

Introduction to Neural Networks

By Imène ALLOUCHE

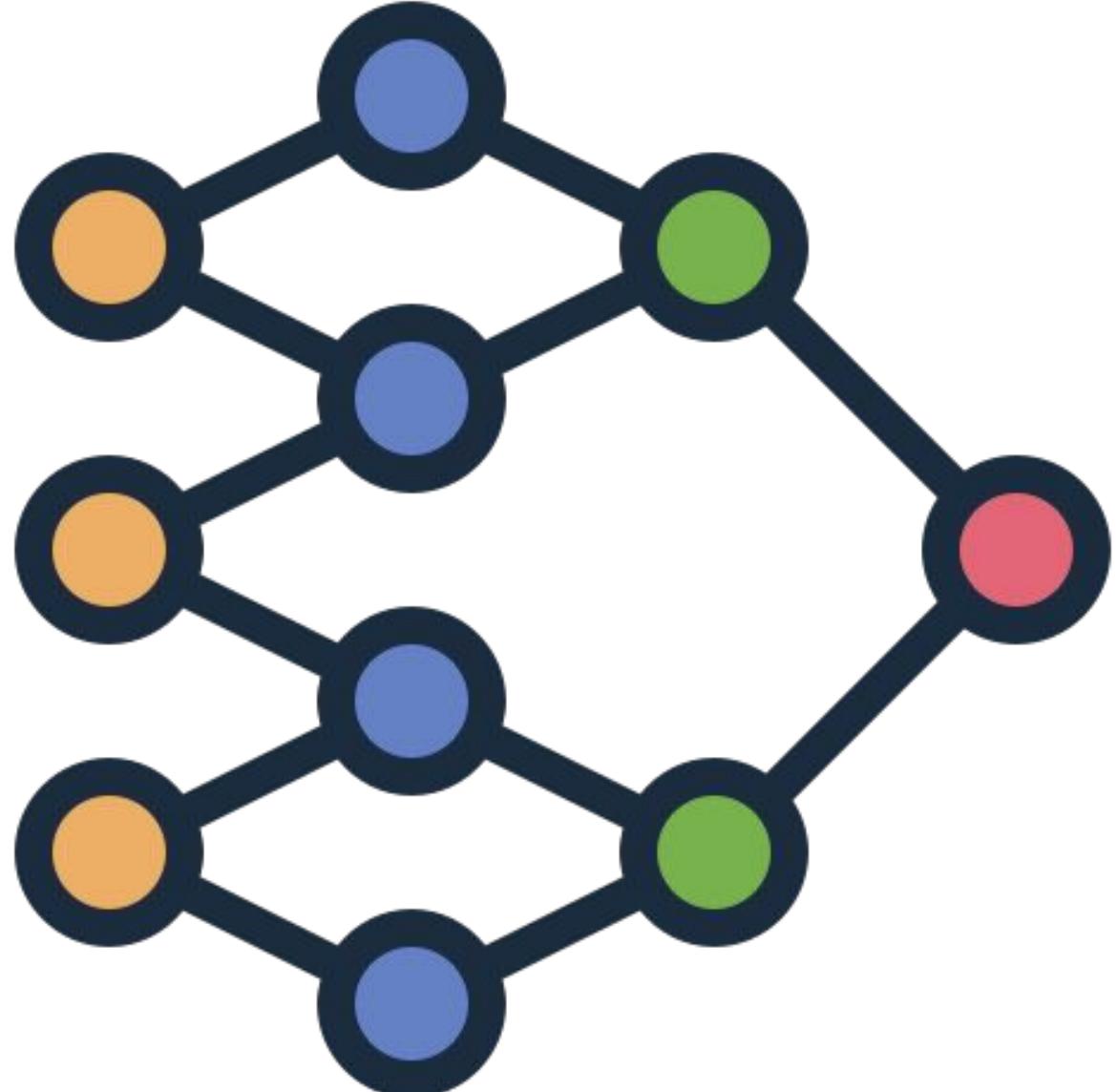


TABLE OF CONTENTS

01

Why and How?

02

Gradient Descent

03

Backpropagation

04

Convolutional Neural Network

05

Recurrent Neural Networks

06

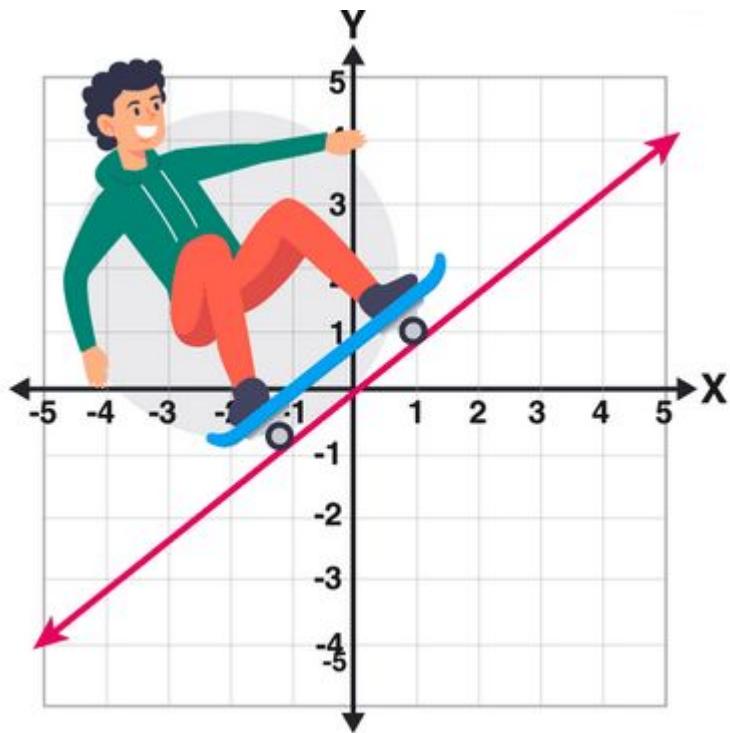
Conclusion

01

Why Neural Network?

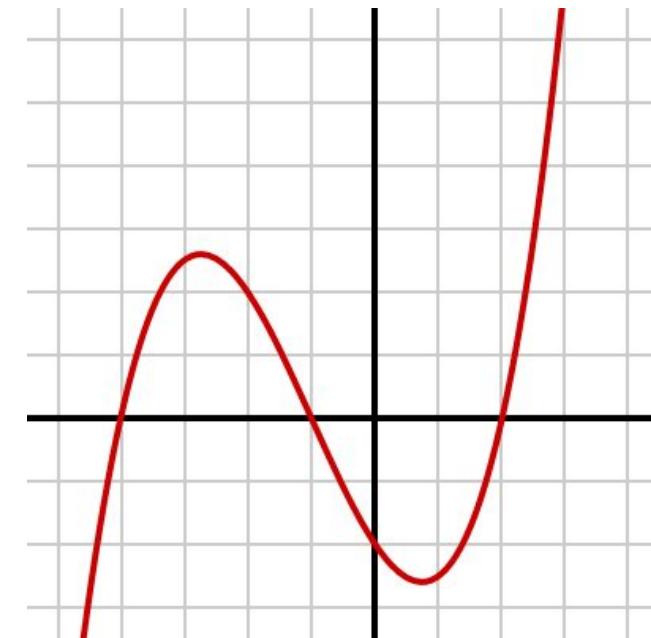
Why Neural Network?

Linear Relations



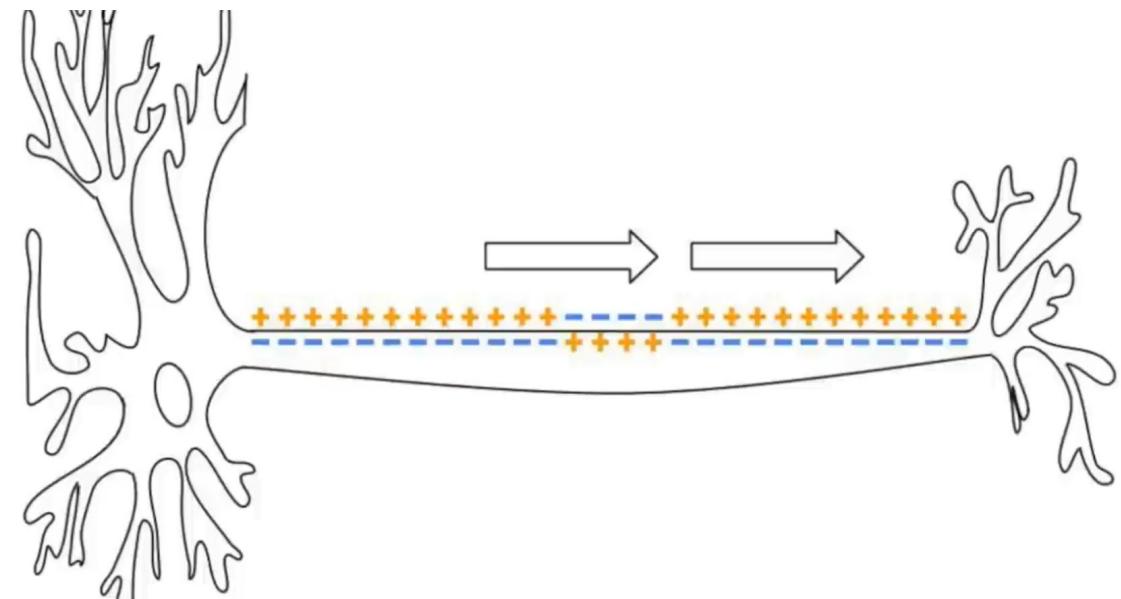
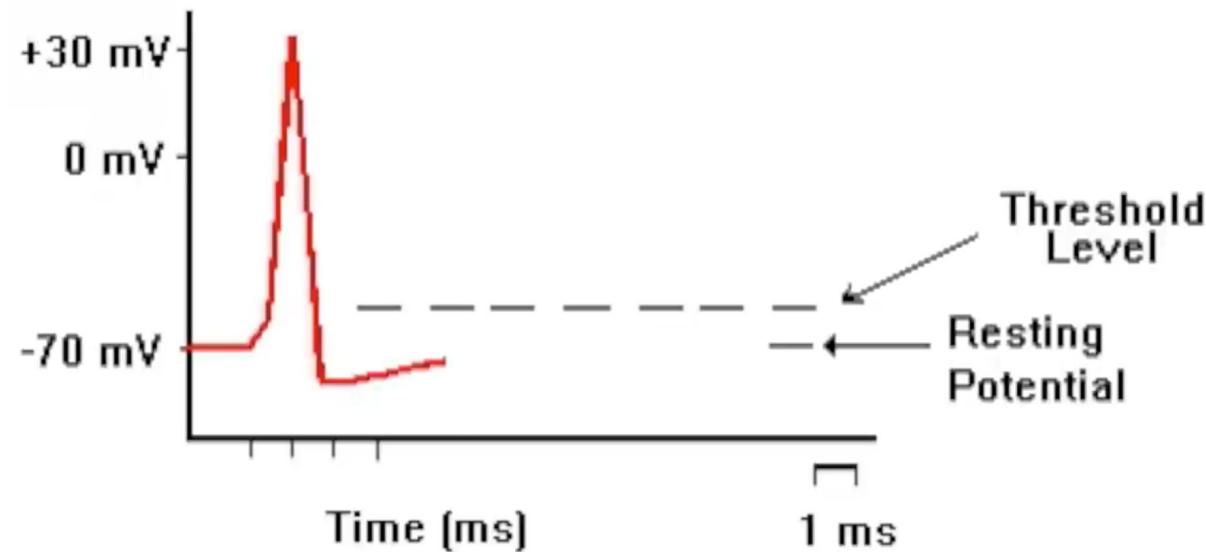
$$y = ax + b$$

Non Linear Relations



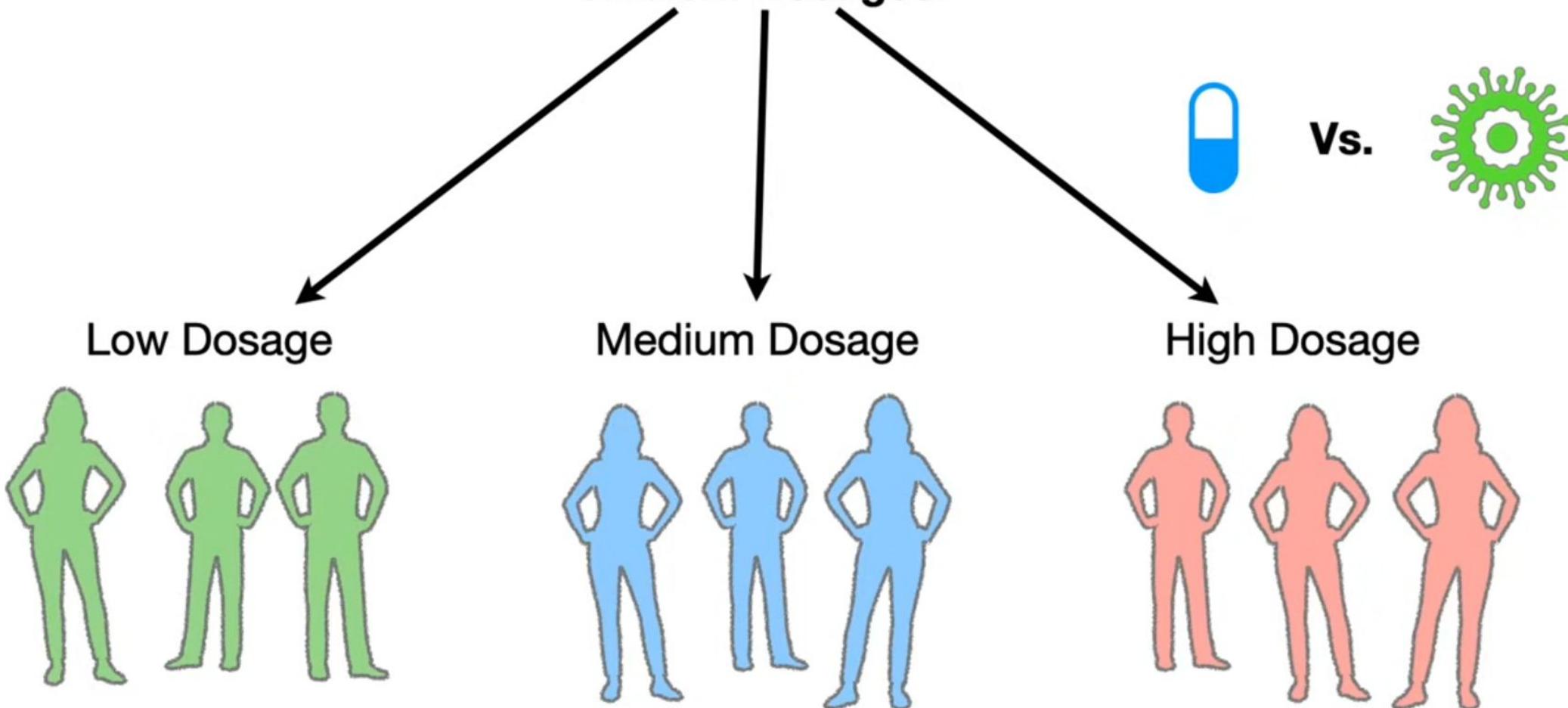
$$y = ?$$

Why Neural Network?

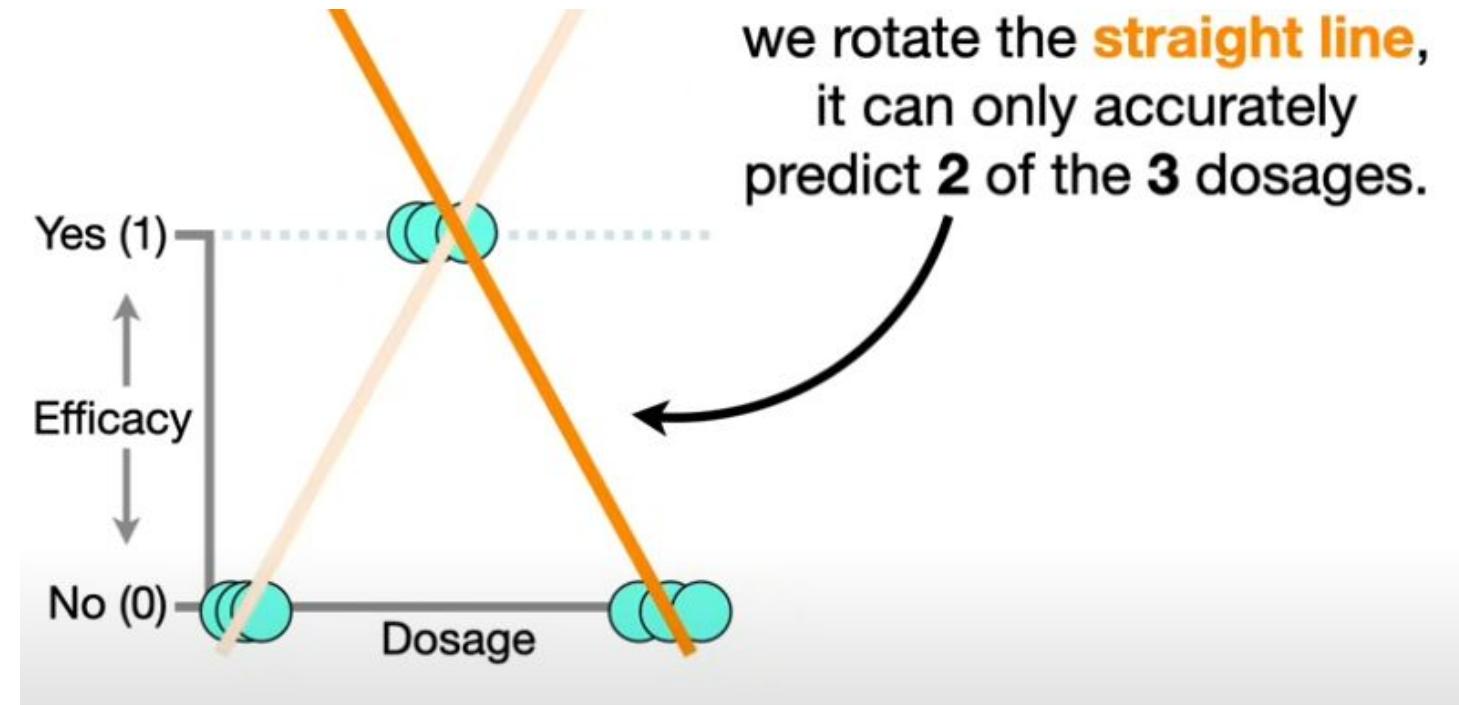
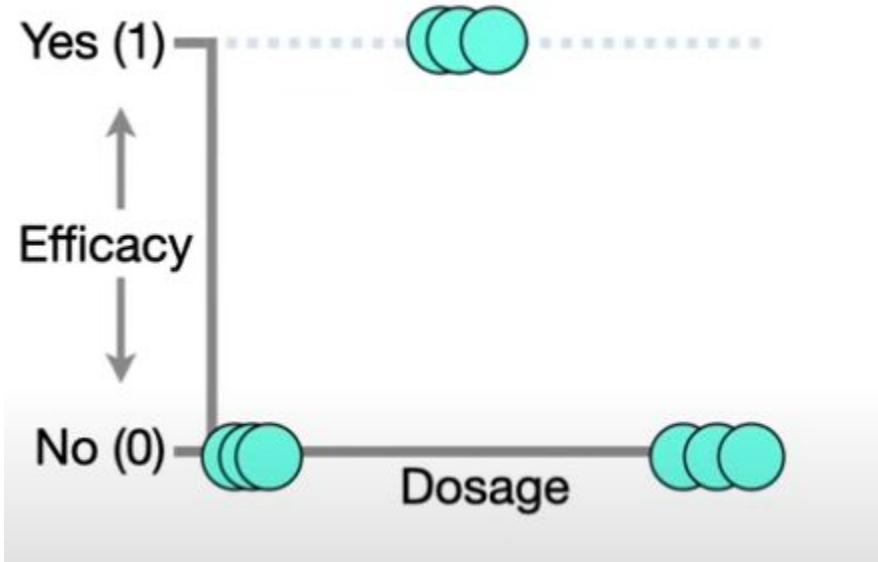


How Neural Network?

