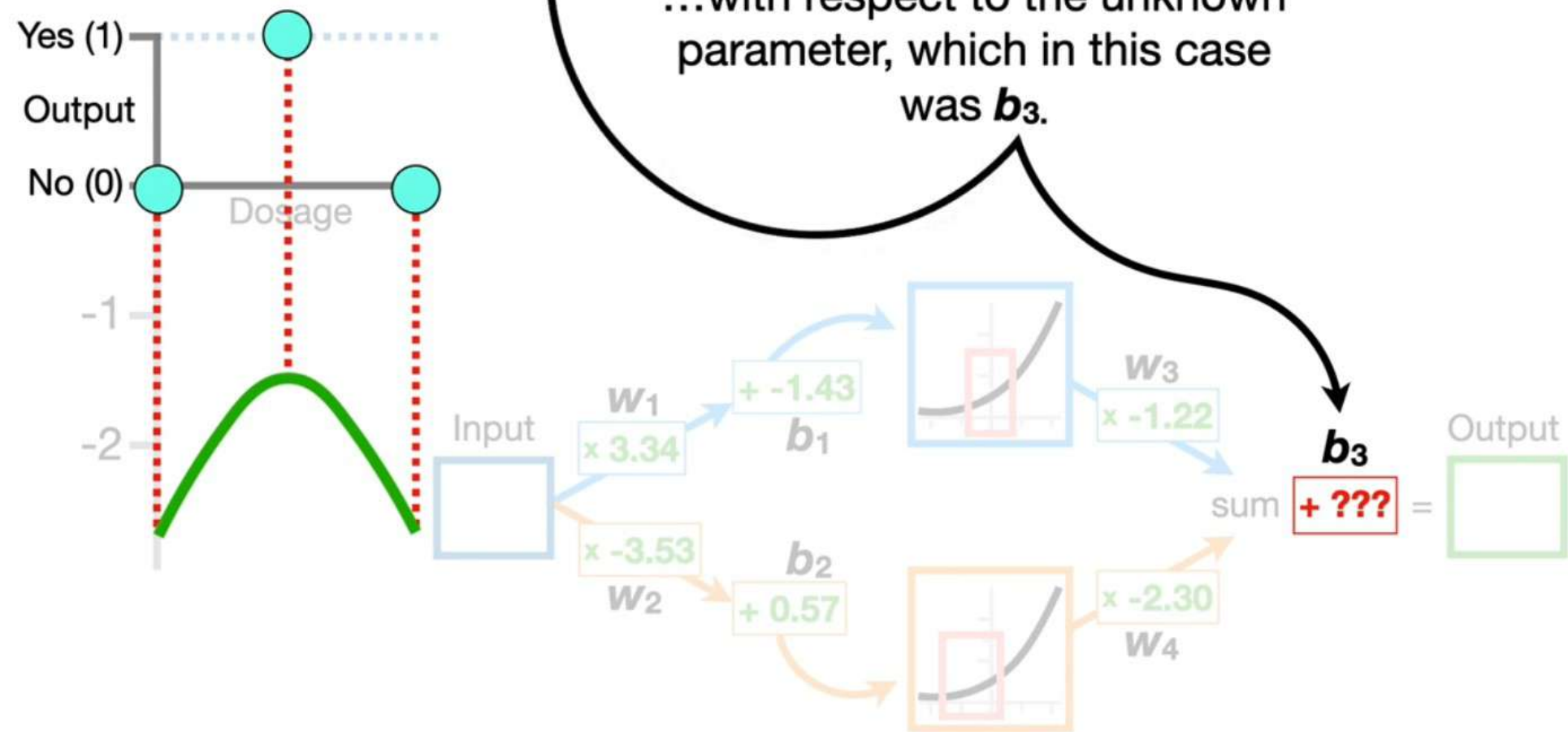...with respect to the unknown parameter, which in this case was $b_3$.

Yes (1)

Output

No (0)

Dosage

-1

-2

Input

$w_1$ × 3.34 → + -1.43 $b_1$

$w_3$ × -1.22

$w_2$ × -3.53 → $b_2$ + 0.57

$w_4$ × -2.30

$b_3$

sum + ??? =

Output

$$\frac{d\ SSR}{d\ b_3} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times \mathbf{1}$$

…with respect to the unknown parameter, which in this case was $b_3$.

