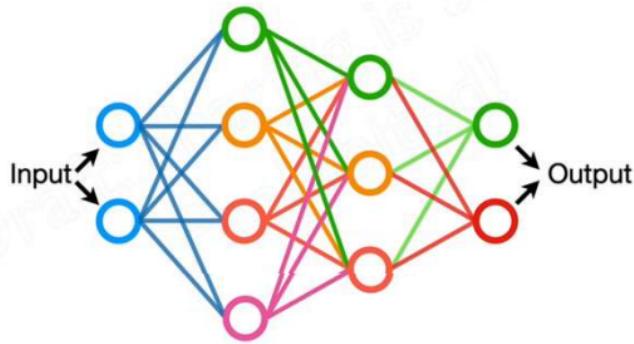
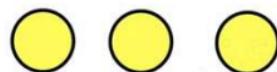




**Neural Networks**, one of the most popular algorithms in **Machine Learning**, cover a broad range of concepts and techniques.





# **Neural Networks...**

So, with that said, let's imagine  
we tested a drug that was  
designed to treat an illness...



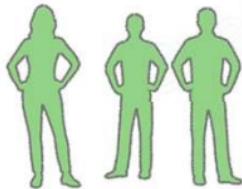
vs.



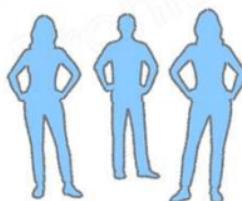
...and we gave the drug to 3 different groups of people with 3 different **Dosages**.



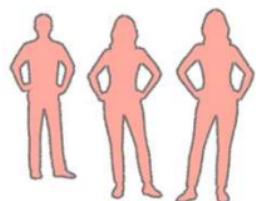
Low Dosage



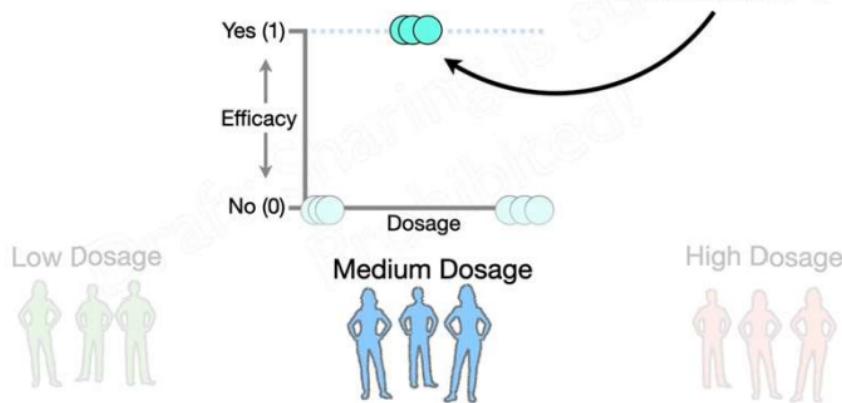
Medium Dosage



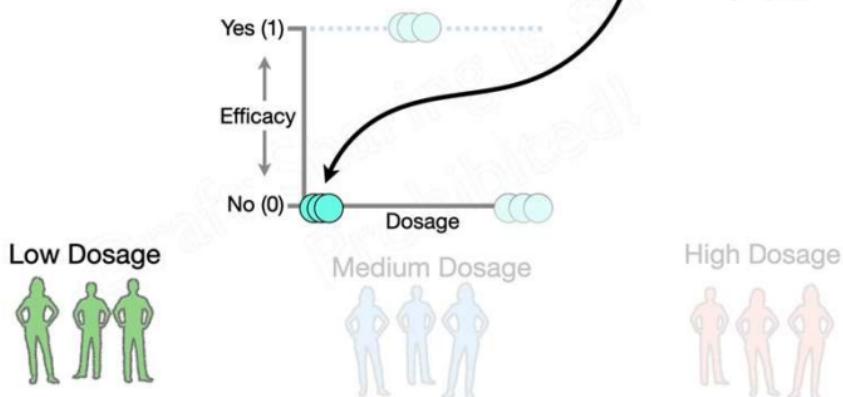
High Dosage

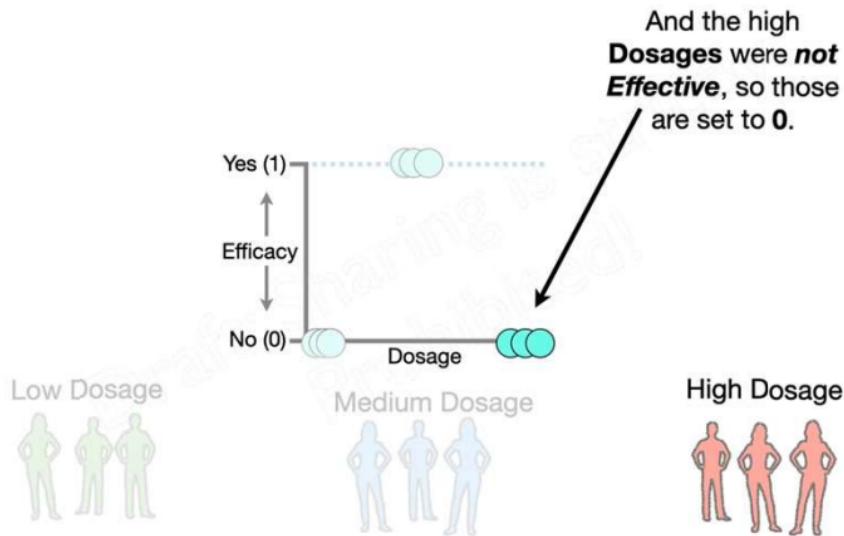


In contrast, the medium **Dosages** were **Effective**, so we set them to **1**.

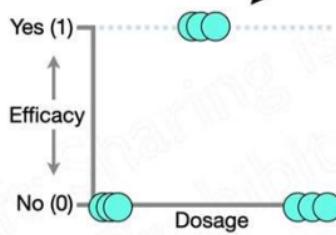


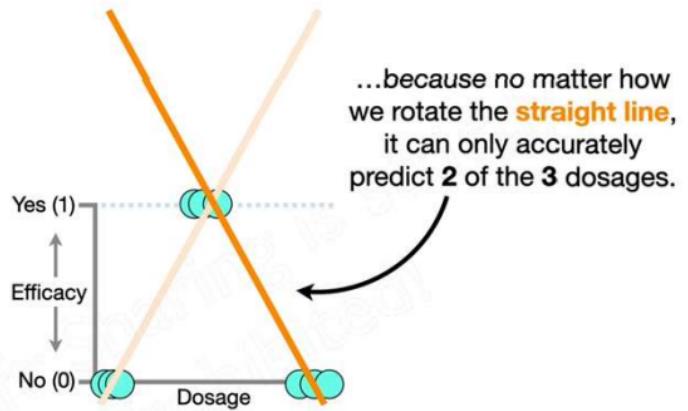
The low Dosages were *not Effective*, so we set them to 0 on this graph.



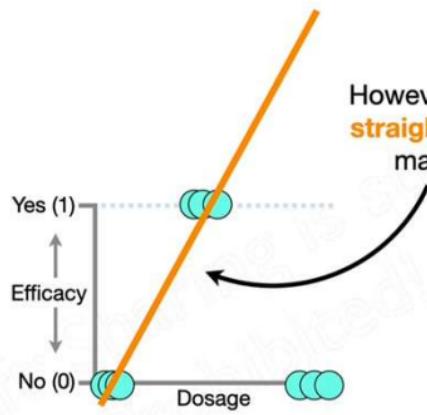


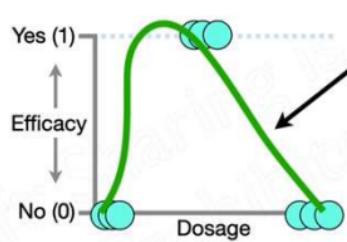
Now that we have this data, we would like to use it to predict whether or not a future **Dosage** will be **Effective**.



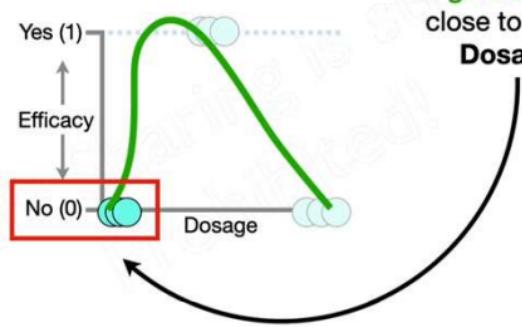


However, we can't just fit a **straight line** to the data to make predictions...

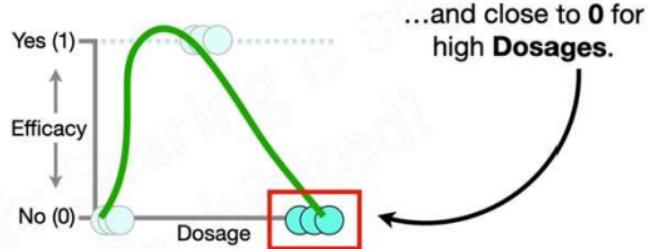


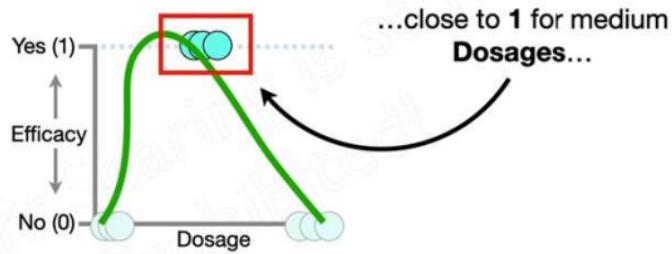


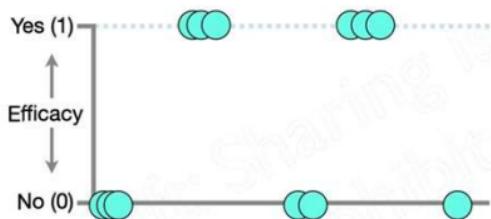
The good news is that a **Neural Network** can fit a squiggle to the data.



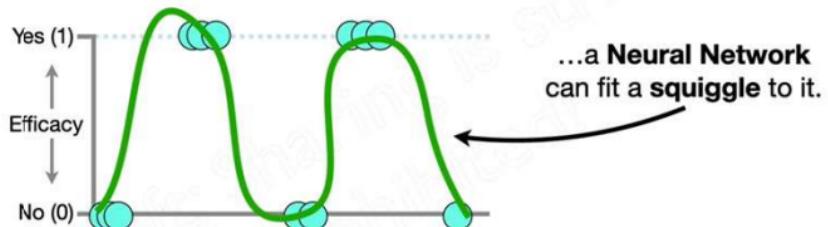
The **green squiggle** is  
close to 0 for low  
**Dosages...**



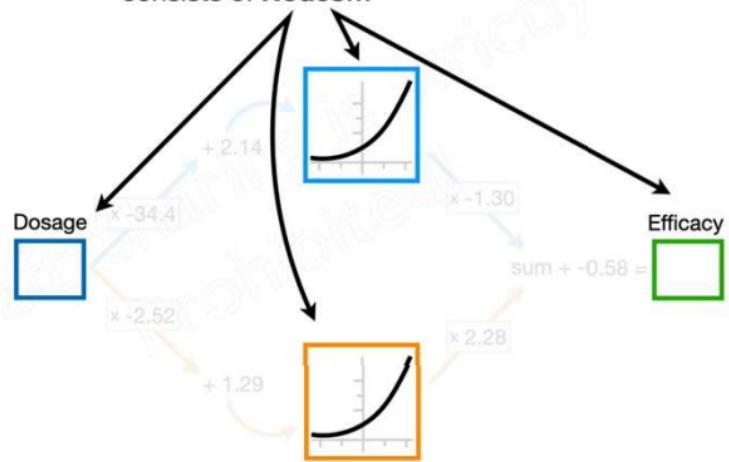




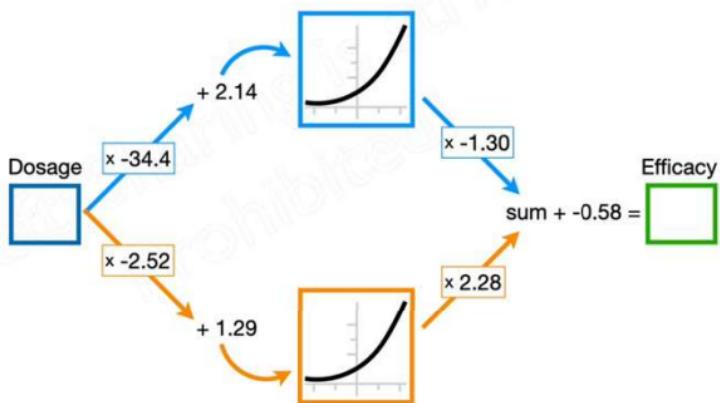
And even if we have  
really complicated  
dataset like this...



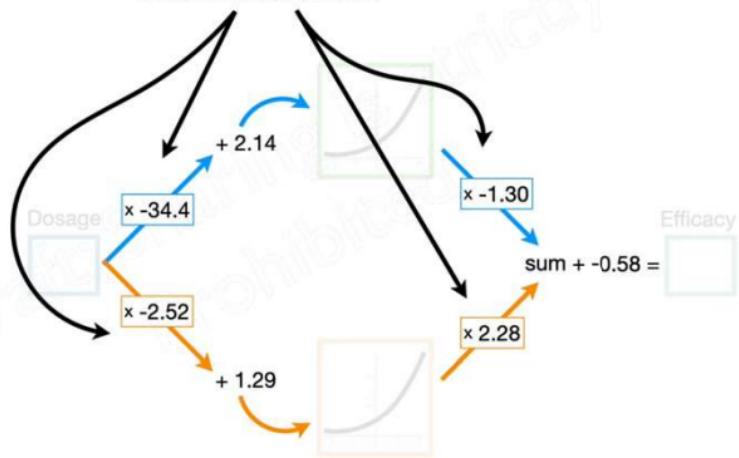
**A Neural Network**  
consists of **Nodes**...



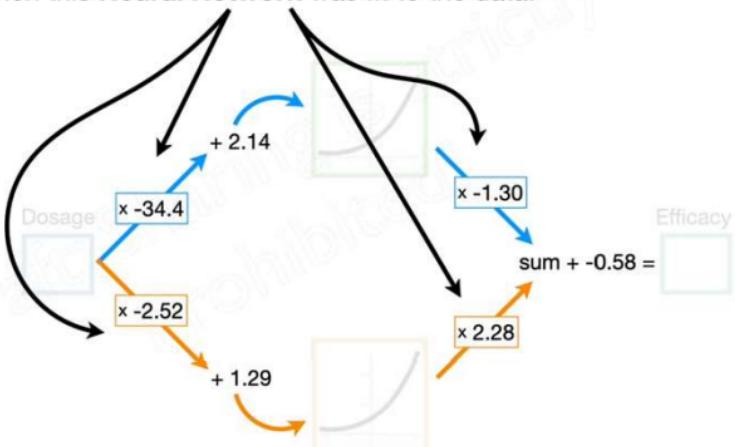
But first let's just talk about  
what a **Neural Network** is.



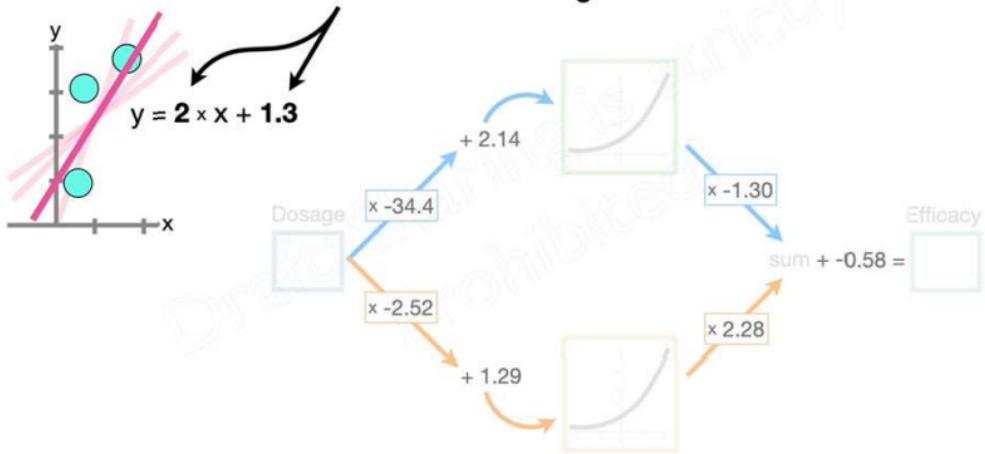
...and **connections**  
between the nodes.



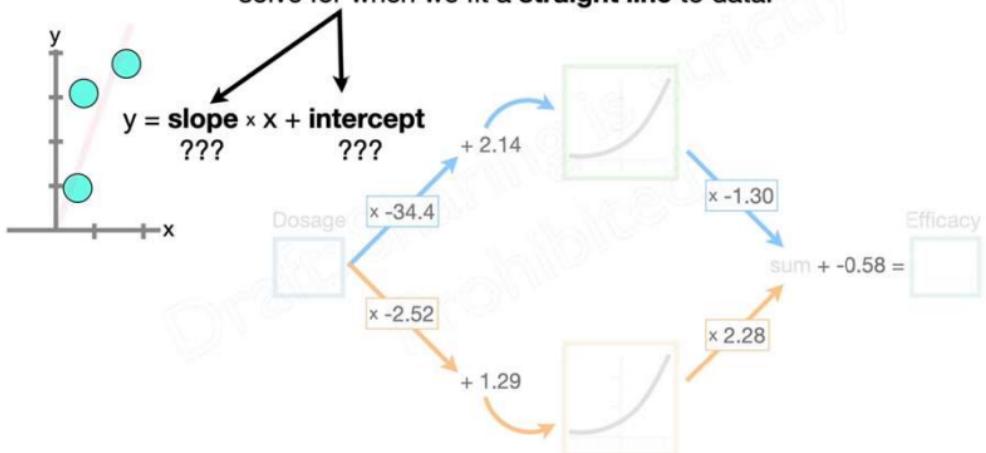
**NOTE:** The numbers along each connection represent parameter values that were estimated when this **Neural Network** was fit to the data.



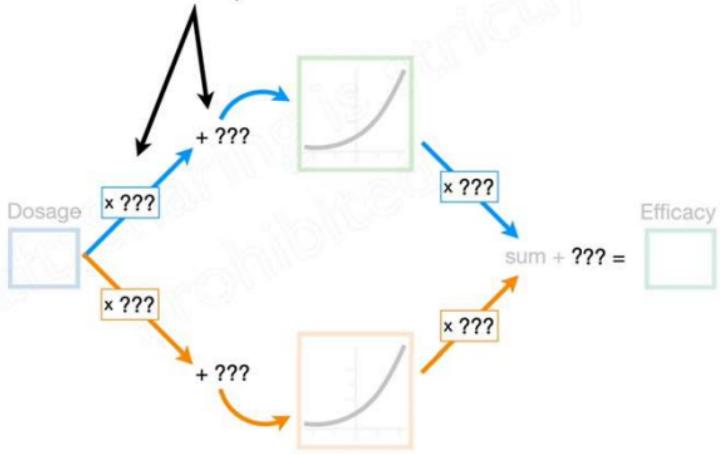
For now, just know that these parameter estimates are analogous to the **slope** and **intercept** values that we solve for when we fit a **straight line** to data.



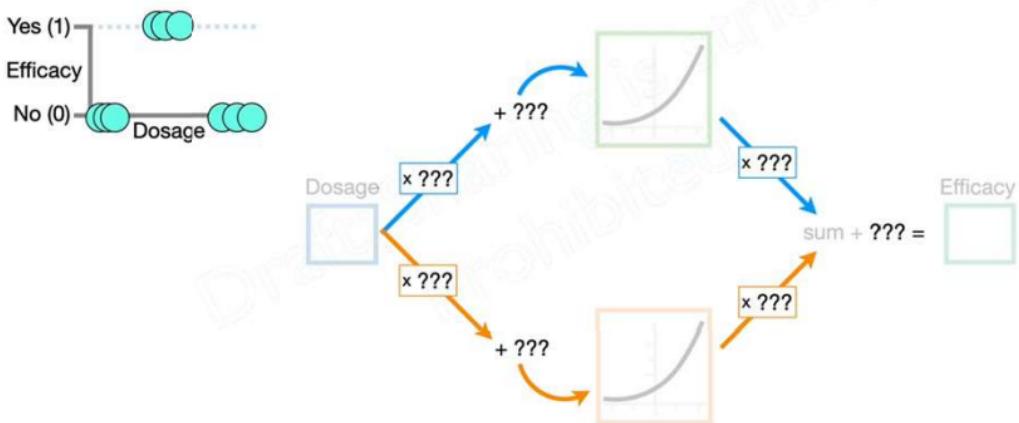
For now, just know that these parameter estimates are analogous to the **slope** and **intercept** values that we solve for when we fit a **straight line** to data.



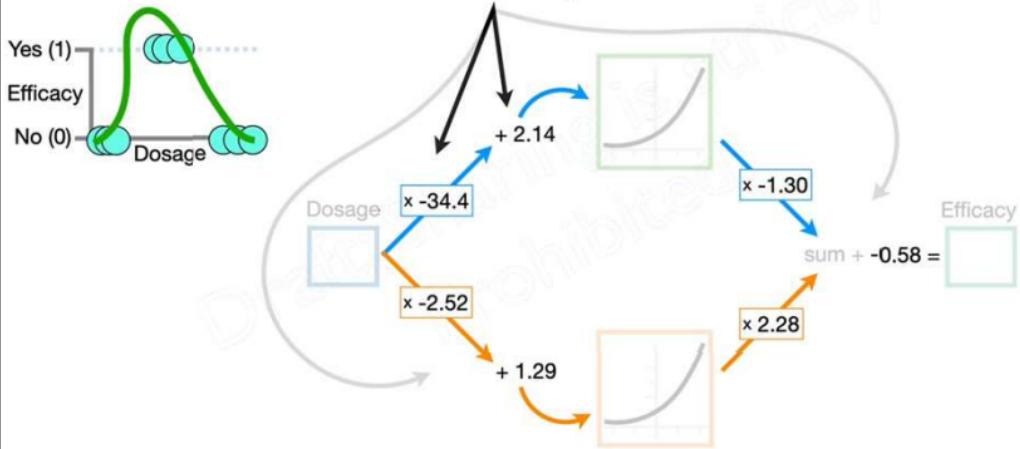
Likewise, a **Neural Network** starts out with unknown parameter values...



...that are estimated when we fit the **Neural Network** to a dataset using a method called **Backpropagation**.



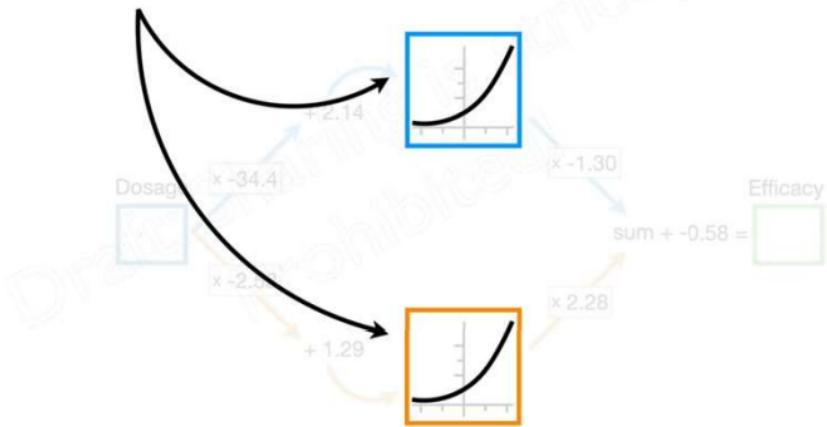
...and that means we have already estimated these parameters.



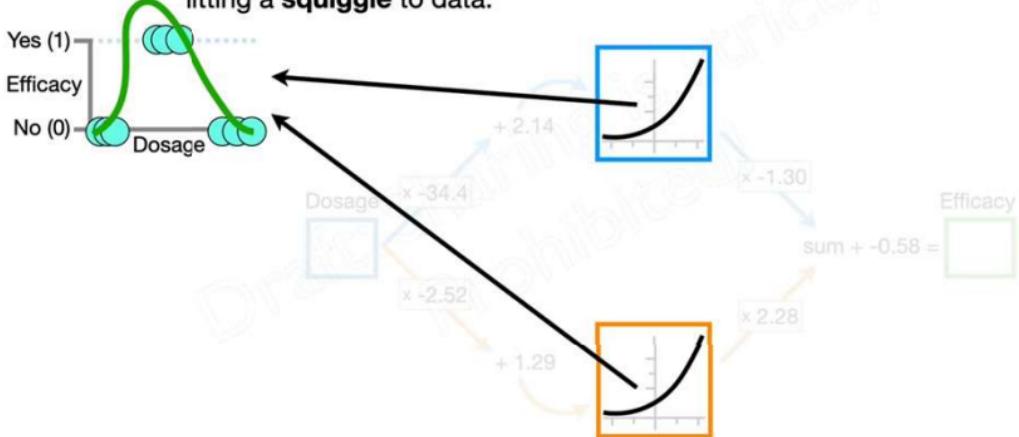
...but for now, just assume that  
we've already fit this **Neural  
Network** to this specific dataset...