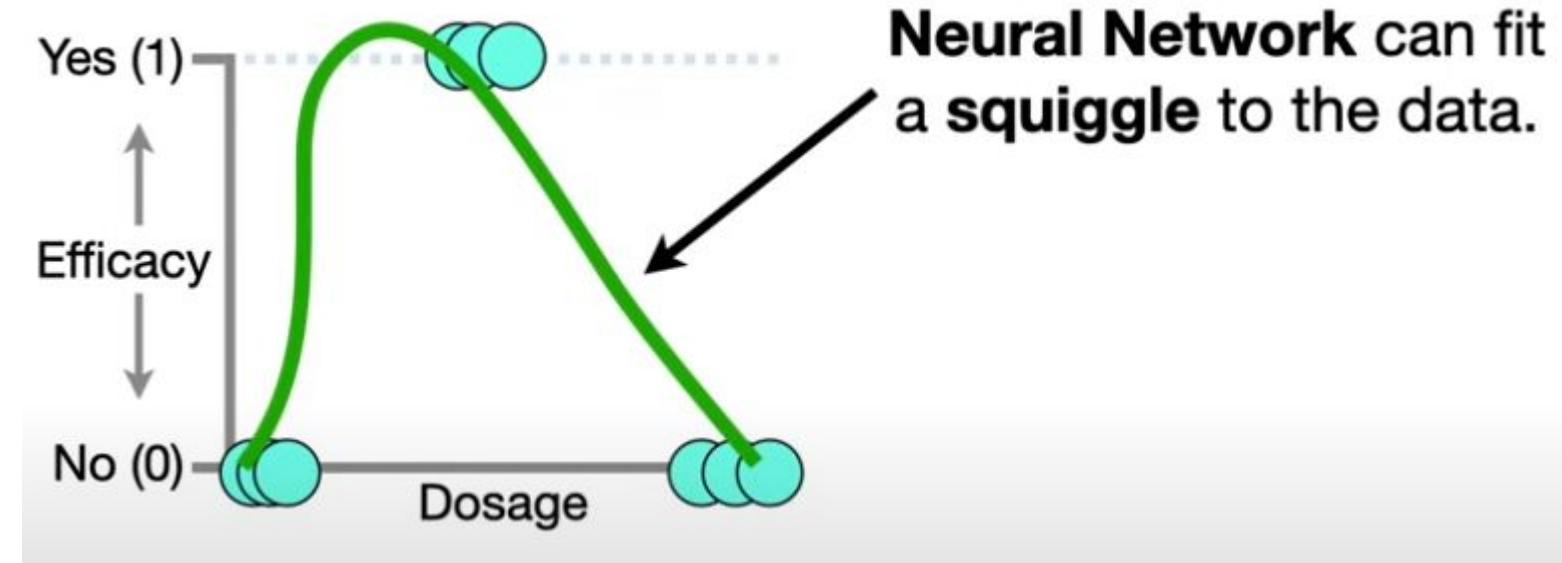
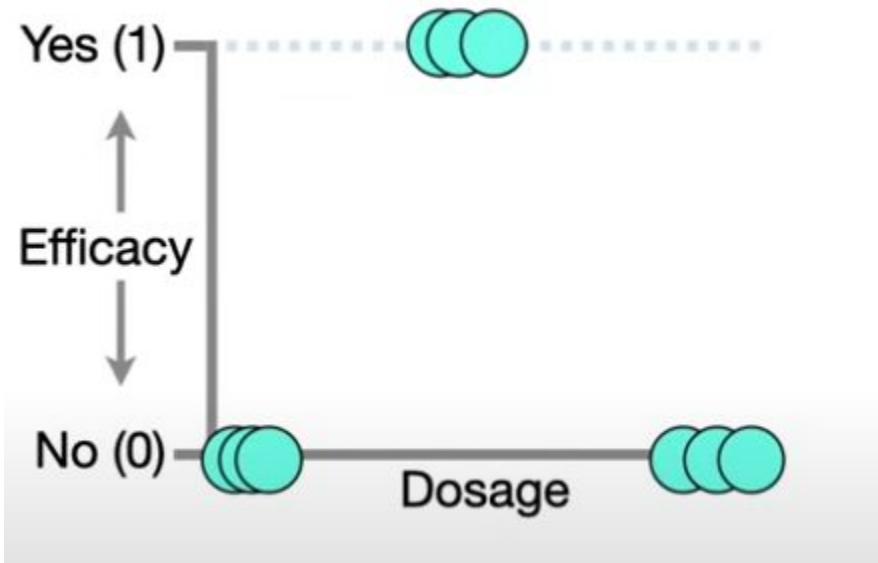
...and we gave the drug to 3 different groups of people with 3 different **Dosages**.



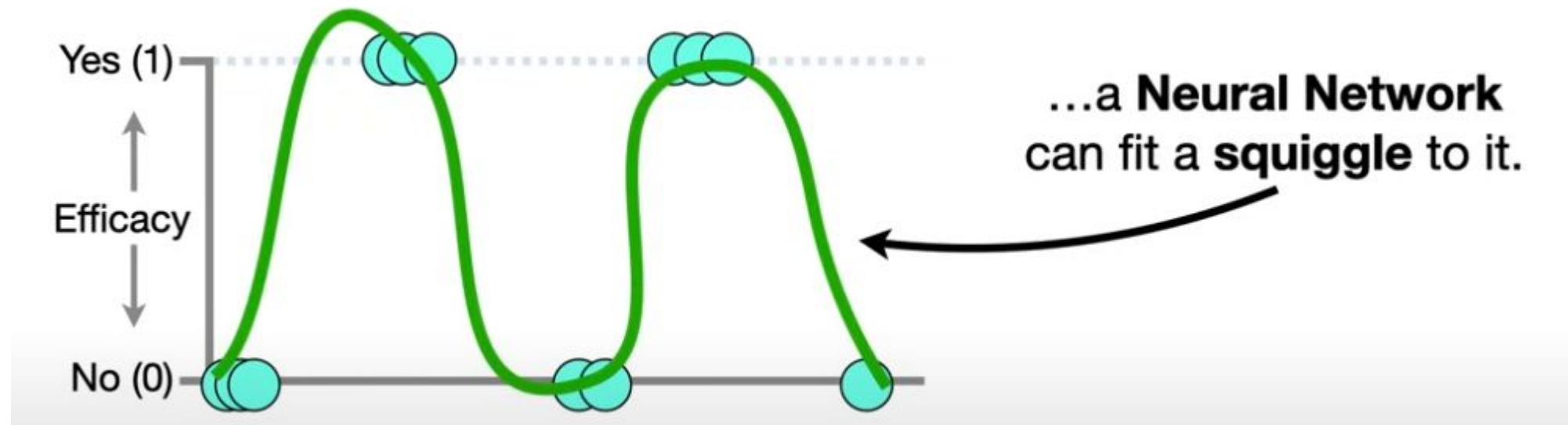
How Neural Network?



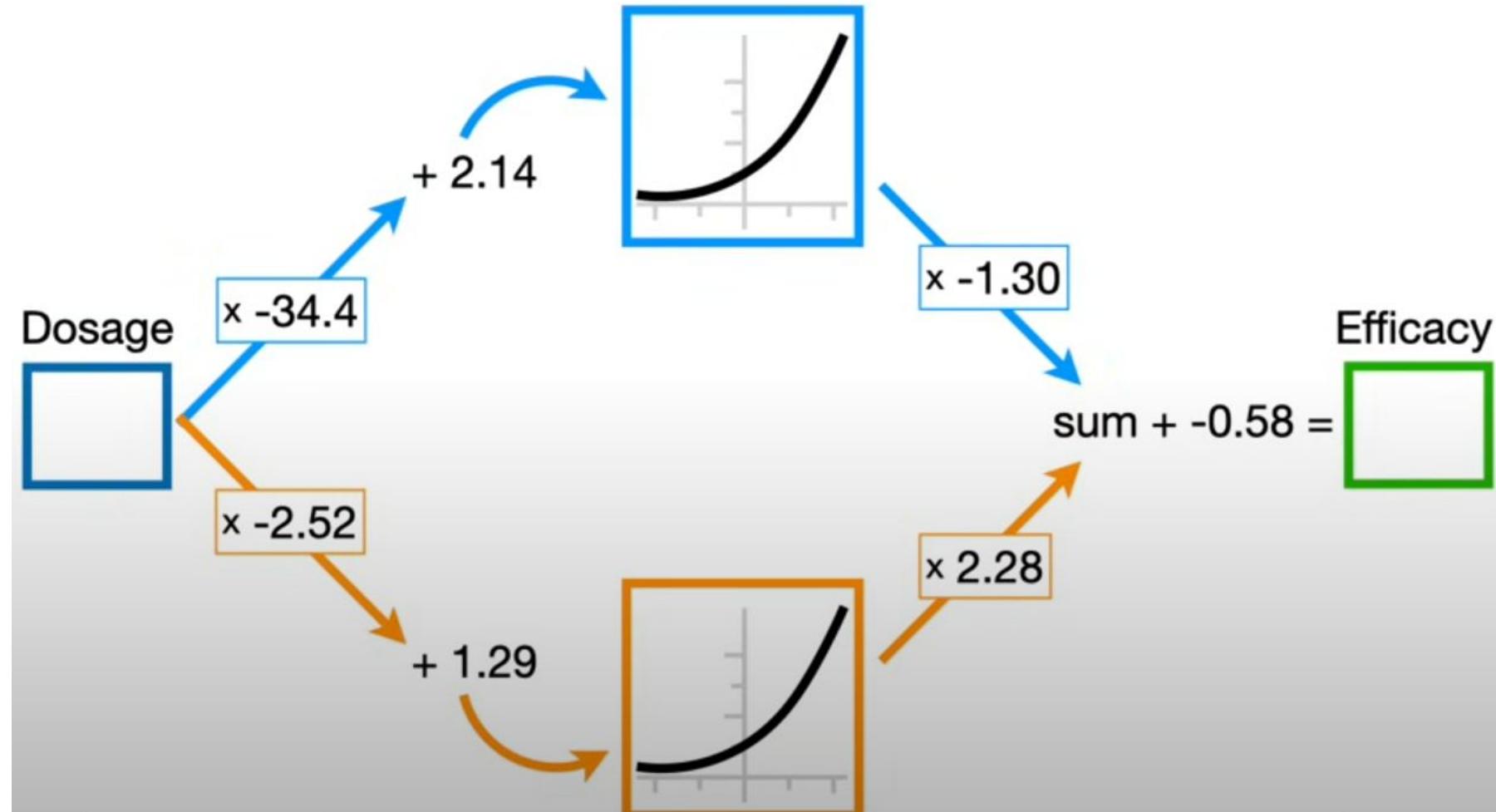
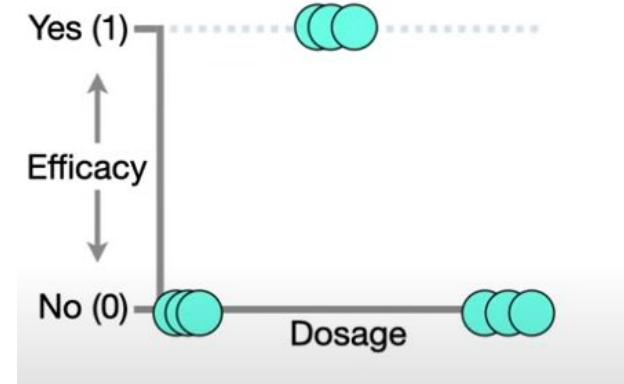
How Neural Network?



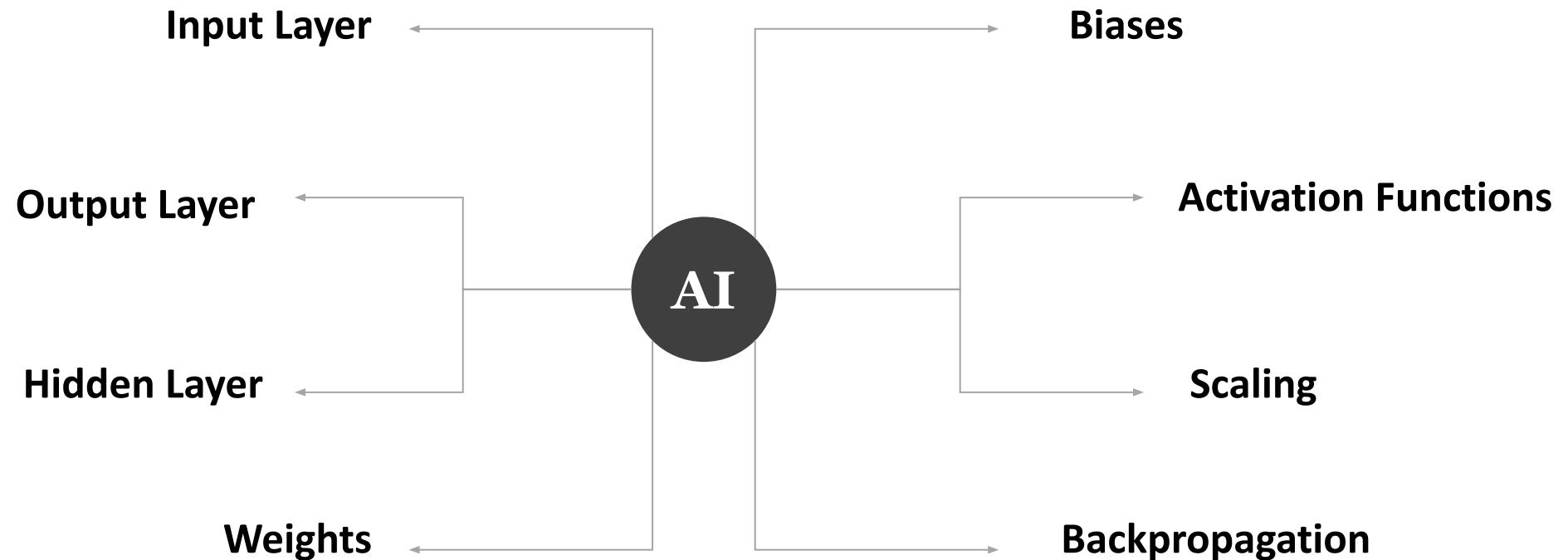
How Neural Network?



How Neural Network?



Terminology



The Math Behind The Hidden Layer

01 Input (Dosages in out example)

02 $X_1 = w_1 * \text{Input} + b_1$

03 $Y_1 = f_1(X_1)$

04 $Y_1 = Y_1 * w_3$

05 $Y = Y_1 + Y_2 + b_3$

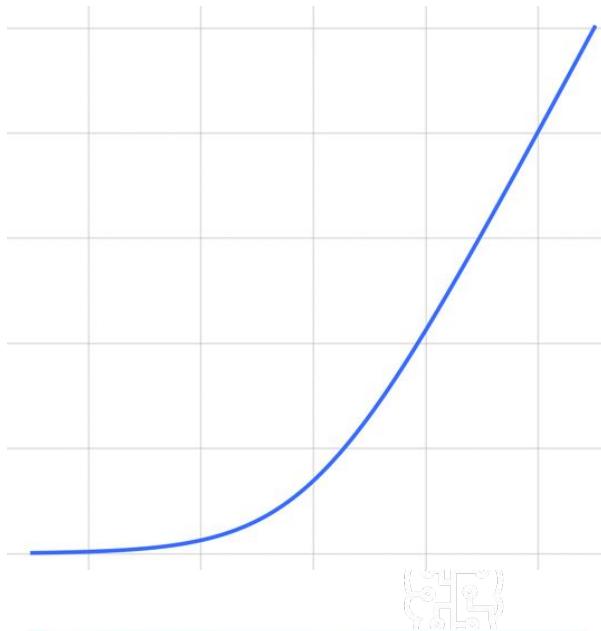
02 $X_2 = w_2 * \text{Input} + b_2$

03 $Y_2 = f_2(X_2)$

04 $Y_2 = Y_2 * w_4$

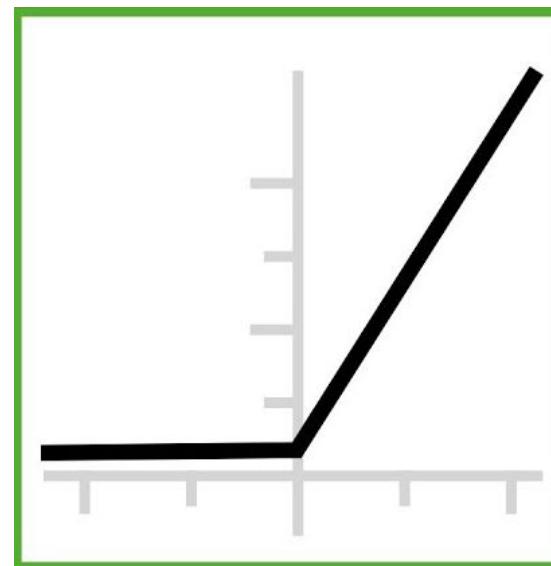
Artificial intelligence development phase

Softplus



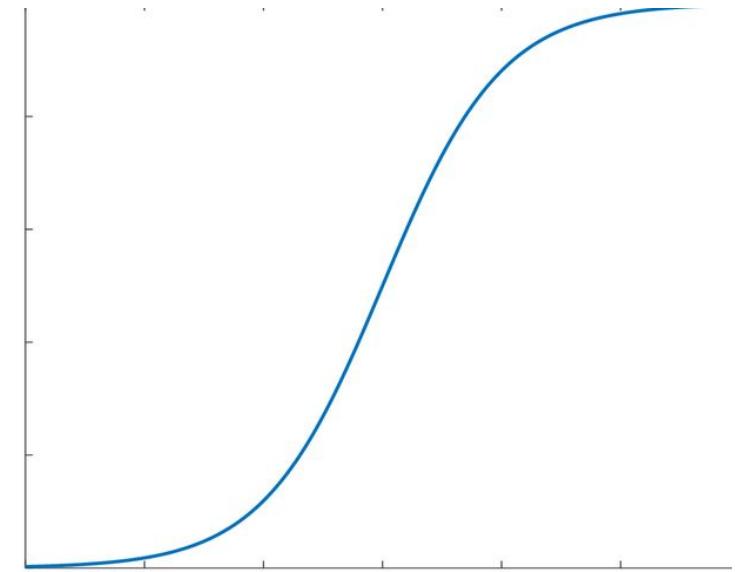
$$f(x) = \log(1 + e^x)$$

ReLU



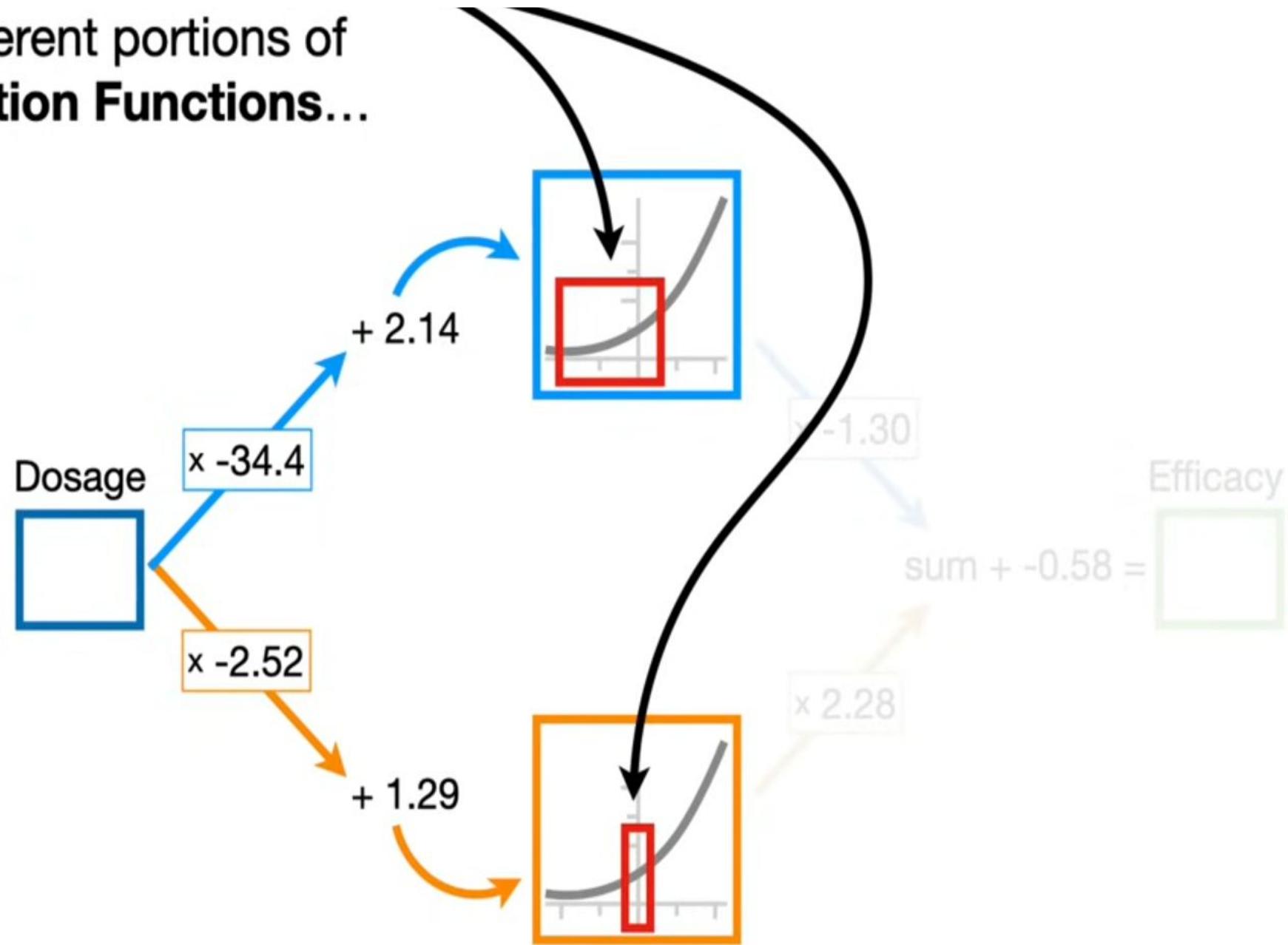
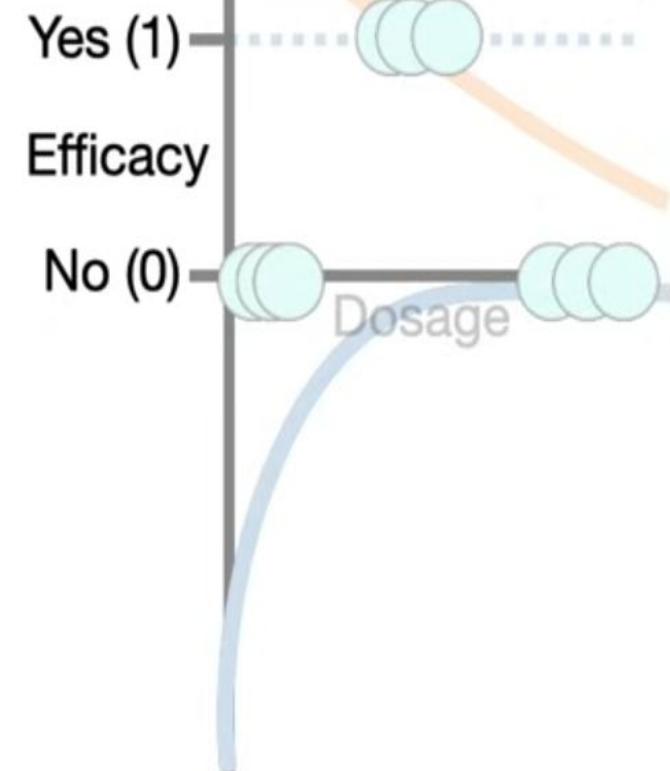
$$f(x) = \max(0, x)$$

