

# Introduction to Machine Learning and Neural Networks by Nicolas Symeou

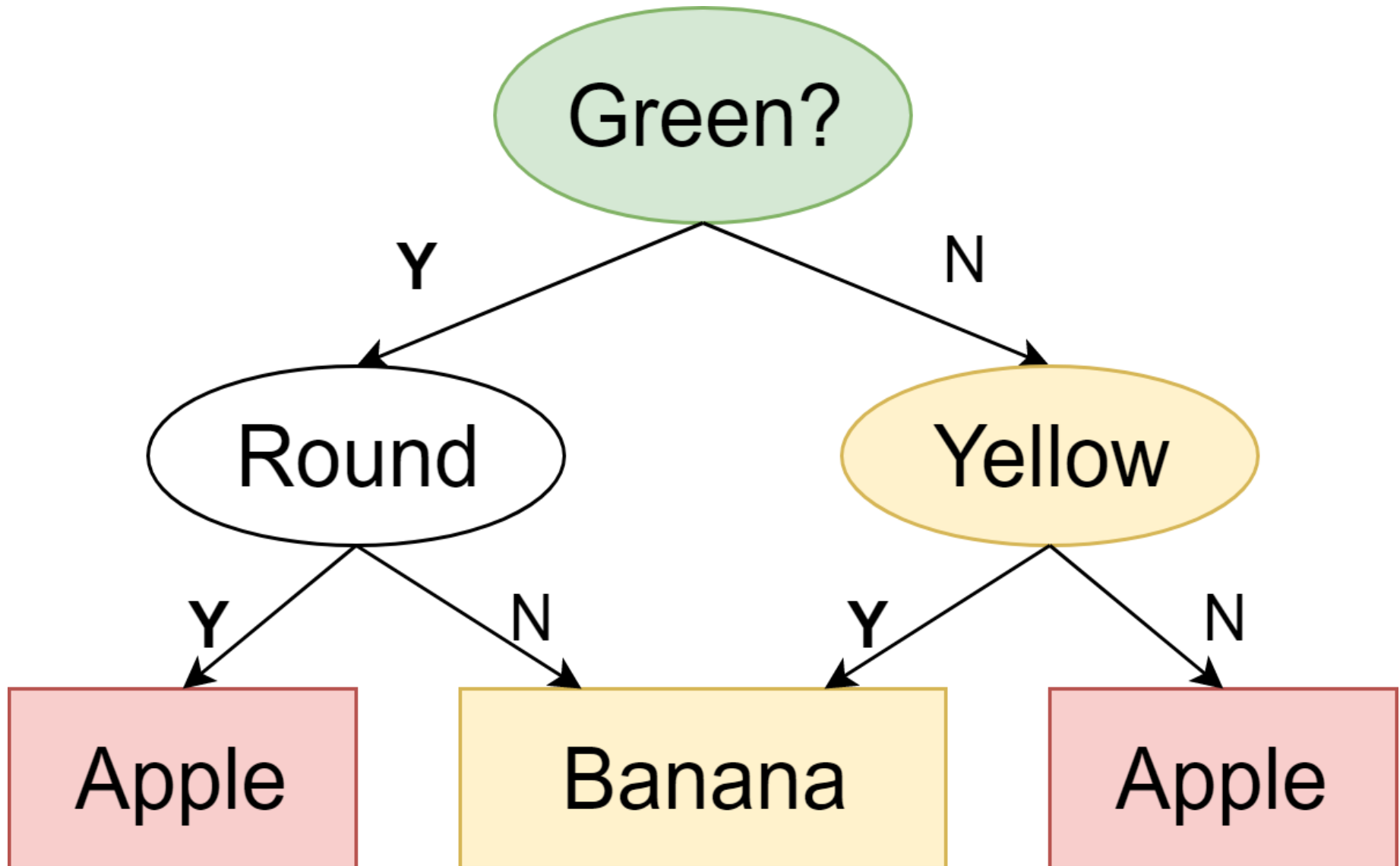
Thanks to

Florin Schwappach  
for correction and feedback