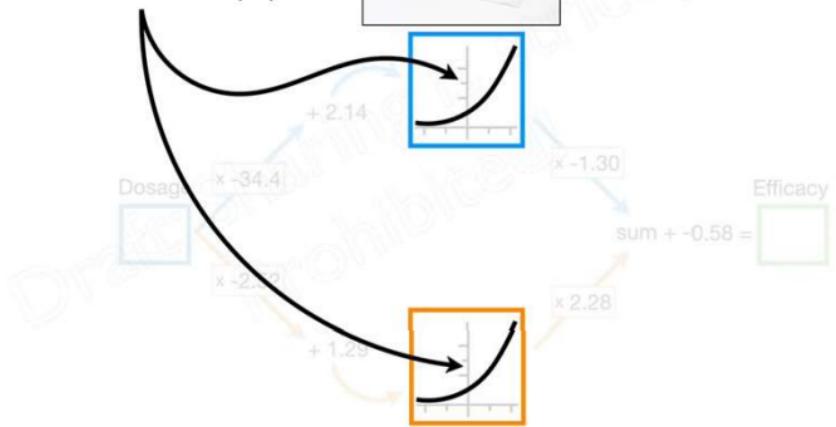
Also, you may have noticed  
that some of the **Nodes** have  
**curved lines** inside of them.



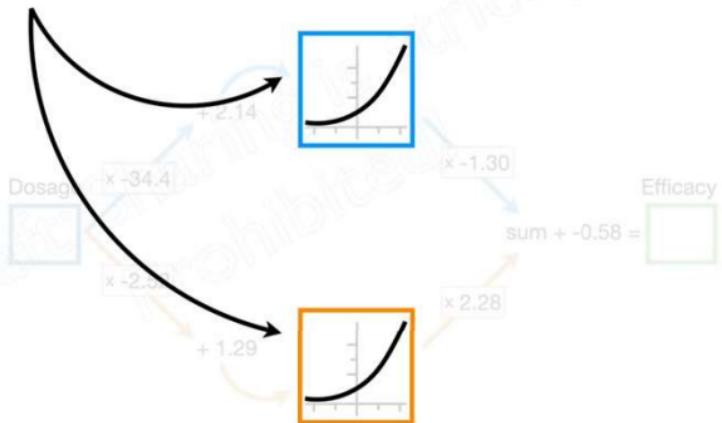
These bent or curved lines  
are the building blocks for  
fitting a squiggle to data.



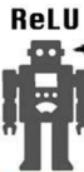
This specific **curved line** is called **softplus**, which sounds like a brand of toilet paper.



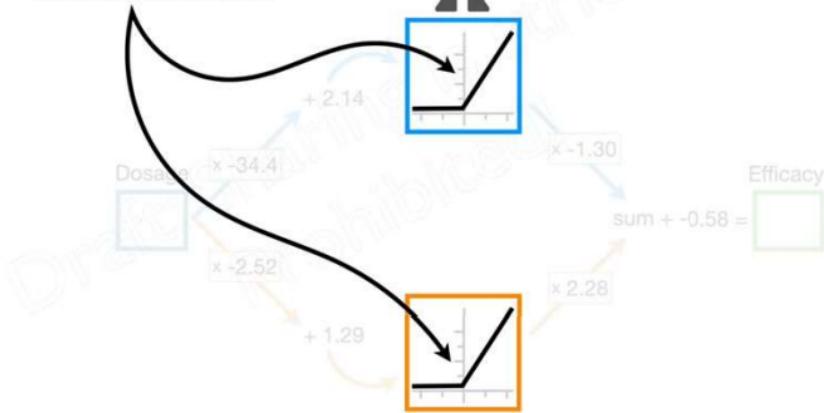
**NOTE:** There are many common bent or curved lines that we can choose for a Neural Network.



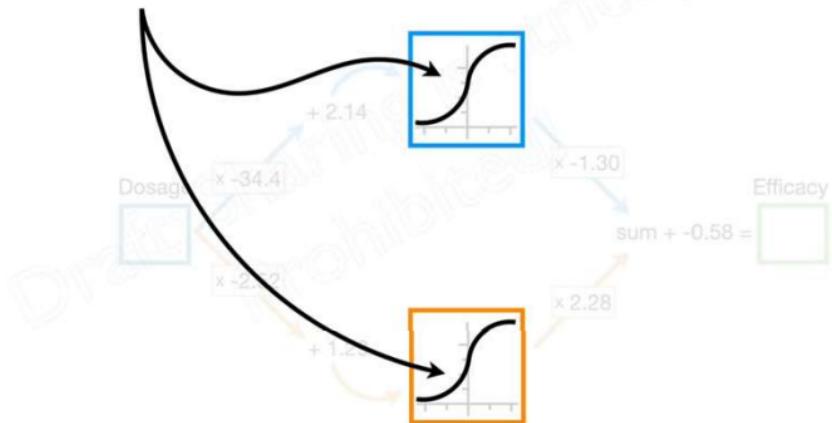
Alternatively, we could use this **bent line**, called, **ReLU**, which is short for **Rectified Linear Unit** and sounds like a robot.



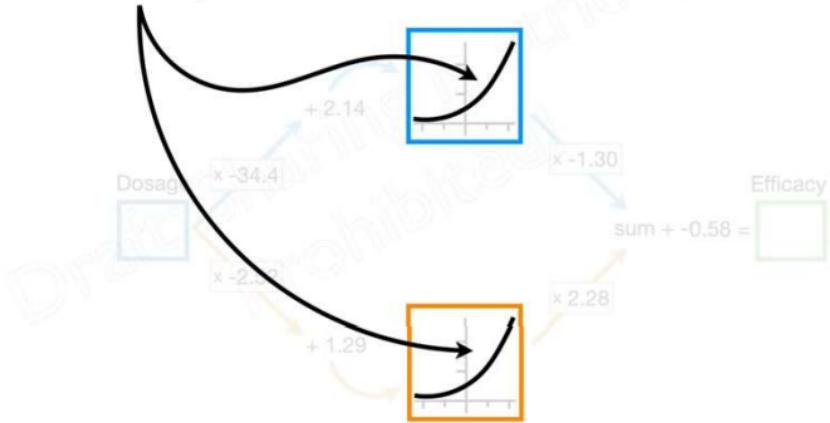
Take me to  
your leader!



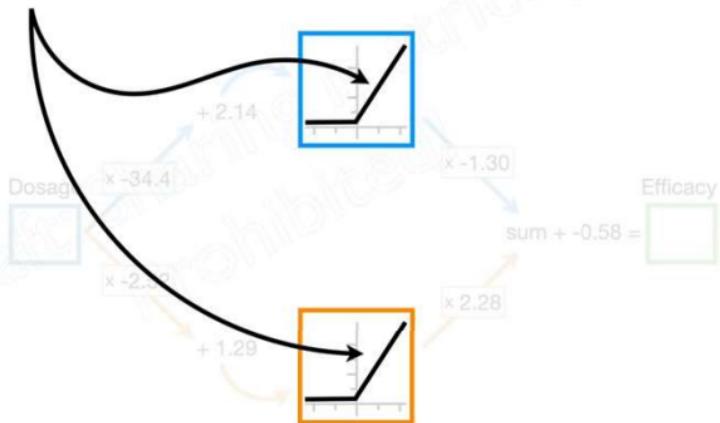
Or we could use a **sigmoid** shape or any other **bent** or **curved line**.



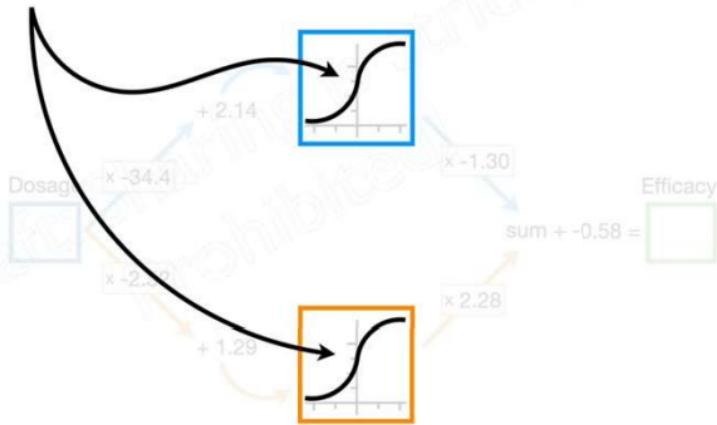
When you build a **Neural Network**, you have to decide which **Activation Function or Functions** you want to use.



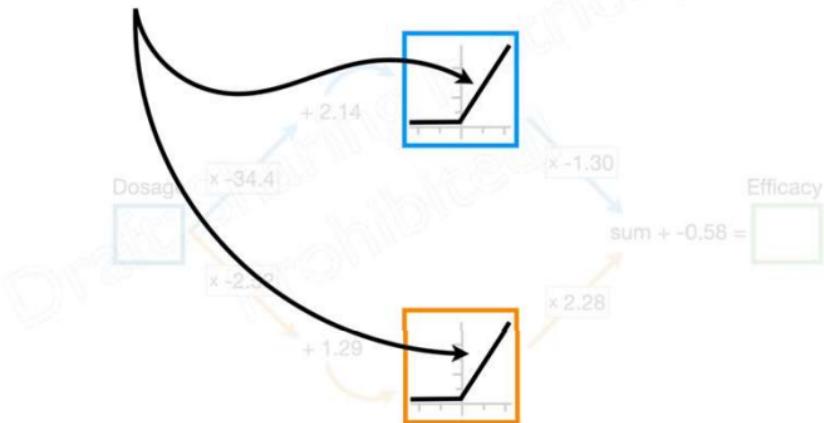
The **curved or bent lines**  
are called **Activation  
Functions.**



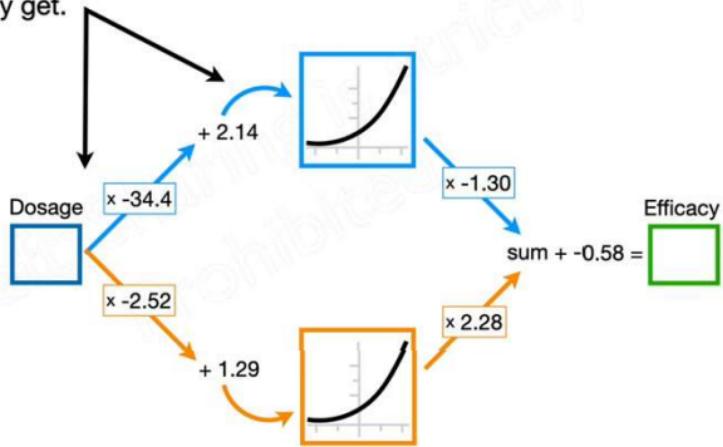
When most people teach **Neural Networks**, they use the **sigmoid Activation Function**.



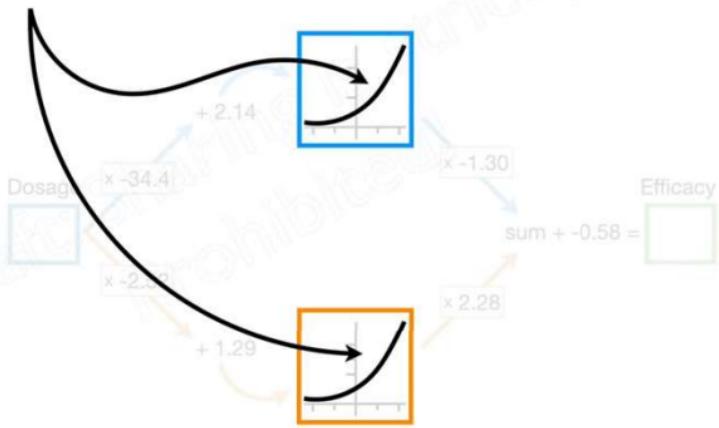
However, in *practice*, it is much more common to use the **ReLU Activation Function**...



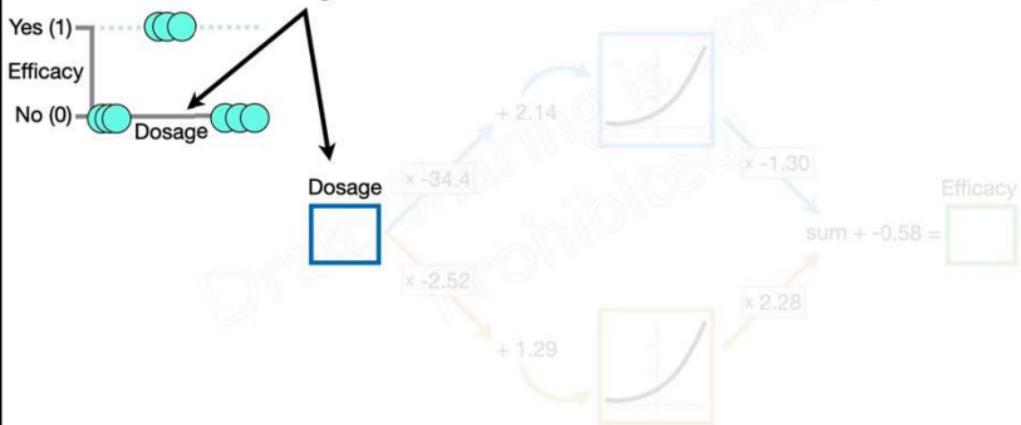
**NOTE:** This specific **Neural Network** is about as simple as they get.



...or the **softplus Activation Function.**



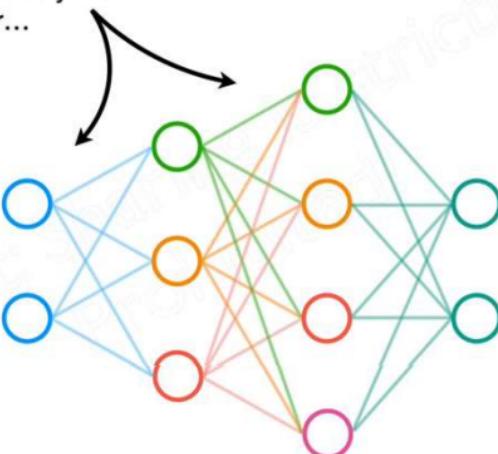
It only has 1 Input Node  
where we plug in the  
**Dosage...**



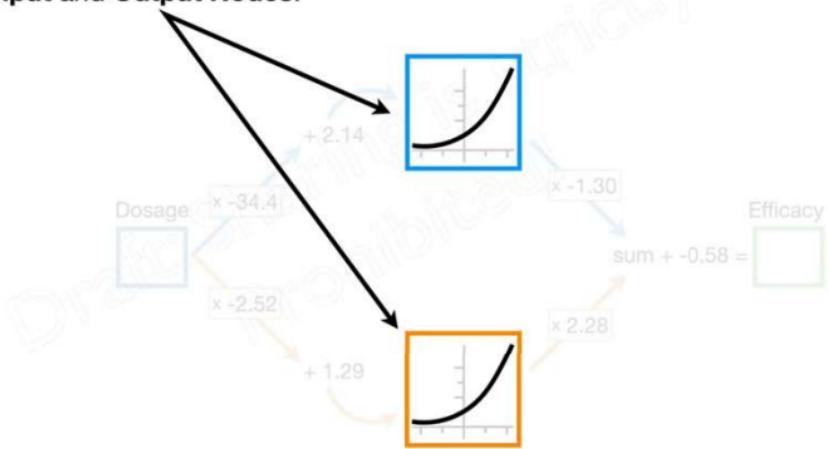
...only 1 Output Node to  
tell us the predicted  
**Effectiveness...**



However, in practice, **Neural Networks** are usually much fancier...



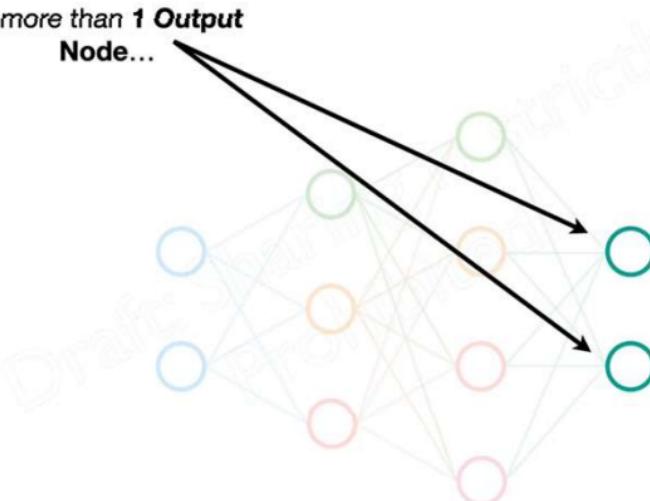
**...and only 2 Nodes between  
the Input and Output Nodes.**



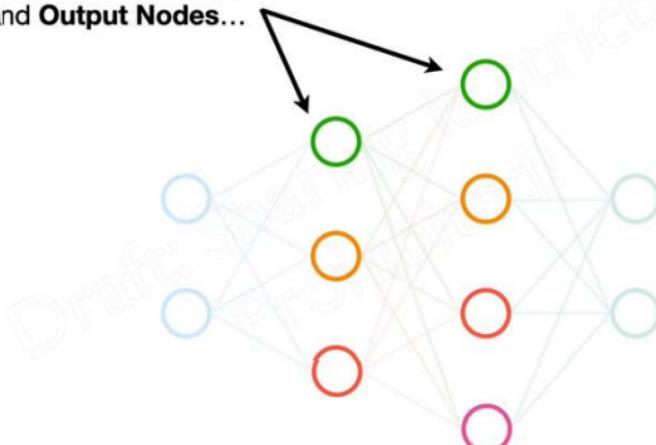
...and have more  
than **1 Input Node**...



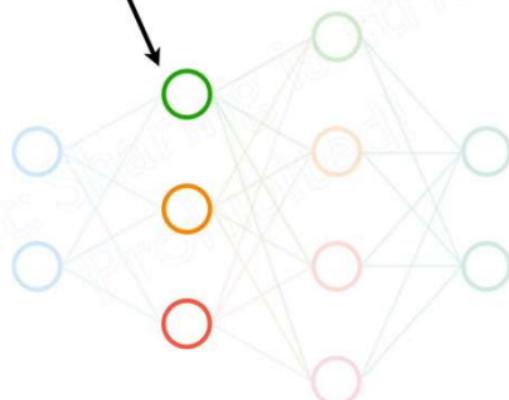
*...more than 1 **Output**  
Node...*



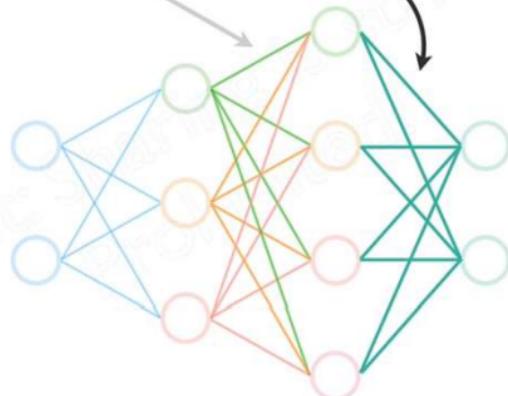
...different layers of  
**Nodes between the Input**  
**and Output Nodes...**



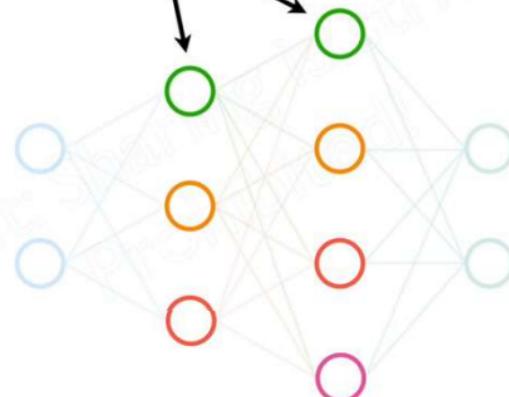
...different layers of  
**Nodes** between the Input  
and Output Nodes...



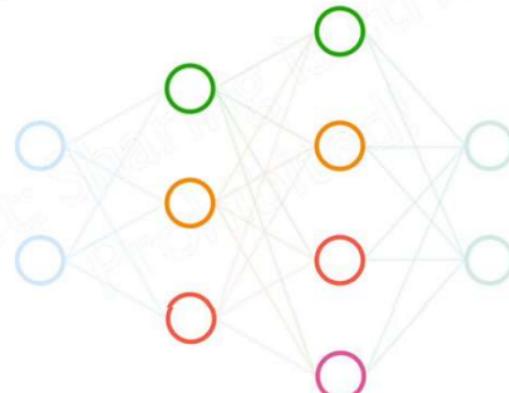
...and a spider web of connections between each layer of **Nodes**.



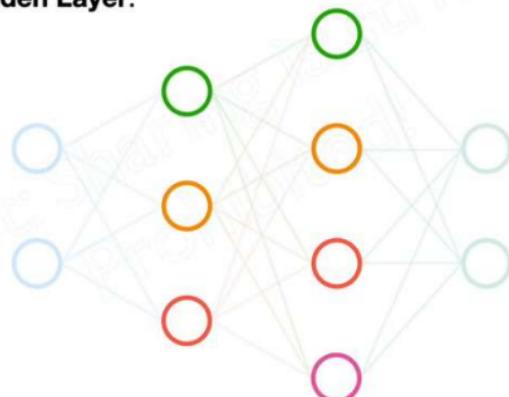
The layers of **Nodes** between  
the **Input and Output Nodes**  
are called **Hidden Layers**.



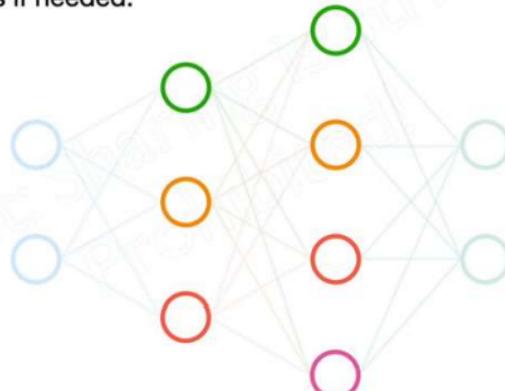
Although there are rules of thumb  
for making decisions about the  
**Hidden Layers...**



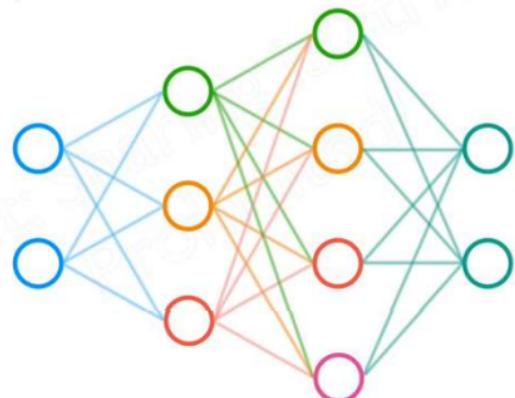
When you build a **Neural Network**,  
one of the first things you do is  
decide how many **Hidden Layers**  
you want, and how many **Nodes**  
go into each **Hidden Layer**.



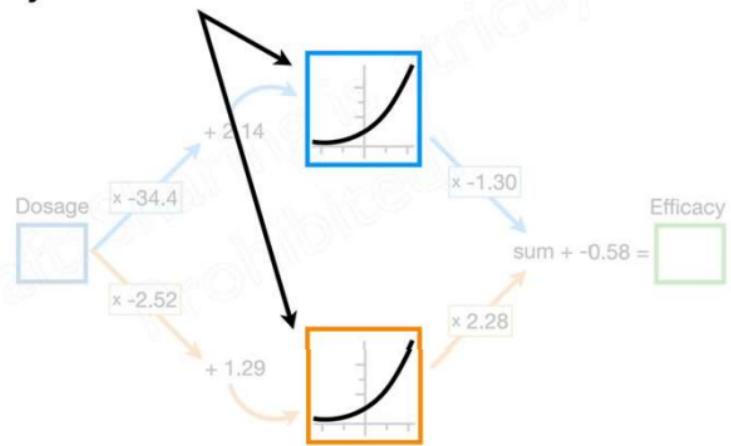
...you essentially make a guess  
and see how well the **Neural**  
**Network** performs, adding more  
layers and nodes if needed.



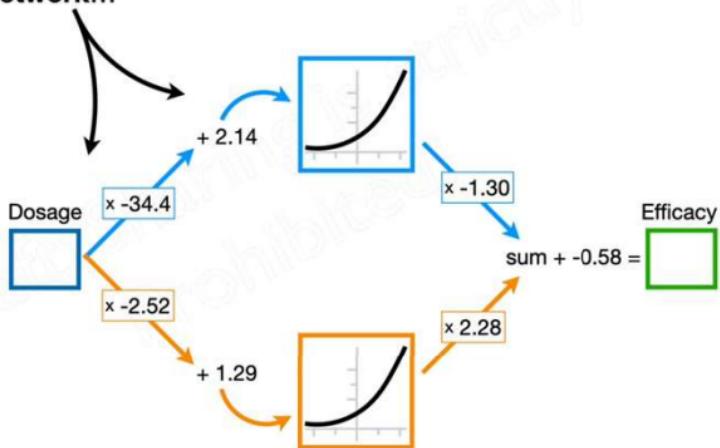
Now, even though this **Neural Network** looks fancy, it is still made from the same parts...