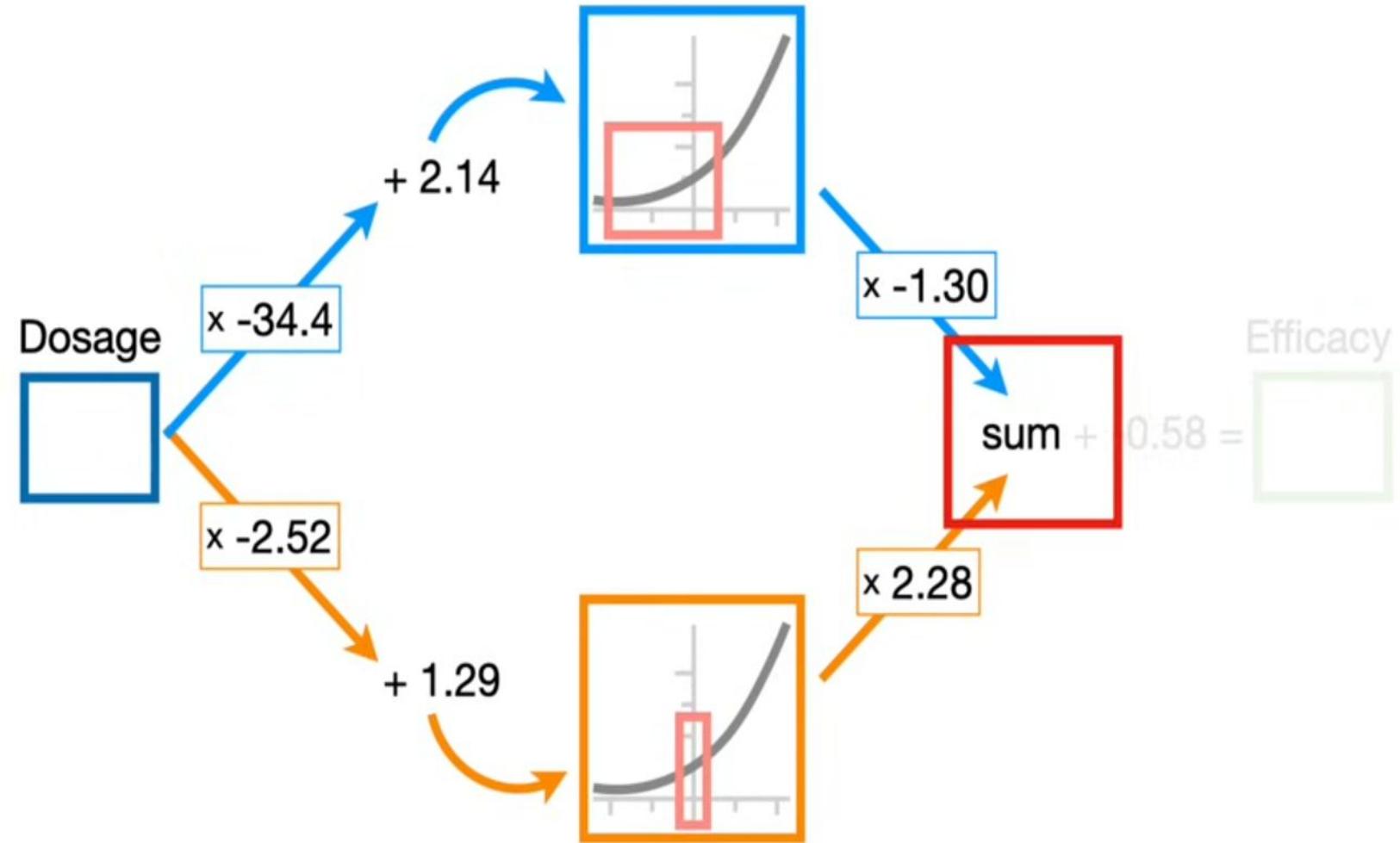
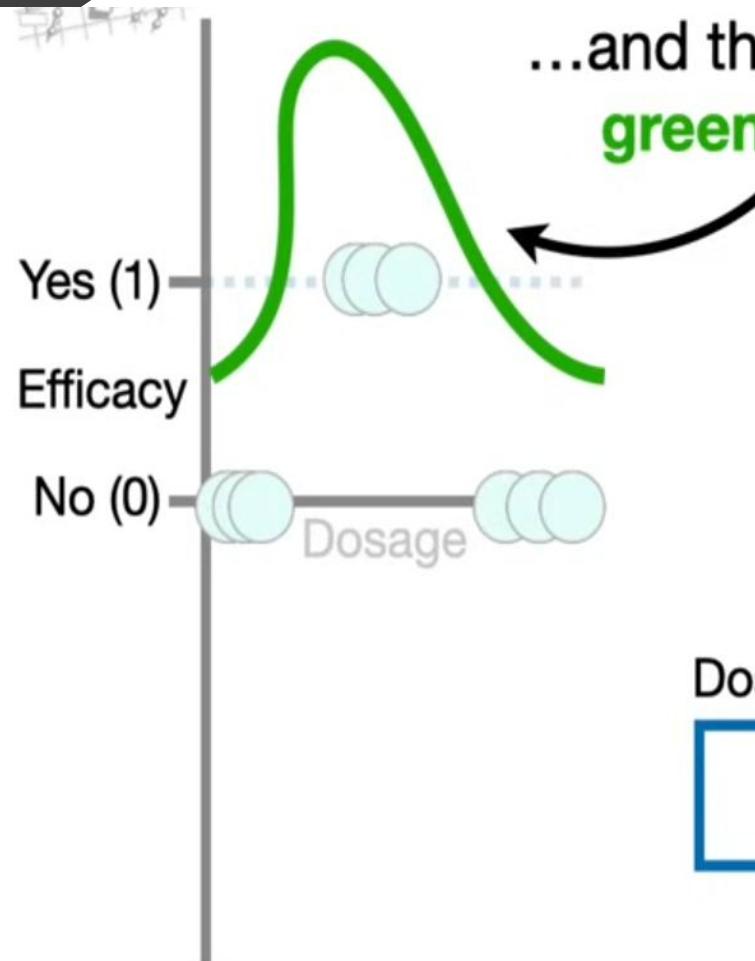
Sigmoid



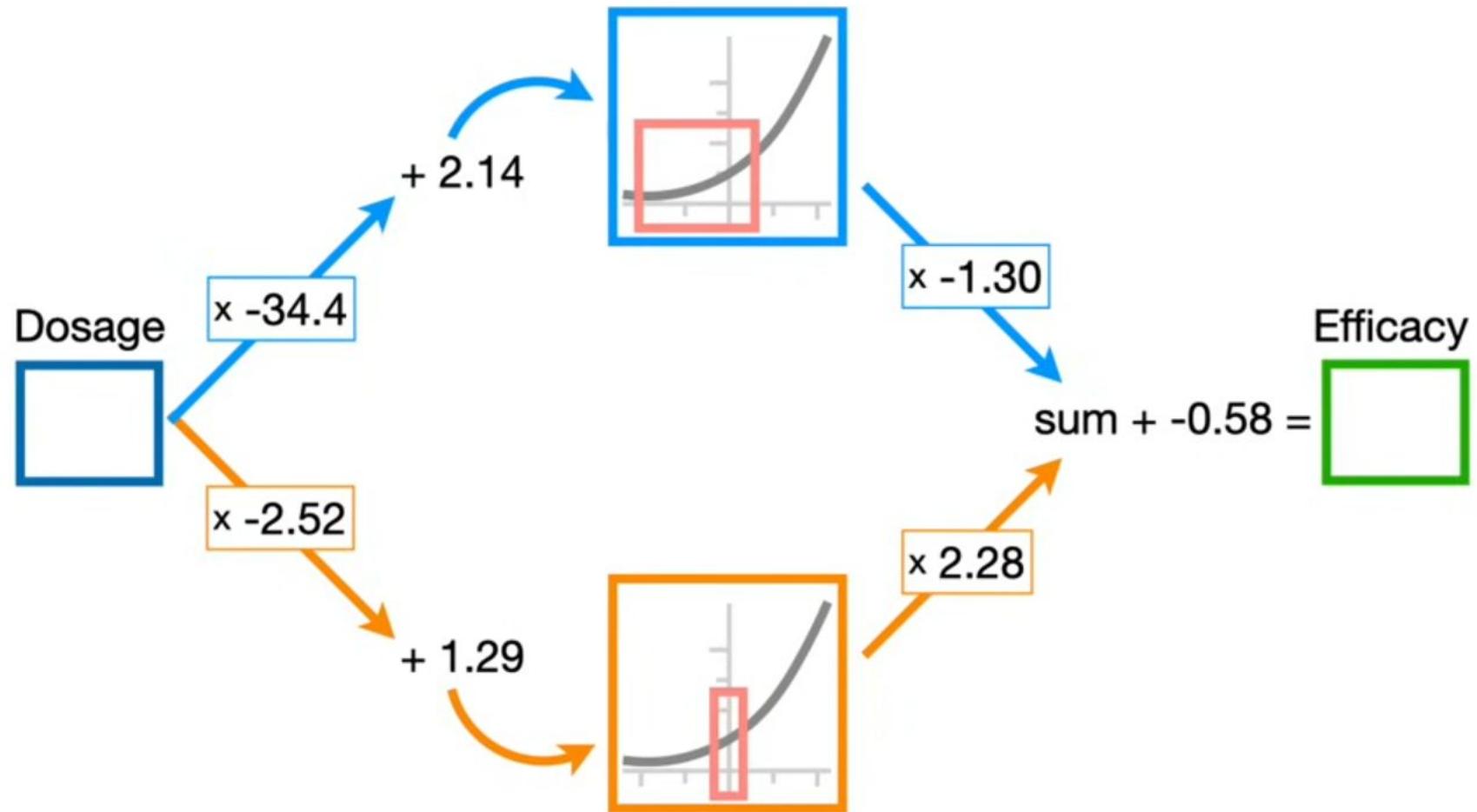
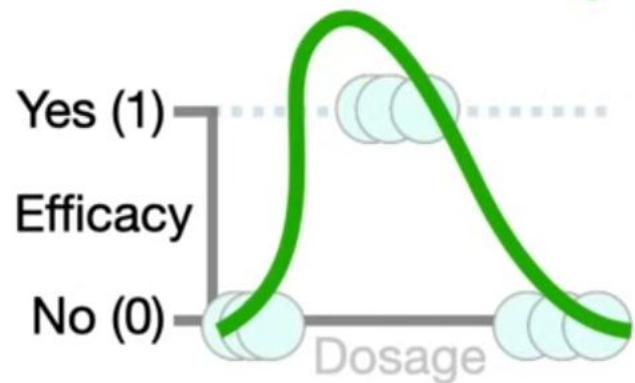
$$f(x) = \frac{e^x}{e^x + 1}$$

using different portions of the **Activation Functions**...





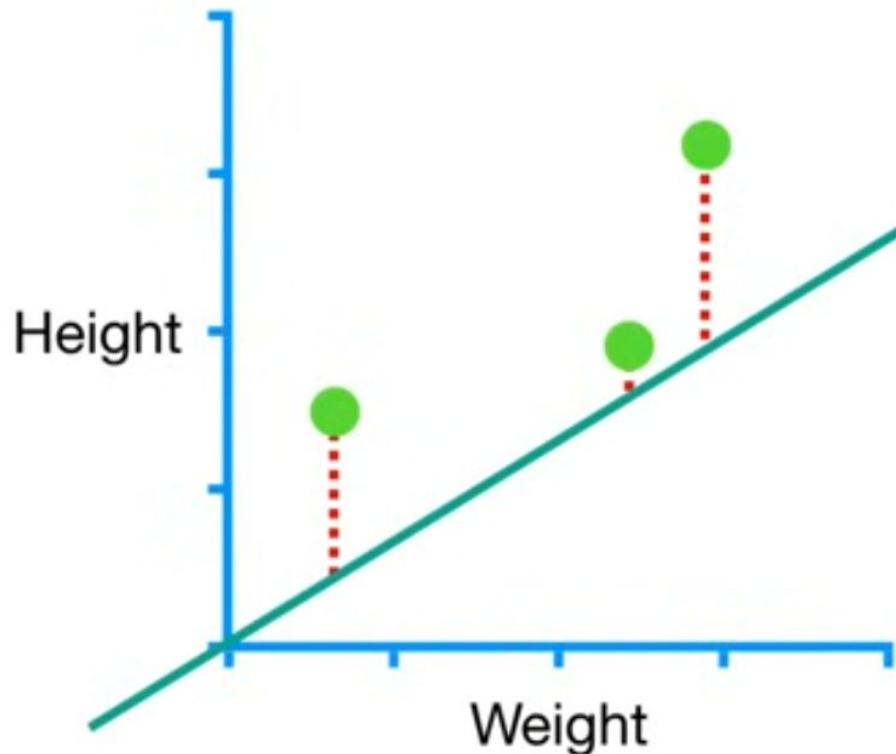
squiggle that fits the data.



02

Gradient Descent

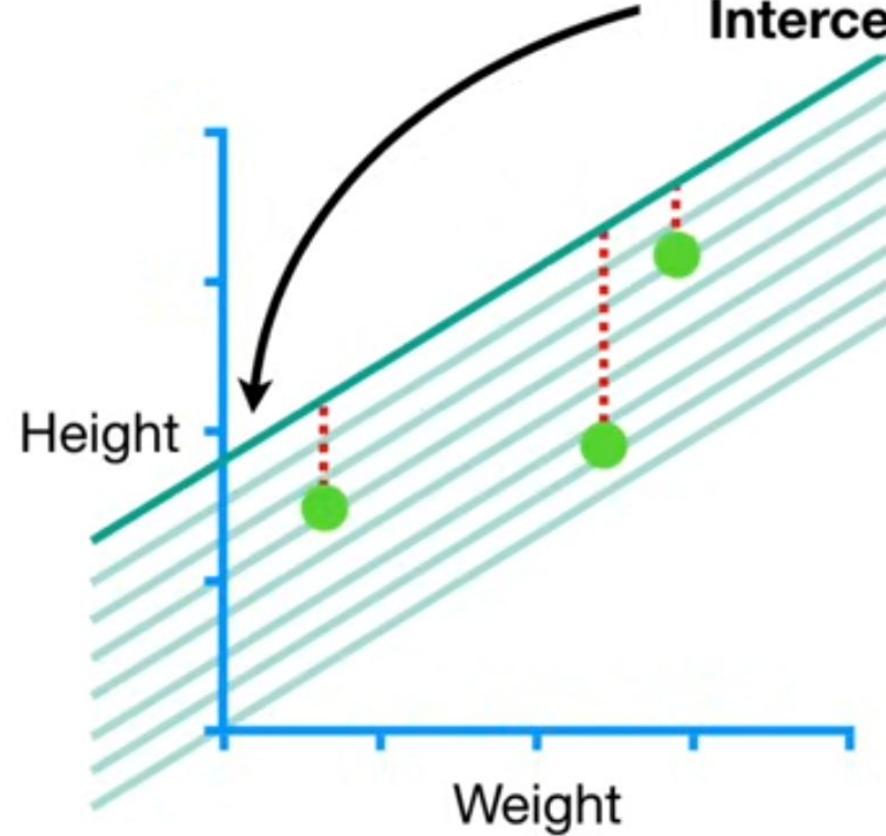
Gradient Descent



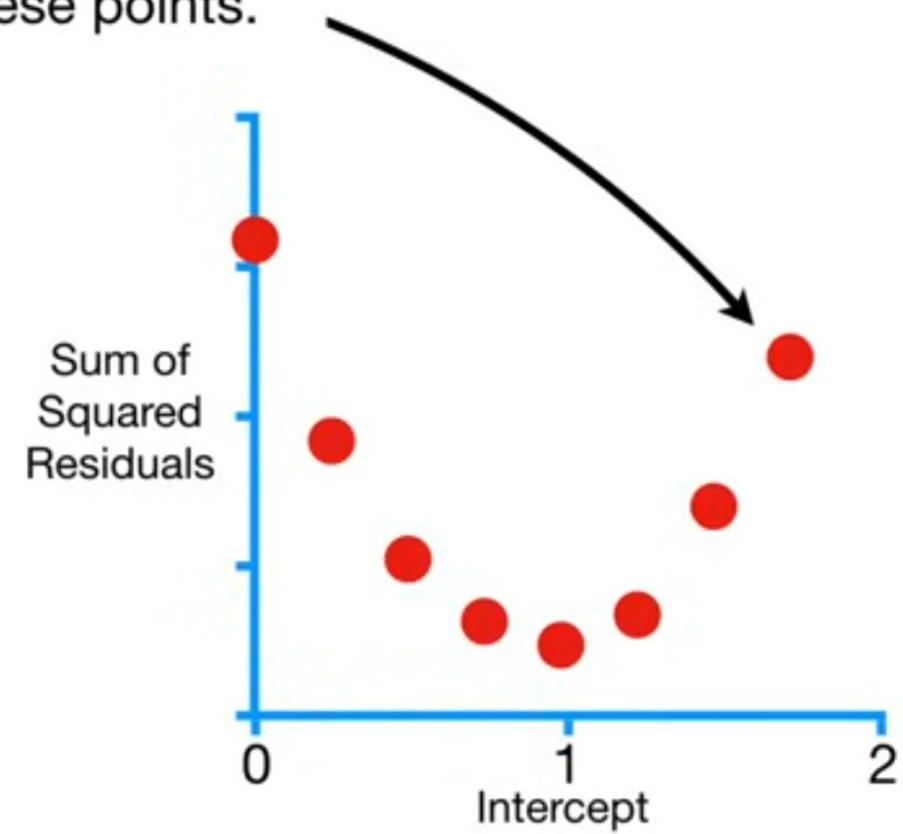
$$\text{Sum of squared residuals} = 1.1^2 + 0.4^2 + 1.3^2 = \boxed{3.1}$$

In the end, **3.1** is the Sum of the Squared Residuals.