Yes (1)

Output

No (0)

Dosage

-1

-2

Input

$w_1$ × 3.34

+ -1.43 $b_1$

$w_3$ × -1.22

$b_3$

$x$ -3.53 $w_2$

$b_2$ + 0.57

$x$ -2.30 $w_4$

sum + ??? =

Output

$$\frac{d\ SSR}{d\ b_3} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times 1$$

Then we initialized the unknown parameter with a number, and in this case we set $b_3 = 0\ldots$

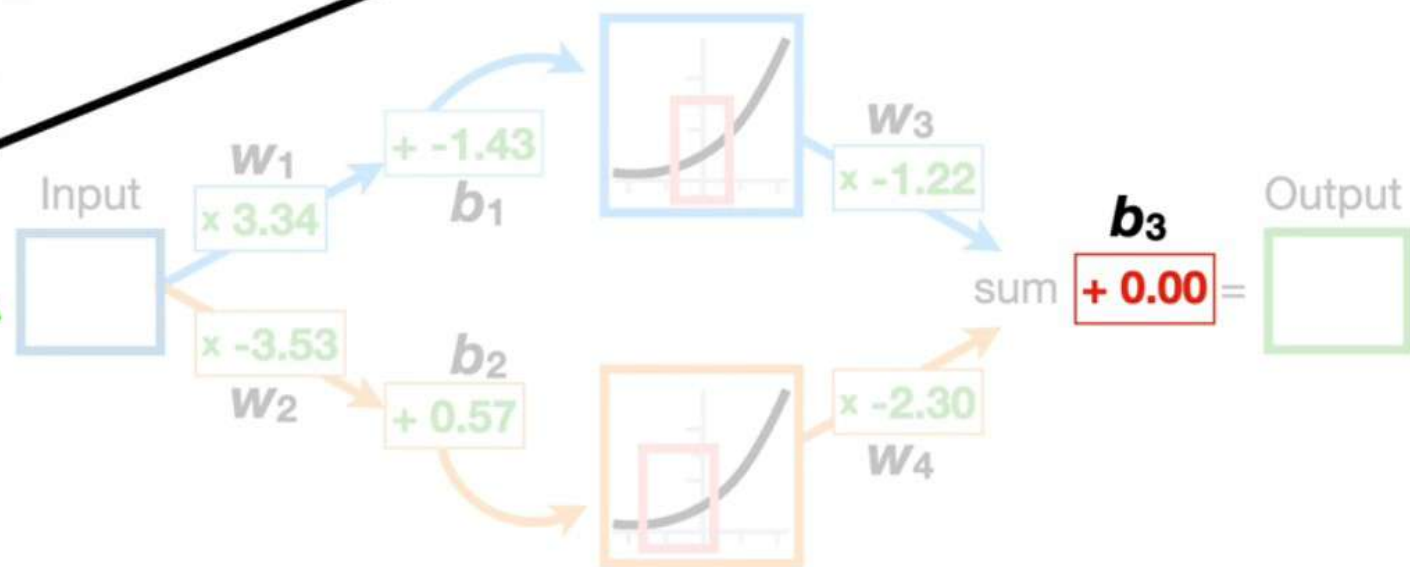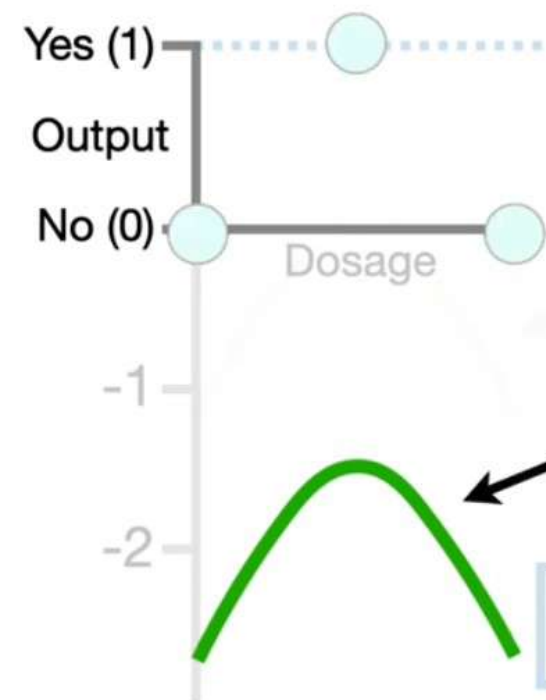$$\frac{d\ SSR}{d\ b_3} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times \mathbf{1}$$

...and used **Gradient Descent** to optimize the unknown parameter.