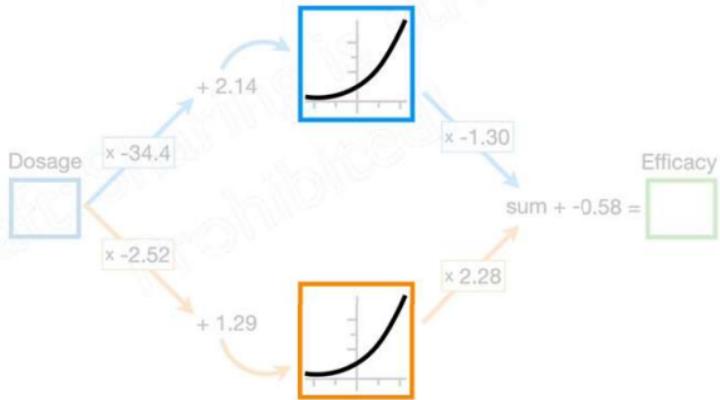
...which only has **1 Hidden Layer** with **2 Nodes**.



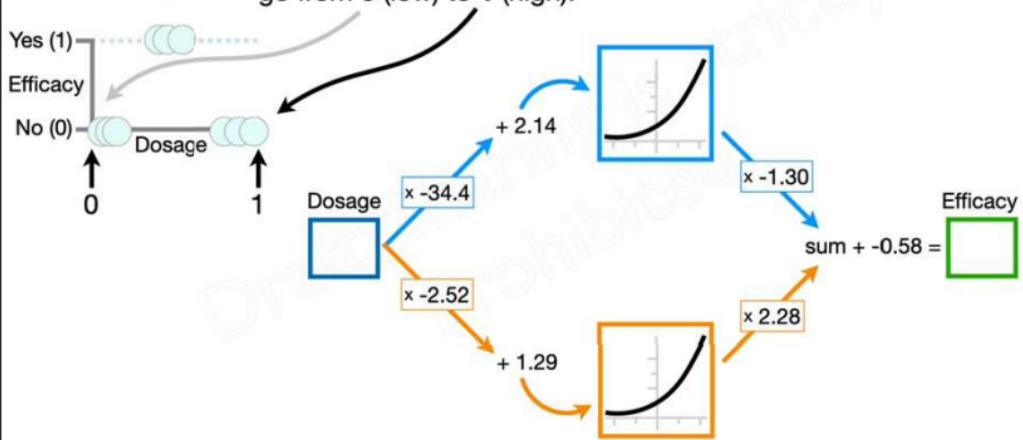
...used in this simple  
**Neural Network**...



So let's learn how this **Neural Network** creates new shapes from the **curved or bent lines** in the **Hidden Layer**...

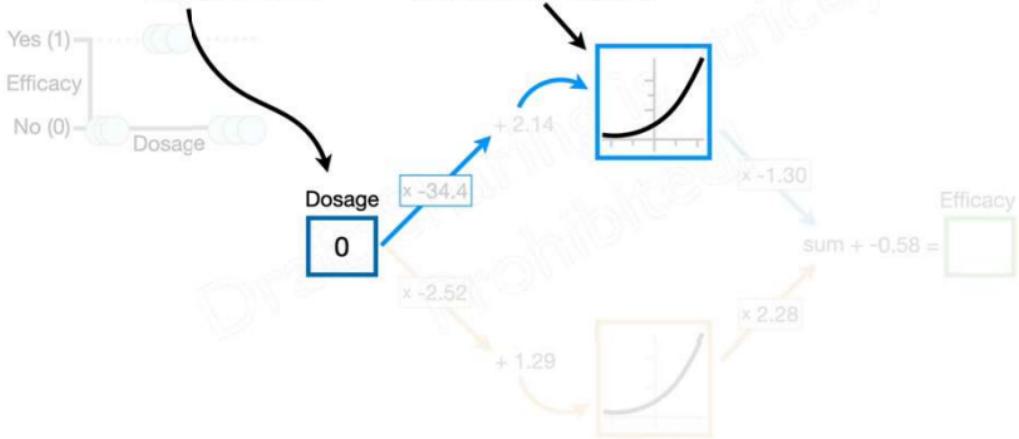


**NOTE:** To keep the math simple, let's assume **Dosages** go from **0** (low) to **1** (high).

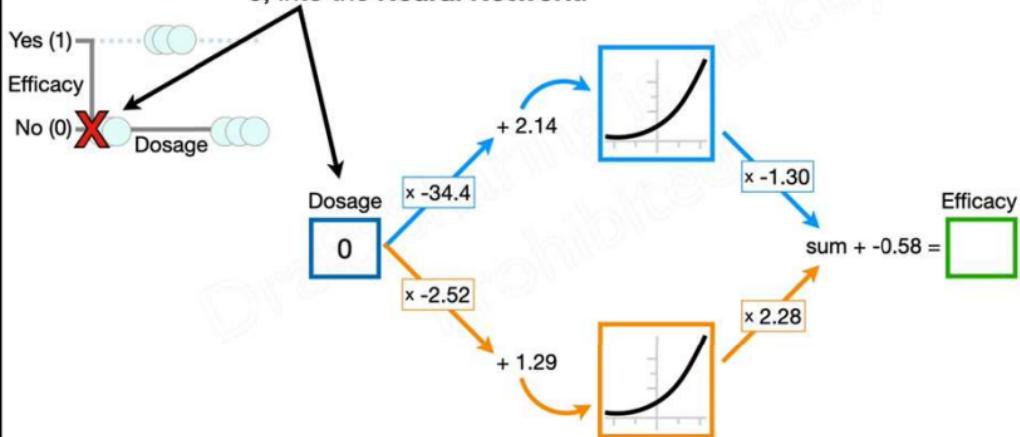


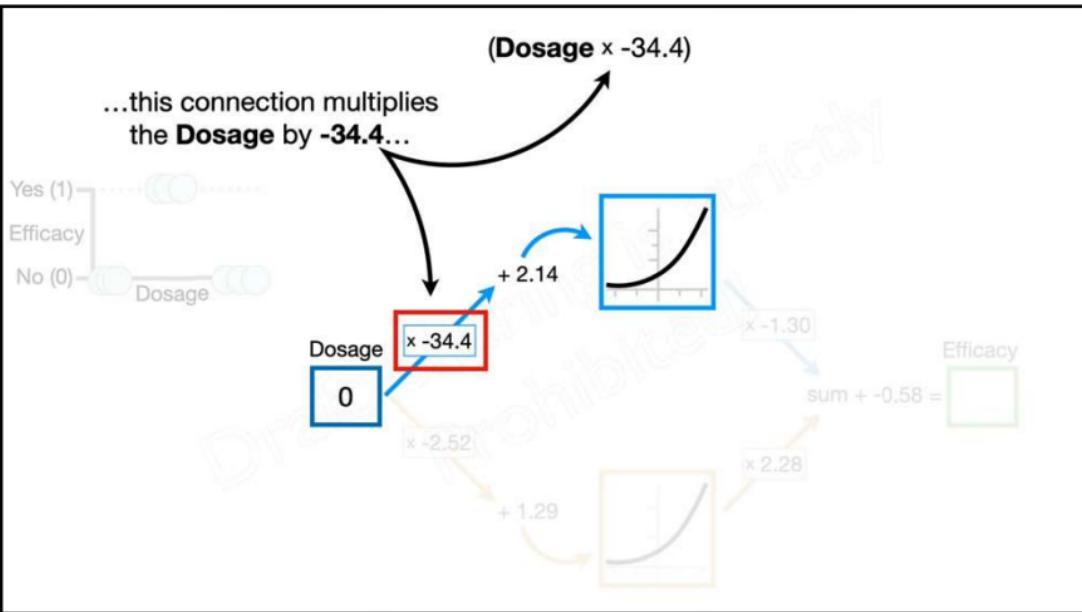
Now, to get from the  
**Input Node**...

...to the top **Node** in  
the **Hidden Layer**...



The first thing we are going to do is plug the lowest **Dosage**, **0**, into the **Neural Network**.





$$(\text{Dosage} \times -34.4) + 2.14 = \text{x-axis coordinate}$$

...and the result is an x-axis coordinate for the **Activation Function**.



Dosage  
0

$x -34.4$

$x -2.52$

+ 1.29



$x -1.30$

sum

+ -0.58 =

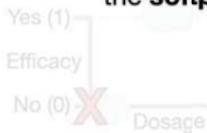
Efficacy



...we plug **2.14** into the **Activation Function**, which, in this case, is the **softplus** function.

$$0 + 2.14 = 2.14$$

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

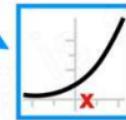


Dosage  
0

$x -34.4$

$x -2.52$

$+ 1.29$



$+ 2.14$

$x -1.30$

$\text{sum} + -0.58 =$

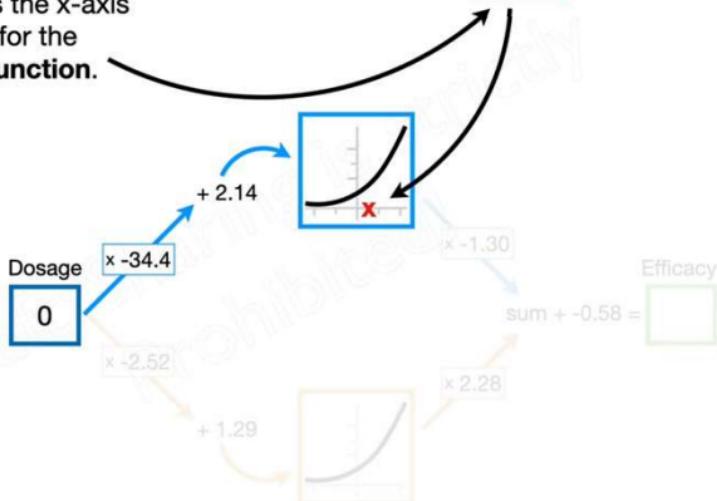


$x 2.28$

$$0 + 2.14 = 2.14$$

...to get **2.14** as the x-axis coordinate for the  
**Activation Function.**

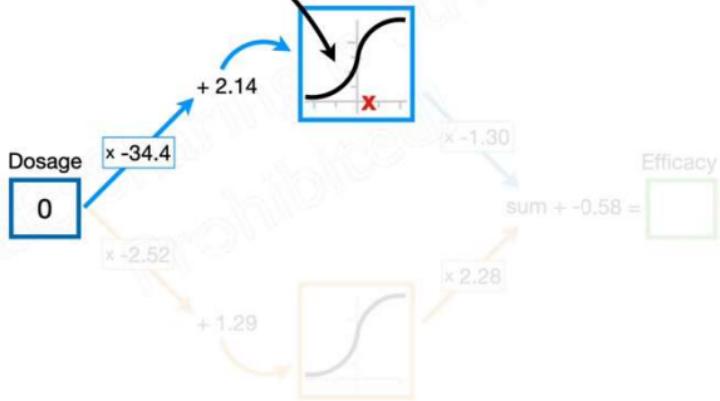
Yes (1)  
Efficacy  
No (0)  Dosage



**NOTE:** If we had chosen the **sigmoid** curve for the **Activation Function**...

$$0 + 2.14 = 2.14$$

Yes (1)  
Efficacy  
No (0) X  
Dosage



...then we would plug  
**2.14** into the equation  
for the **sigmoid** curve...

$$0 + 2.14 = 2.14$$

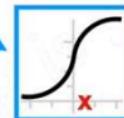
$$f(x) = \frac{e^x}{e^x + 1}$$



Dosage  
0

$$x - 34.4$$

$$+ 2.14$$



$$x - 1.30$$

Efficacy  
sum + -0.58 =

$$x - 2.52$$

$$+ 1.29$$



$$x 2.28$$

...then we would plug  
**2.14** into the **ReLU**  
equation.

$$0 + 2.14 = 2.14$$

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

Yes (1)  
Efficacy  
No (0) X  
Dosage

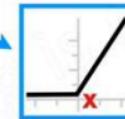
Dosage  
0

x -34.4

x -2.52

+ 1.29

+ 2.14



x -1.30

sum + -0.58 =

x 2.28

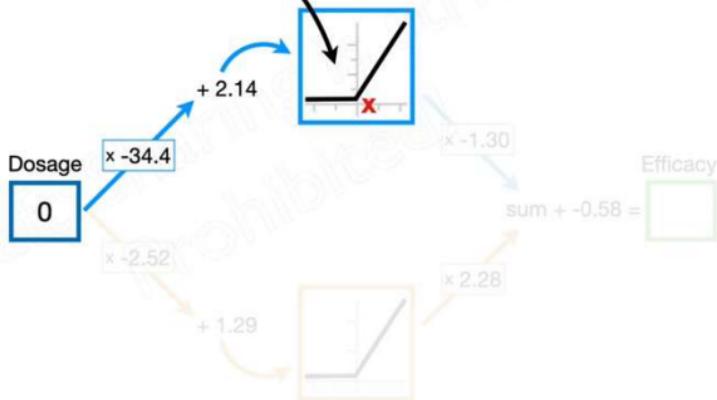
Efficacy  
□



$$0 + 2.14 = 2.14$$

...and if had chosen the  
**ReLU bent line** for the  
**Activation Function...**

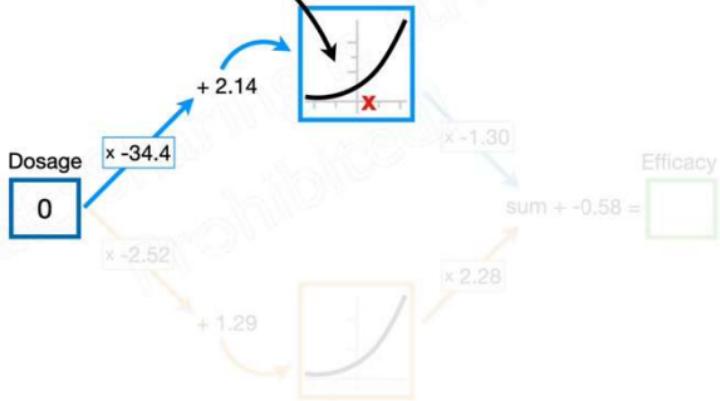
Yes (1)  
Efficacy  
No (0) ~~X~~  
Dosage



$$0 + 2.14 = 2.14$$

But since we are using  
**softplus** for the  
Activation Function...

Yes (1)  
Efficacy  
No (0) X  
Dosage



...we plug **2.14** into the  
**softplus** equation...

$$0 + 2.14 = 2.14$$

$$f(2.14) = \log(1 + e^{2.14})$$

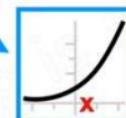
Yes (1)  
Efficacy  
No (0) X  
Dosage

Dosage  
0

$$x - 34.4$$

$$x - 2.52$$

$$+ 1.29$$



$$x - 1.30$$

sum

+ -0.58 =

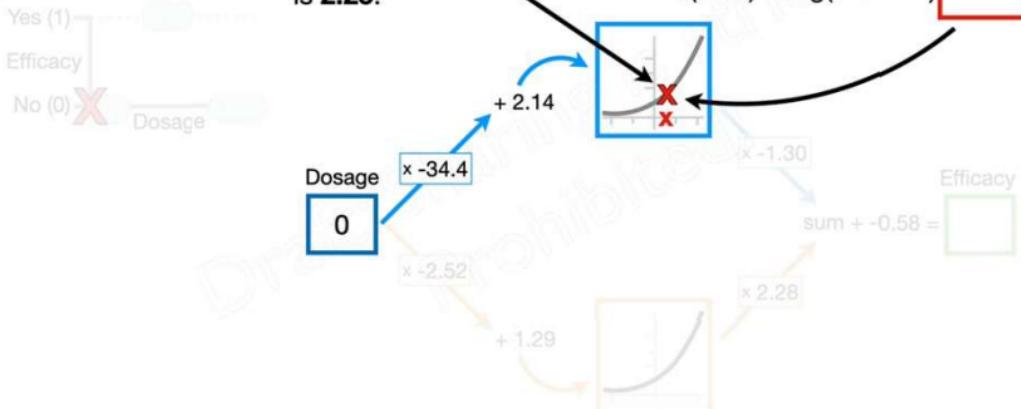
Efficacy  
□

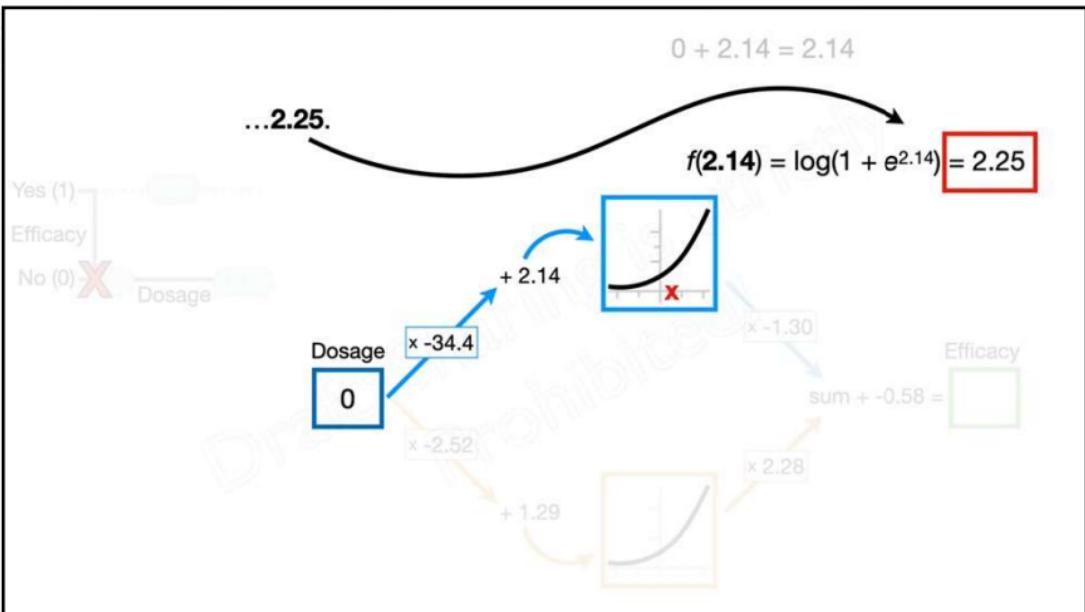


$$x 2.28$$

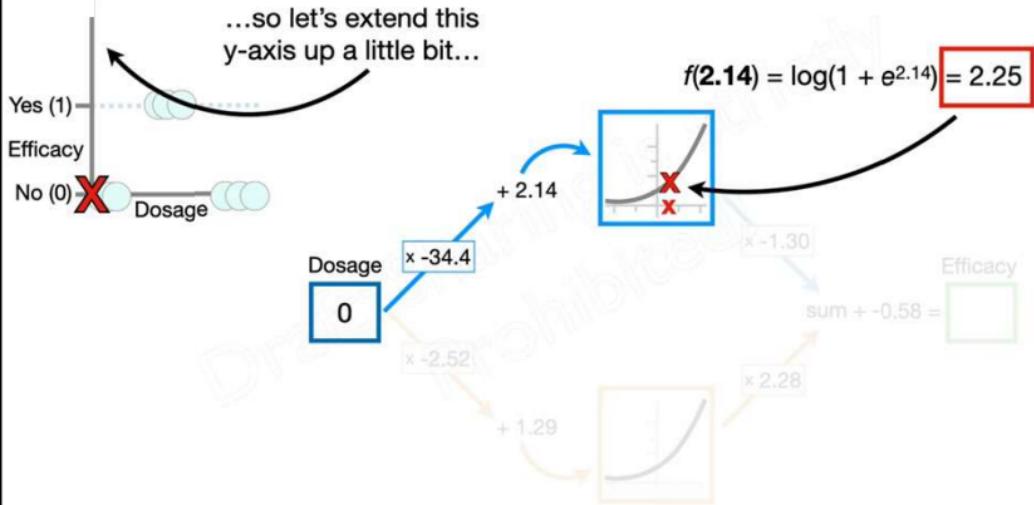
Anyway, the y-axis coordinate for the **Activation Function** is 2.25.

$$f(2.14) = \log(1 + e^{2.14}) = 2.25$$

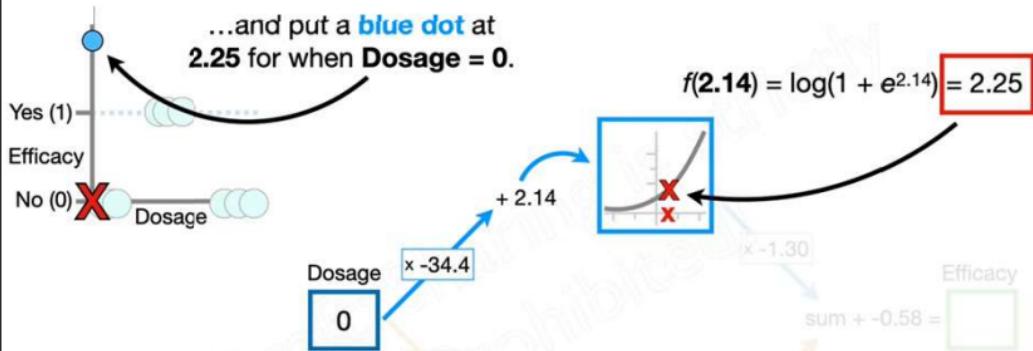




$$(0 \times -34.4) + 2.14 = 2.14$$

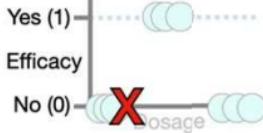
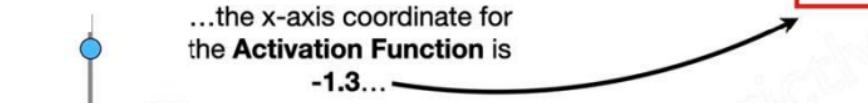