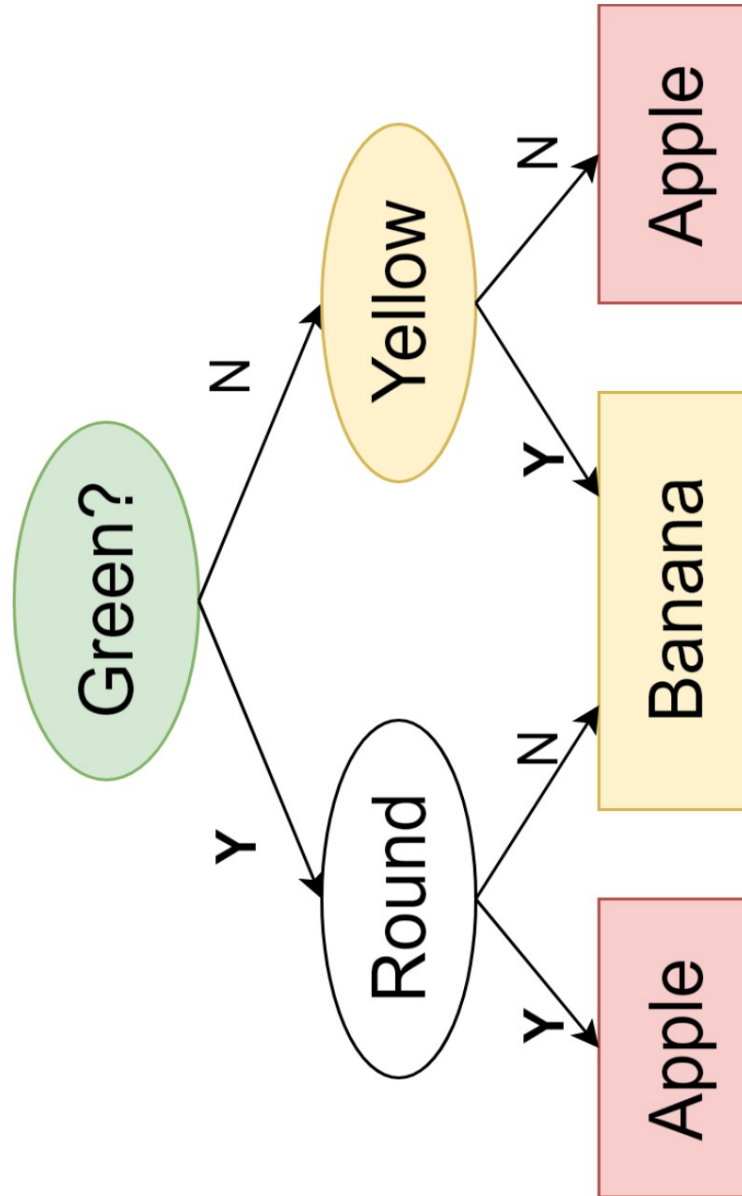
Machine Learning,  
An introduction 1:

The BIG question

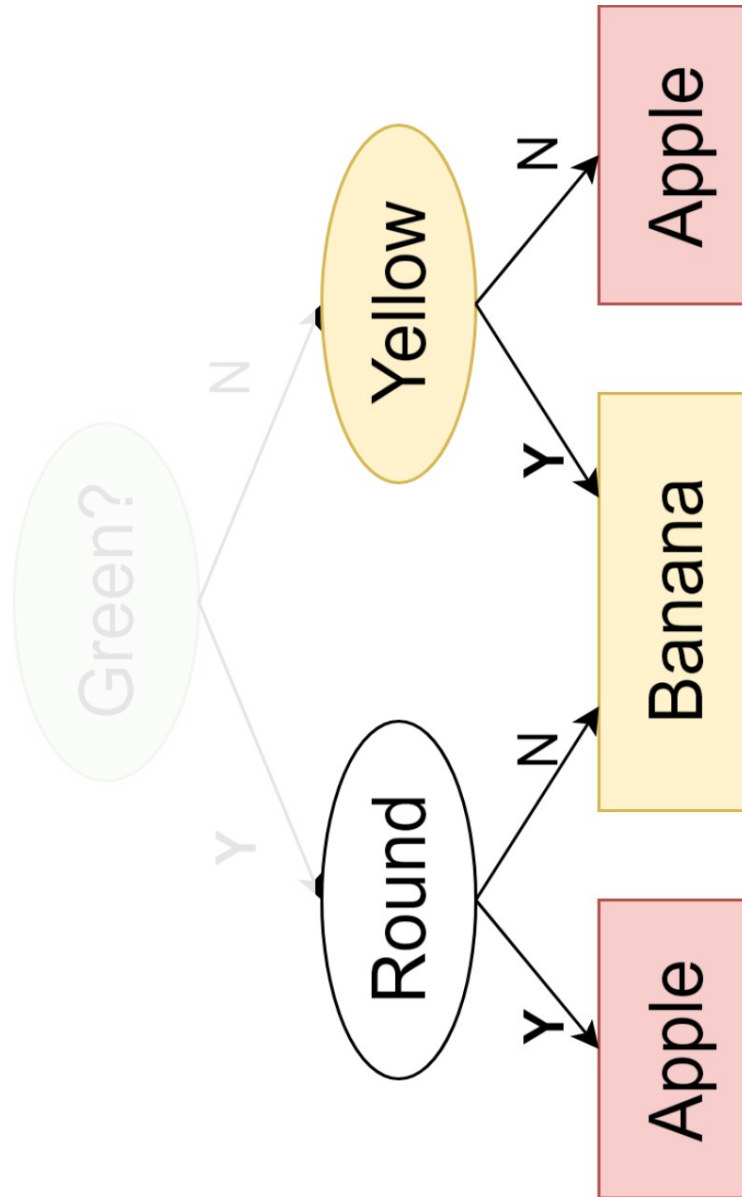
# Banana or Apple?