SSR

Bias$_3$ ($b_3$)

Yes (1)

Output

No (0)

Dosage

-1

-2

Input

$w_1$
× 3.34

+ -1.43
$b_1$

$w_3$
× -1.22

$b_3$
+ 0.00 =

Output

× -3.53
$w_2$

$b_2$
+ 0.57

× -2.30
$w_4$

sum

$$\frac{d\ SSR}{d\ b_3} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times 1$$

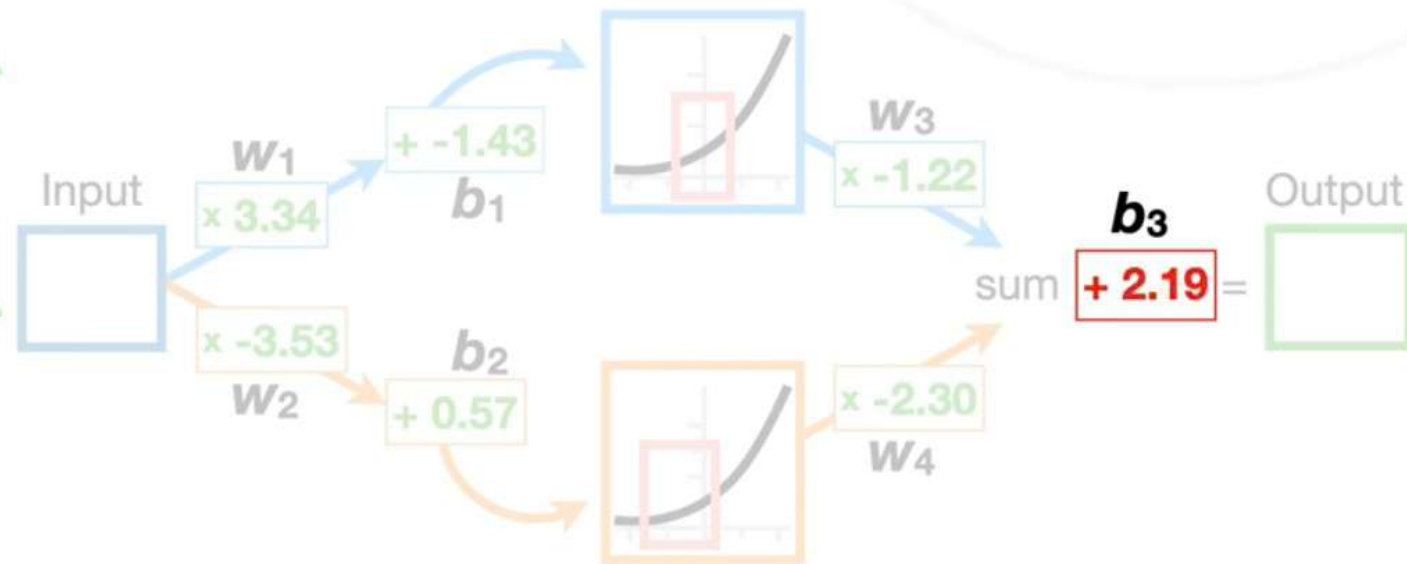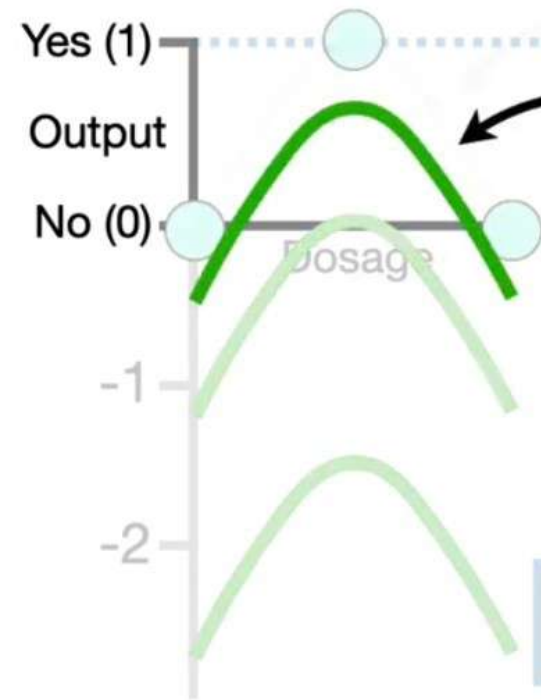…and used **Gradient Descent** to optimize the unknown parameter.
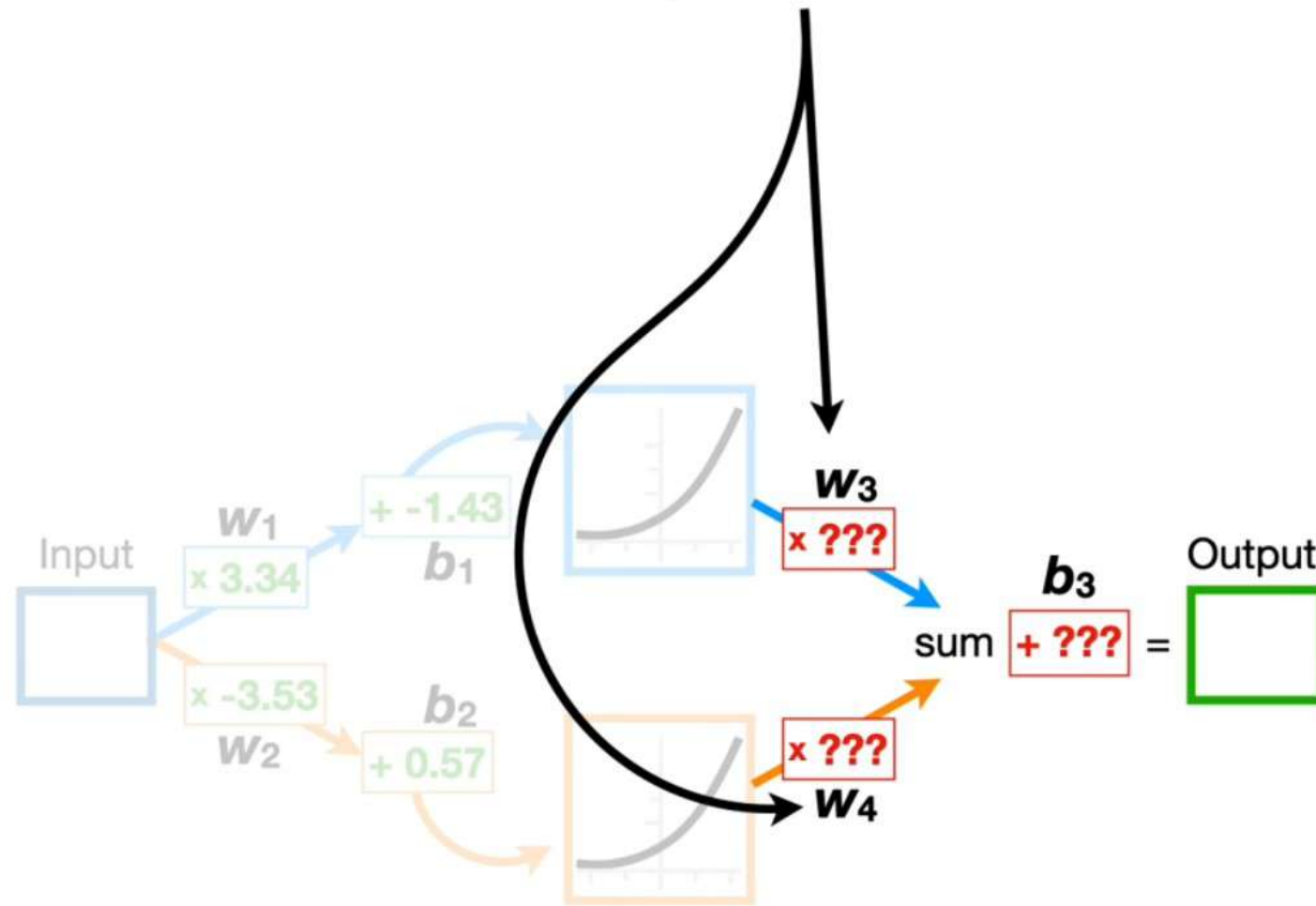
SSR