$$(0 \times -34.4) + 2.14 = 2.14$$



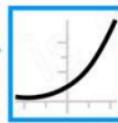
$$(0.1 \times -34.4) + 2.14 = -1.3$$

...the x-axis coordinate for  
the **Activation Function** is  
**-1.3...**



Dosage  
**0.1**

$$\begin{matrix} x -34.4 \\ + 2.14 \end{matrix}$$



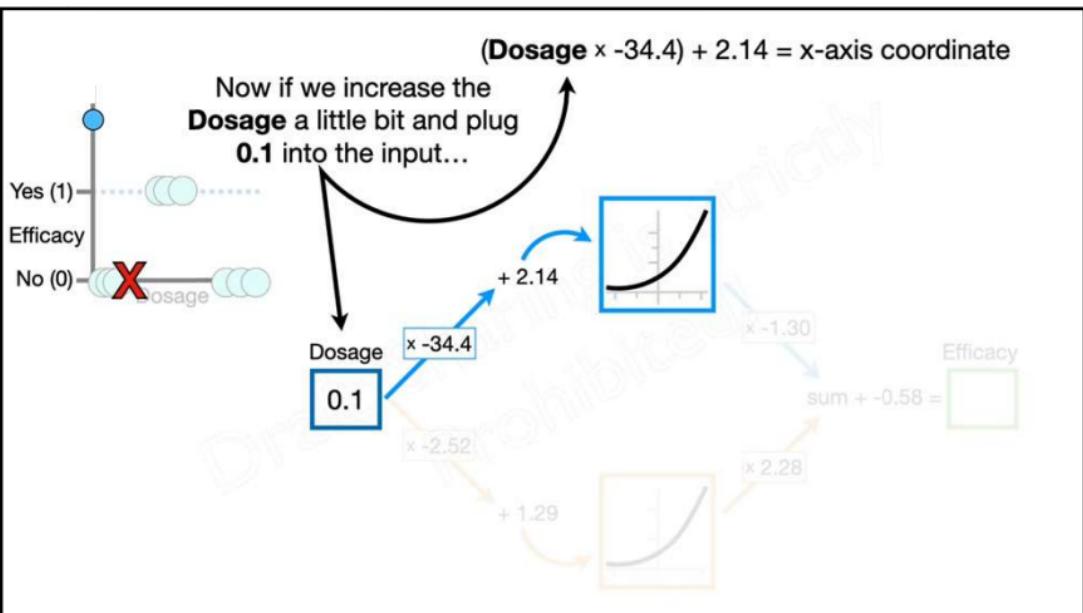
$$x -1.30$$

$$\text{sum} + -0.58 =$$

Efficacy

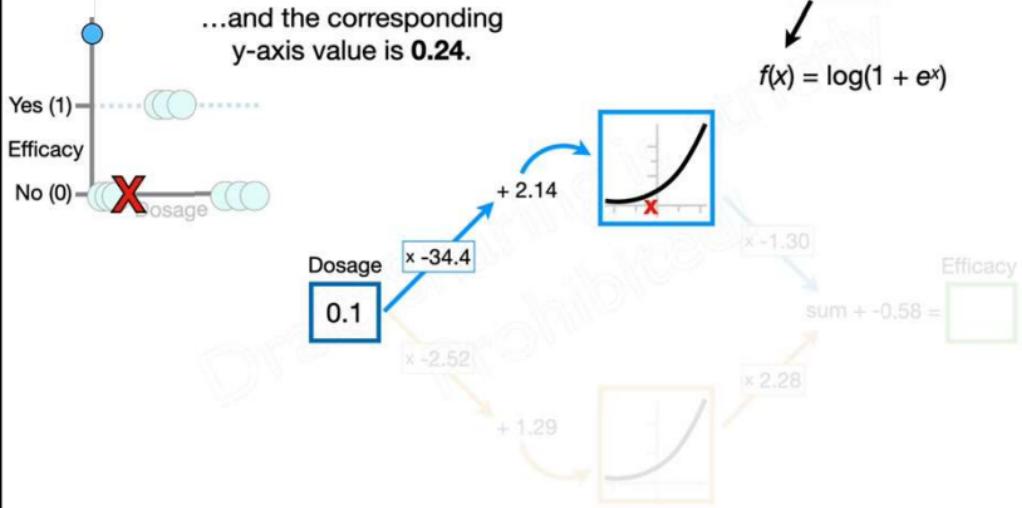


$$x 2.28$$



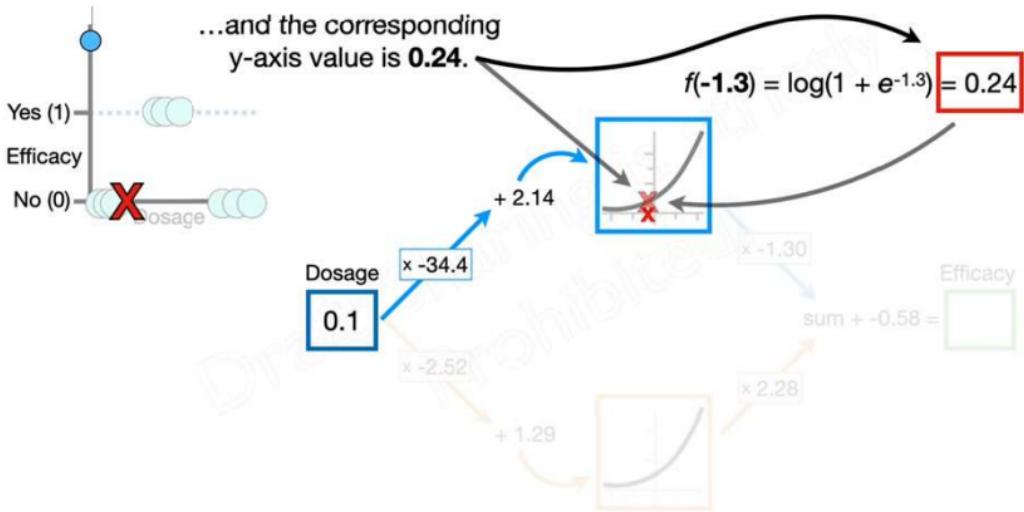
$$(0.1 \times -34.4) + 2.14 = -1.3$$

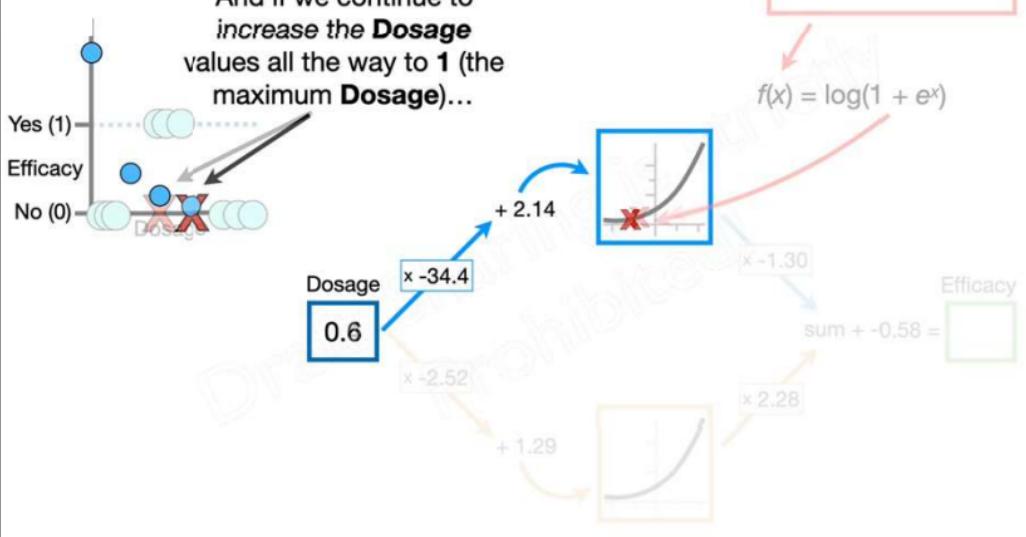
$$f(x) = \log(1 + e^x)$$



$$(0.1 \times -34.4) + 2.14 = -1.3$$

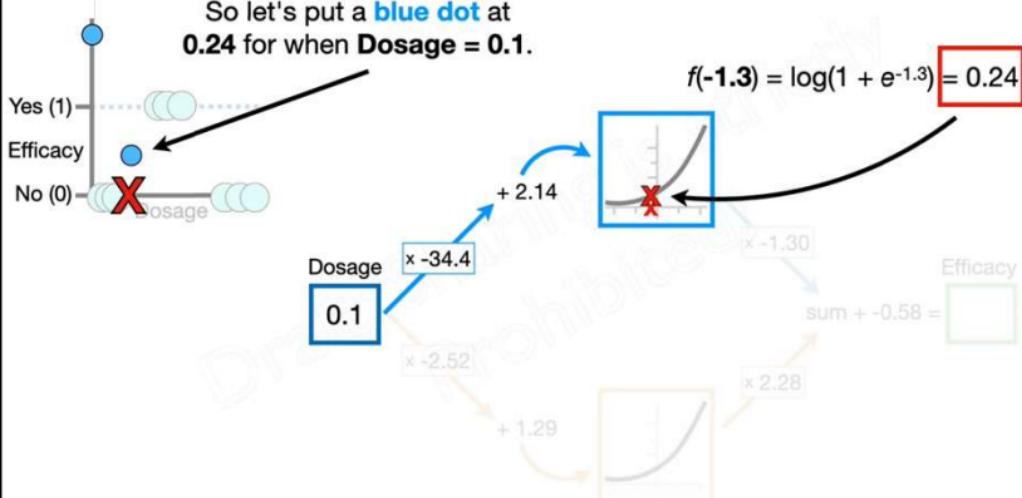
...and the corresponding  
y-axis value is **0.24**.

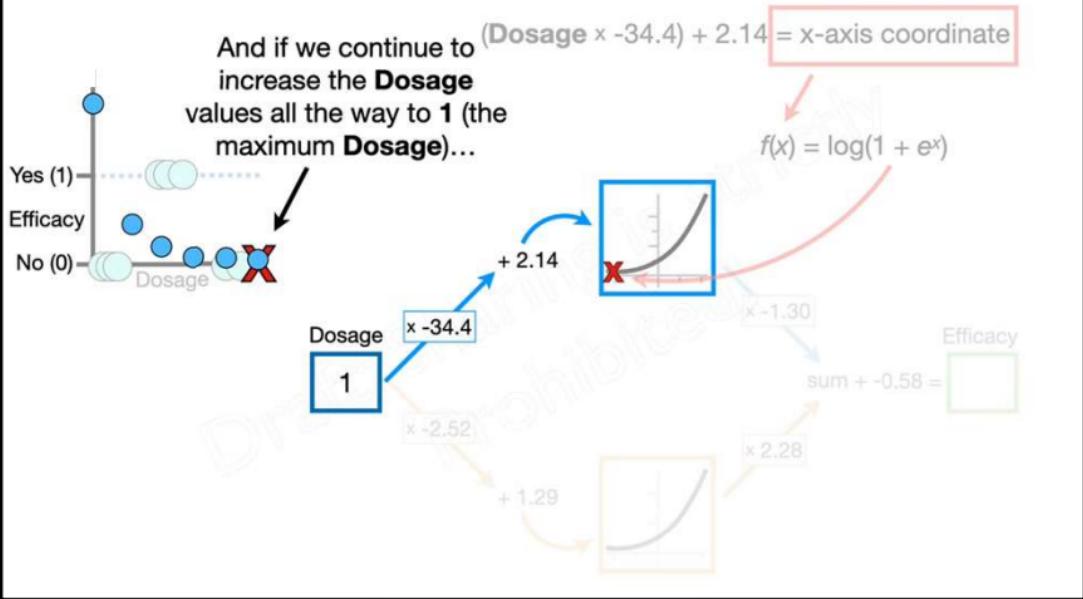


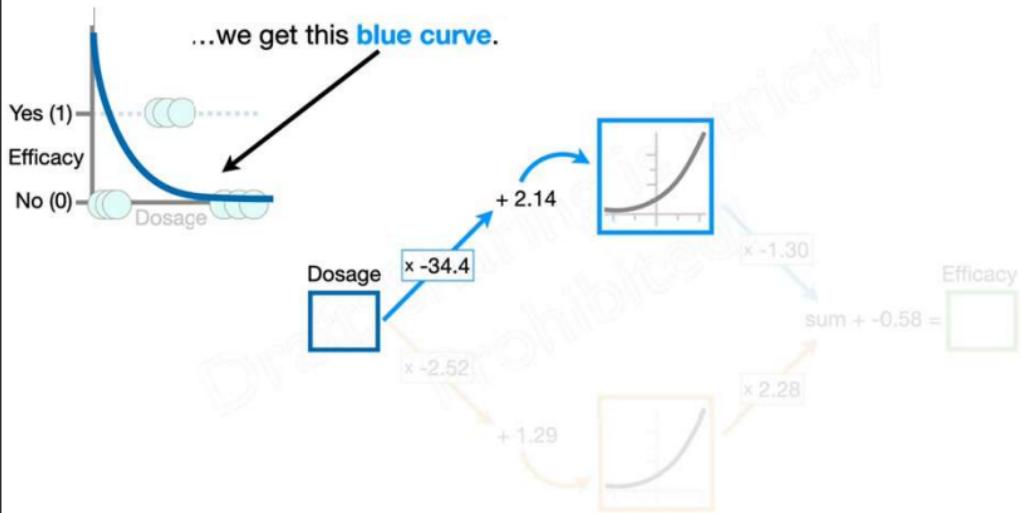


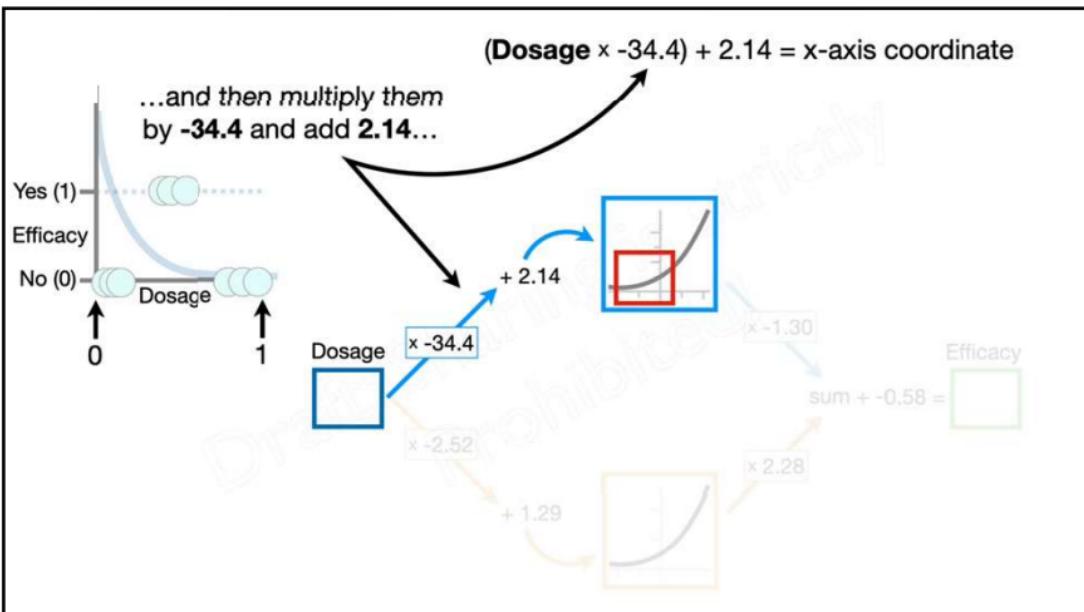
$$(0.1 \times -34.4) + 2.14 = -1.3$$

So let's put a **blue dot** at  
**0.24** for when **Dosage = 0.1**.

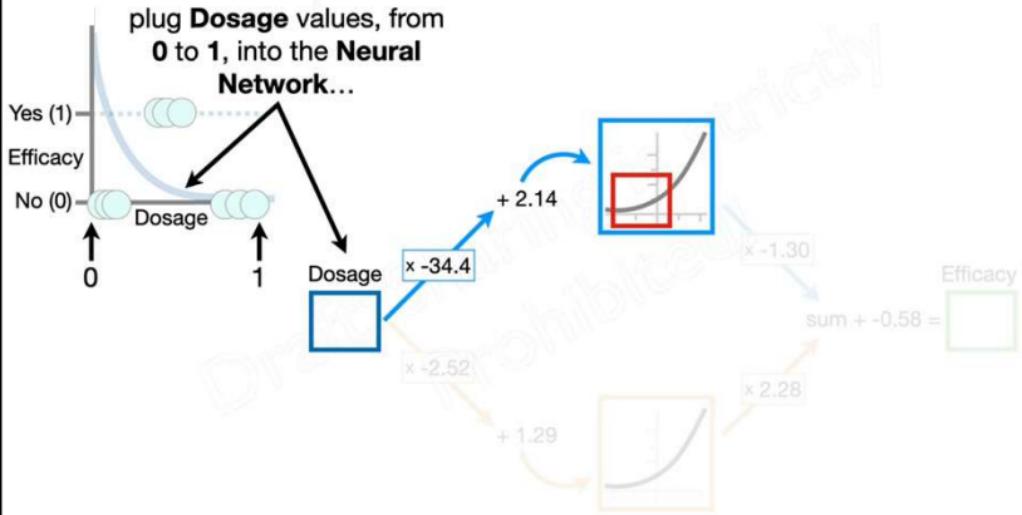






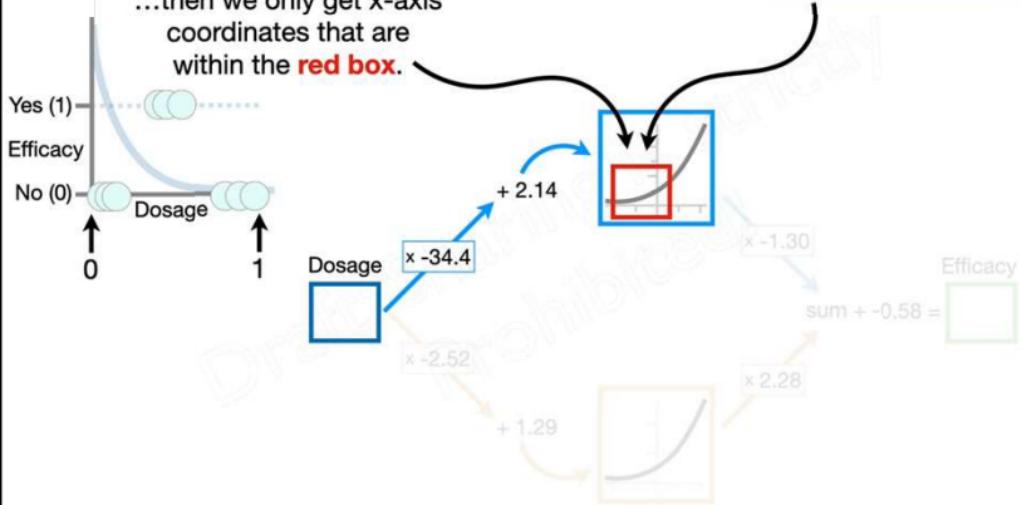


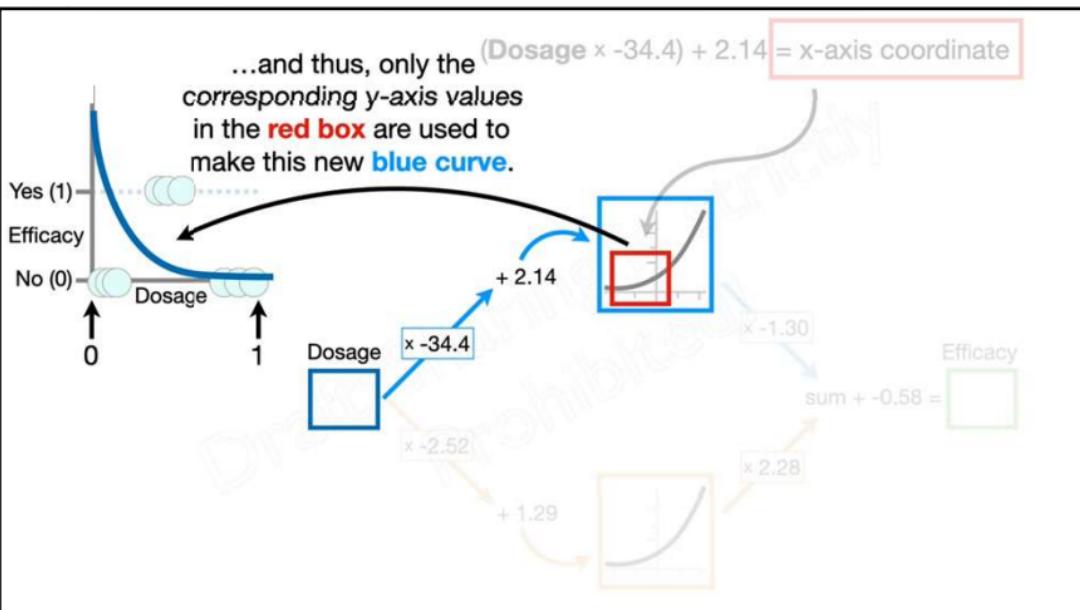
In other words, when we plug **Dosage** values, from **0 to 1**, into the **Neural Network**...



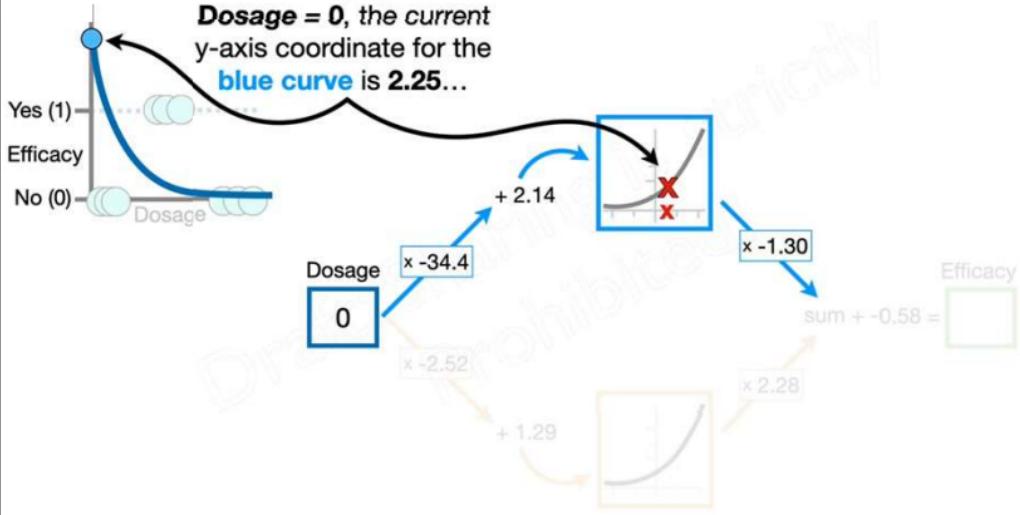
$$(\text{Dosage} \times -34.4) + 2.14 = \text{x-axis coordinate}$$

...then we only get x-axis coordinates that are within the red box.

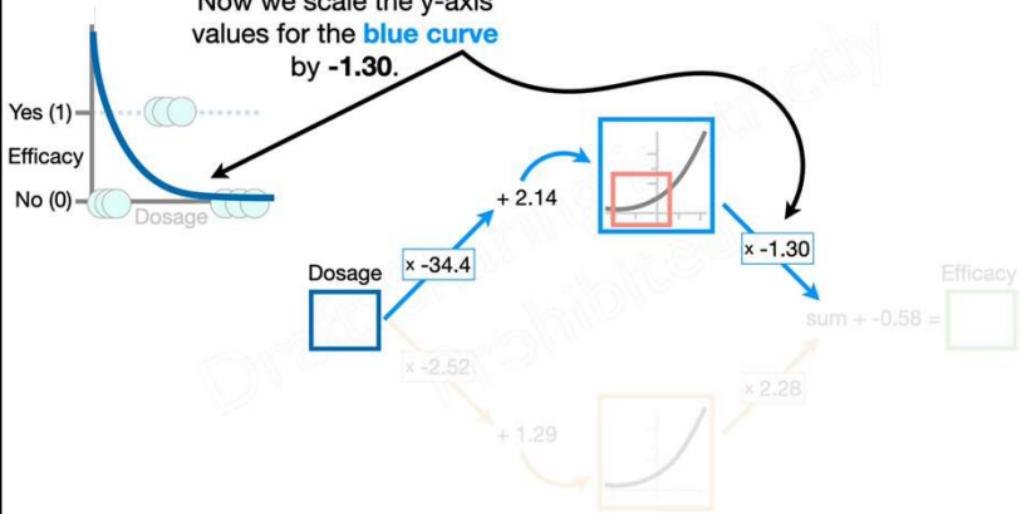


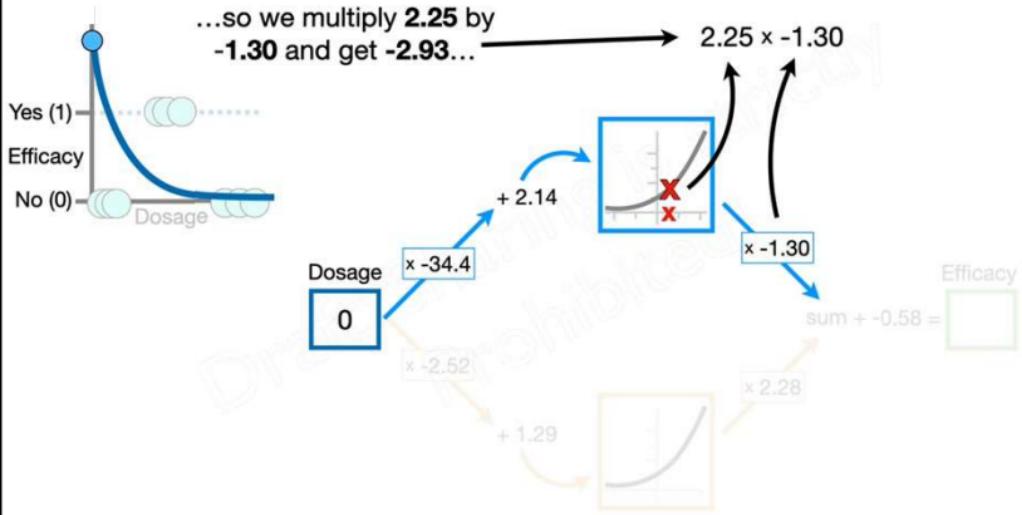


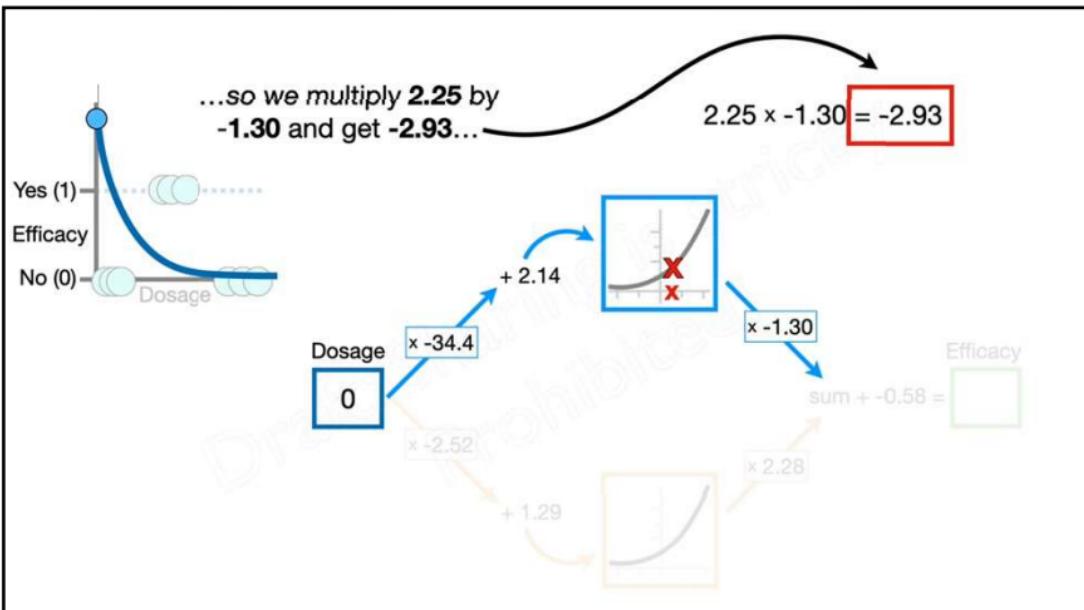
For example, when  
**Dosage = 0**, the current  
y-axis coordinate for the  
**blue curve** is 2.25...