Gradient Descent



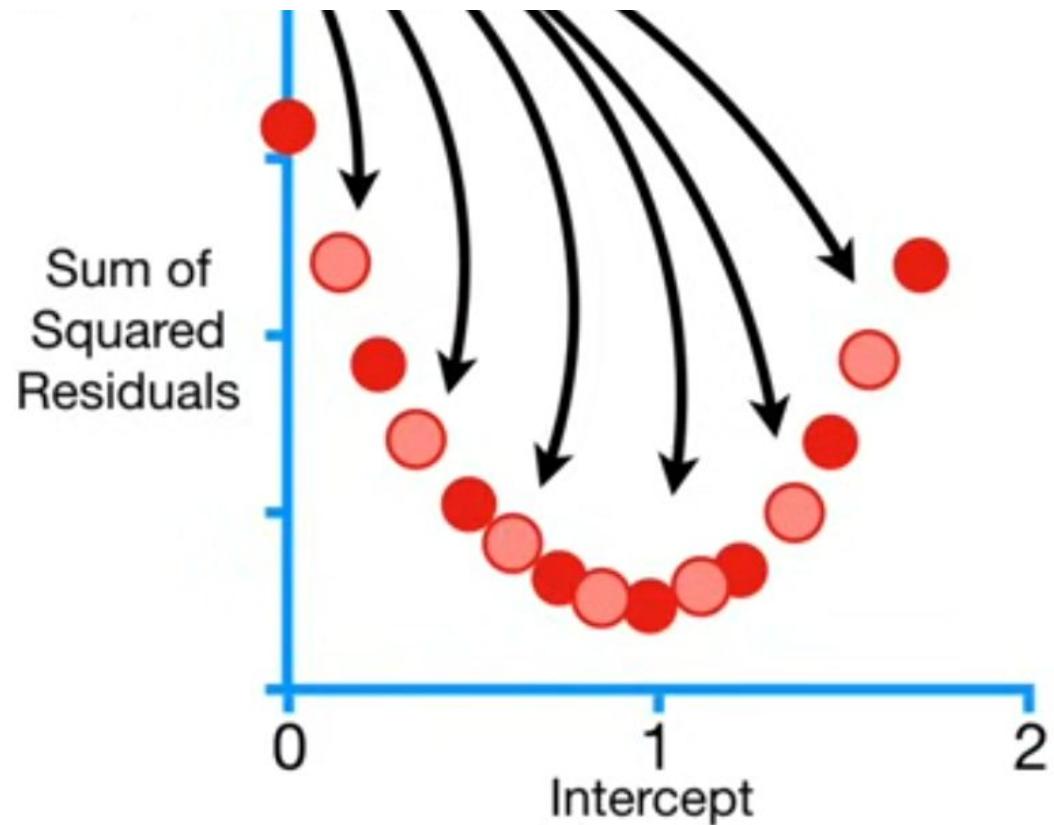
And for increasing values for the **Intercept**, we get these points.



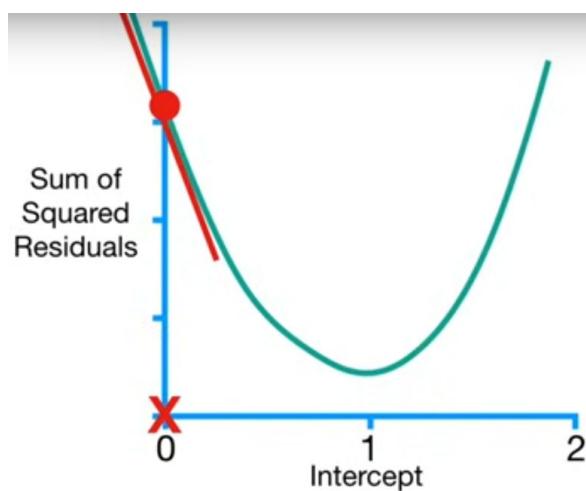
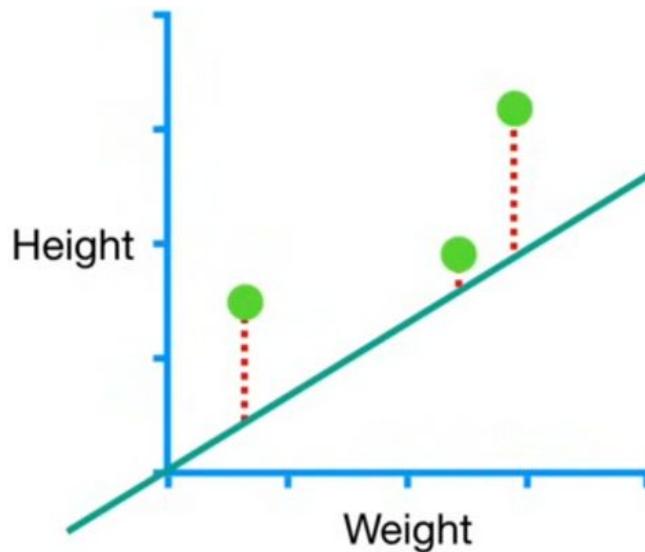
Gradient Descent

A slow and painful method for finding minimal sum of the squared Residuals is to plug a bunch of values of the Intercept

$$Y = \text{Slope} * x + \text{Intercept}$$



Gradient Descent



$$\begin{aligned} \text{Sum of squared residuals} = & (1.4 - (\text{intercept} + 0.64 \times 0.5))^2 \\ & + (1.9 - (\text{intercept} + 0.64 \times 2.3))^2 \\ & + (3.2 - (\text{intercept} + 0.64 \times 2.9))^2 \end{aligned}$$

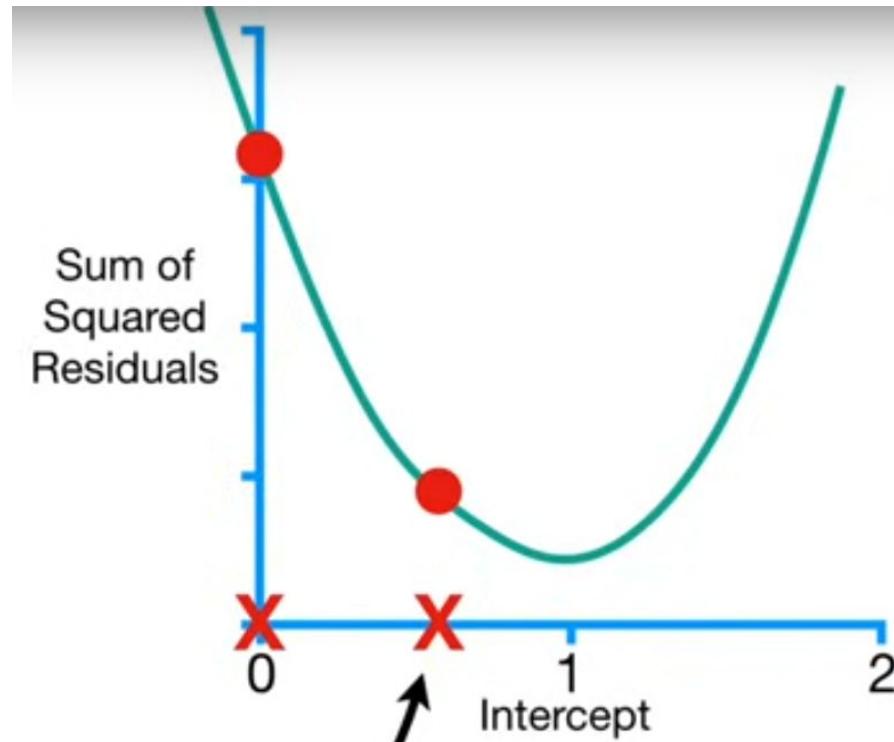
$$\begin{aligned} \frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} = & -2(1.4 - (0 + 0.64 \times 0.5)) \\ & + -2(1.9 - (0 + 0.64 \times 2.3)) \\ & + -2(3.2 - (0 + 0.64 \times 2.9)) \\ = & -5.7 \end{aligned}$$

Step Size = $-5.7 \times 0.1 = -0.57$

Step Size = $-5.7 * \text{Learning Rate}$

When the **Intercept** = 0, the
Step Size = **-0.57**.

Gradient Descent



New Intercept = Old Intercept - Step Size

$$\text{New Intercept} = 0 - (-0.57) = \boxed{0.57}$$

...and the the **New Intercept = 0.57.**



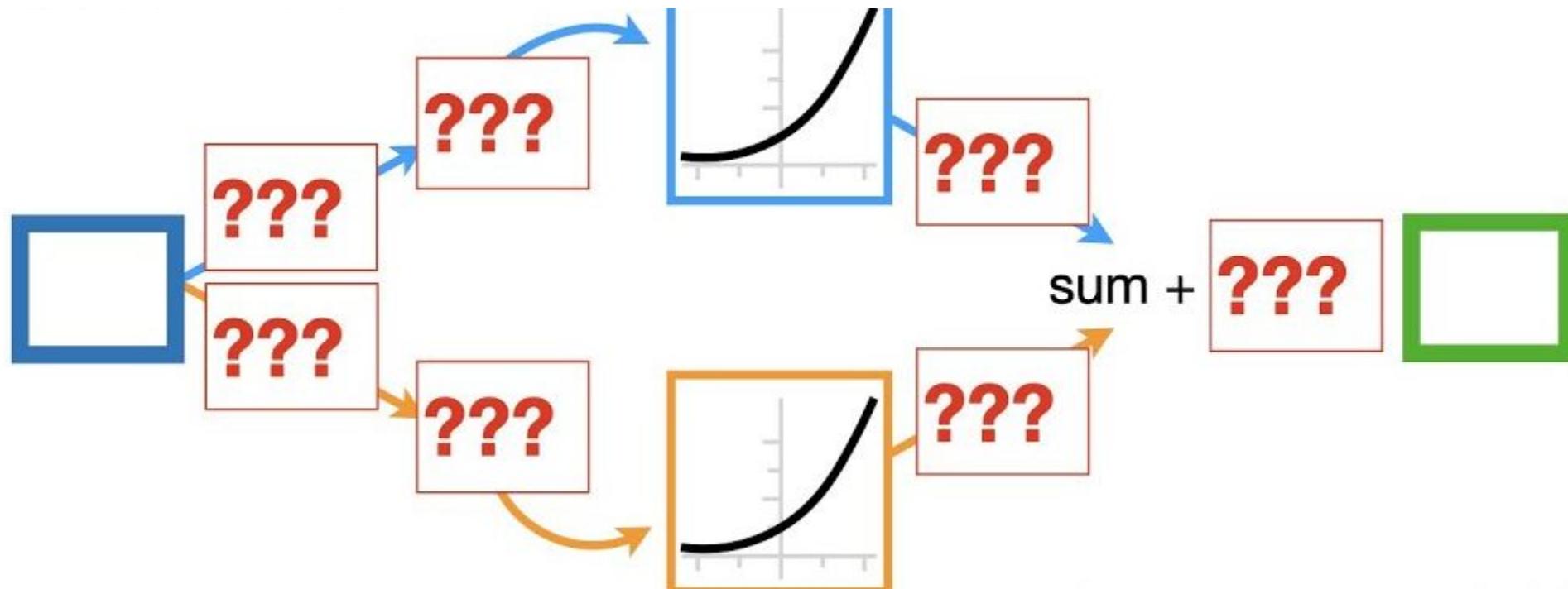
When do we stop?

- when reaching an SSR ≤ 0.001
- Could be changed by the developer
- When reaching a maximum number of steps is reached
- To be fixed by the developer

03

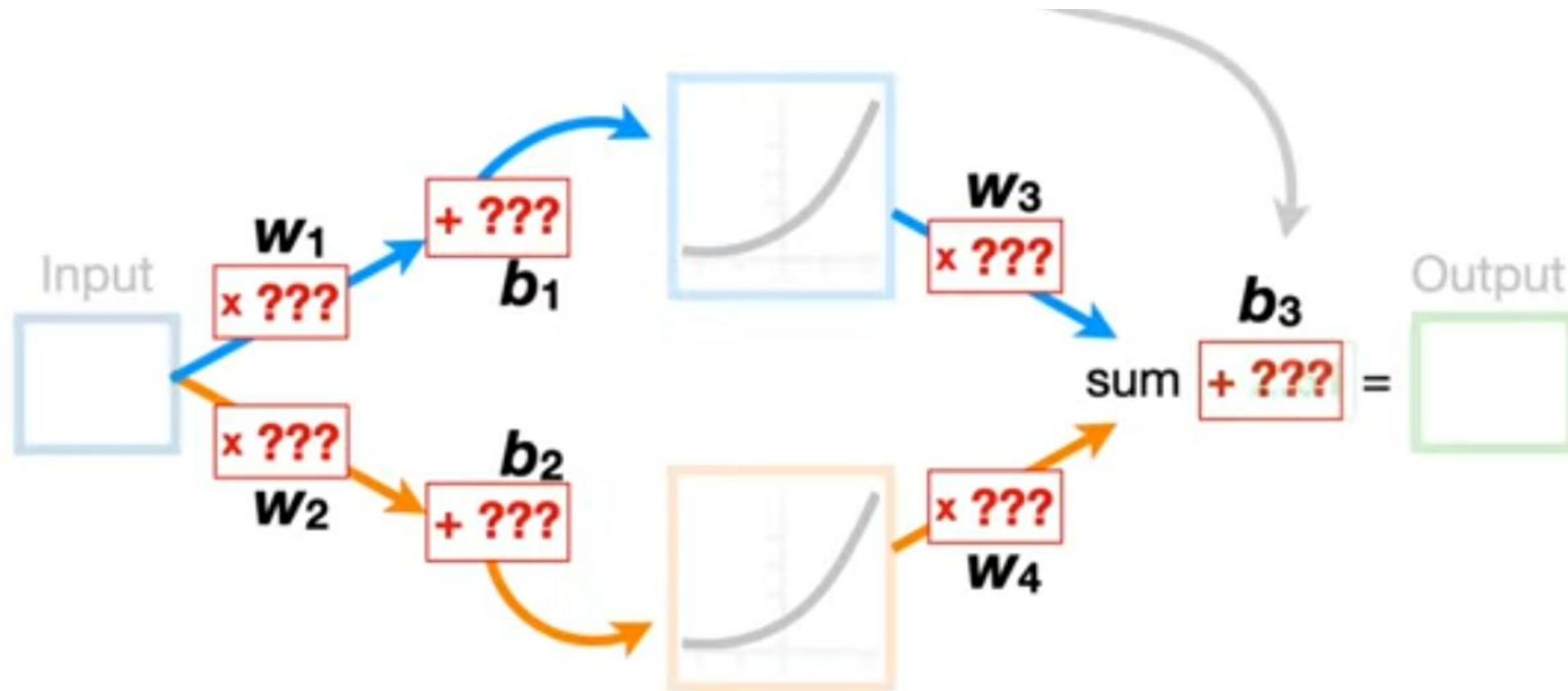
Backpropagation

How to find the best Weights and biases?



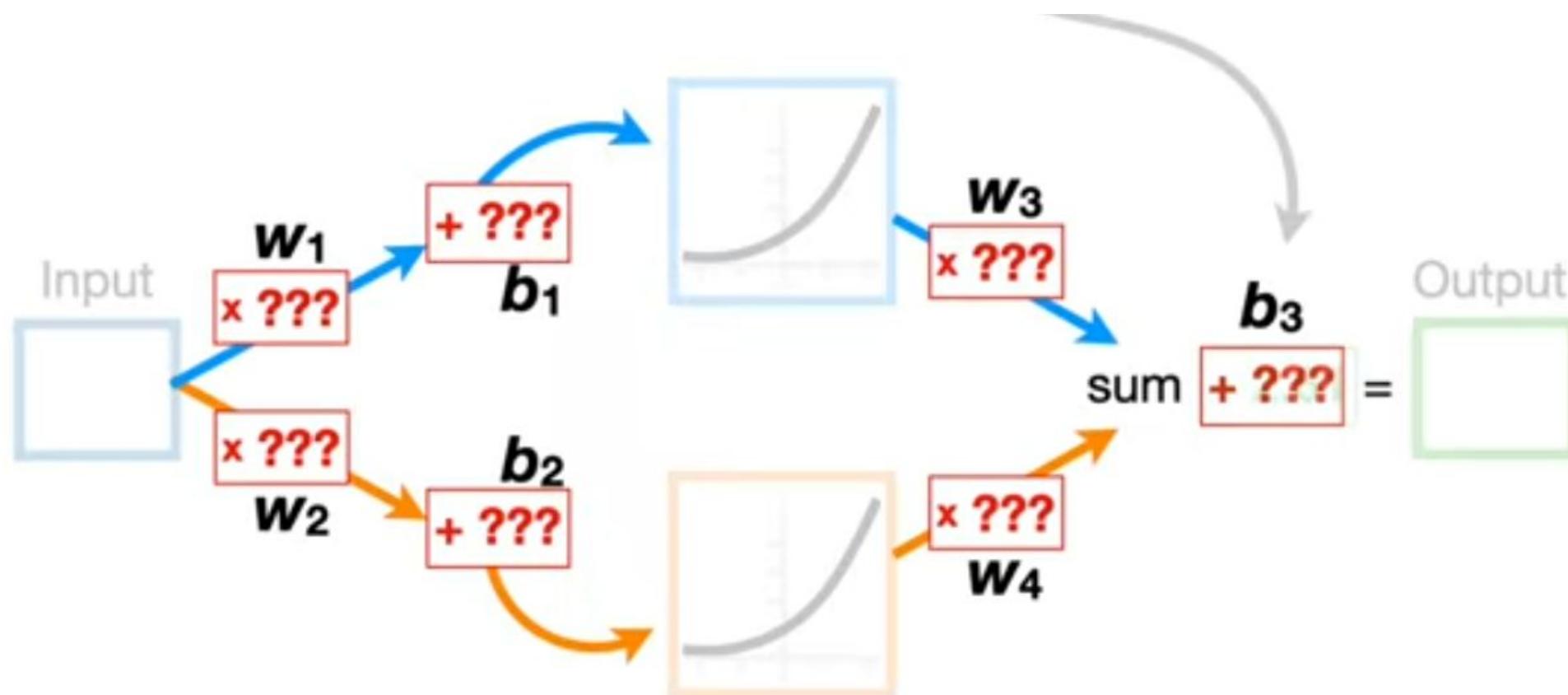
$$\frac{d \text{SSR}}{d w_1} = \frac{d \text{SSR}}{d \text{Predicted}} \times \frac{d \text{Predicted}}{d y_1} \times \frac{d y_1}{d x_1} \times \frac{d x_1}{d w_1}$$

$$\frac{d \text{SSR}}{d w_1} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times w_3 \times \frac{e^x}{1 + e^x} \times \text{Input}_i$$



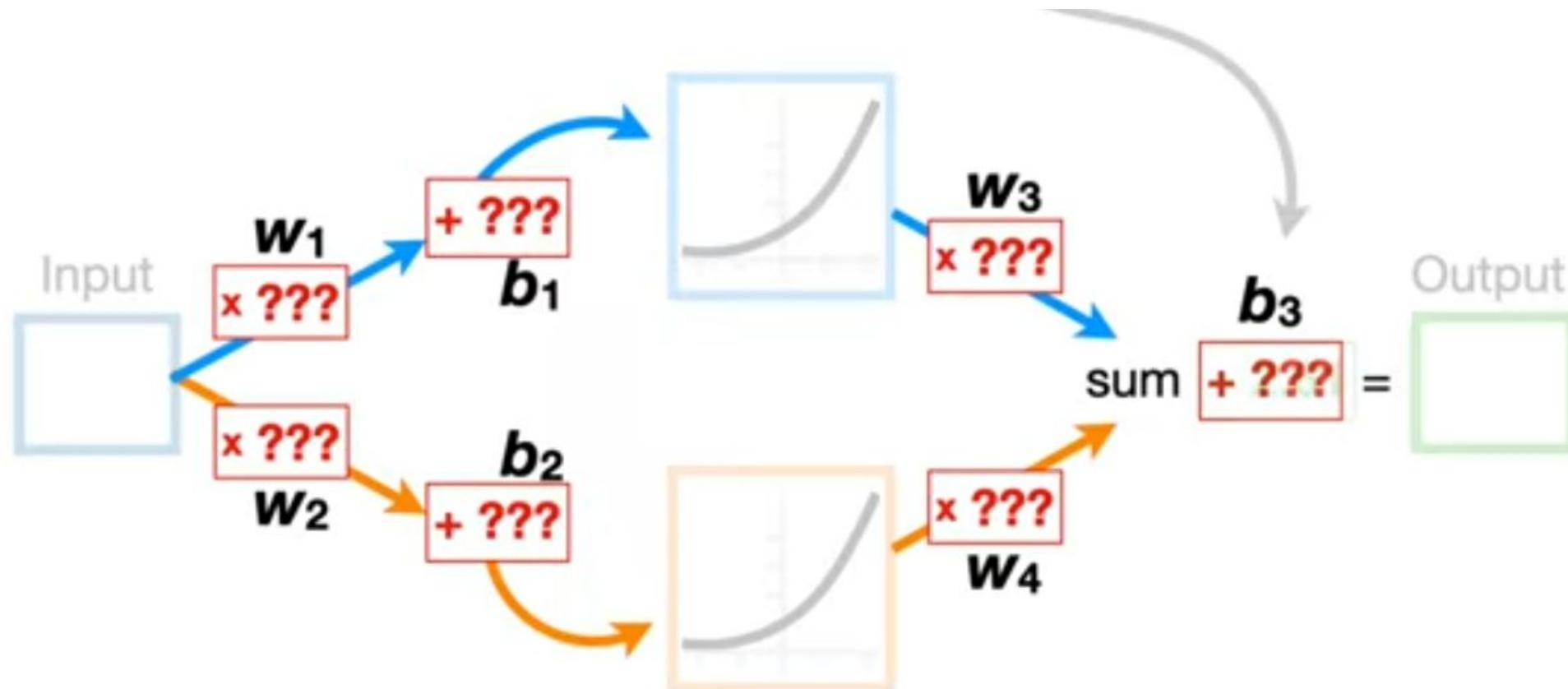
$$\frac{d \text{SSR}}{d b_1} = \frac{d \text{SSR}}{d \text{Predicted}} \times \frac{d \text{Predicted}}{d y_1} \times \frac{d y_1}{d x_1} \times \frac{d x_1}{d b_1}$$

$$\frac{d \text{SSR}}{d b_1} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times w_3 \times \frac{e^x}{1 + e^x} \times 1$$