Bias₃ ($b_3$)

Yes (1)

Output

No (0)

Dosage

-1

-2

Input

$w_1$
$\times 3.34$

$+ -1.43$
$b_1$

$w_3$
$\times -1.22$

$\times -3.53$
$b_2$
$+ 0.57$
$w_2$

$\times -2.30$
$w_4$

$b_3$
sum $+ 2.19$ =

Output

The first thing we do is initialize
the **Weights**, $w_3$ and $w_4$, with
random starting values…

Input

$w_1$
$\times 3.34$

$+ -1.43$
$b_1$

$w_3$
$\times \text{???}$

sum $+ \text{???}$ $=$

$b_3$

Output

$\times -3.53$

$w_2$

$b_2$
$+ 0.57$

$\times \text{???}$

$w_4$

…and, in this example, that means we randomly select **2** values from a **Standard Normal Distribution** (mean = 0, standard deviation = 1).

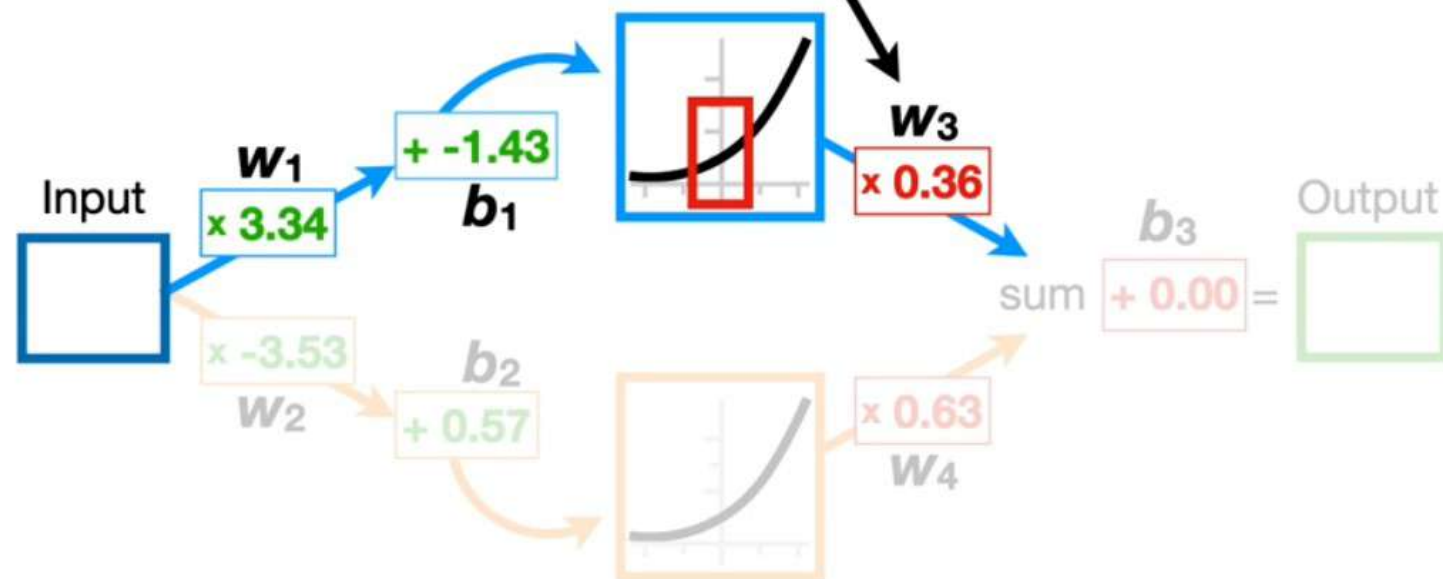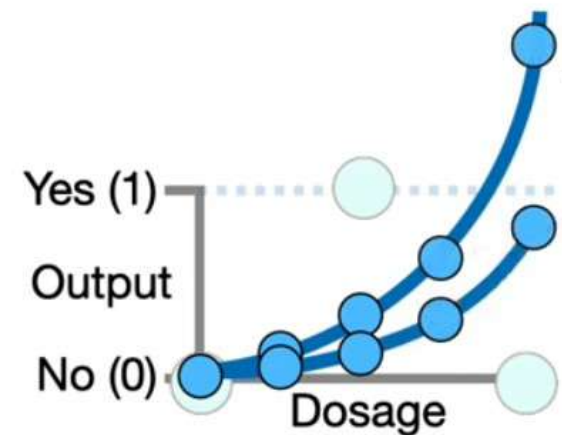Now we multiply the y-axis coordinates on the **blue curve** by $w_3$, which starts out with the random value **0.36**…

Yes (1)

Output

No (0)