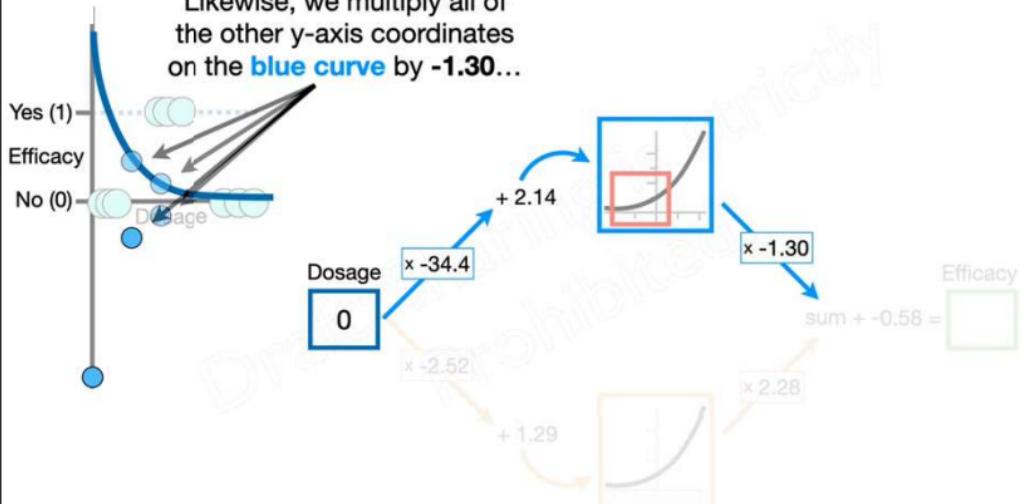
Now we scale the y-axis values for the **blue curve** by -1.30.

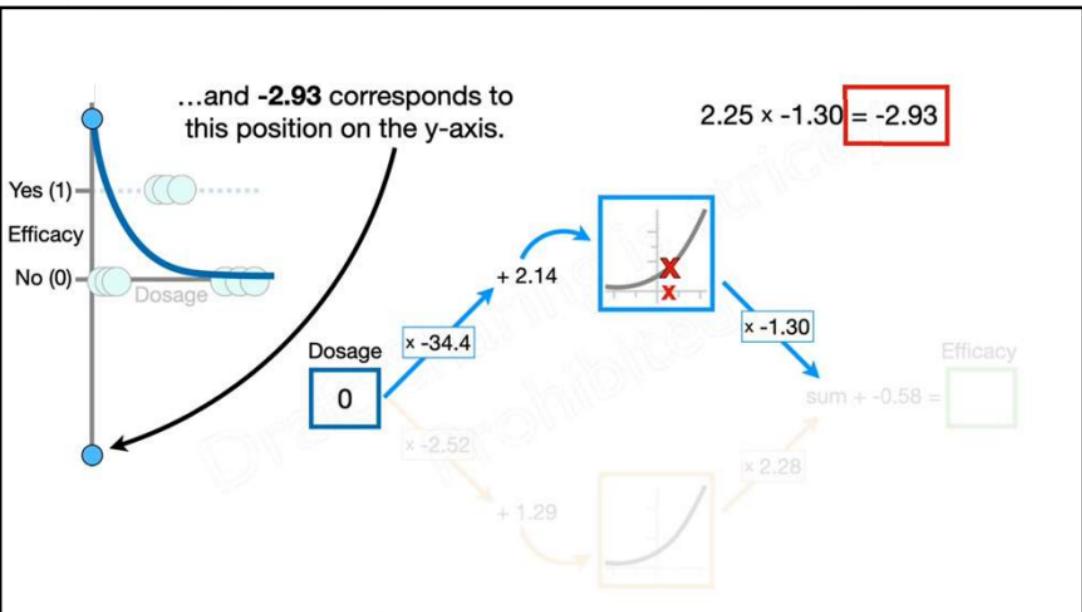




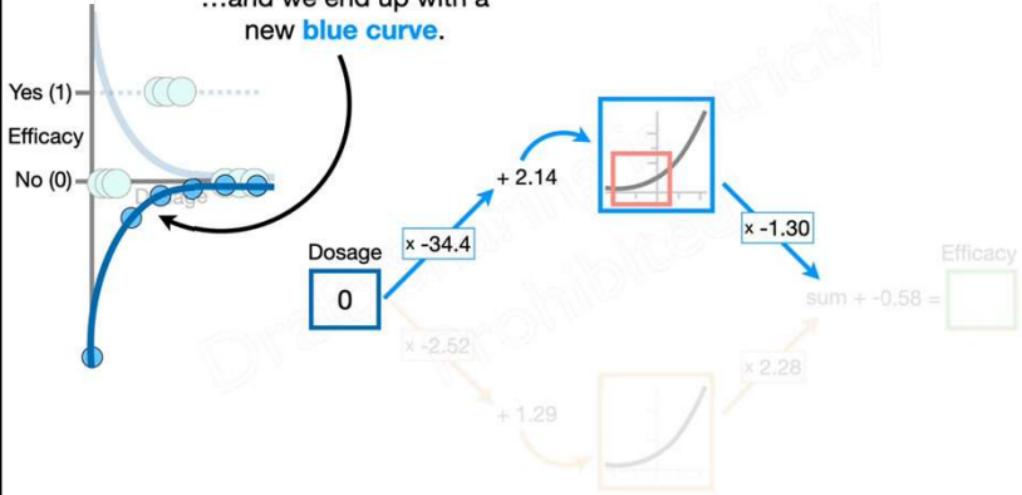


Likewise, we multiply all of the other y-axis coordinates on the **blue curve** by **-1.30...**

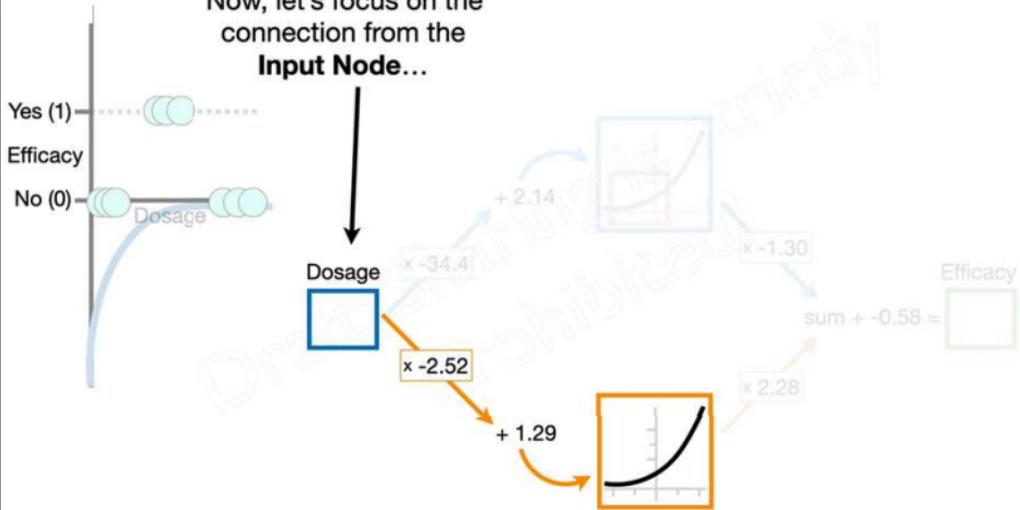




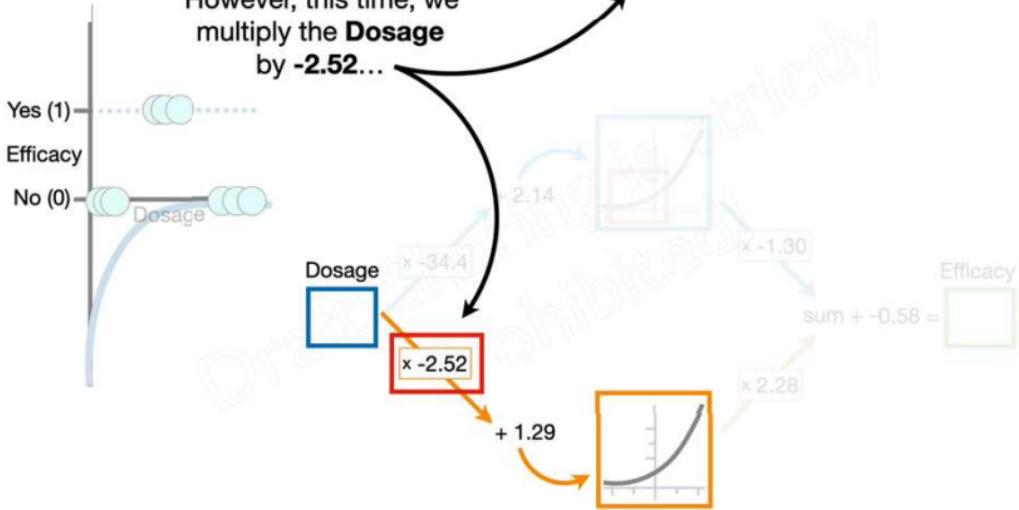
...and we end up with a new **blue curve**.



Now, let's focus on the connection from the **Input Node...**

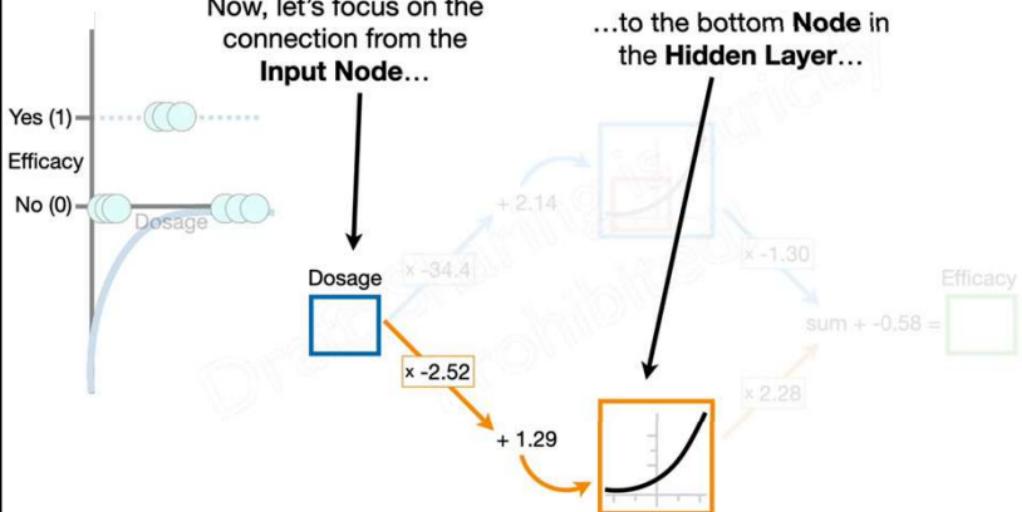


However, this time, we multiply the **Dosage** by **-2.52...**



Now, let's focus on the connection from the **Input Node**...

...to the bottom **Node** in the **Hidden Layer**...



(Dosage  $\times$  -2.52)

... instead of -34.4...

Yes (1) 

Efficacy

No (0)  Dosage

Dosage

$\times$  -34.4

$\times$  -2.52

+ 2.14

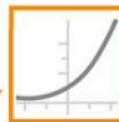
+ 1.29

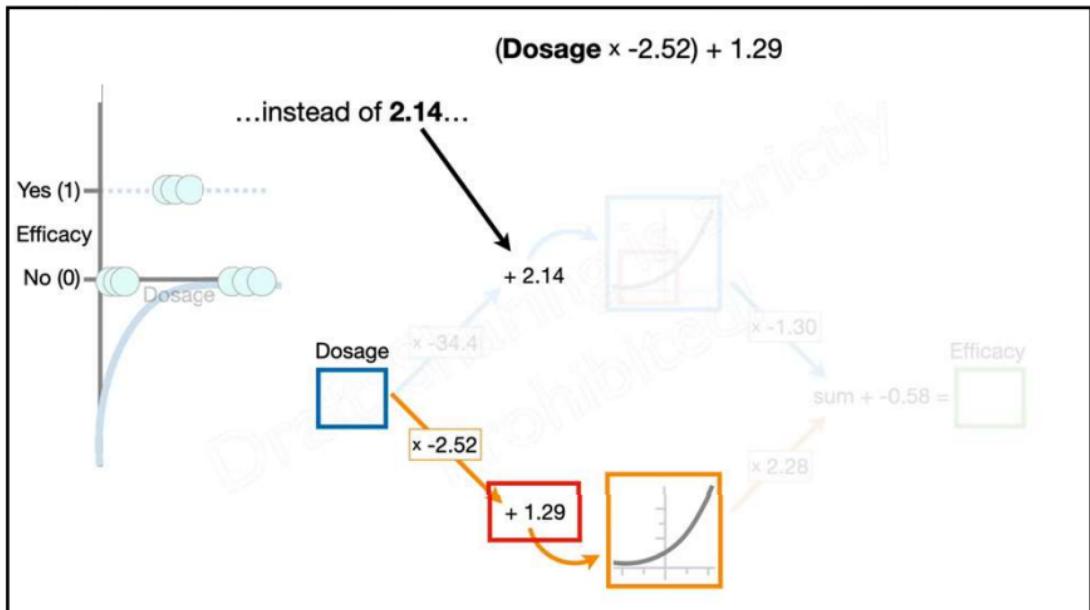
$\times$  -1.30

sum + -0.58 =

Efficacy 

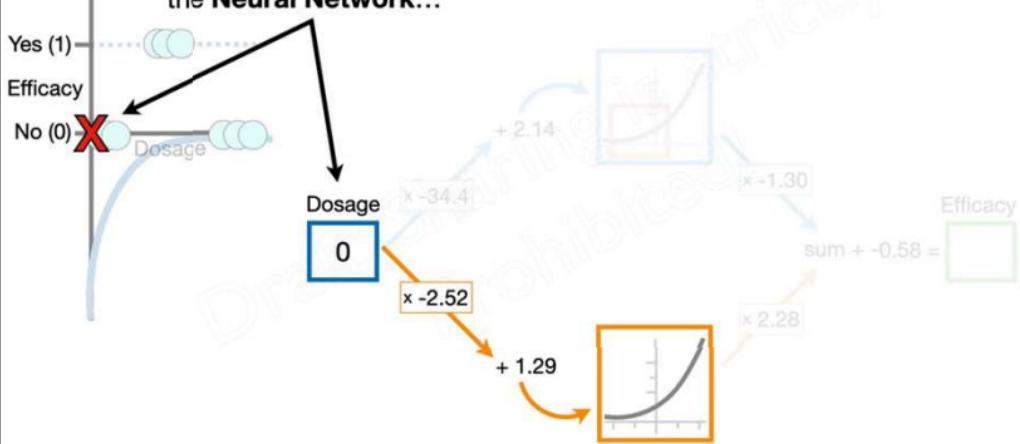
$\times$  2.28

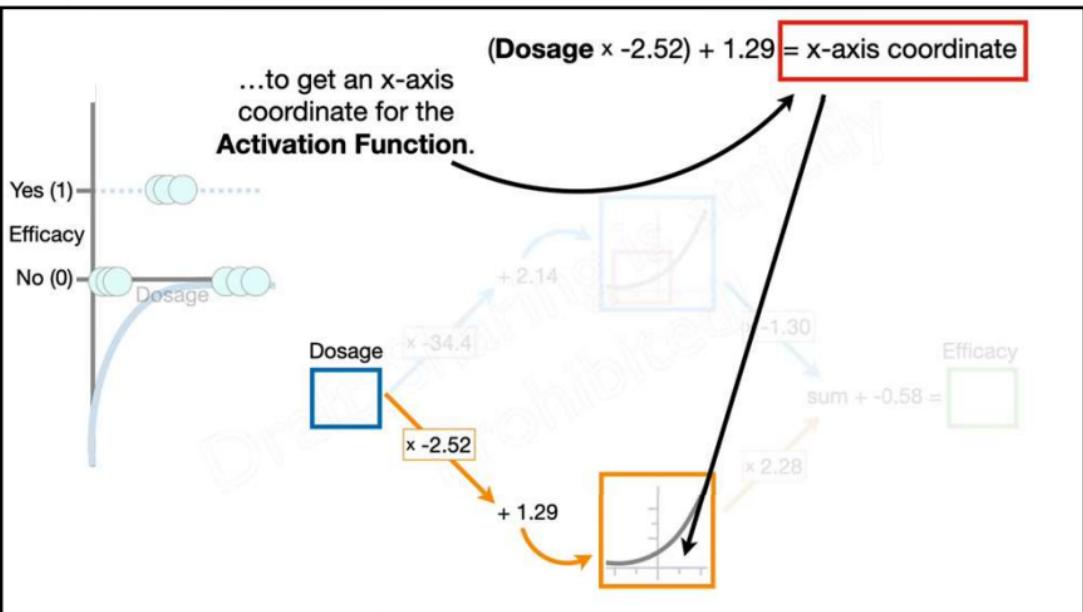




$$(\text{Dosage} \times -2.52) + 1.29 = \text{x-axis coordinate}$$

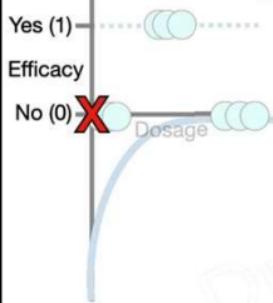
Now, if we plug the lowest **Dosage**, 0, into the **Neural Network**...





$$(0 \times -2.52) + 1.29 = 1.29$$

...then the x-axis coordinate for the **Activation Function** is **1.29**.



Dosage  
0

$$x - 2.52$$

$$+ 1.29$$



$$+ 2.14$$

$$x - 34.4$$

$$\text{sum} + -0.58 =$$

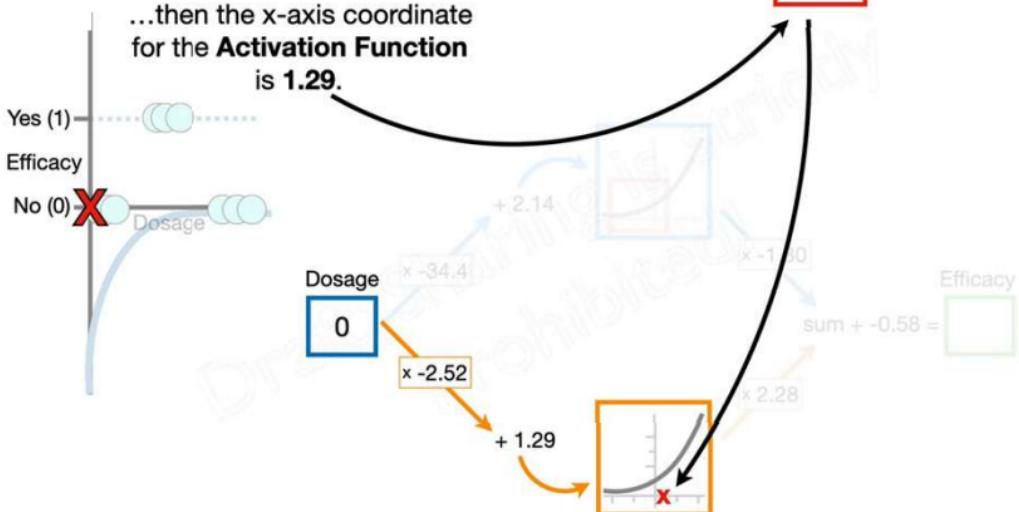
$$x - 1.30$$

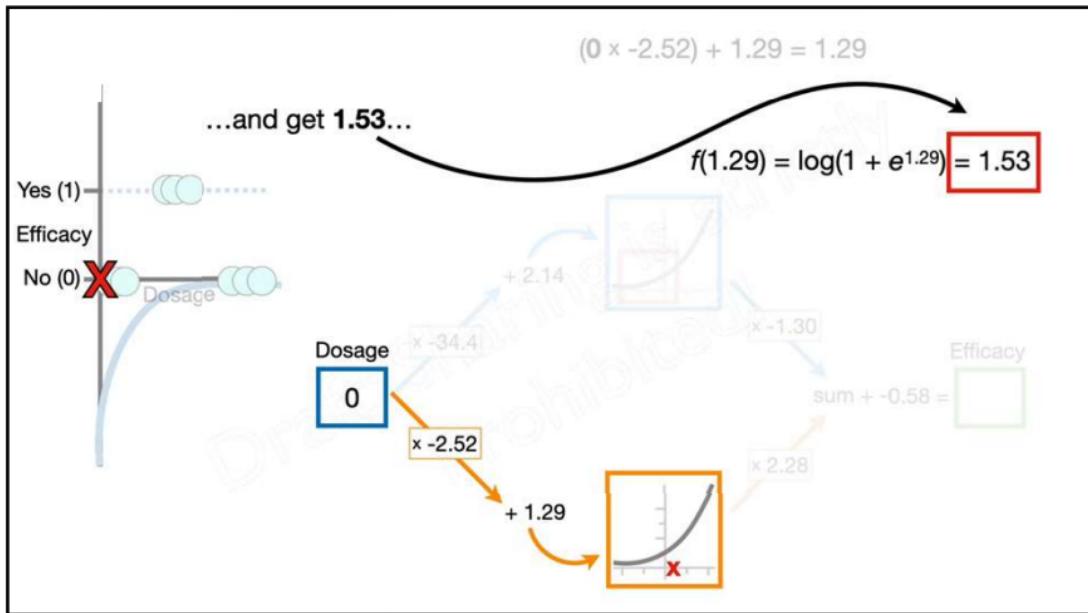
$$x 2.28$$

Efficacy

$$(0 \times -2.52) + 1.29 = 1.29$$

...then the x-axis coordinate  
for the **Activation Function**  
is **1.29**.

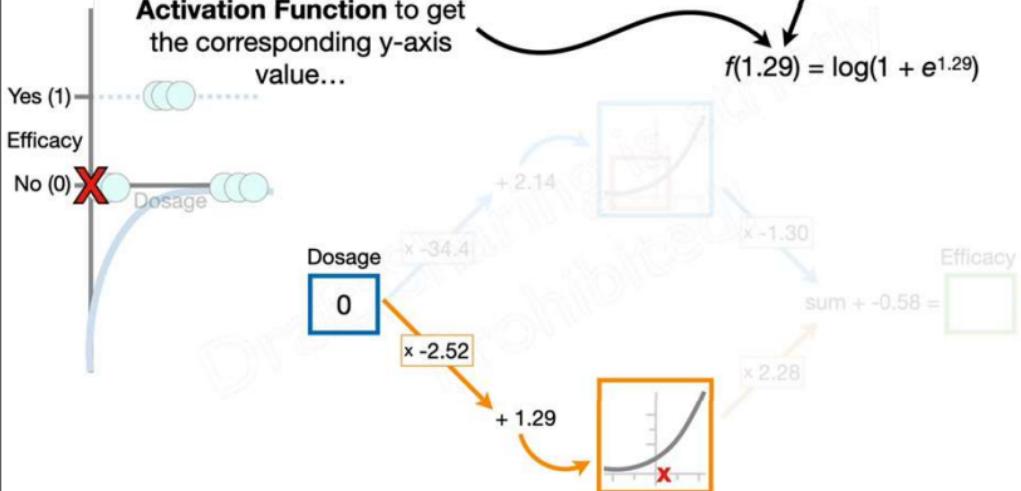




Now we plug **1.29** into the **Activation Function** to get the corresponding y-axis value...

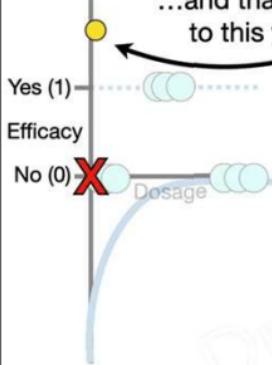
$$(0 \times -2.52) + 1.29 = 1.29$$

$$f(1.29) = \log(1 + e^{1.29})$$



$$(0 \times -2.52) + 1.29 = 1.29$$

...and that corresponds to this yellow dot.