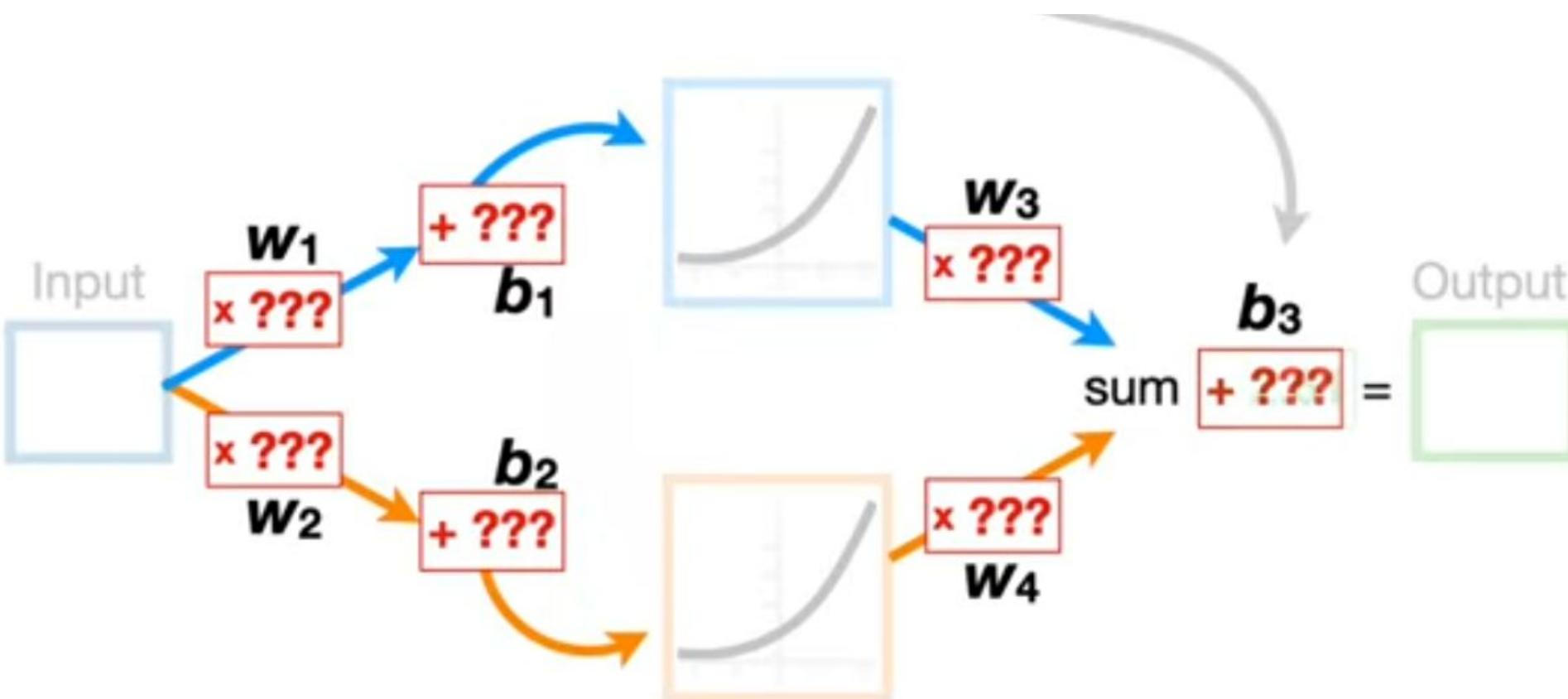


$$\frac{d \text{SSR}}{d w_2} = \frac{d \text{SSR}}{d \text{Predicted}} \times \frac{d \text{Predicted}}{d y_2} \times \frac{d y_2}{d x_2} \times \frac{d x_2}{d w_2}$$

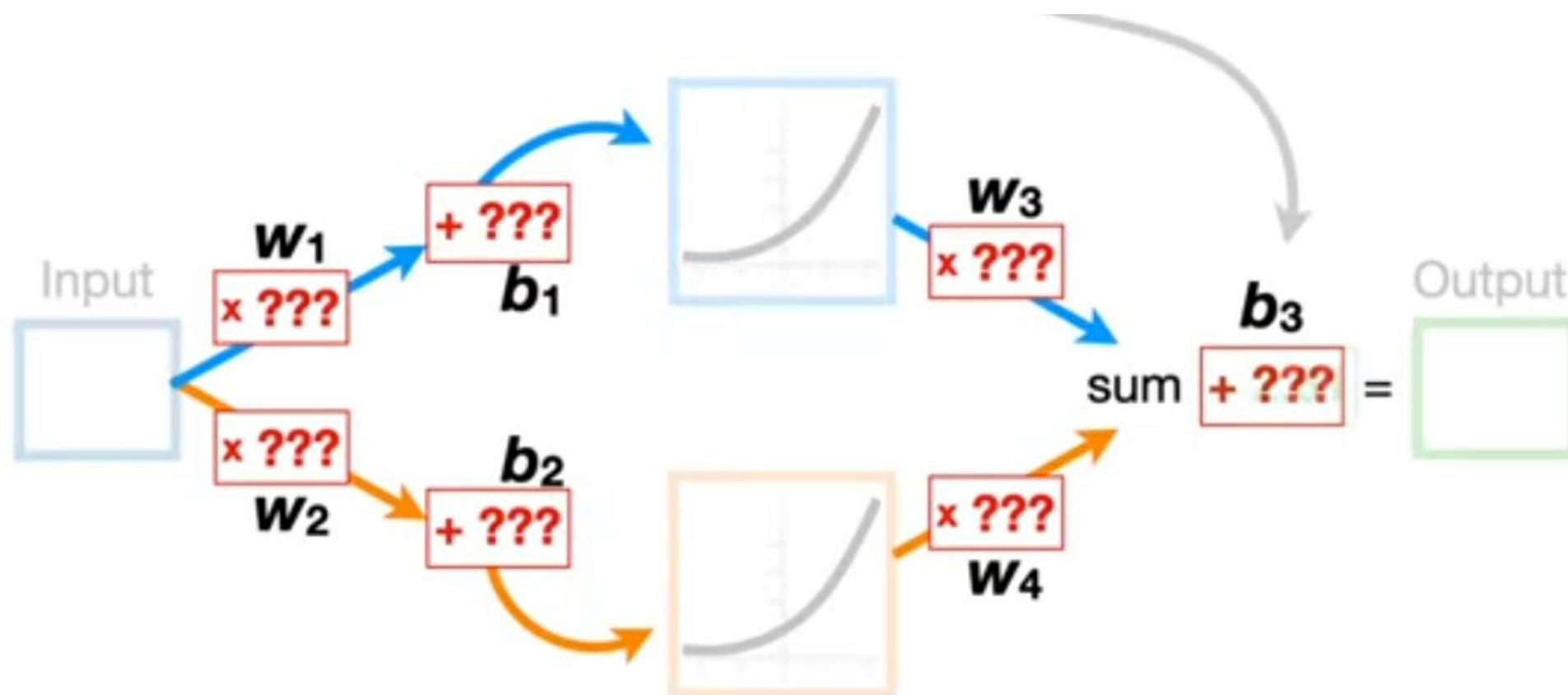
$$\frac{d \text{SSR}}{d w_2} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times w_4 \times \frac{e^x}{1 + e^x} \times \text{Input}_i$$



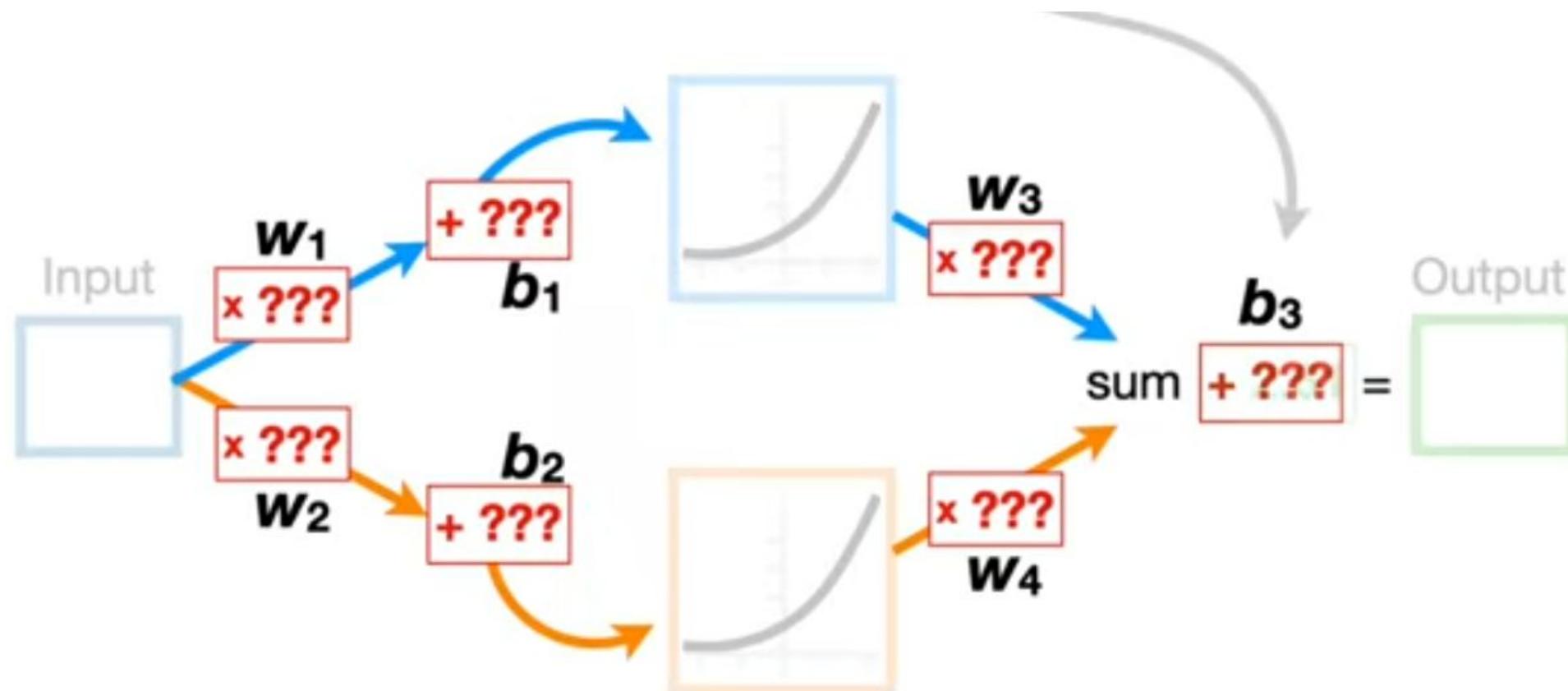
$$\frac{d \text{SSR}}{d b_2} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times w_4 \times \frac{e^x}{1 + e^x} \times 1$$



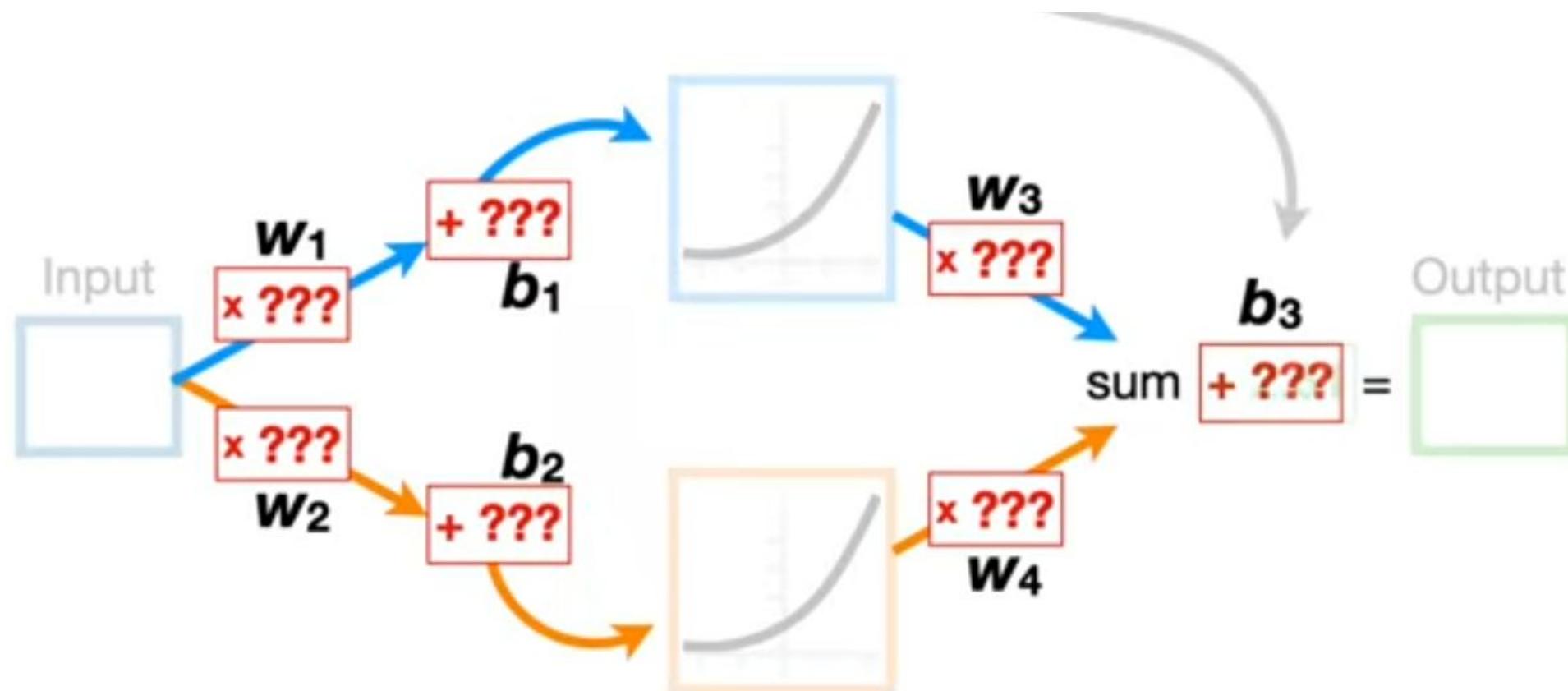
$$\frac{d \text{SSR}}{d w_3} = \frac{d \text{SSR}}{d \text{Predicted}} \times \frac{d \text{Predicted}}{d w_3} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times y_{1,i}$$



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



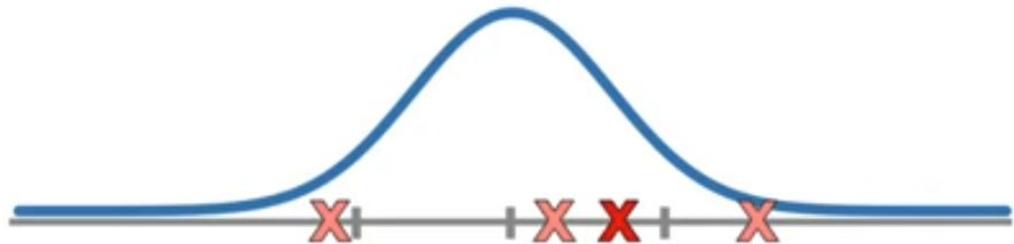
$$\frac{d \text{SSR}}{d b_3} = \frac{d \text{SSR}}{d \text{Predicted}} \times \frac{d \text{Predicted}}{d b_3} = \sum_{i=1}^{n=3} -2 \times (\text{Observed}_i - \text{Predicted}_i) \times 1$$



How do we pick the initial weights and biases

Weights

we pick initial numbers from a standard normal Distribution



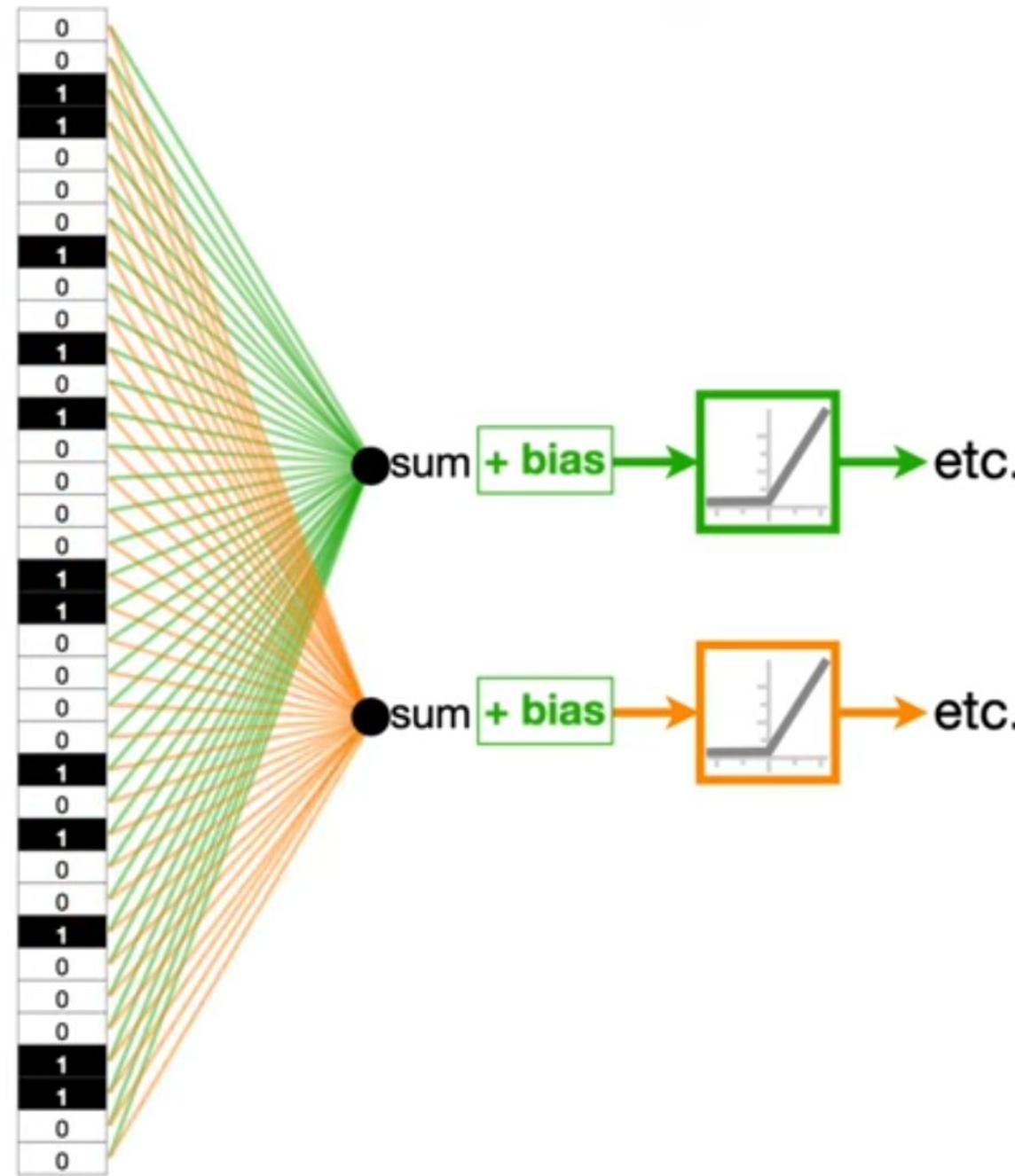
Biases

We initialize them to 0 by standard

04

Convolutional Neural Networks

0	0	1	1	0	0
0	1	0	0	1	0
1	0	0	0	0	1
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0