Dosage

Input

$w_1$
$\times$ 3.34

$+$ -1.43
$b_1$

$w_3$
$\times$ 0.36

$b_3$
sum $+$ 0.00 $=$

Output
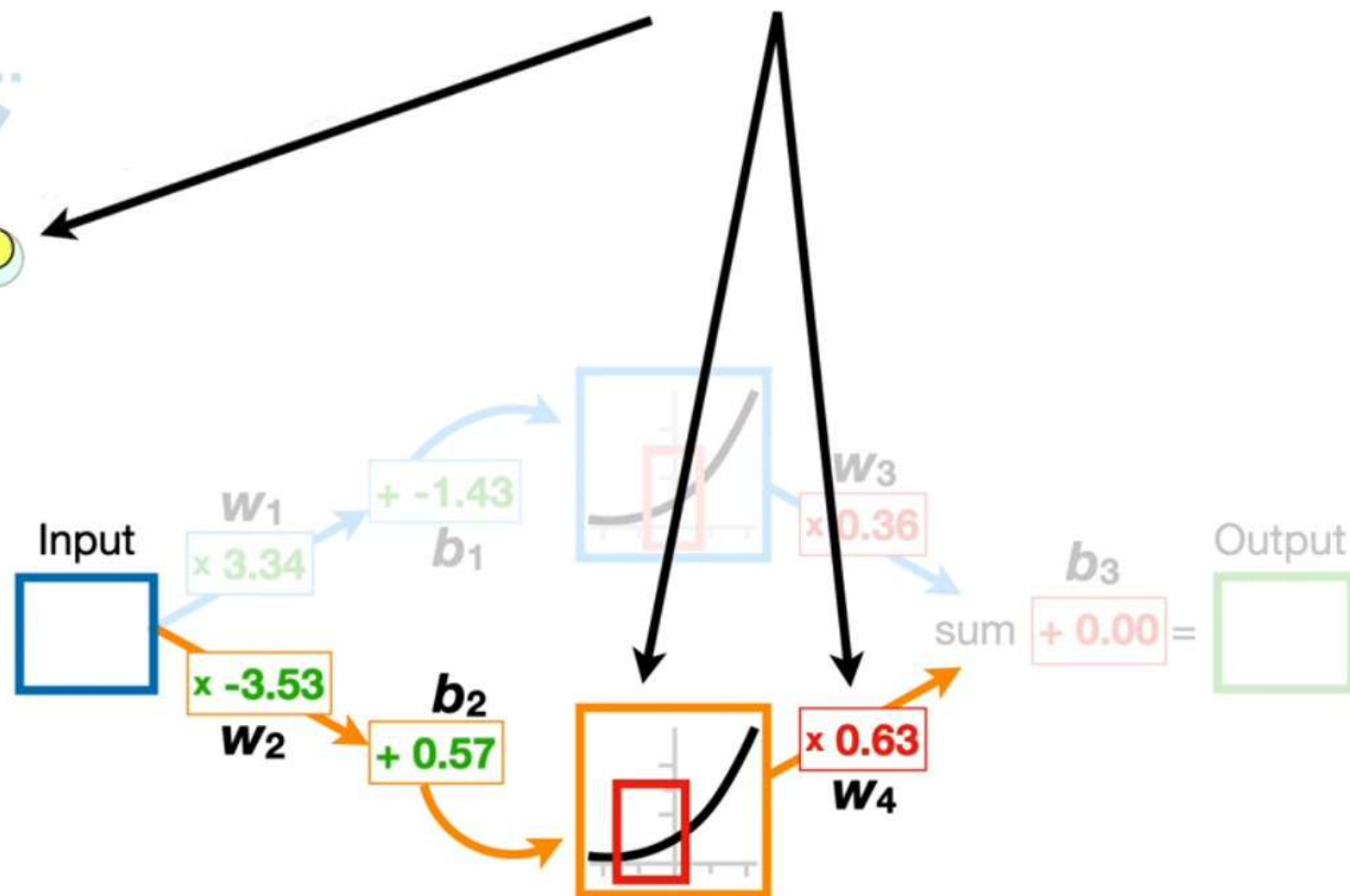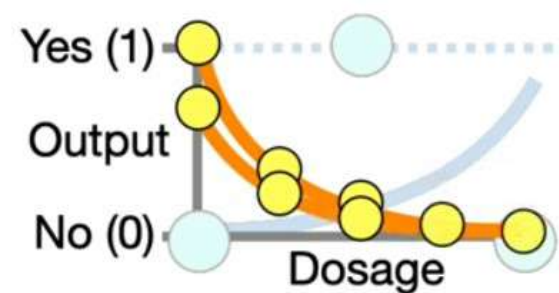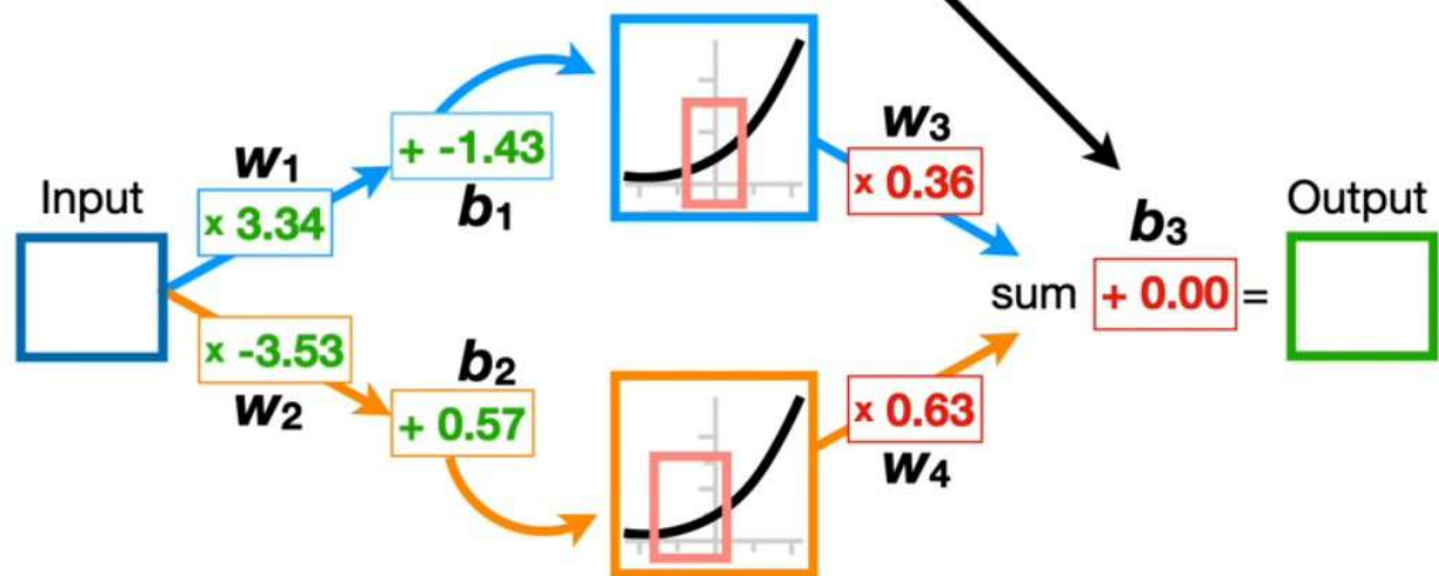
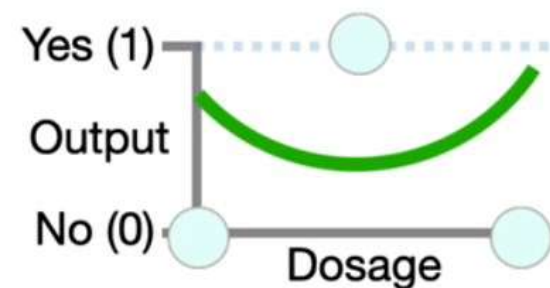$\times$ -3.53
$w_2$