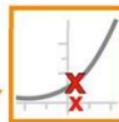


$$f(1.29) = \log(1 + e^{1.29}) = 1.53$$

Dosage  
0

$x -2.52$

$+ 1.29$



$+ 2.14$

$x -34.4$

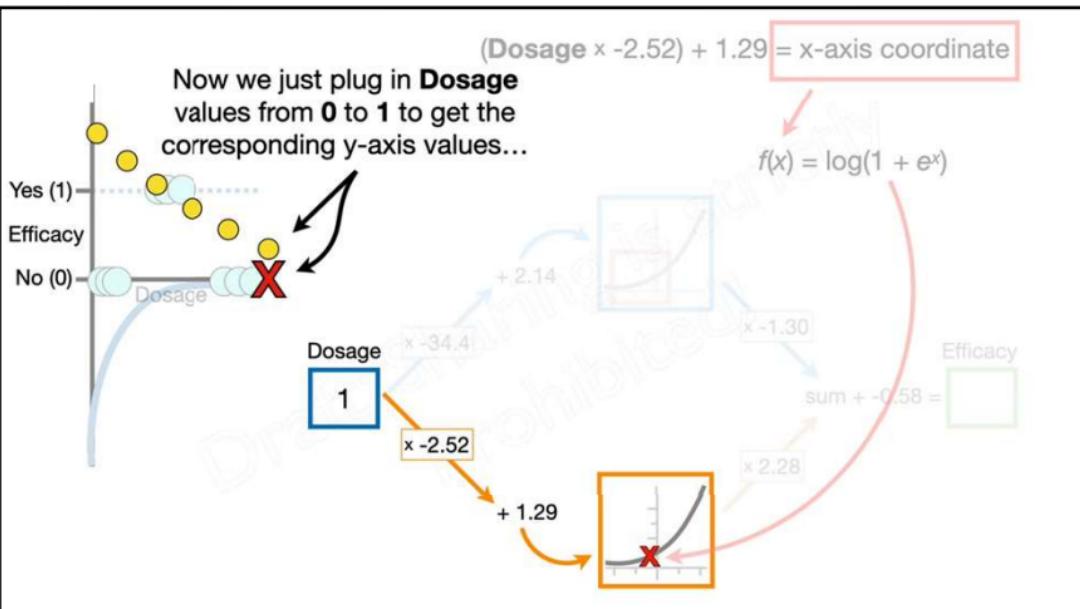
$x -1.30$

sum  $+ -0.58 =$

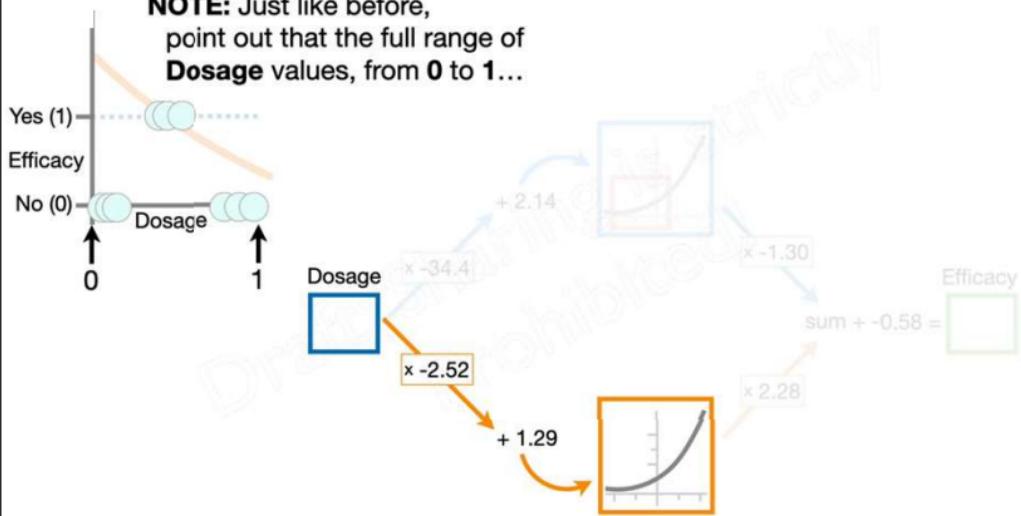
Efficacy

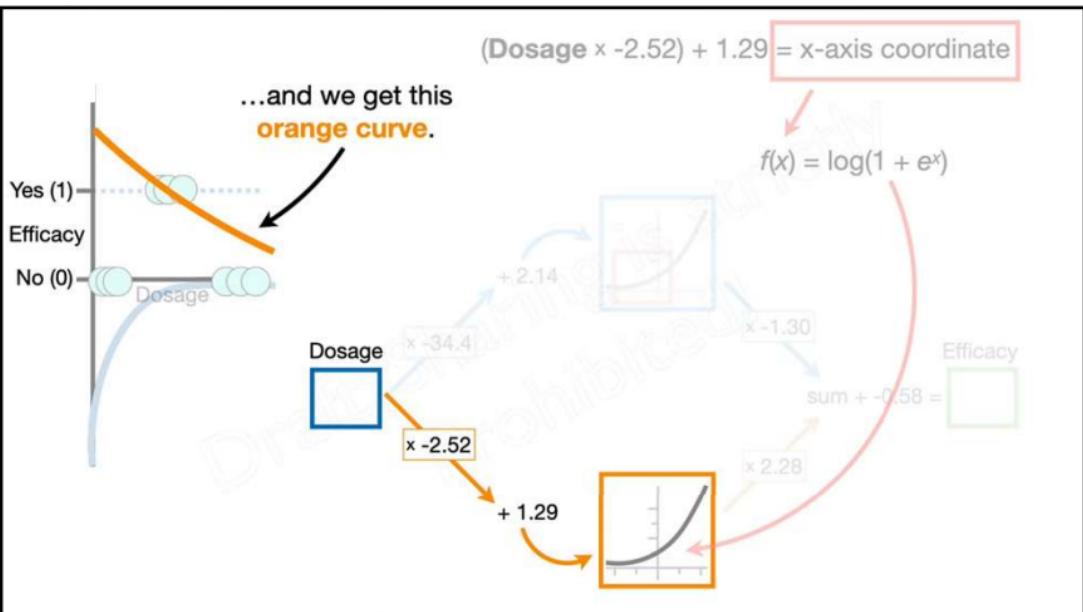


$x 2.28$

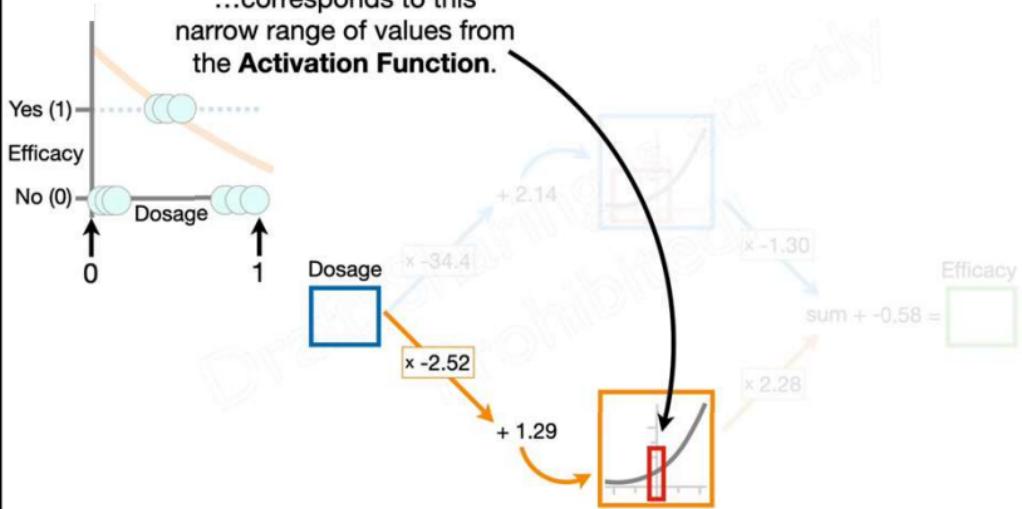


**NOTE:** Just like before,  
point out that the full range of  
**Dosage** values, from **0** to **1**...

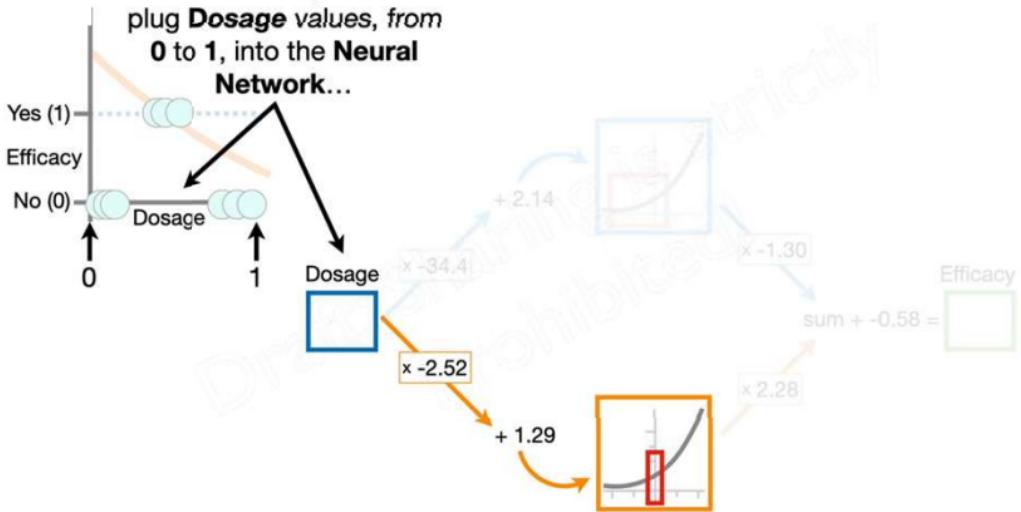




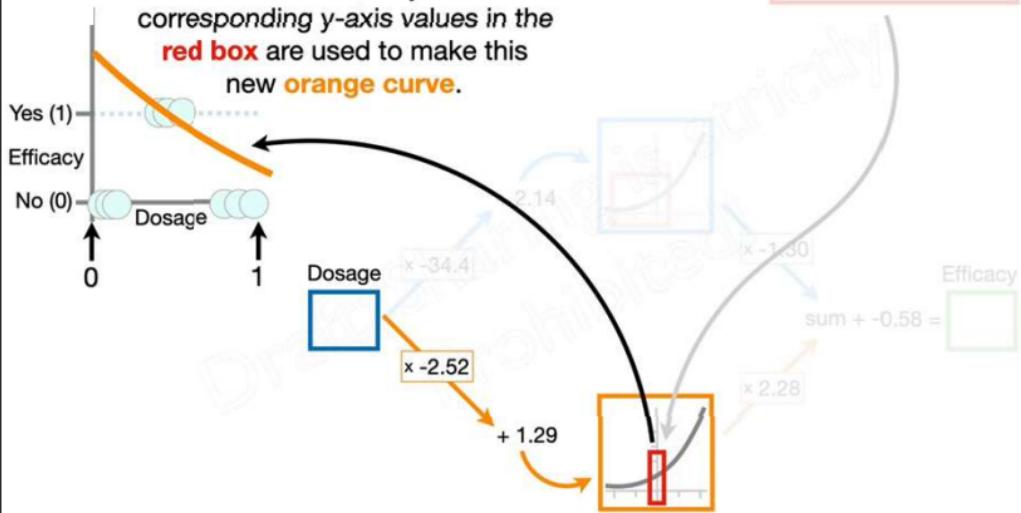
...corresponds to this narrow range of values from the **Activation Function**.



In other words, when we plug **Dosage** values, from 0 to 1, into the **Neural Network**...



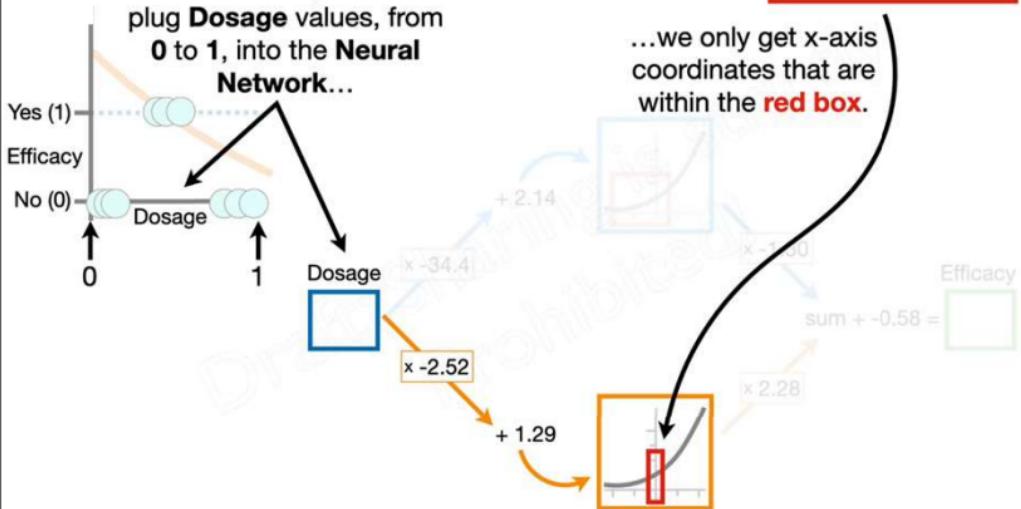
...and thus, only the corresponding y-axis values in the red box are used to make this new orange curve.



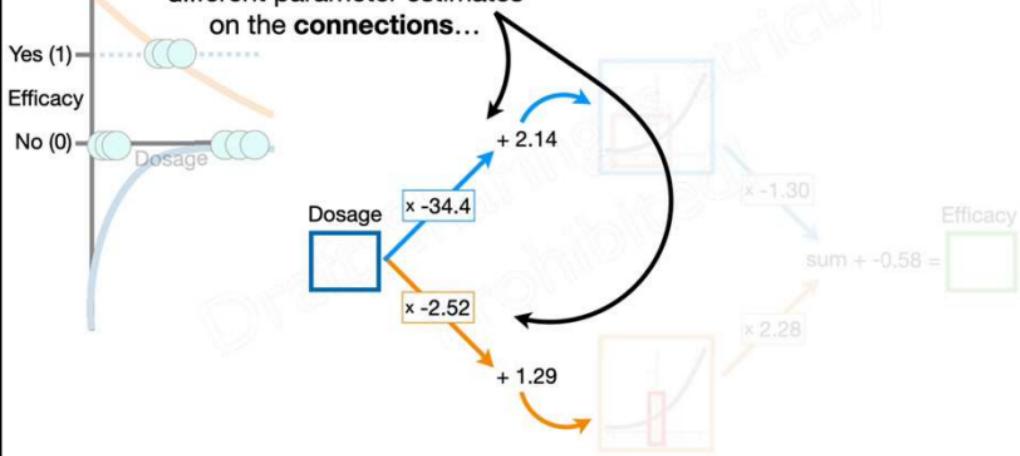
In other words, when we plug **Dosage** values, from **0** to **1**, into the **Neural Network**...

$$(\text{Dosage} \times -2.52) + 1.29 = \text{x-axis coordinate}$$

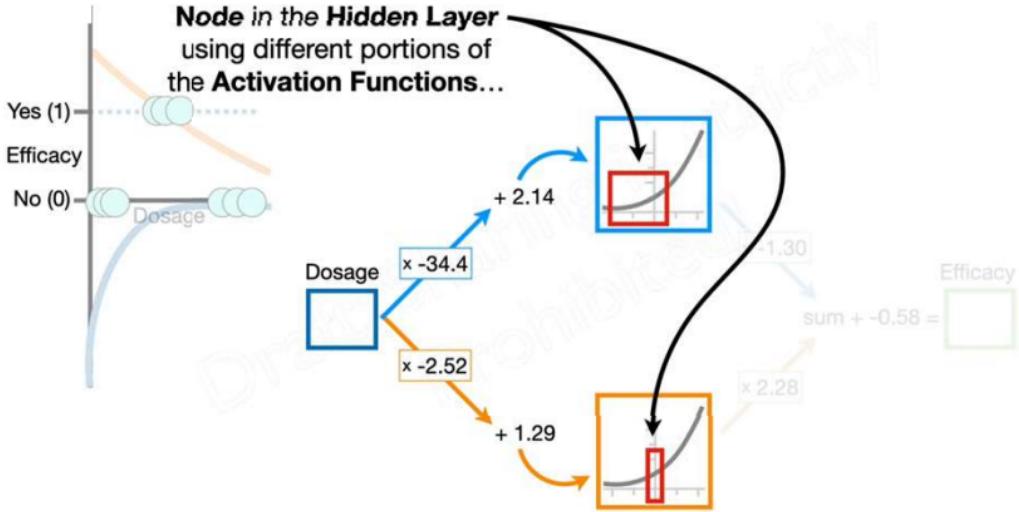
...we only get x-axis coordinates that are within the **red box**.



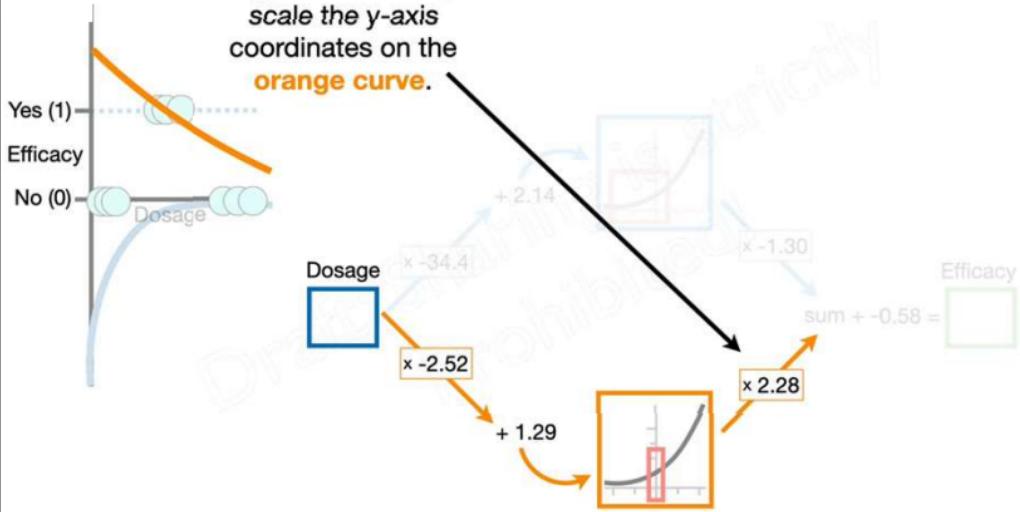
So we see that fitting a **Neural Network** to data gives us different parameter estimates on the **connections**...



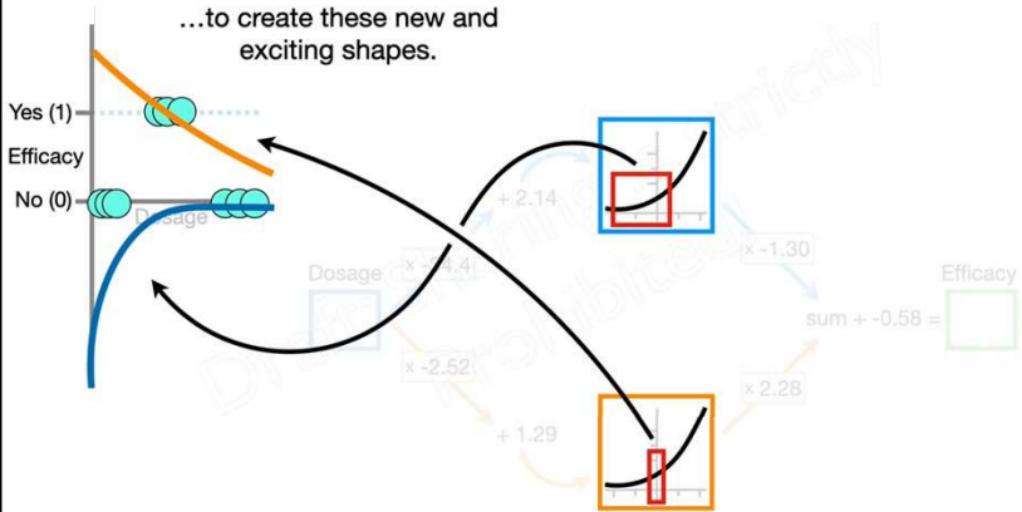
...and that results in each  
**Node in the Hidden Layer**  
using different portions of  
the **Activation Functions**...



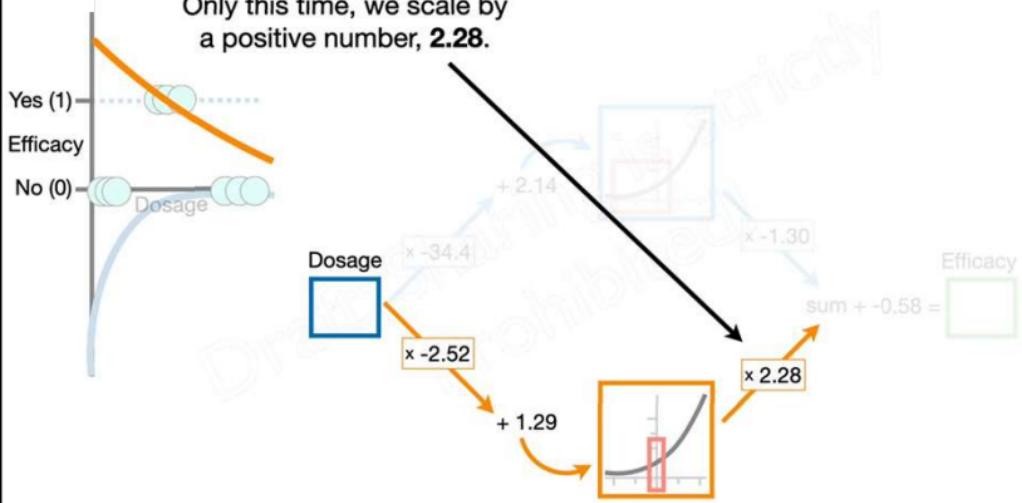
Now, just like before, we scale the y-axis coordinates on the **orange curve**.

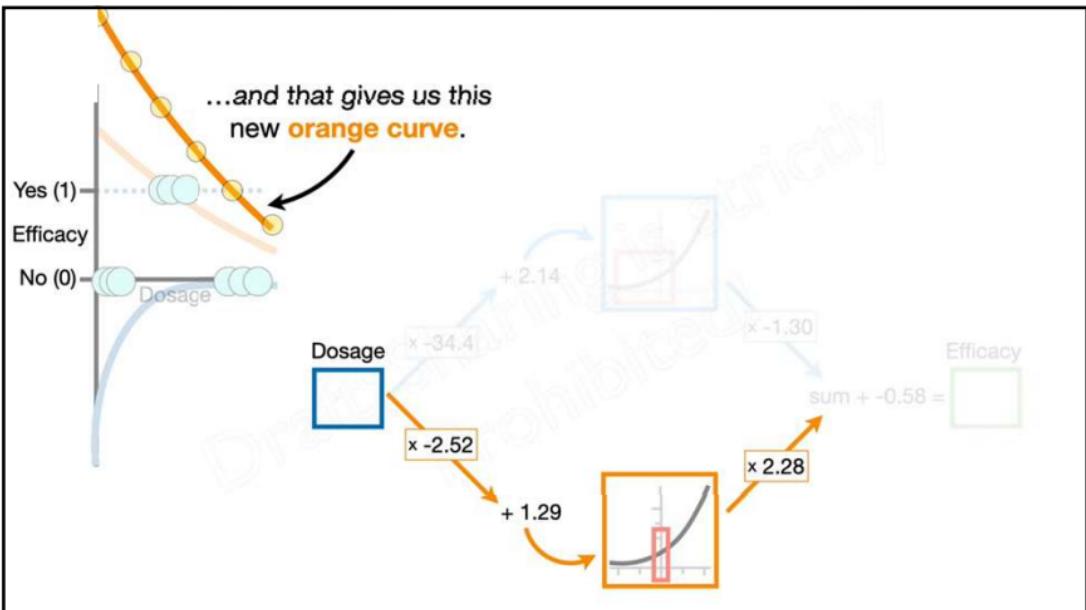


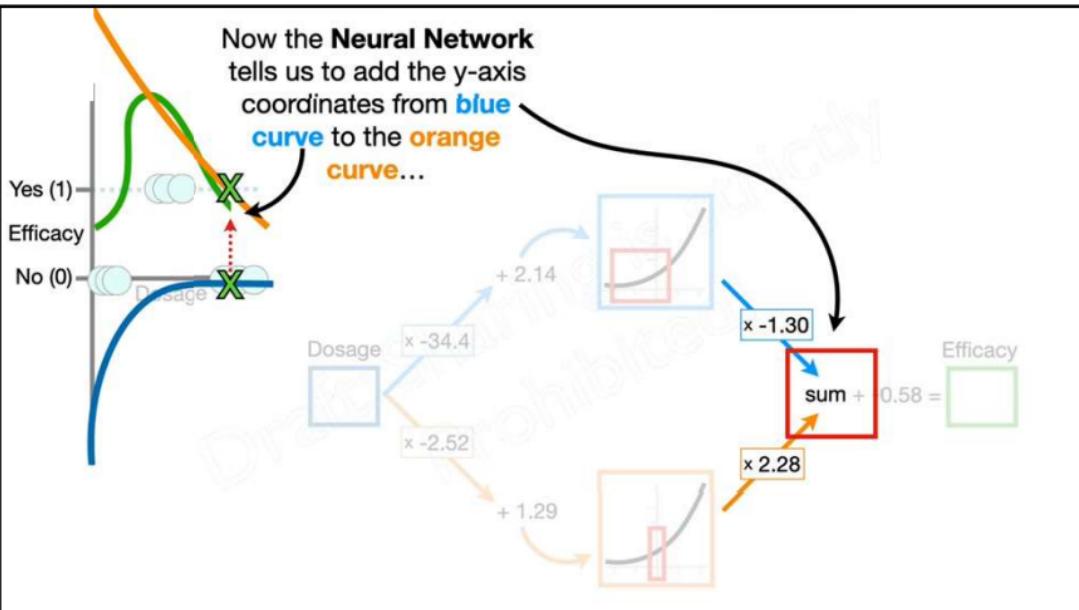
...to create these new and exciting shapes.