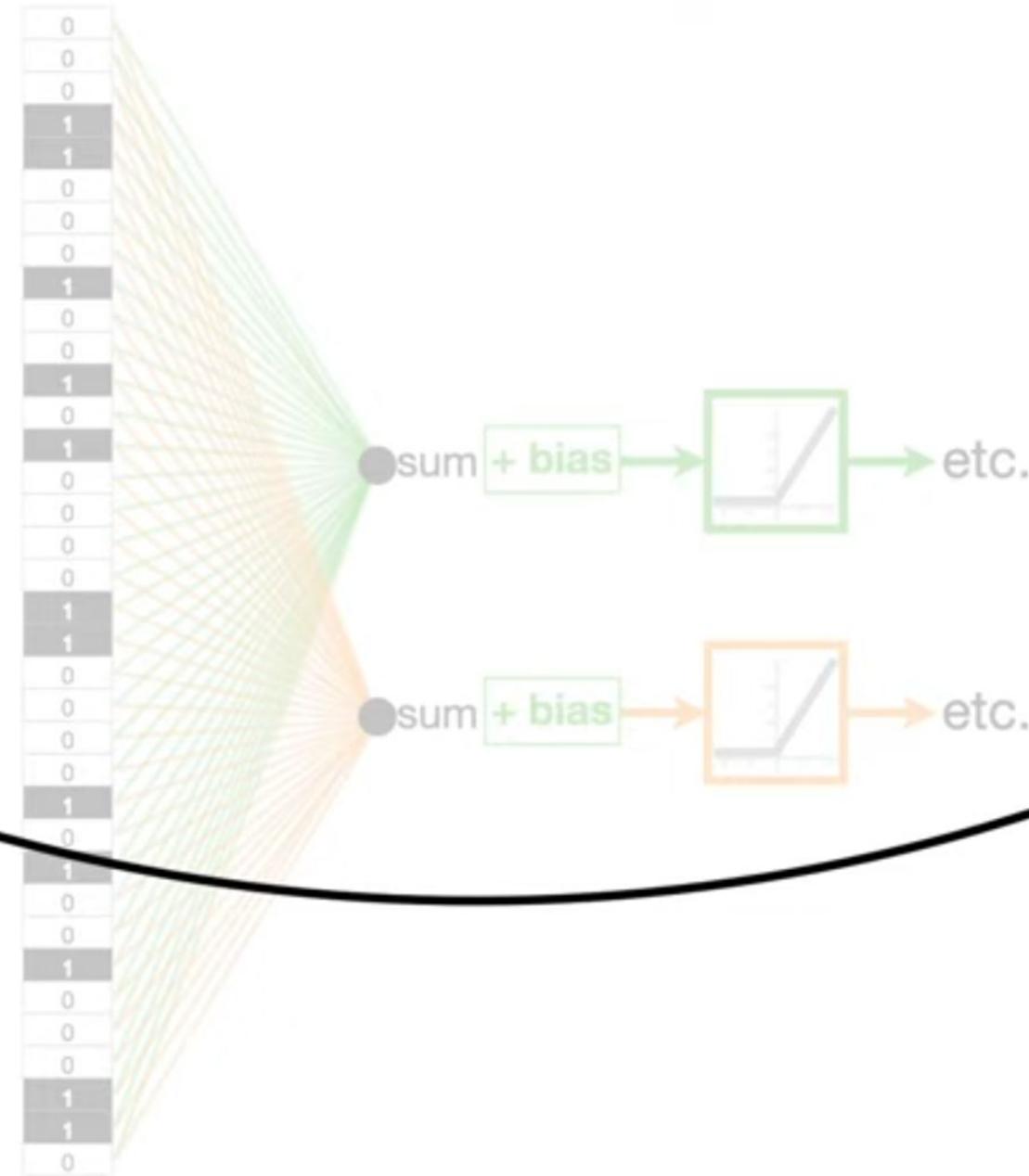


0	0	1	1	0	0
0	1	0	0	1	0
1	0	0	0	0	1
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0



0	0	0	1	1	0
0	0	1	0	0	1
0	1	0	0	0	0
0	1	0	0	0	0
0	0	1	0	0	1
0	0	0	1	1	0

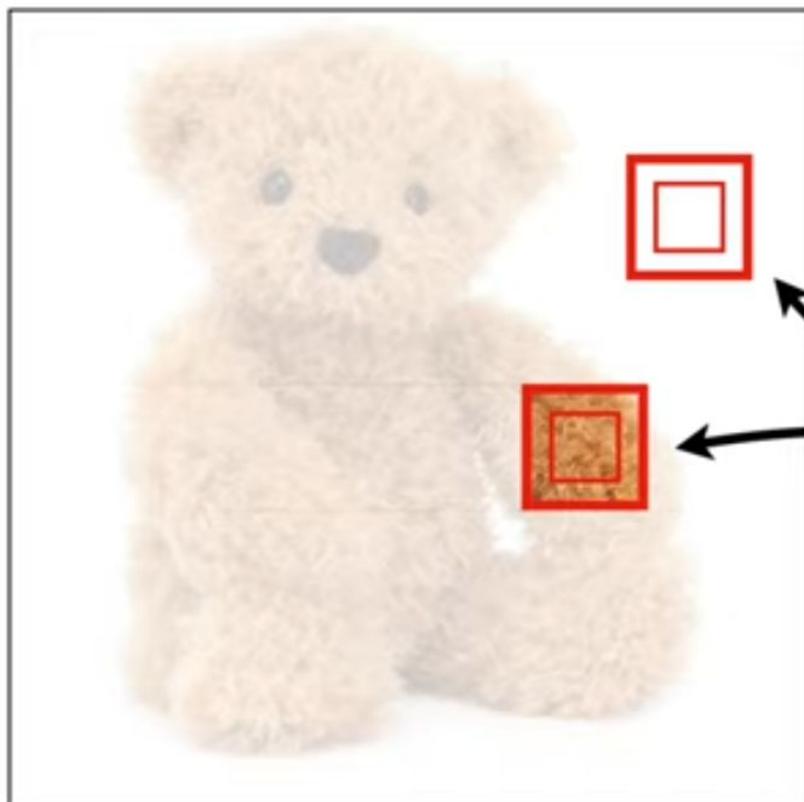
Shifted to the right 1 pixel.



...then it is not clear
that the **Neural
Network** will still
recognize this letter
“O” correctly if each
pixel is shifted to the
right by one.

Convolutional Neural Networks

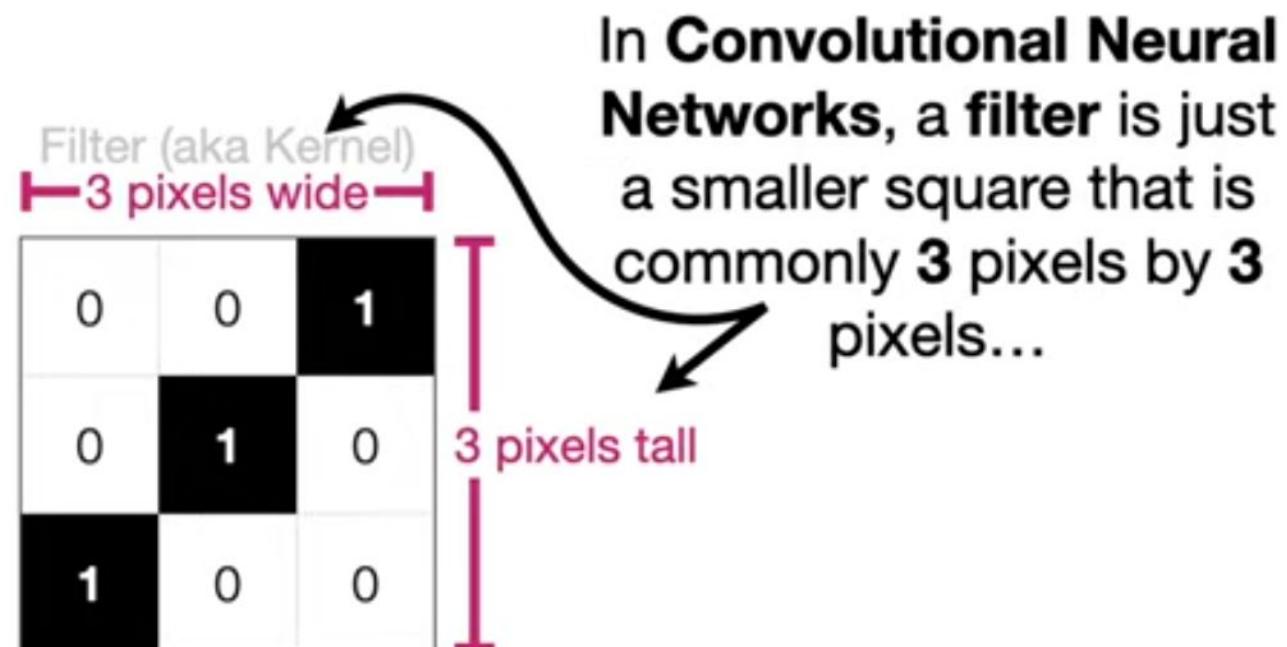
do 3 things to make image classification practical:

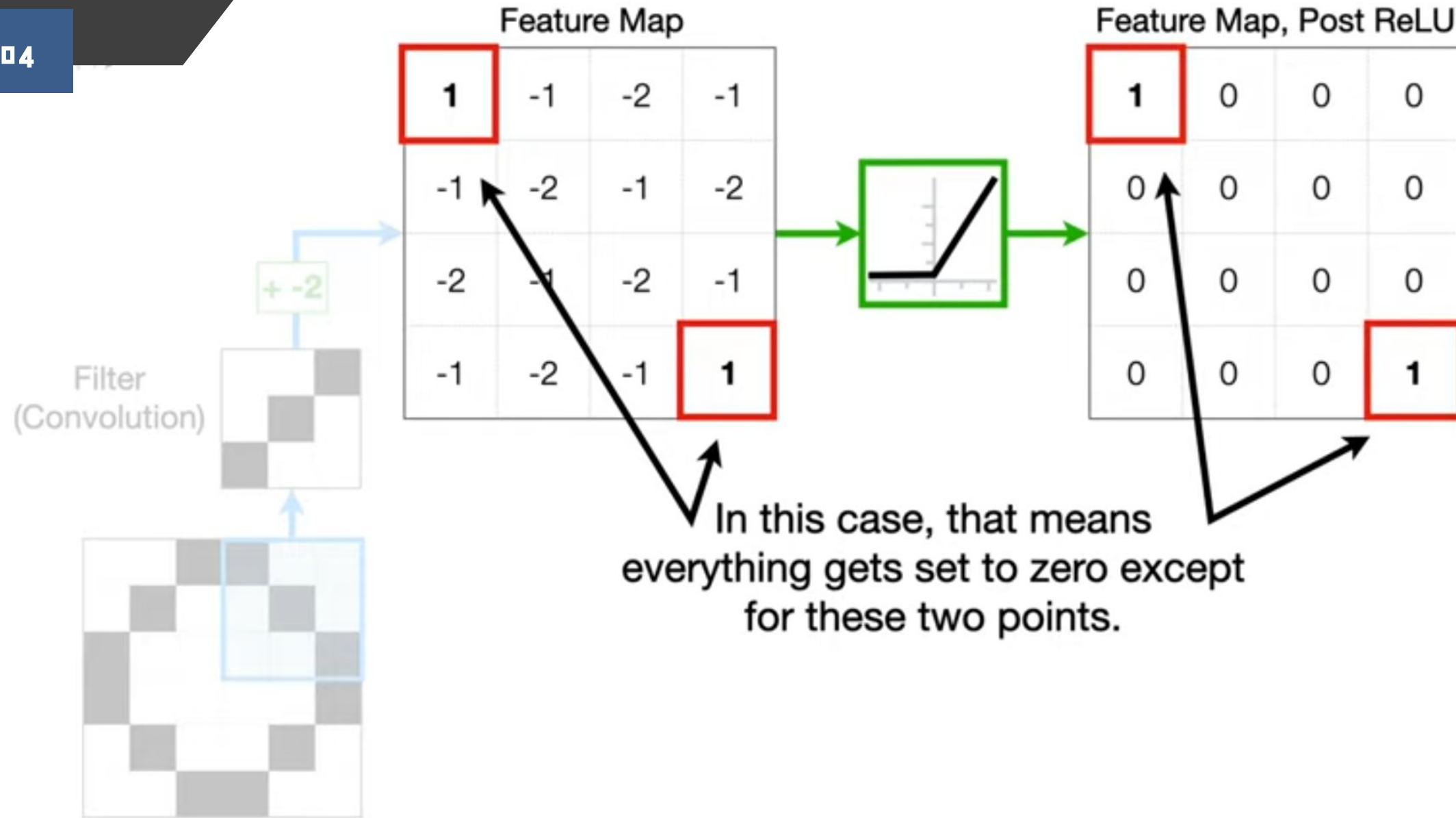


- 1) Reduce the number of input nodes.
- 2) Tolerate small shifts in where the pixels are in the image.
- 3) Take advantage of the correlations that we observe in complex images.

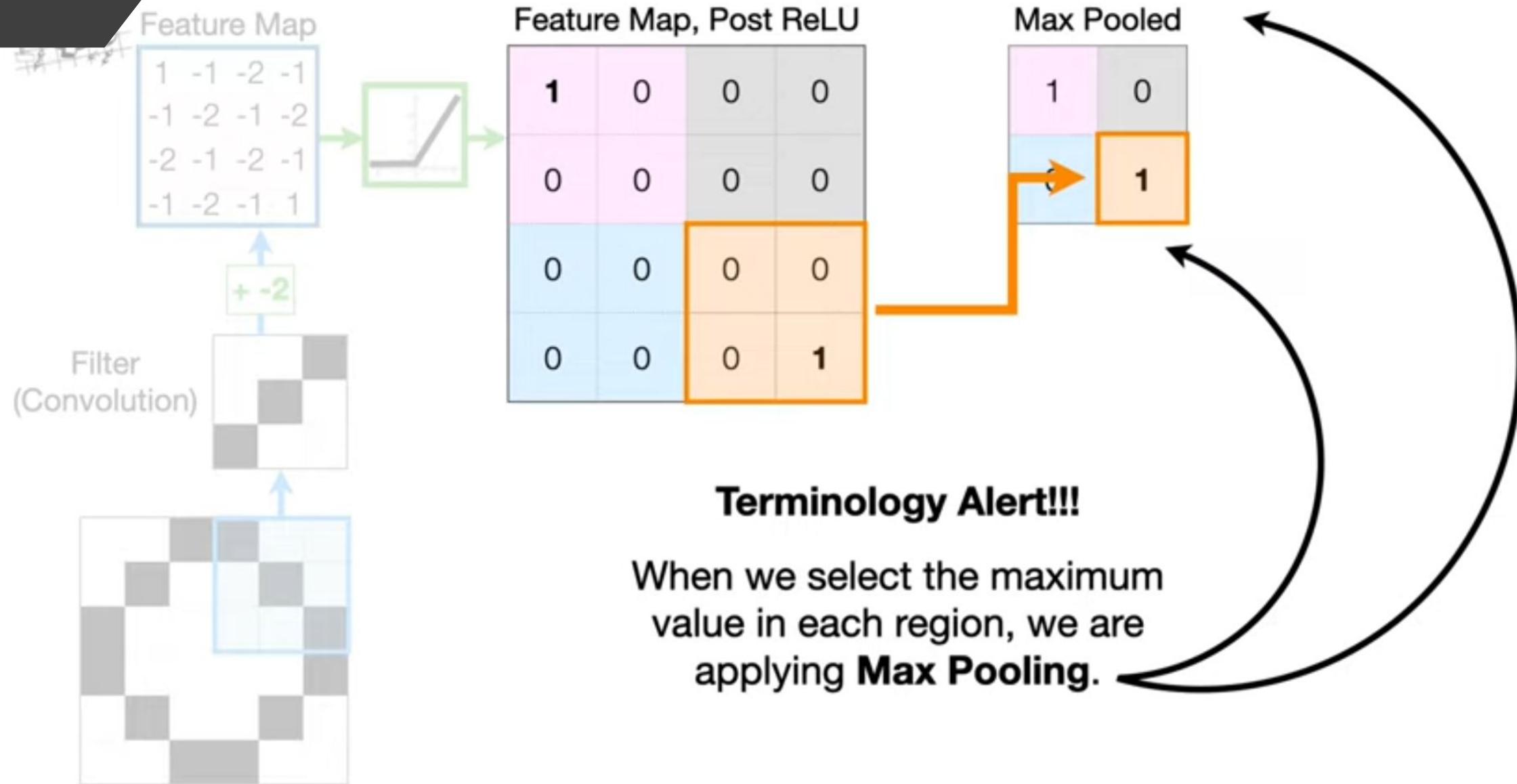
Input Image

0	0	1	1	0	0
0	1	0	0	1	0
1	0	0	0	0	1
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0