# Banana or Apple?



# Banana or Apple?



# Banana or Apple?

A yellow oval with a thin black border, containing the word "Yellow" in black text.

Yellow

A white oval with a thin black border, containing the word "Round" in black text.

Round

# Banana or Apple?

A yellow oval with a thin brown border.

Yellow

A pink rectangle with a thin brown border.

Apple

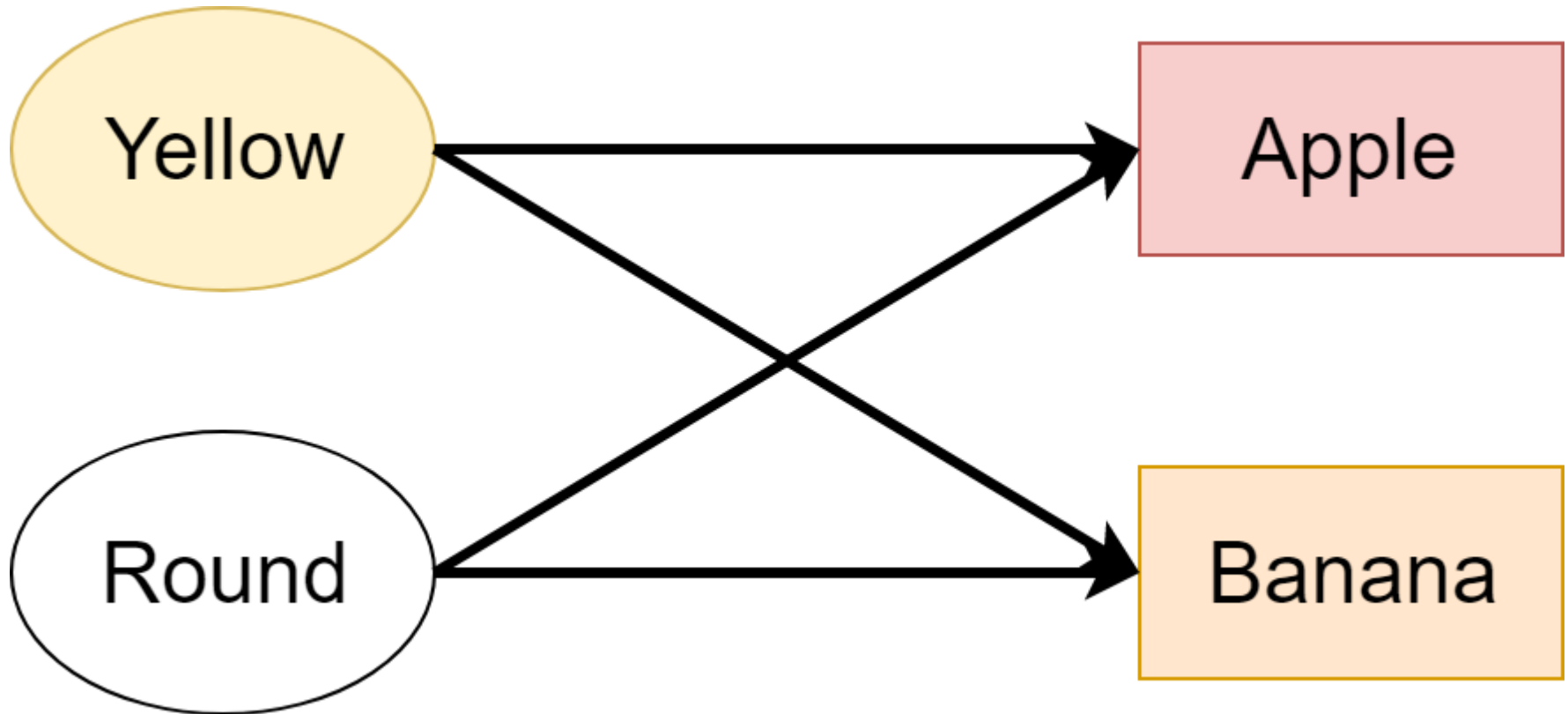
A white oval with a thin black border.