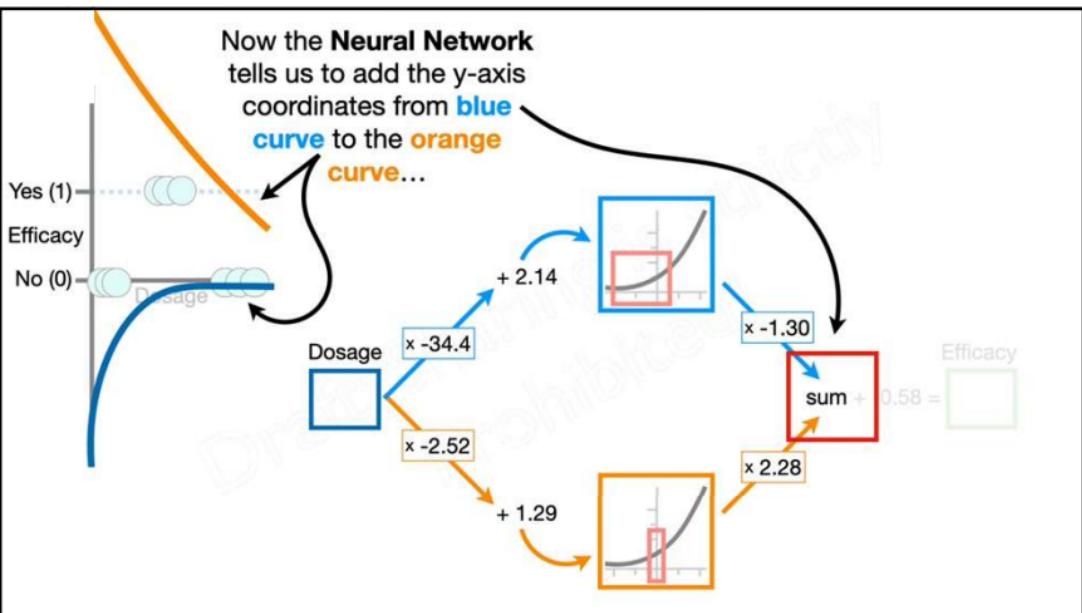


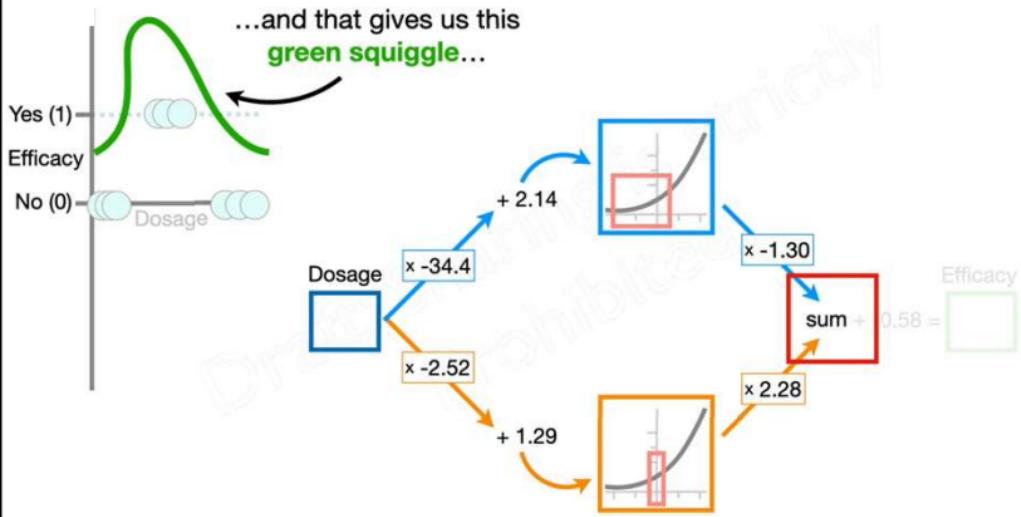
Only this time, we scale by a positive number, **2.28**.

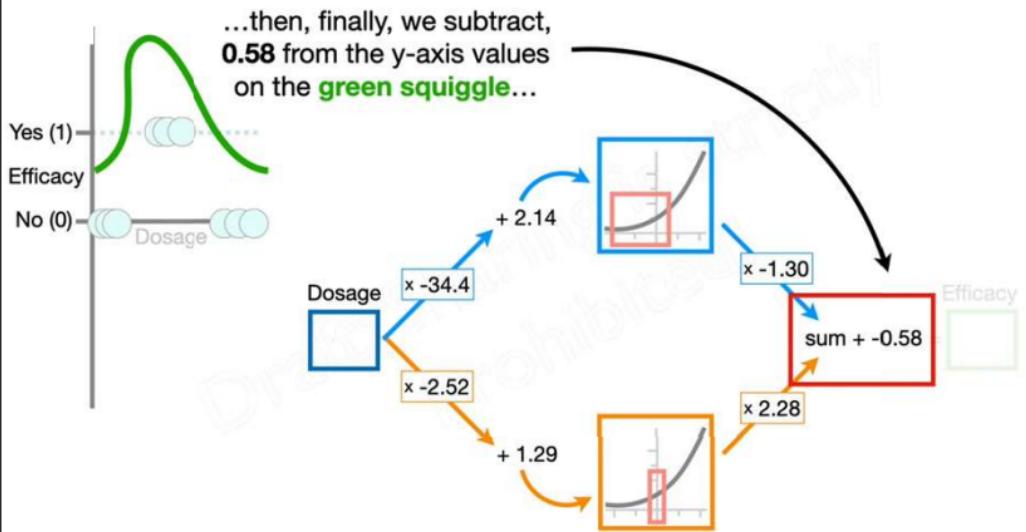




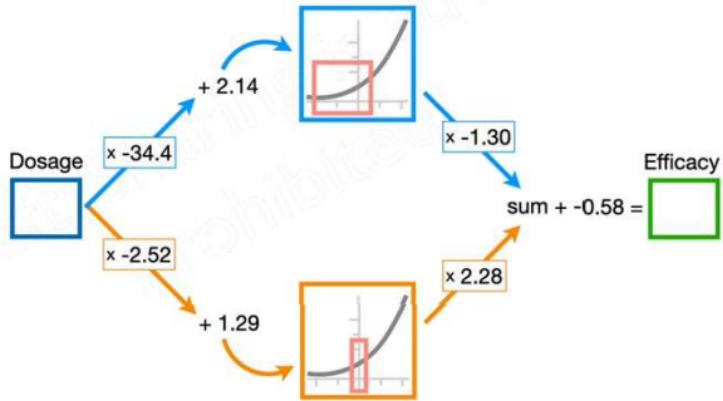
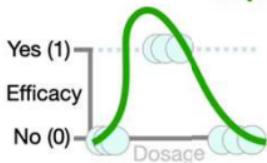


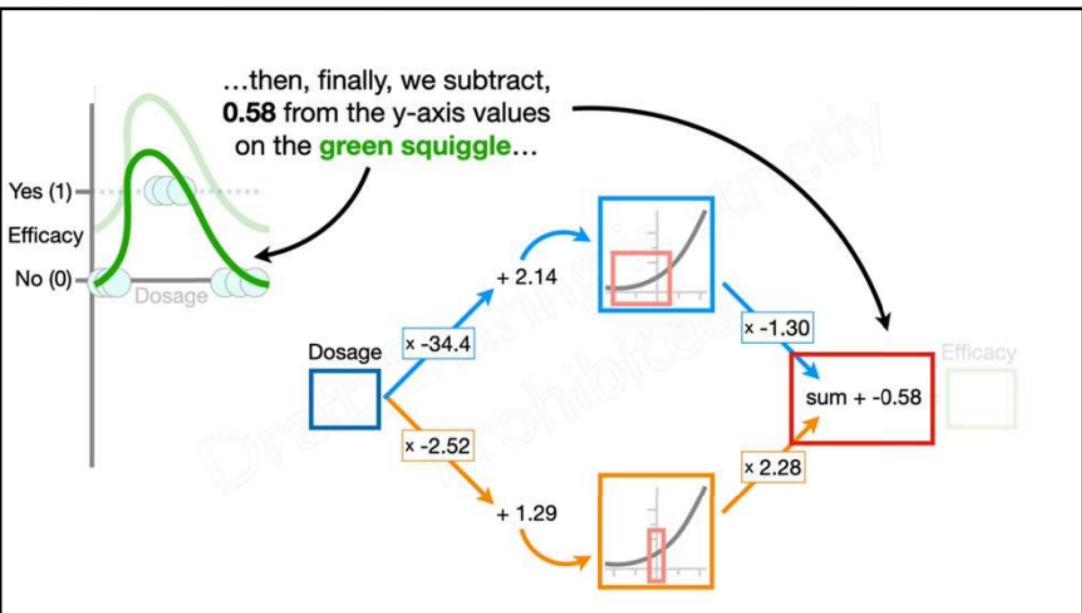




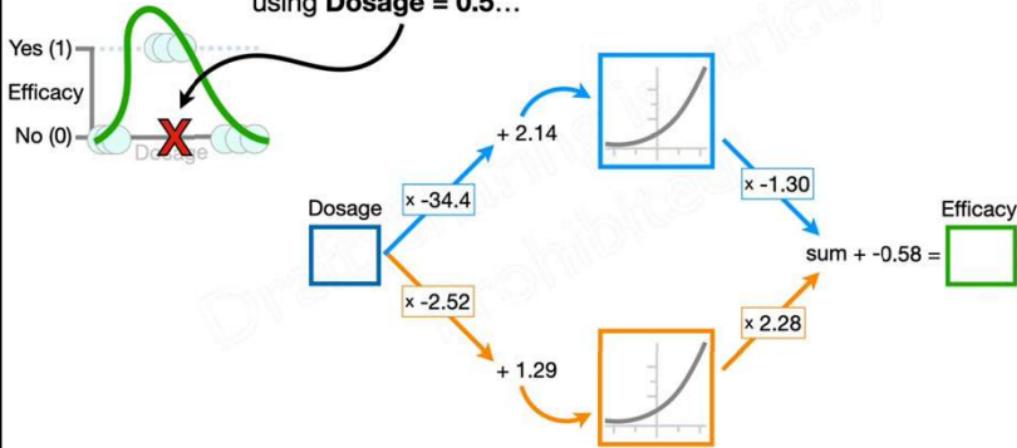


...and we have a **green** squiggle that fits the data.

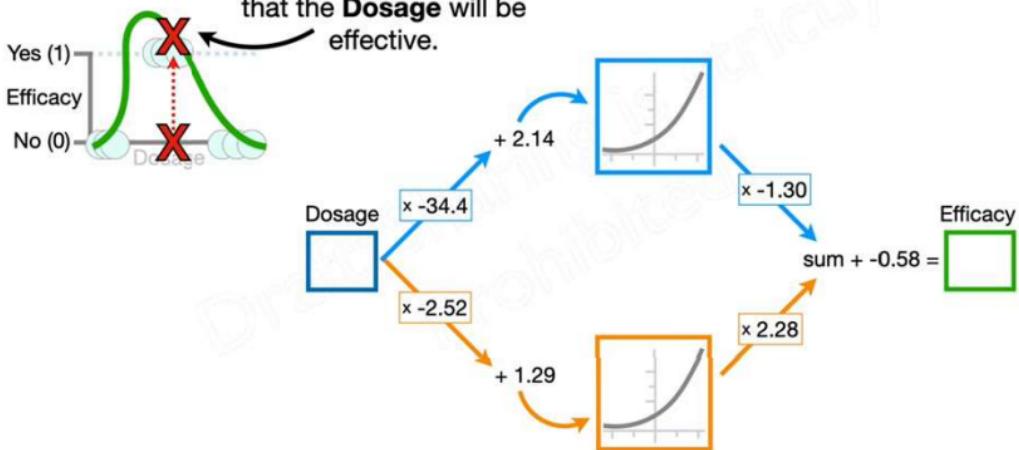




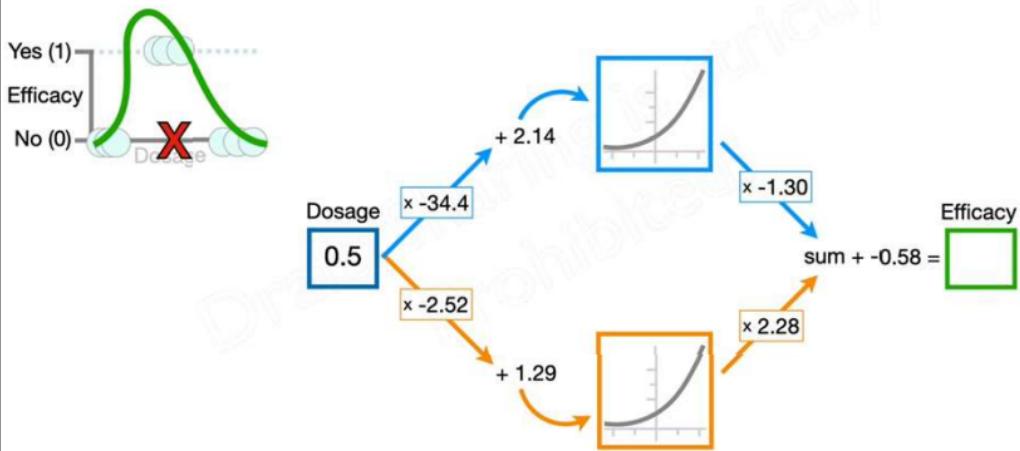
Now, if someone comes along and says that they are using **Dosage = 0.5...**



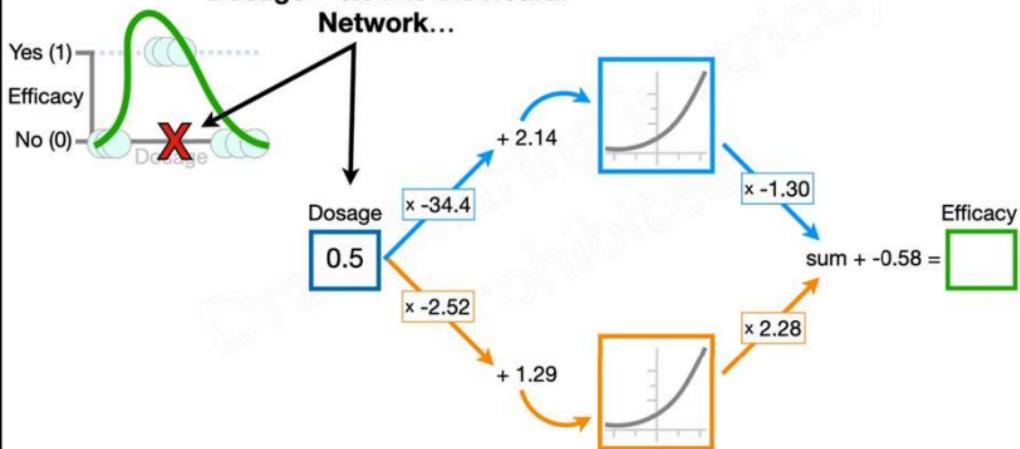
...we can look at the corresponding y-axis coordinate on the **green squiggle** and see that the **Dosage** will be effective.



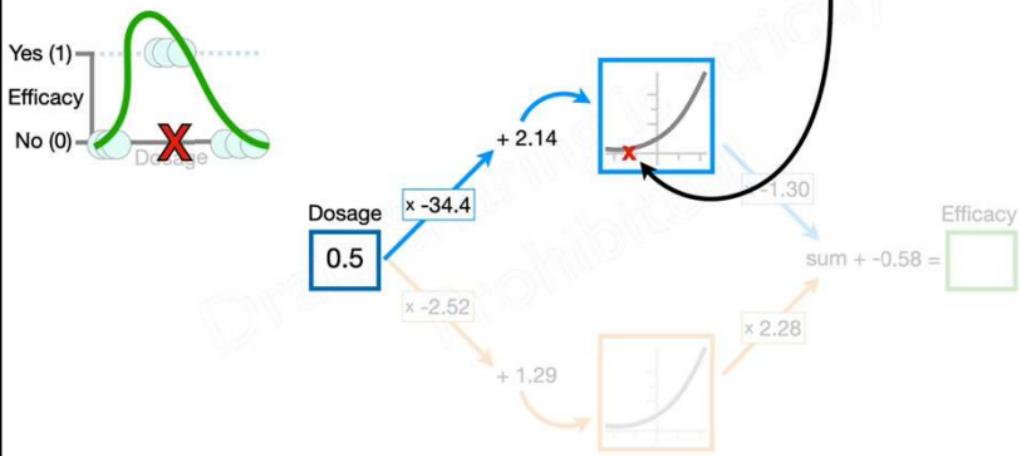
...and do the math...



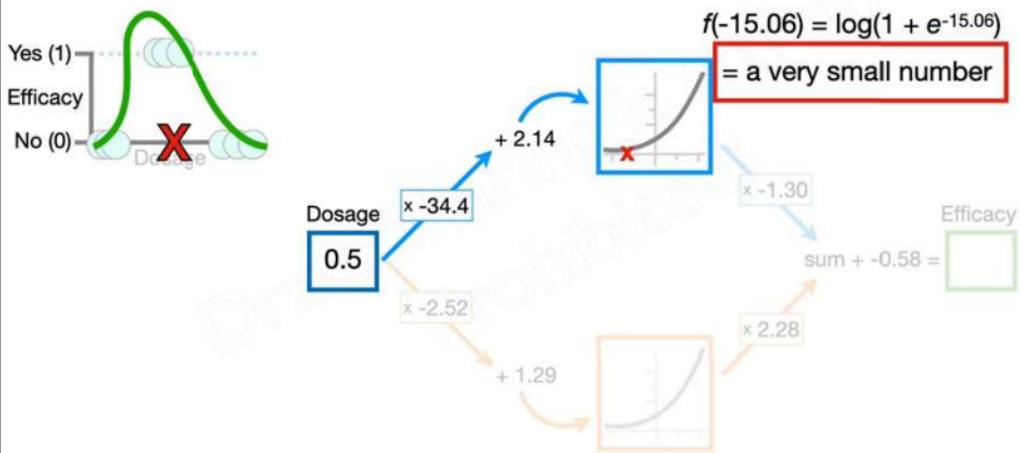
Or we can solve for the y-axis  
coordinate by plugging  
**Dosage = 0.5** into the **Neural  
Network...**



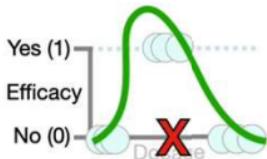
$$(0.5 \times -34.4) + 2.14 = -15.06$$



$$(0.5 \times -34.4) + 2.14 = -15.06$$



$$(\text{Dosage} \times -2.52) + 1.29 = \text{x-axis coordinate}$$



Dosage

0.5

$\times -2.52$

+ 1.29

$\times -34.4$

+ 2.14

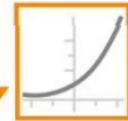
$\times -1.30$

sum + -0.58 =

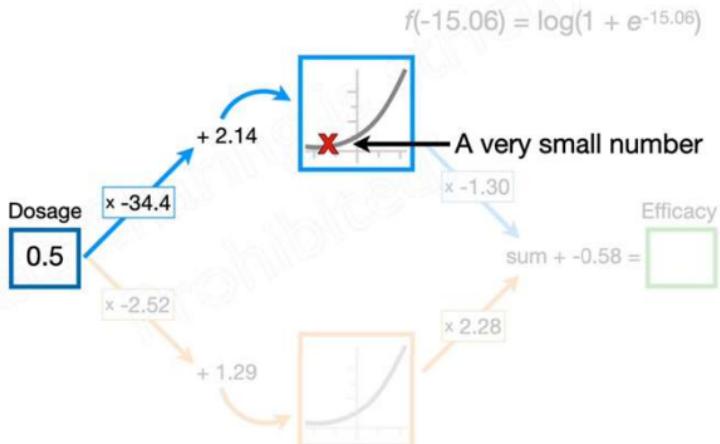
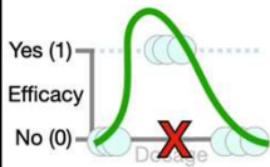
Efficacy



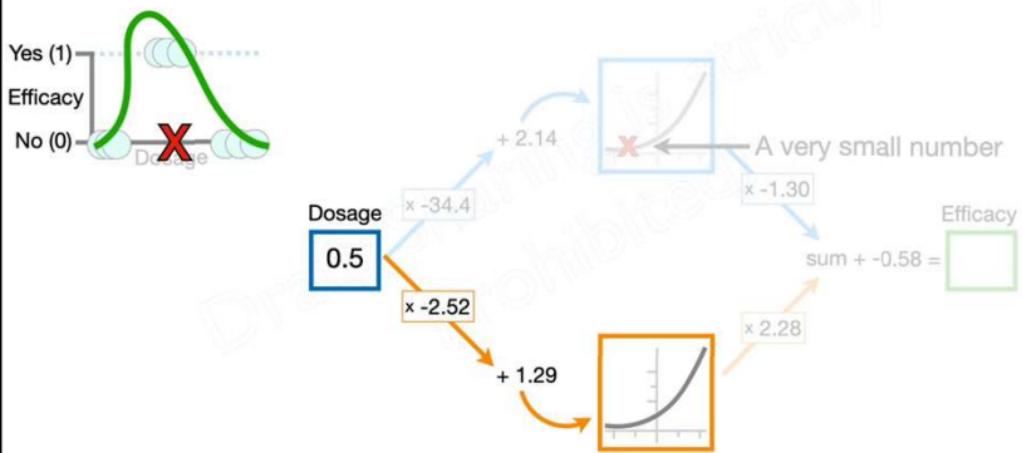
$\times 2.28$



$$(0.5 \times -34.4) + 2.14 = -15.06$$

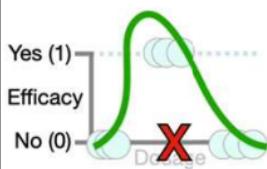


$$(0.5 \times -2.52) + 1.29 = 0.3$$



$$(0.5 \times -2.52) + 1.29 = 0.3$$

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



Dosage  
**0.5**

$\times -34.4$

$\times -2.52$

$+ 2.14$

$+ 1.29$

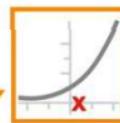


$\times -1.30$

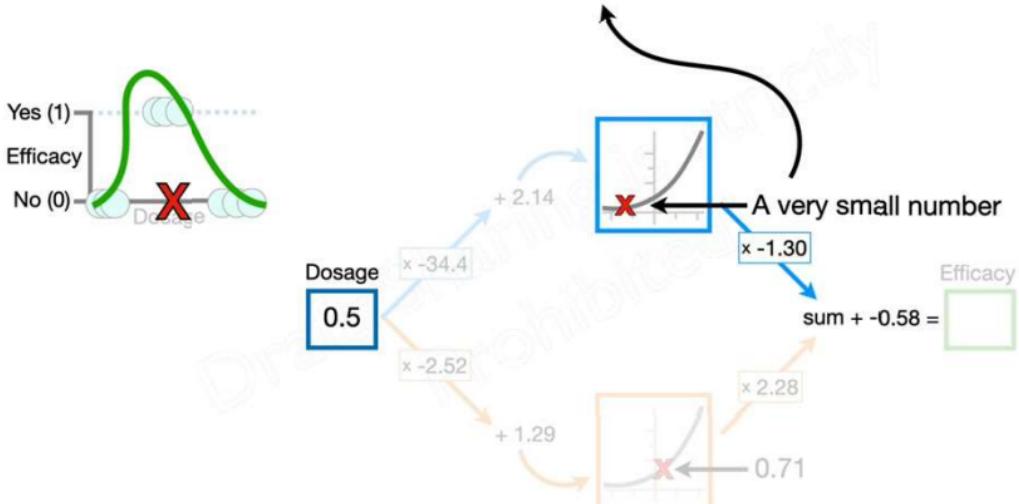
sum +

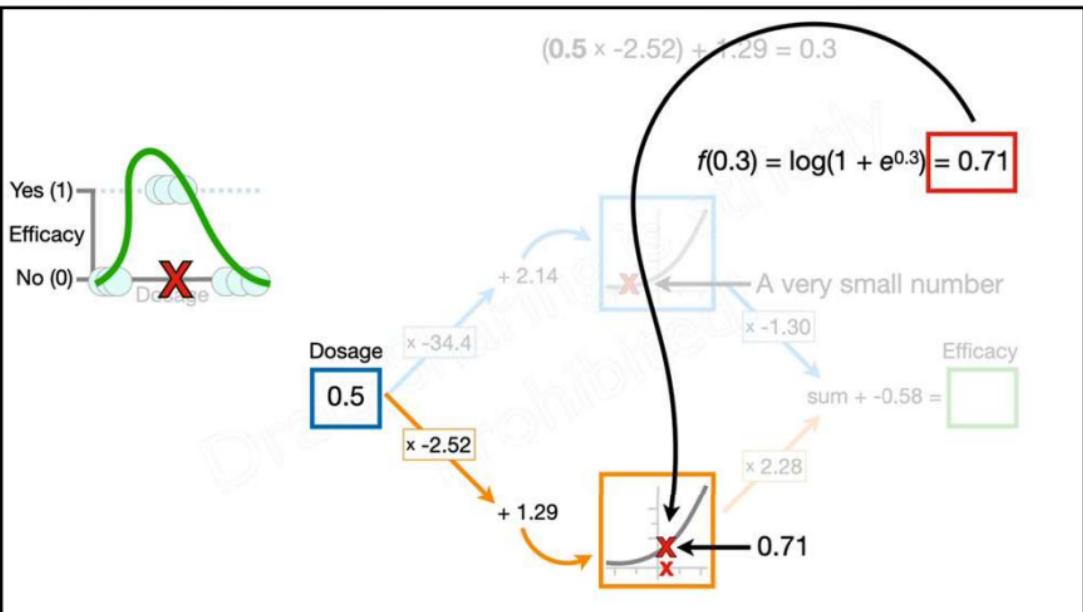
-0.58 =

Efficacy

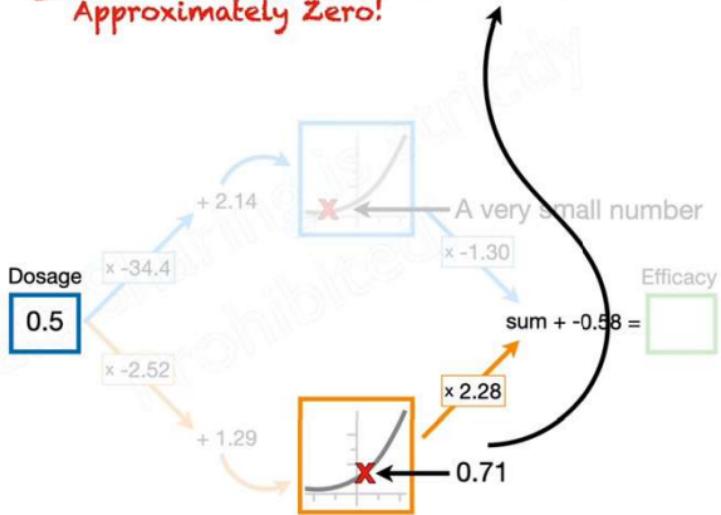
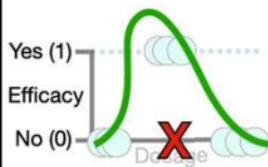


(A really small number  $\times$  -1.30)

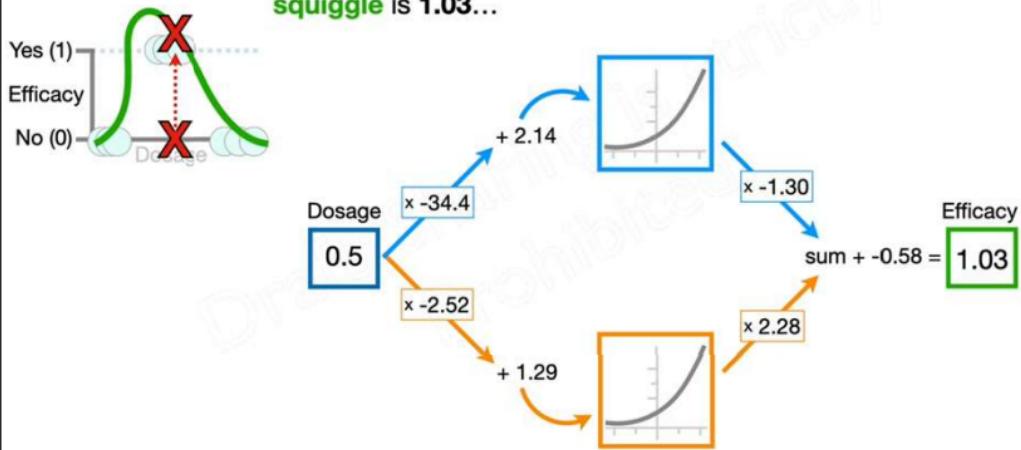




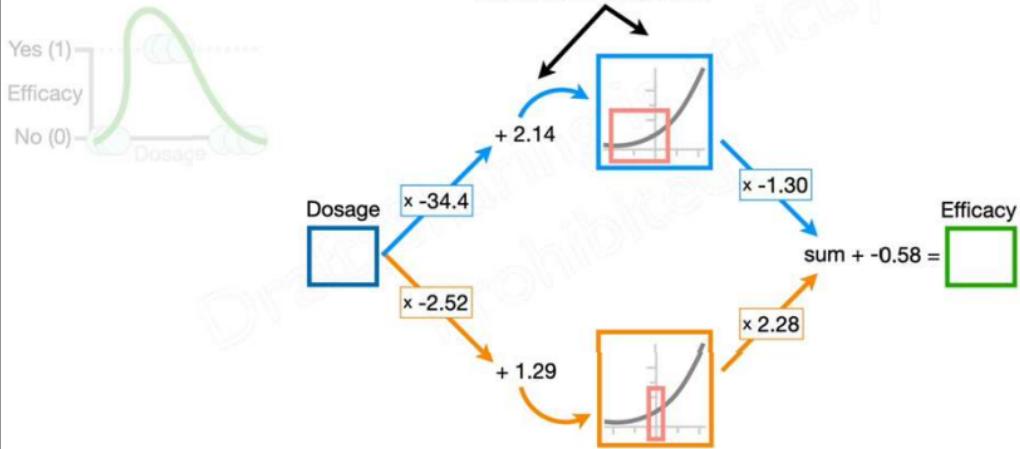
~~(A really small number  $x - 1.30$ )~~  $(0.71 \times 2.28)$   
Approximately Zero!