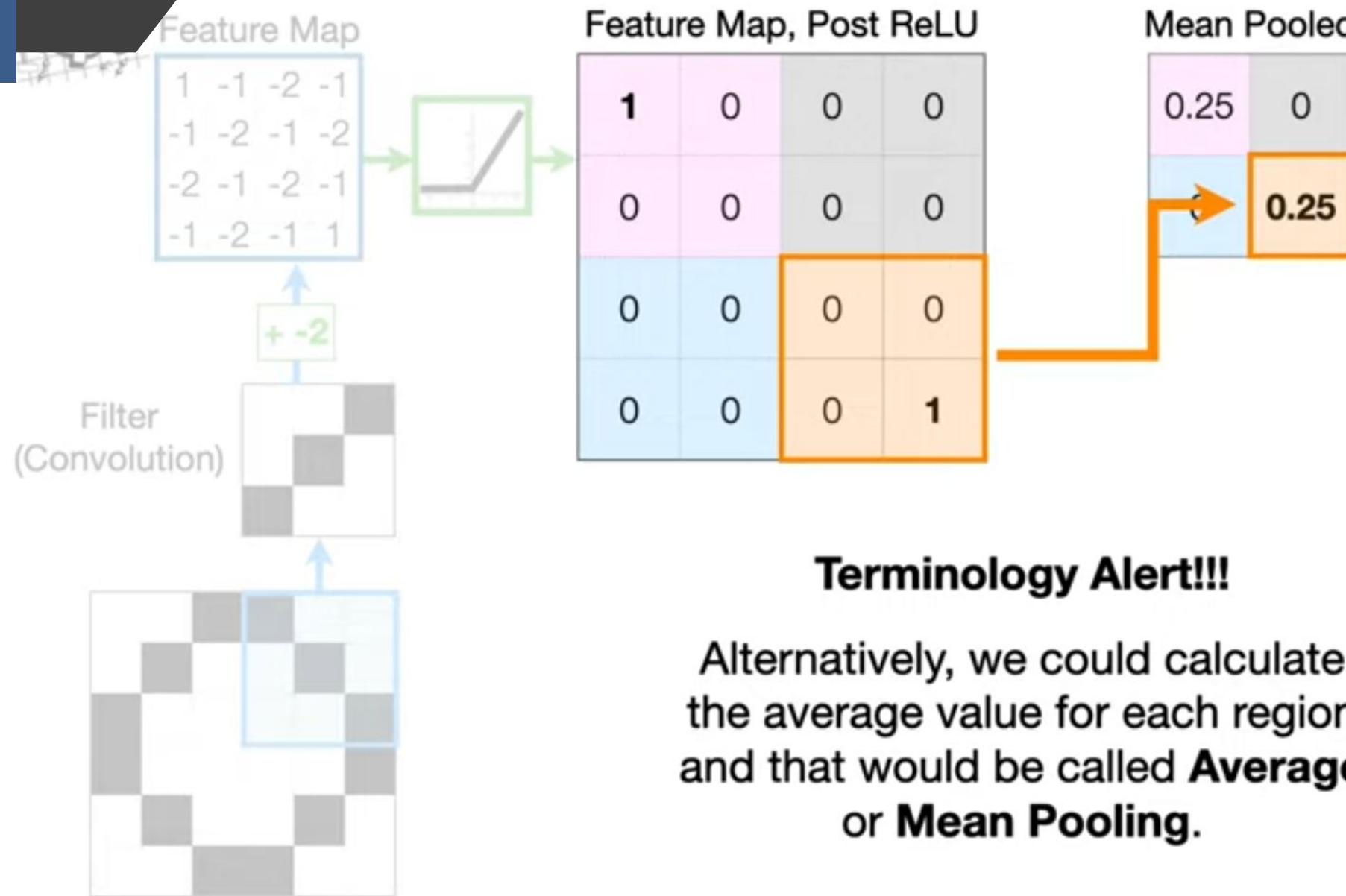


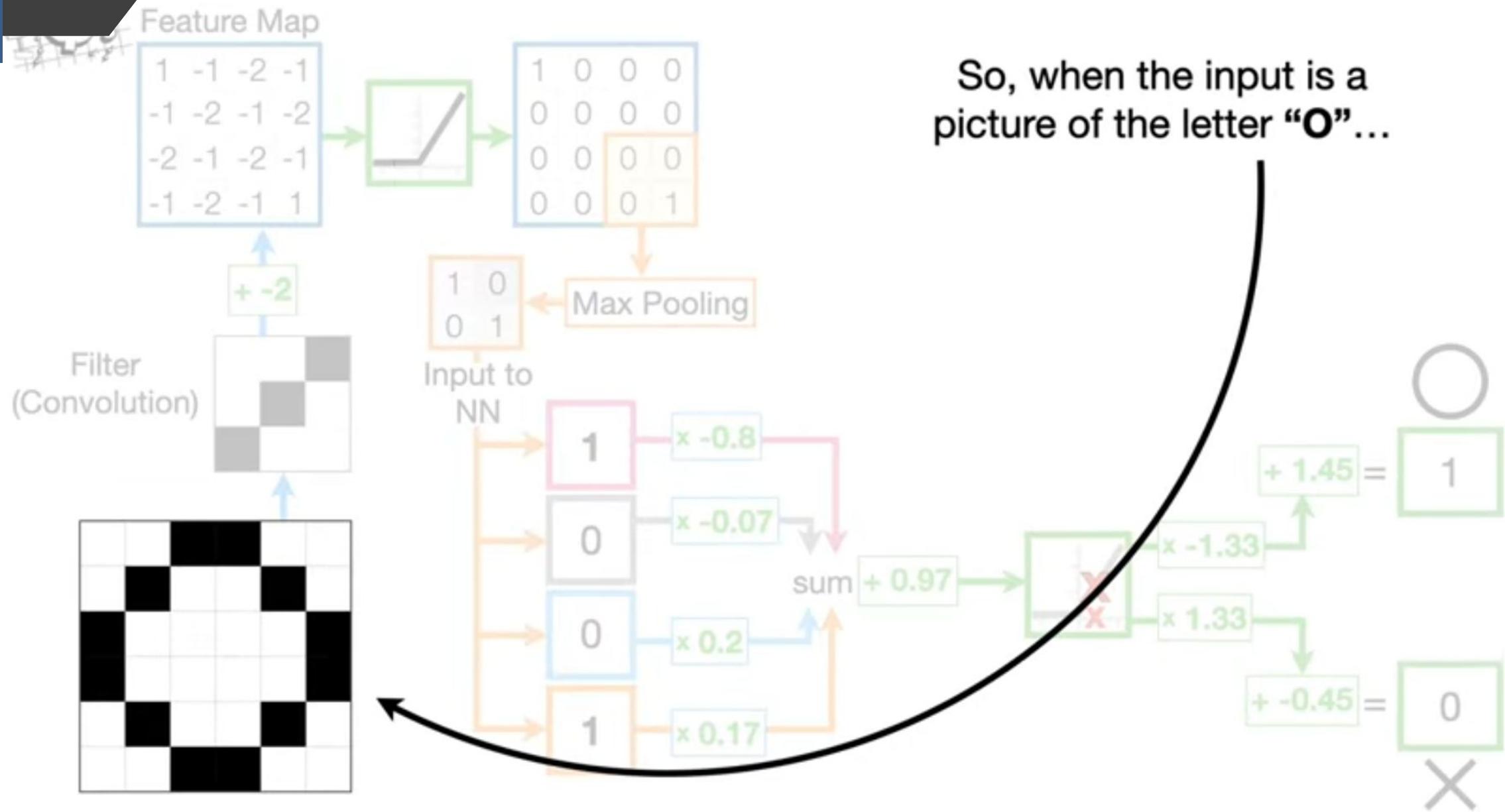
In this case, that means
everything gets set to zero except
for these two points.

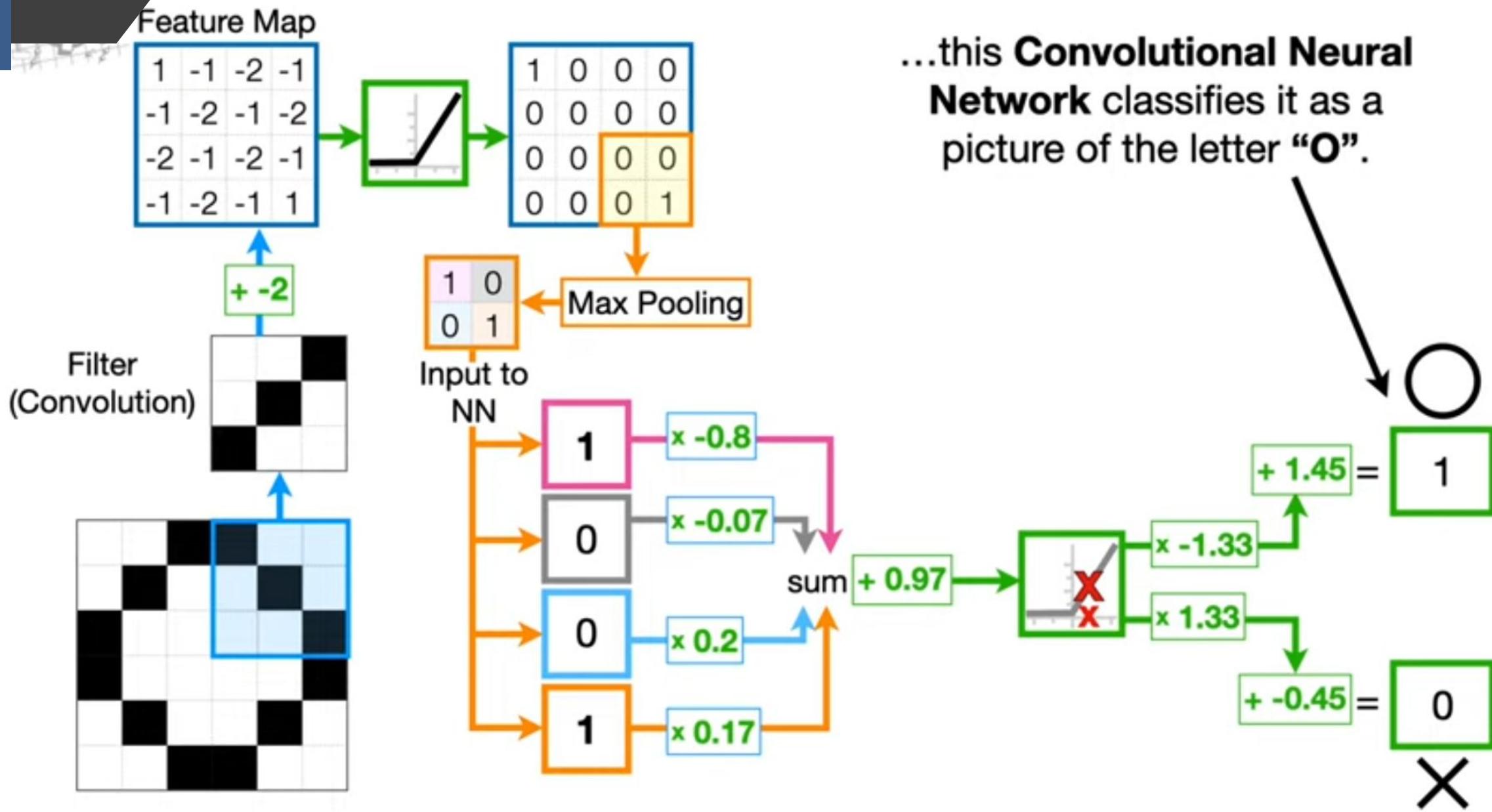


Terminology Alert!!!

When we select the maximum value in each region, we are applying **Max Pooling**.



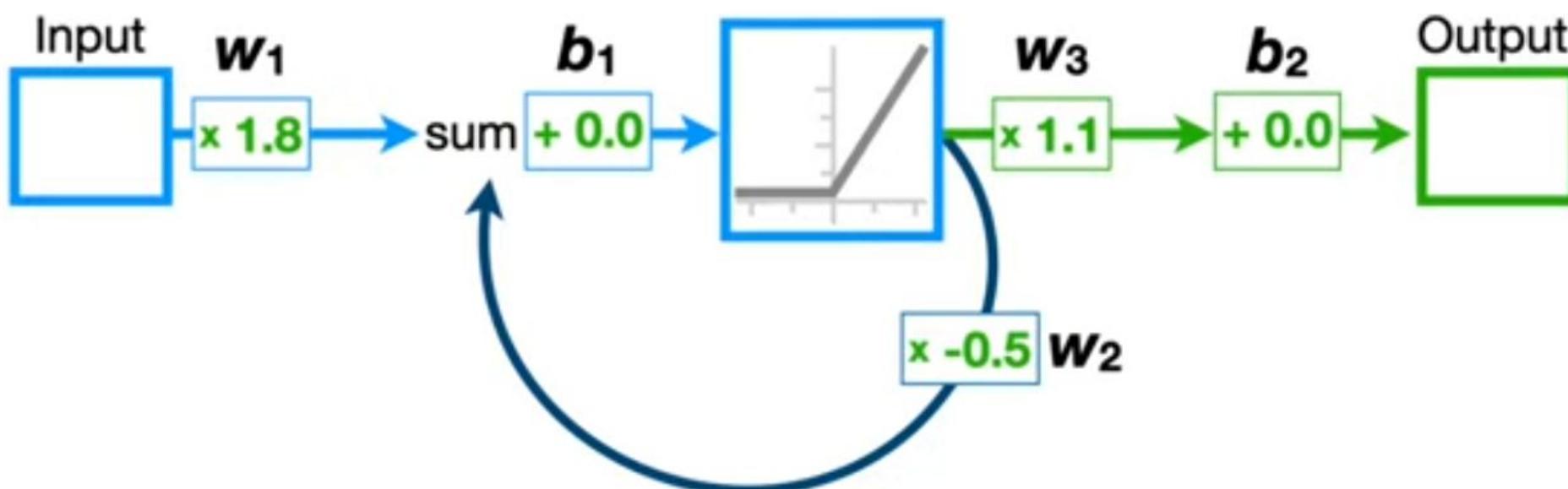


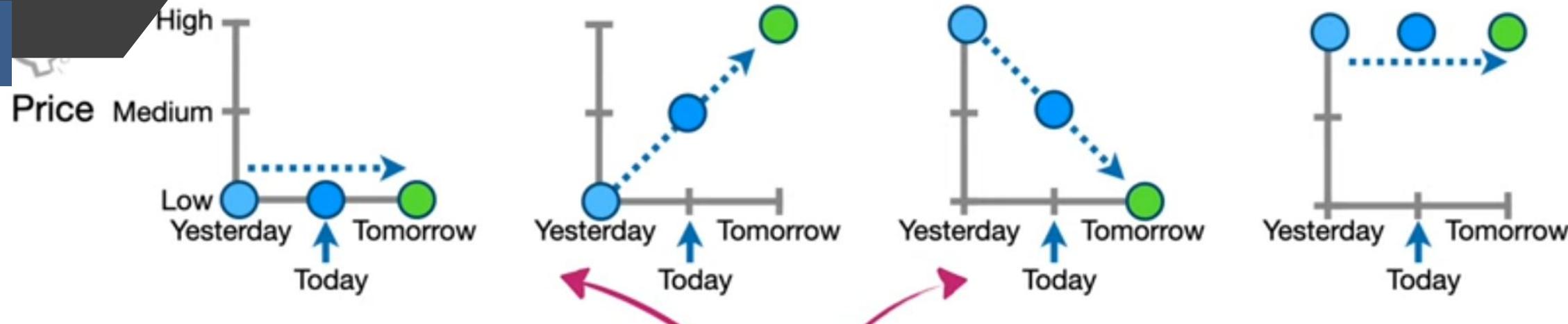




Recurrent Neural Networks

Just like the other neural networks that we've seen before, **Recurrent Neural Networks have weights, biases,**

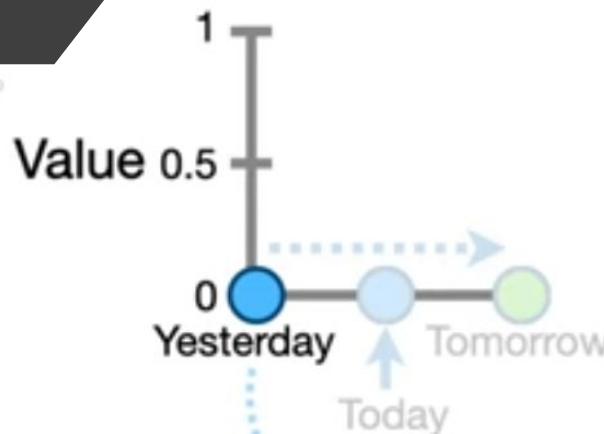




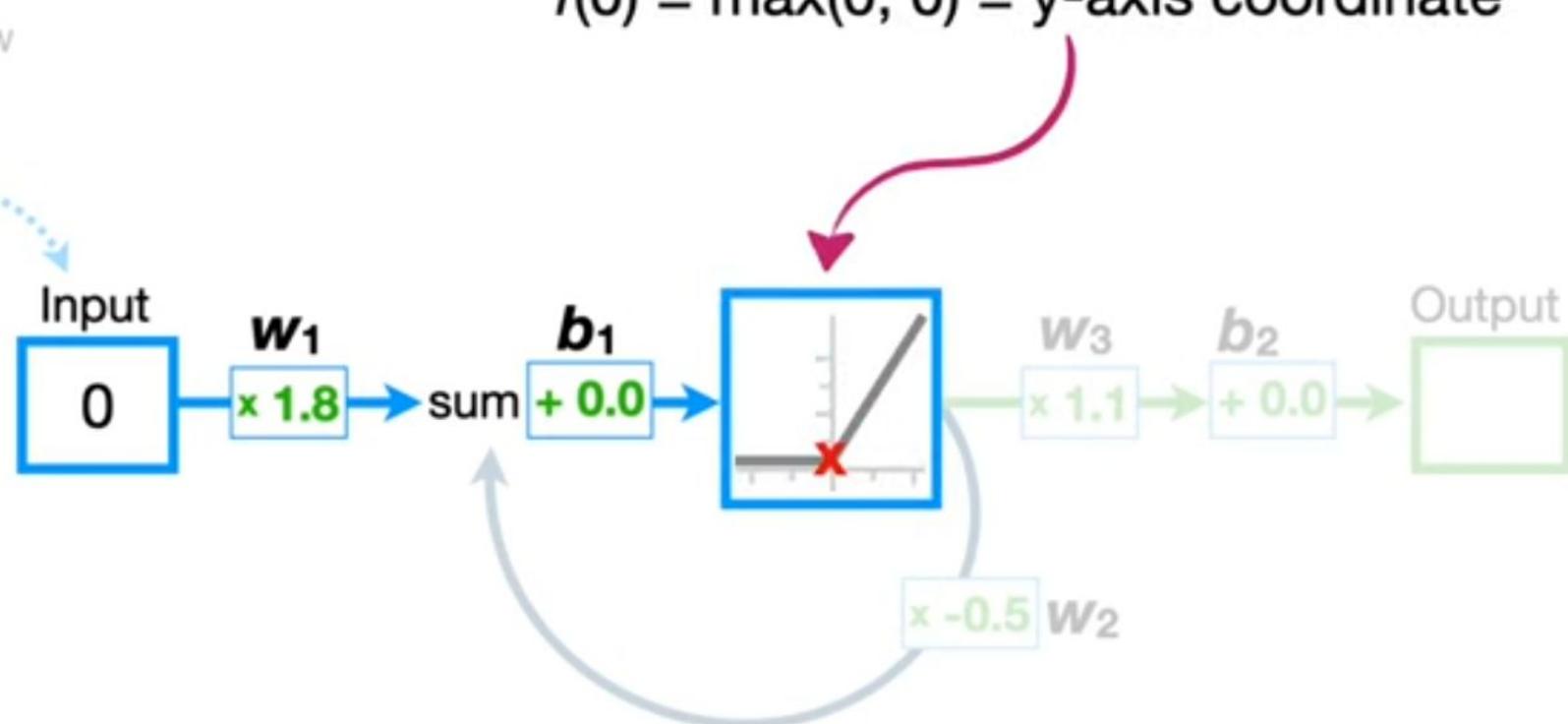


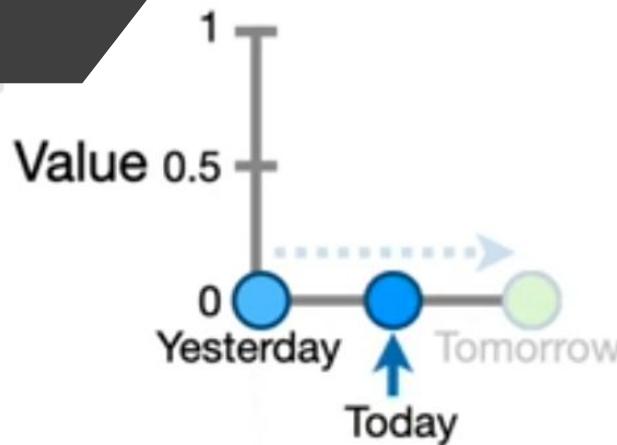
Yesterday $\times w_1 + b_1$ = x-axis coordinate

$$0 \times 1.8 + 0.0 = 0$$

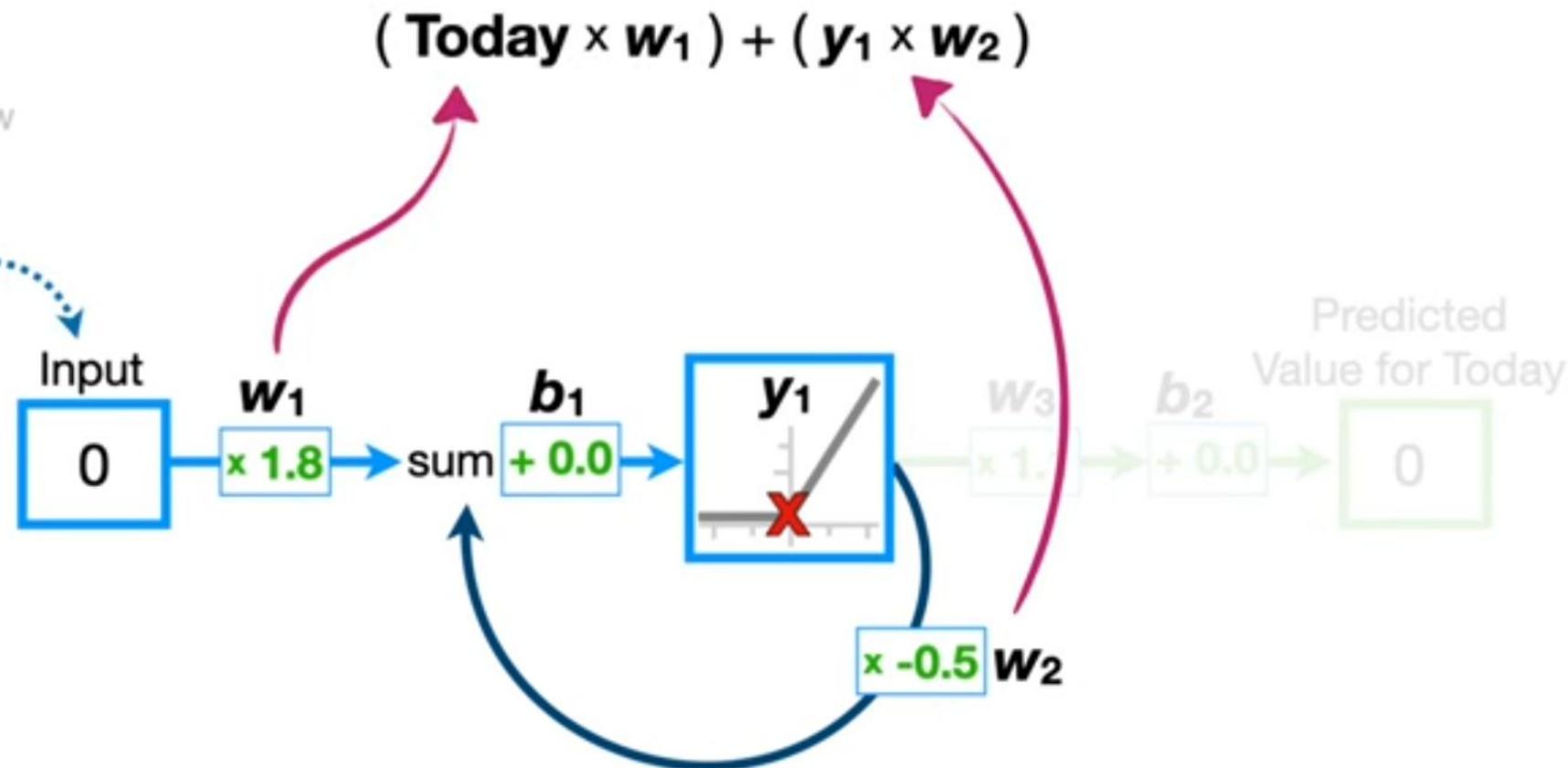


Now we can do the math just like we would for any other neural network.