$b_2$
$+$ 0.57

$\times$ 0.63
$w_4$

Now we multiply the y-axis coordinates on the **orange curve** by $w_4$, which starts with the random value **0.63**…

Lastly, since the initial value for $b_3$ is **0**,
adding it to the y-axis values on the **green
squiggle** does not change anything.

And just like before, if we change $b_3$ then we will change the **SSR**…

$SSR = (0 - -0.28)^2$
$+ (1 - -0.54)^2$
$+ (0 - -0.22)^2 = 2.5$

SSR

Bias$_3$ ($b_3$)

Yes (1)

Output

No (0)

Dosage

Input

$w_1$
$\times 3.34$
$+ -1.43$
$b_1$

$w_3$
$\times 0.36$

$\times -3.53$
$w_2$
$b_2$
$+ 0.57$

$\times 0.63$
$w_4$

$b_3$

sum $+ -1.0$ =

Output

Now let's talk about how to calculate
the derivatives of the **SSR** with
respect to the **Weights $w_3$ and $w_4$.**

Yes (1)

Output

No (0)

Dosage

Input

$w_1$
× 3.34

+ -1.43
$b_1$

$w_3$
× 0.36

$b_3$
sum + 0.00 =

Output

× -3.53
$w_2$

$b_2$
+ 0.57

× 0.63
$w_4$

And that means we can plug
$y_{1,i}$ times $w_3$ into the equation
for the **Predicted** values.

Yes (1)

Output

No (0)

Dosage

Predicted$_i$ = **green squiggle**$_i$ = $y_{1,i}w_3$ + **orange** + $b_3$

$y_{1,i}$

$w_3$

Input

$w_1$

× 3.34

+ -1.43

$b_1$

$x_{1,i}$

× 0.36

$b_3$

+ 0.00

sum

Output

=

× -3.53

$w_2$

+ 0.57

$b_2$

$y_{2,i}$

$x_{2,i}$

× 0.63

$w_4$

And that means we can plug
$y_{2,i}$ times $w_4$ into the equation
for the **Predicted** values.

Yes (1)

Output

No (0)

Dosage

$Predicted_i = $ **green squiggle**$_i = y_{1,i}w_3 + y_{2,i}w_4 + b_3$

$y_{1,i}$

Input

$w_1$
× 3.34

+ -1.43
$b_1$

× 0.36

$b_3$
+ 0.00 =

Output

sum

× -3.53
$w_2$

$b_2$
+ 0.57

$y_{2,i}$

× 0.63
$w_4$

$x_{2,i}$

...then the **SSR** are linked to $w_3$ and $w_4$...

$$\boxed{\text{SSR}} = \sum_{i=1}^{n=3} (\text{Observed}_i - \text{Predicted}_i)^2$$

Yes (1)

Output

No (0)

Dosage

$\text{Predicted}_i = \textbf{\textcolor{green}{green squiggle}}_i = \textcolor{blue}{y_1,} \boxed{\textcolor{blue}{w_3}} + \textcolor{orange}{y_2,} \boxed{\textcolor{orange}{w_4}} + \boldsymbol{b_3}$