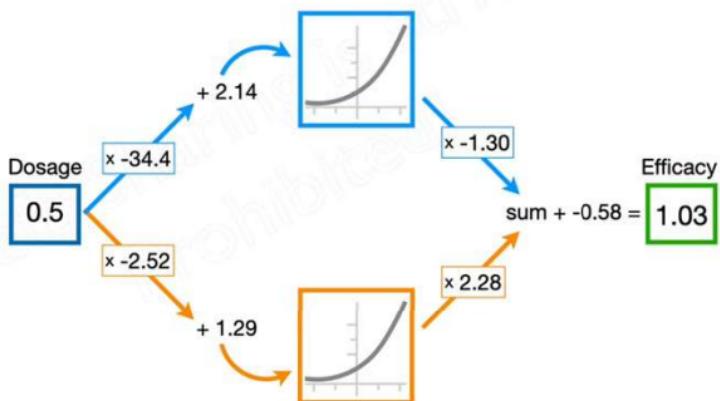
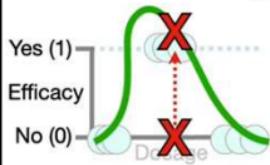
...and we see that the y-axis coordinate on the green squiggle is 1.03...



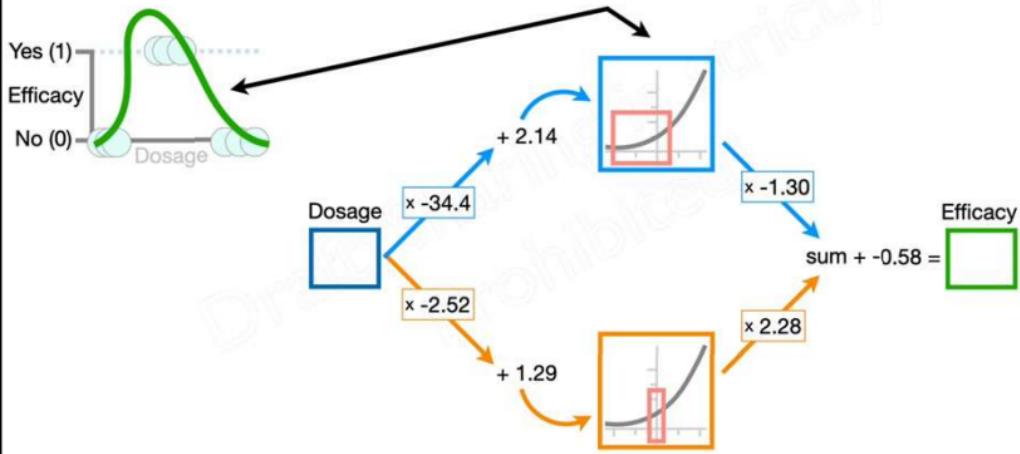
Now, if you've made it this far, you may  
be wondering why this is called a  
**Neural Network...**



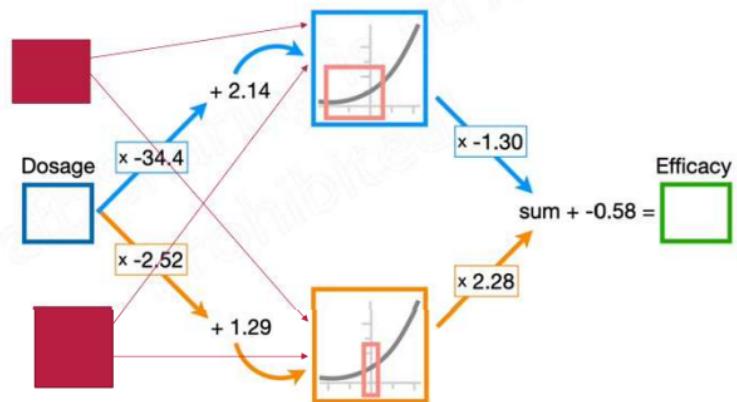
...and since **1.03** is closer to **1** than **0**, we will conclude that a **Dosage = 0.5** is effective.



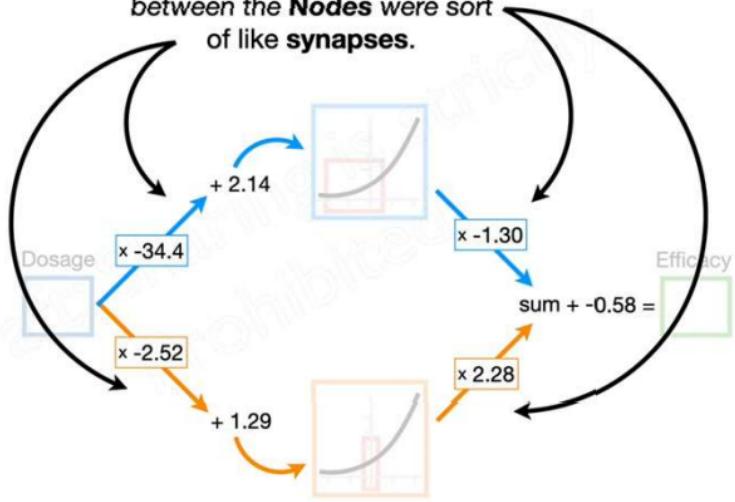
...instead of a **Big Fancy Squiggle**  
**Fitting Machine.**



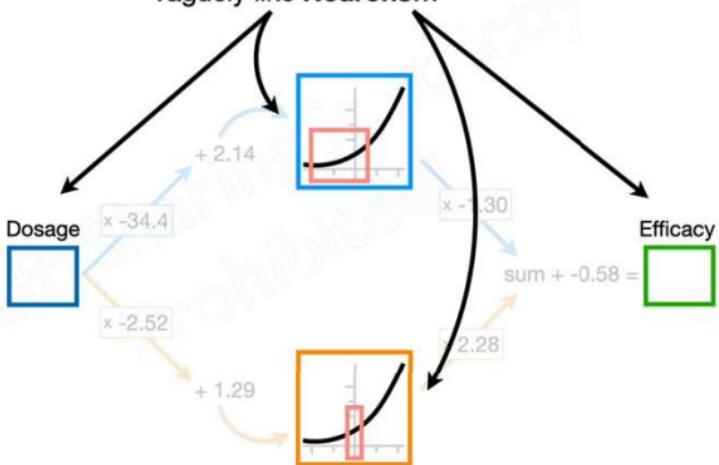
The reason is that way back in the  
1940s and 50s when **Neural  
Networks** were invented...



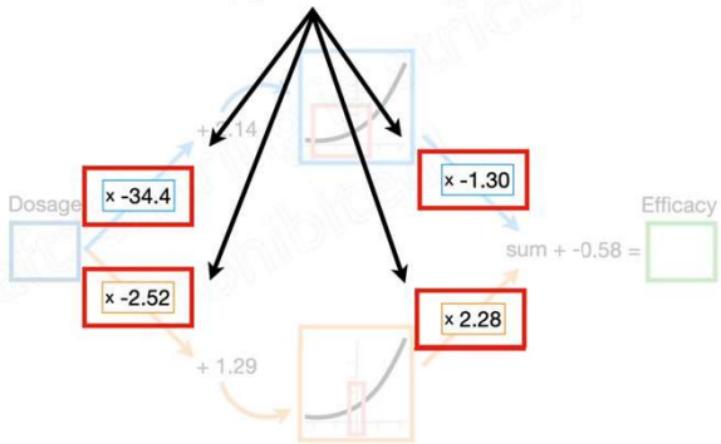
...and the **connections**  
between the **Nodes** were sort  
of like **synapses**.



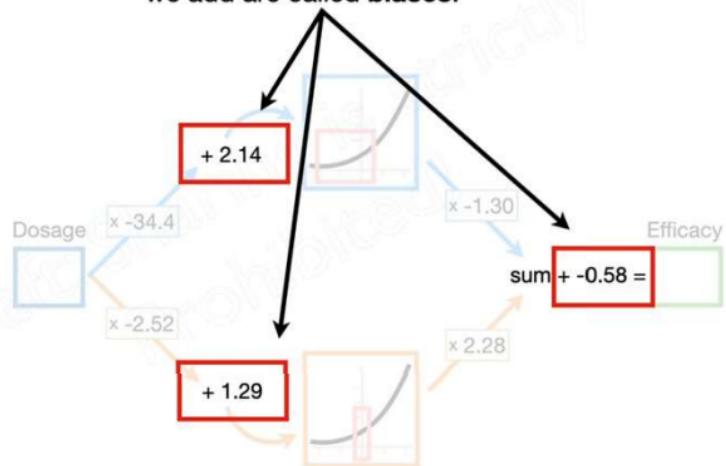
...they thought the **Nodes** were  
vaguely like **Neurons**...



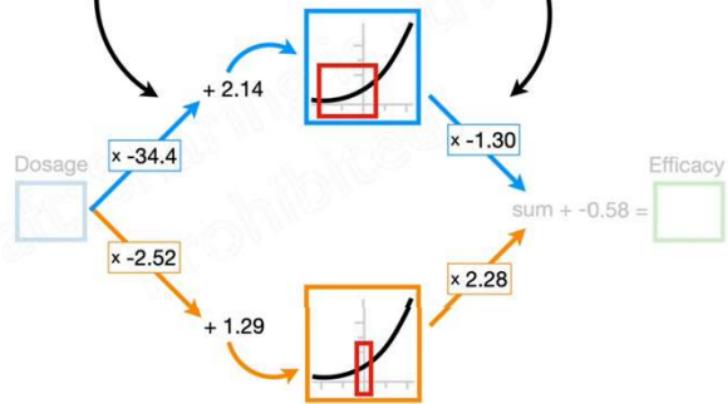
**NOTE:** Whether or not you call it a **Squiggle Fitting Machine**, the parameters that we multiply are called **weights**...



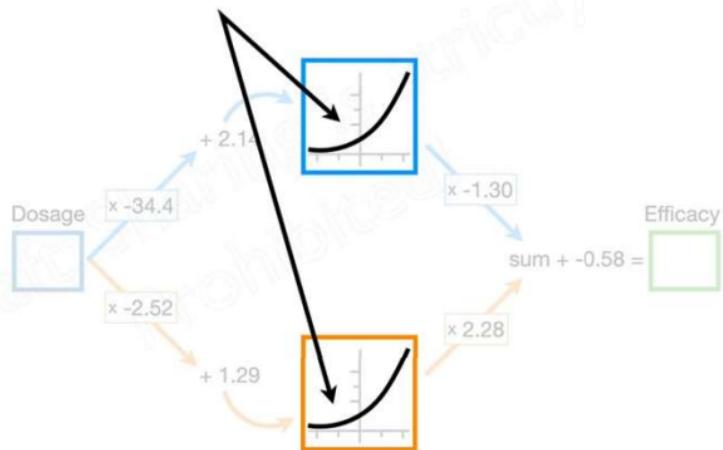
...and the parameters that we add are called **biases**.



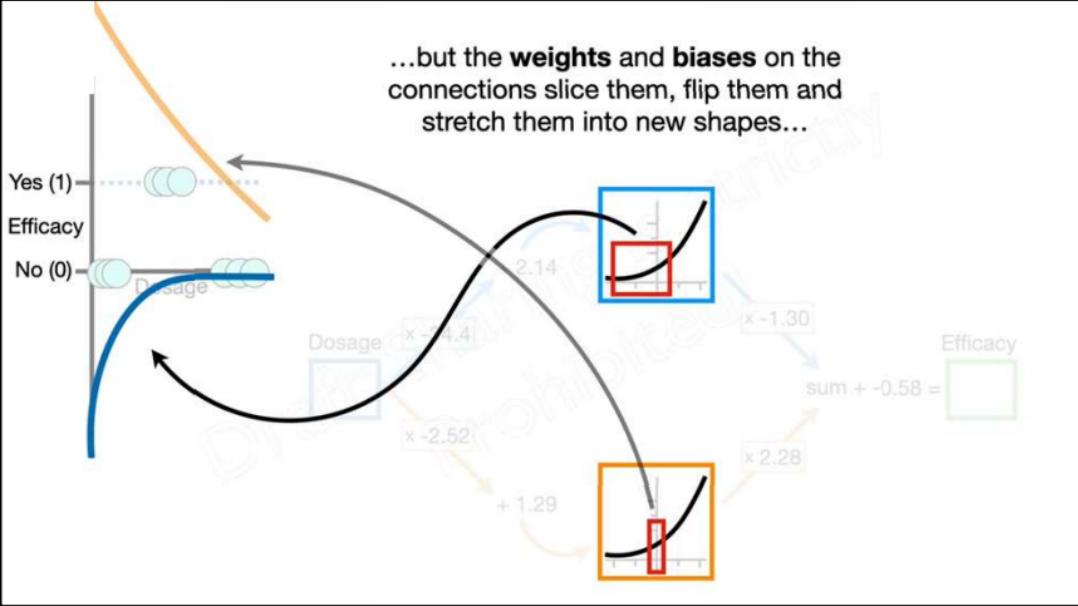
...but the **weights** and **biases** on the connections slice them, flip them and stretch them into new shapes...



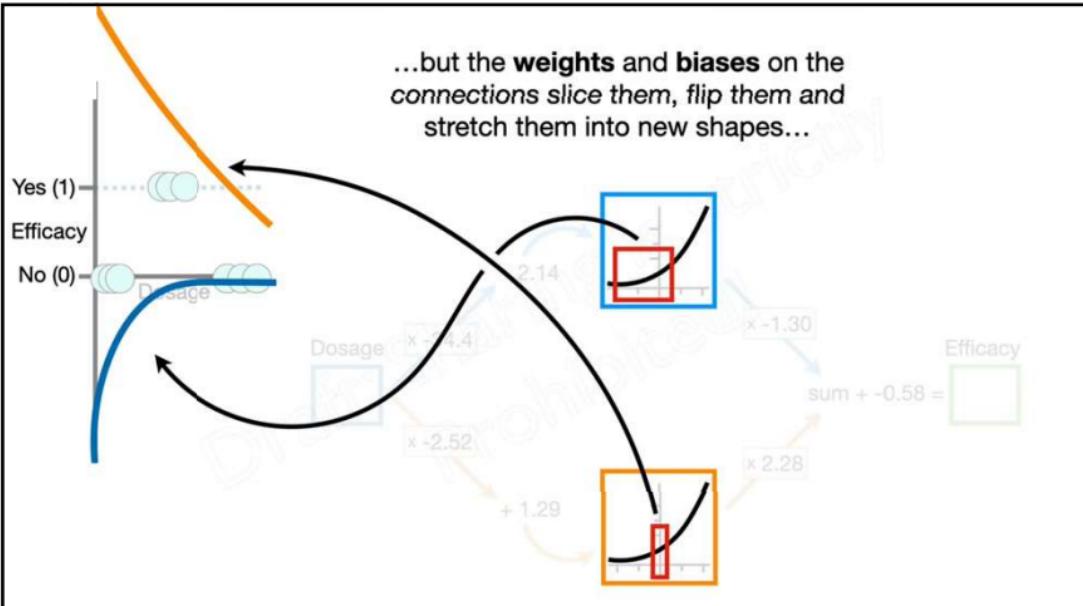
**NOTE:** This Neural Network starts with  
two identical Activation Functions...

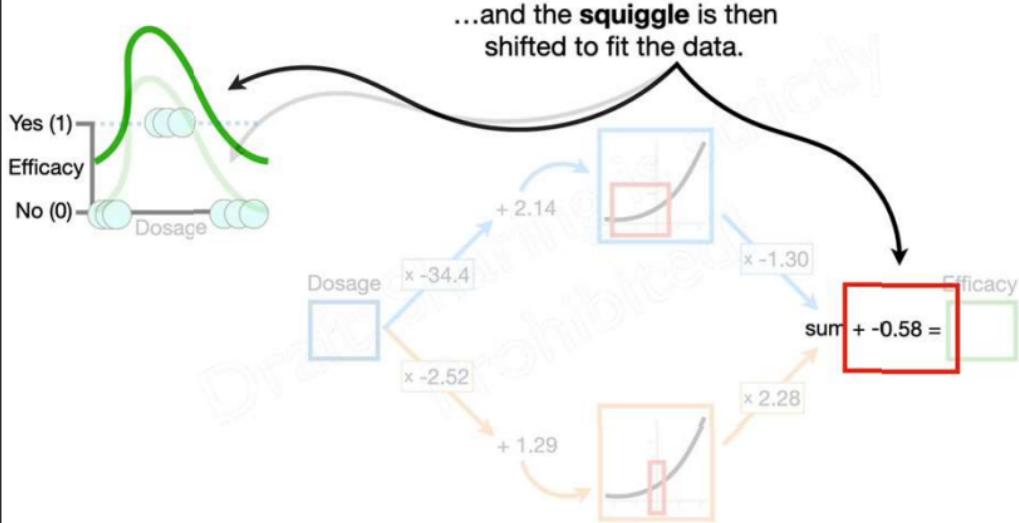


...but the **weights** and **biases** on the connections slice them, flip them and stretch them into new shapes...

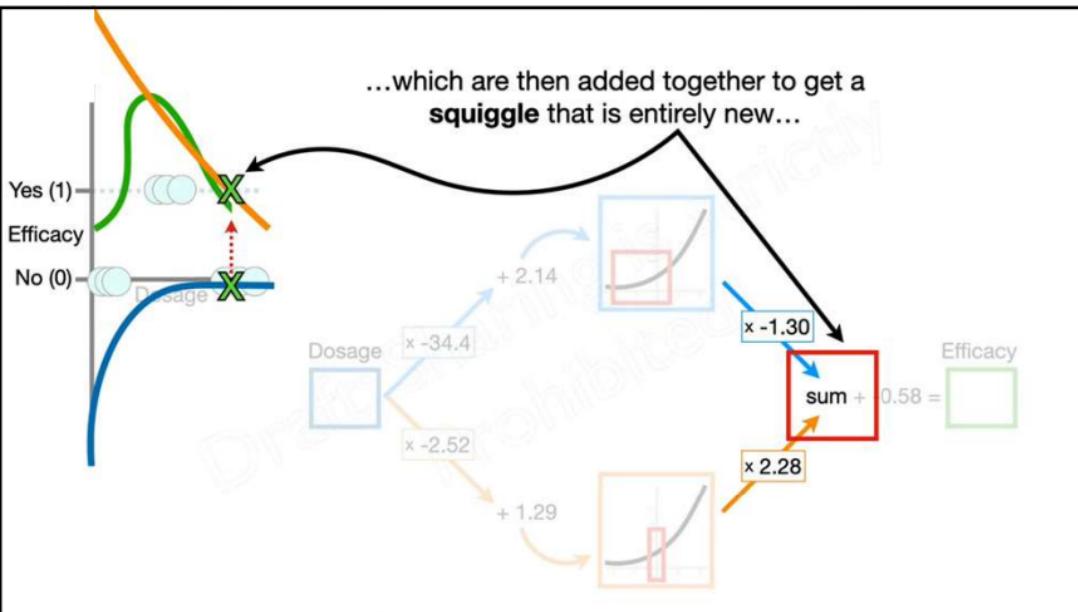


...but the **weights** and **biases** on the connections slice them, flip them and stretch them into new shapes...

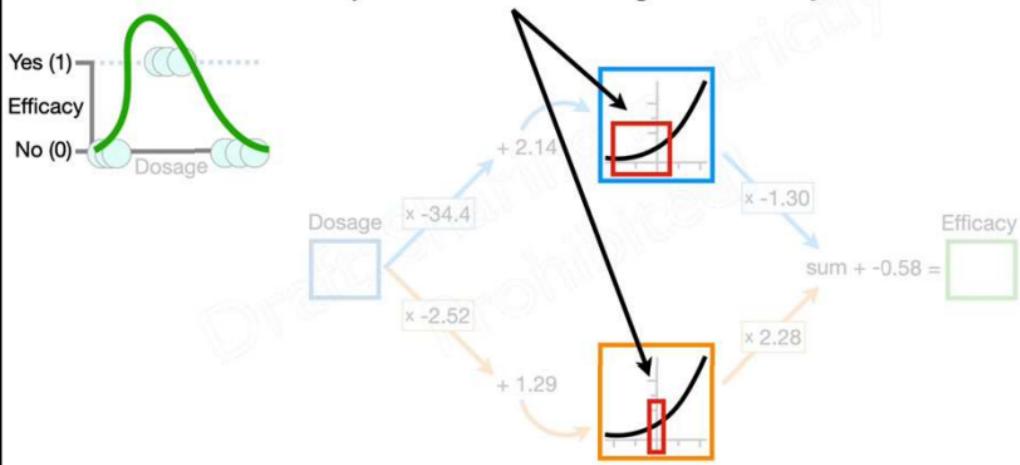




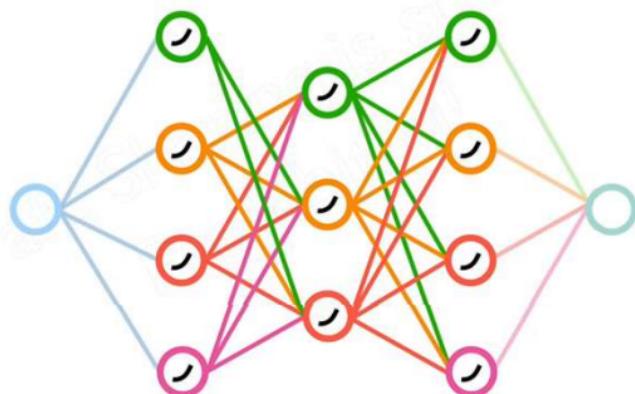
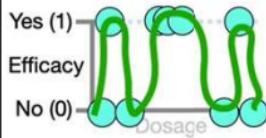
...which are then added together to get a squiggle that is entirely new...



Now, if we can create this **green squiggle** with just two **Nodes** in a single **Hidden Layer**...



...just imagine what types of **green squiggles**  
we could fit with more **Hidden Layers** and  
more **Nodes** in each **Hidden Layer**.



**THANK YOU!**

In theory, **Neural Networks** can fit a **green squiggle** to just about any dataset, no matter how complicated!