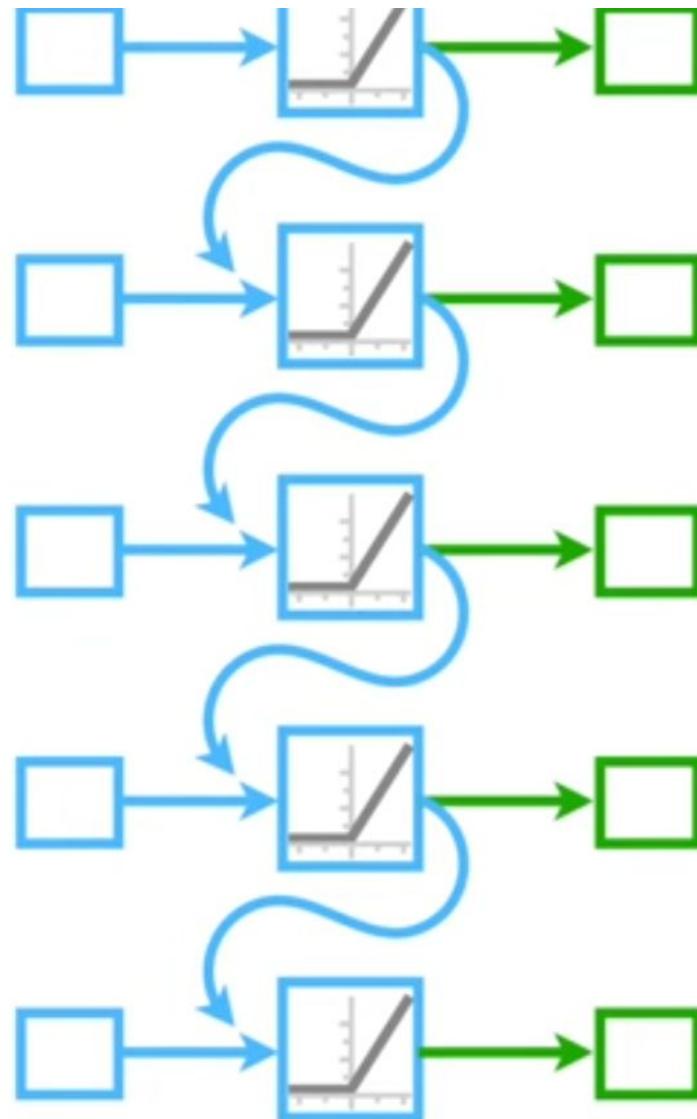




In other words, the **feedback loop** allows both **yesterday** and **today's** values to influence the prediction.



This problem is
called **The
Vanishing/Exploding
Gradient Problem.**





...we multiply the input value by w_2 , which is 2, raised to the number of times we **unrolled**, which is 4.

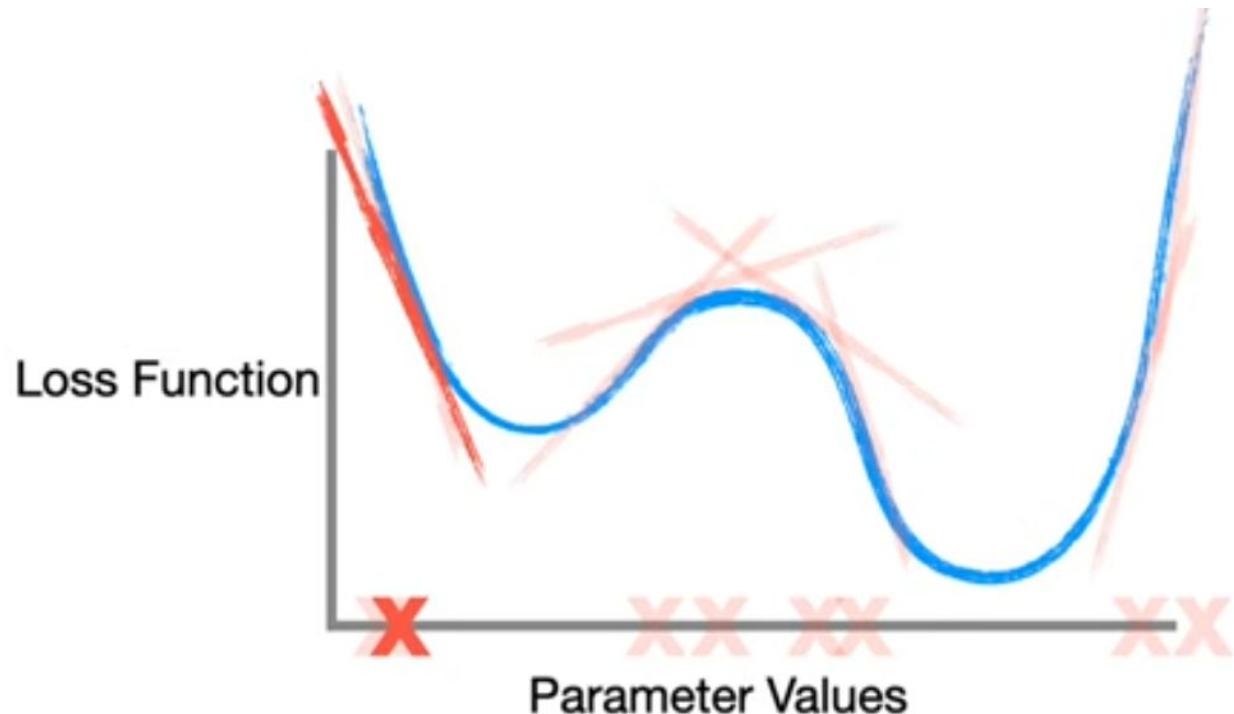
$$\text{Input}_1 \times 2 \times 2 \times 2 \times 2$$

$$= \text{Input}_1 \times 2^4$$

$$= \text{Input}_1 \times w_2^{\text{Num. Unroll}}$$



One way to prevent
**The Exploding
Gradient Problem**
would be to limit w_2
to values < 1 .

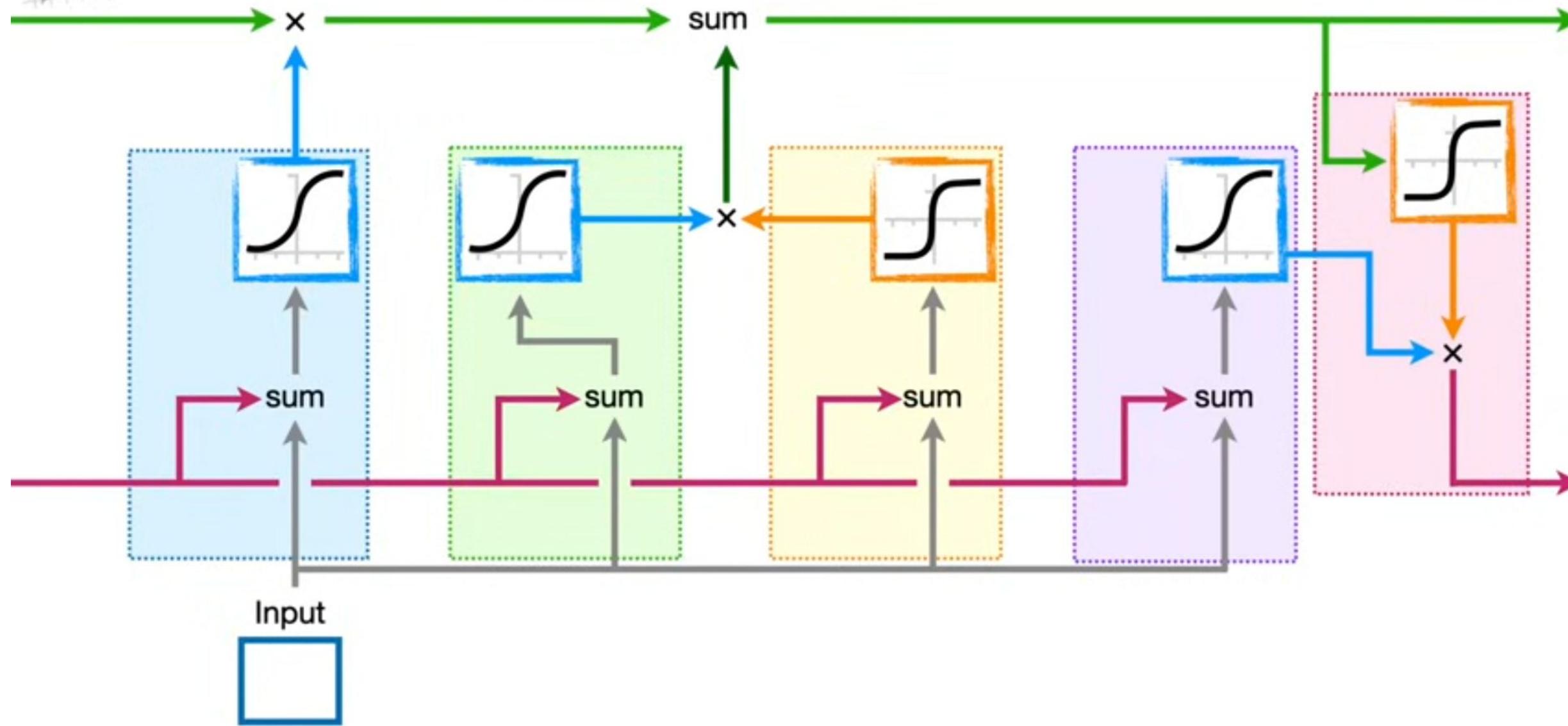


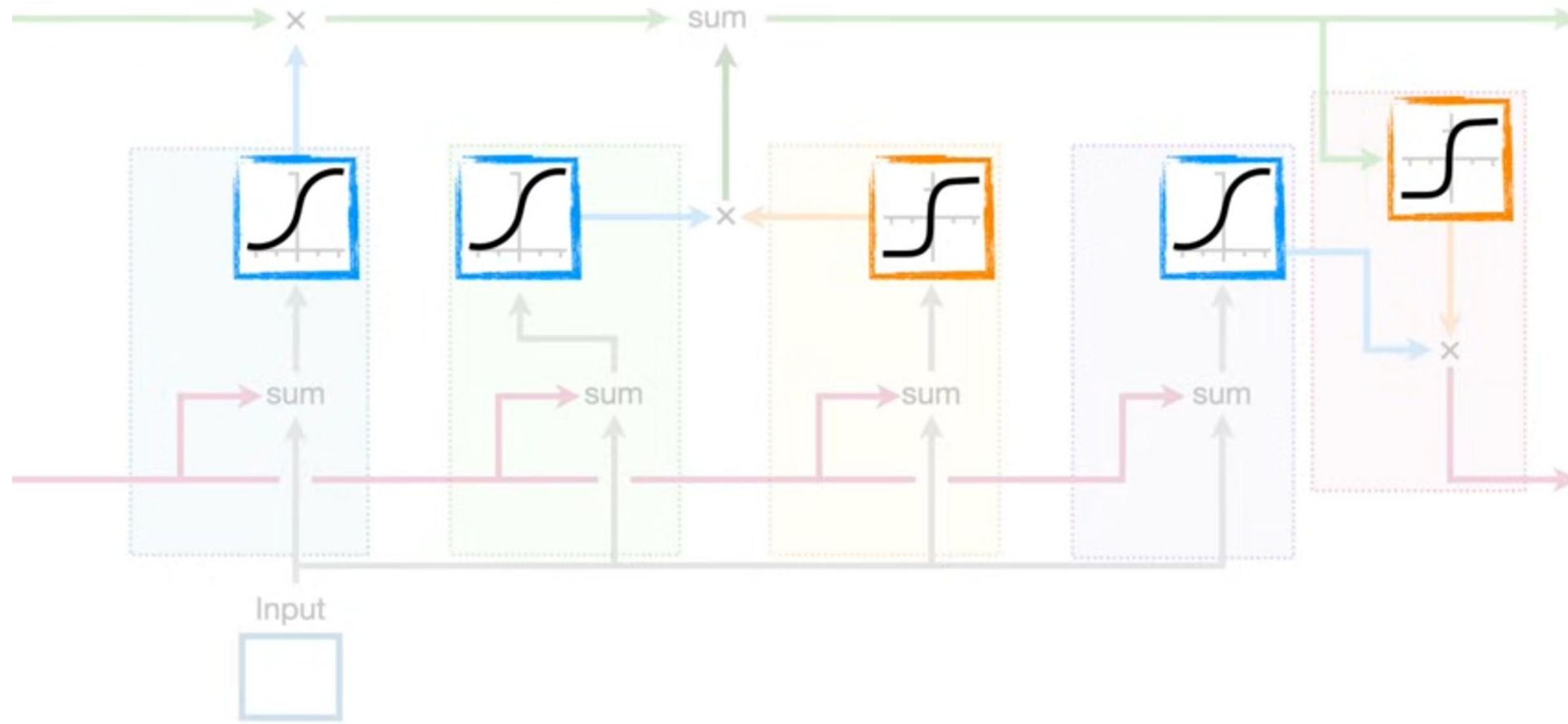
Now, just like before, we multiply the first input by w_2 raised to the number of times we **unroll** the network.

$\text{Input}_1 \times w_2^{\text{Num. Unroll}}$

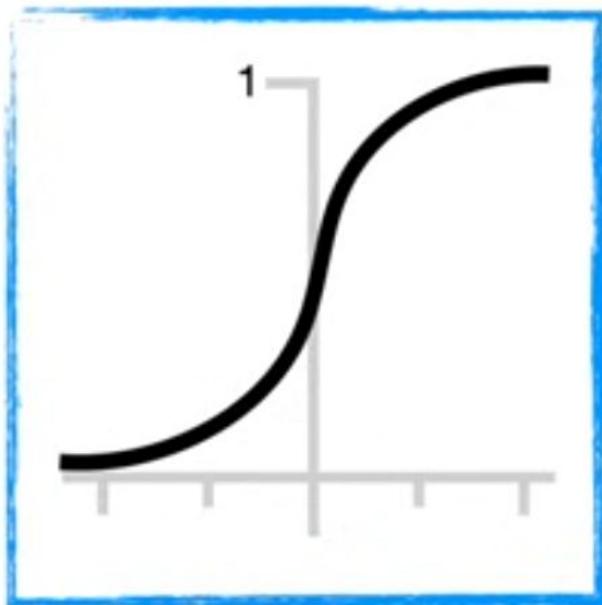


...Long Short-Term Memory is based
on a much more complicated unit.

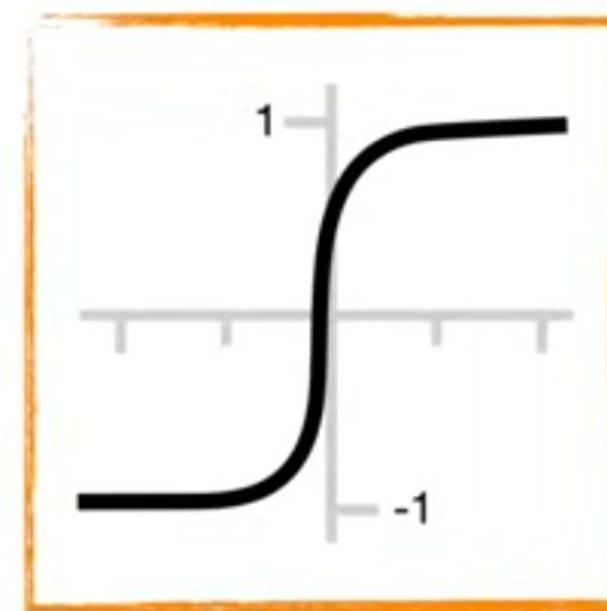


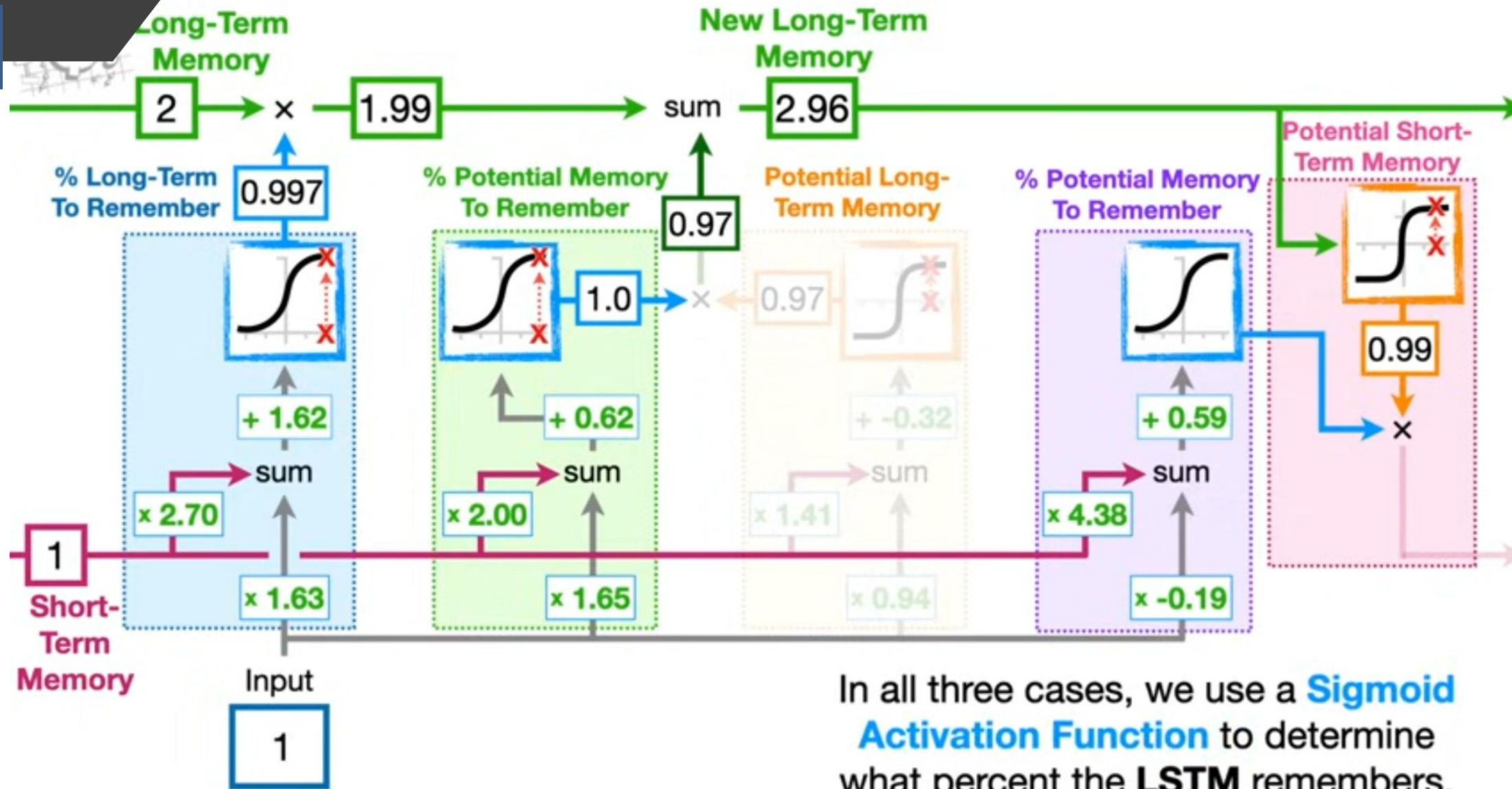


So, now that we know that the **Sigmoid Activation Function** turns any input into a number between **0** and **1**...



...and the **Tanh Activation Function** turns any input into a number between **-1** and **1**...





In all three cases, we use a **Sigmoid Activation Function** to determine what percent the **LSTM** remembers.

06

Conclusion



Thank You