Input

$w_1$
× 3.34

+ -1.43
$b_1$

$y_{1,i}$
$x_{1,i}$

$w_3$
× 0.36

$b_3$
sum + 0.00 =

Output

× -3.53
$w_2$

$b_2$
+ 0.57

$y_{2,i}$
$x_{2,i}$

× 0.63
$w_4$

$$\frac{d\ SSR}{d\ w_3} = \frac{d\ SSR}{d\ Predicted} \times \frac{d\ Predicted}{d\ w_3}$$

$$\frac{d\ SSR}{d\ w_4} = \frac{d\ SSR}{d\ Predicted} \times \frac{d\ Predicted}{d\ w_4}$$

...times the derivative of the **Predicted** values with respect to $w_4$.

$$SSR = \sum_{i=1}^{n=3} (Observed_i - Predicted_i)^2$$
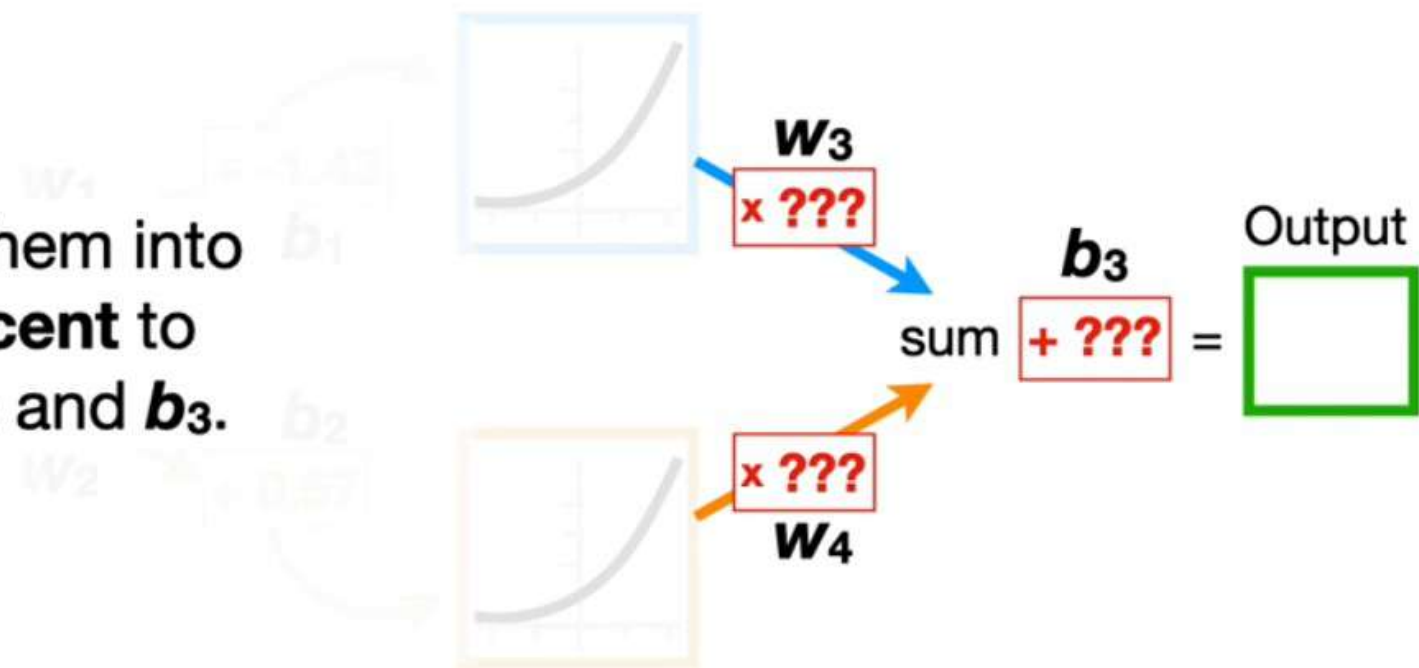
$$Predicted_i = \text{green squiggle}_i = y_{1,i}\,w_3 + y_{2,i}\,w_4 + b_3$$

$$\frac{d \ SSR}{d \ w_3} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times y_{1,i}$$

$$\frac{d \ SSR}{d \ w_4} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times y_{2,i}$$

$$\frac{d \ SSR}{d \ b_3} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times 1$$

…we can plug them into **Gradient Descent** to optimize $w_3$, $w_4$ and $b_3$.

$$\frac{d\ SSR}{d\ w_3} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times y_{1,i}$$

$$\frac{d\ SSR}{d\ w_4} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times y_{2,i}$$

$$\frac{d\ SSR}{d\ b_3} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times 1$$

Now we repeat that process until the **Predictions** no longer improve very much, or we reach a maximum number of steps or we meet some other criteria.