Round

An orange rectangle with a thin brown border.

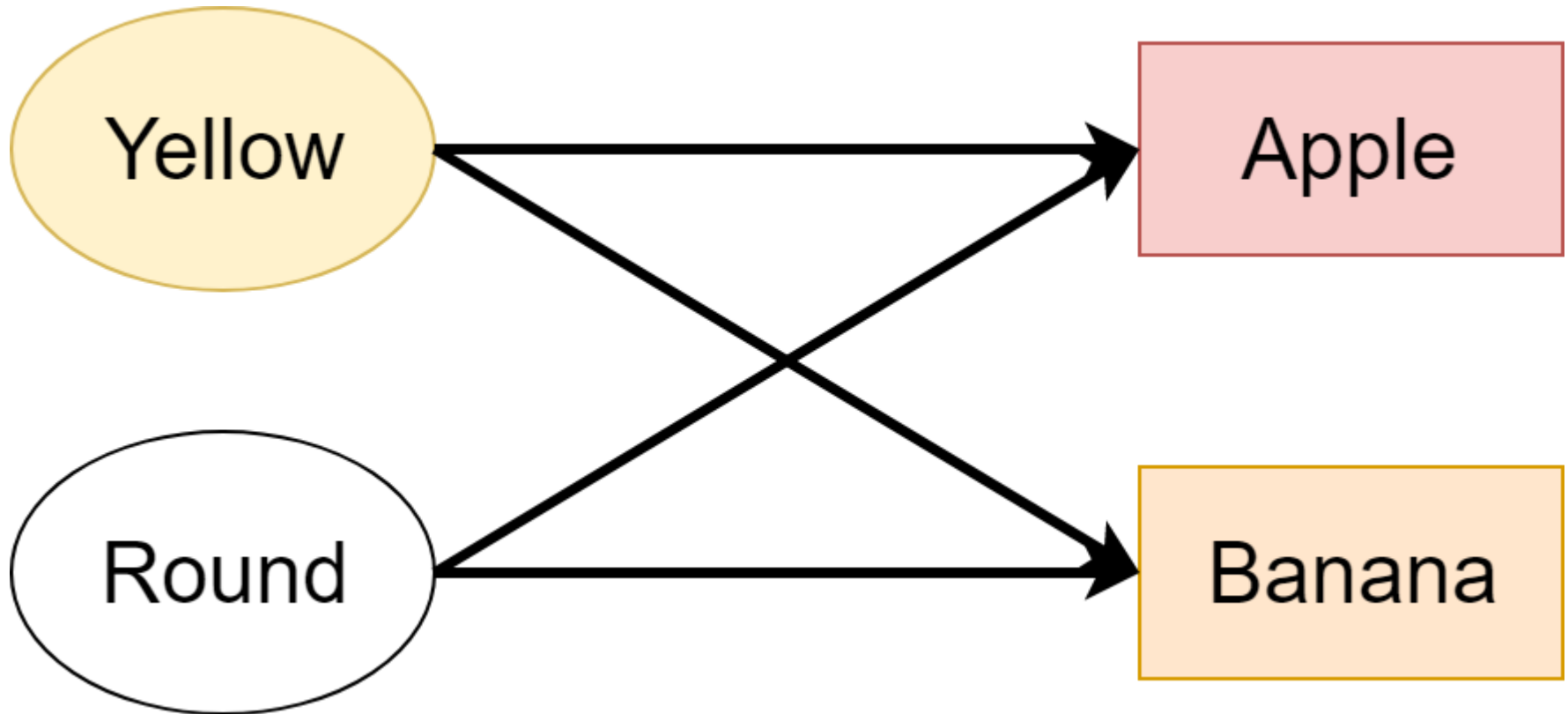
Banana

# Banana or Apple?



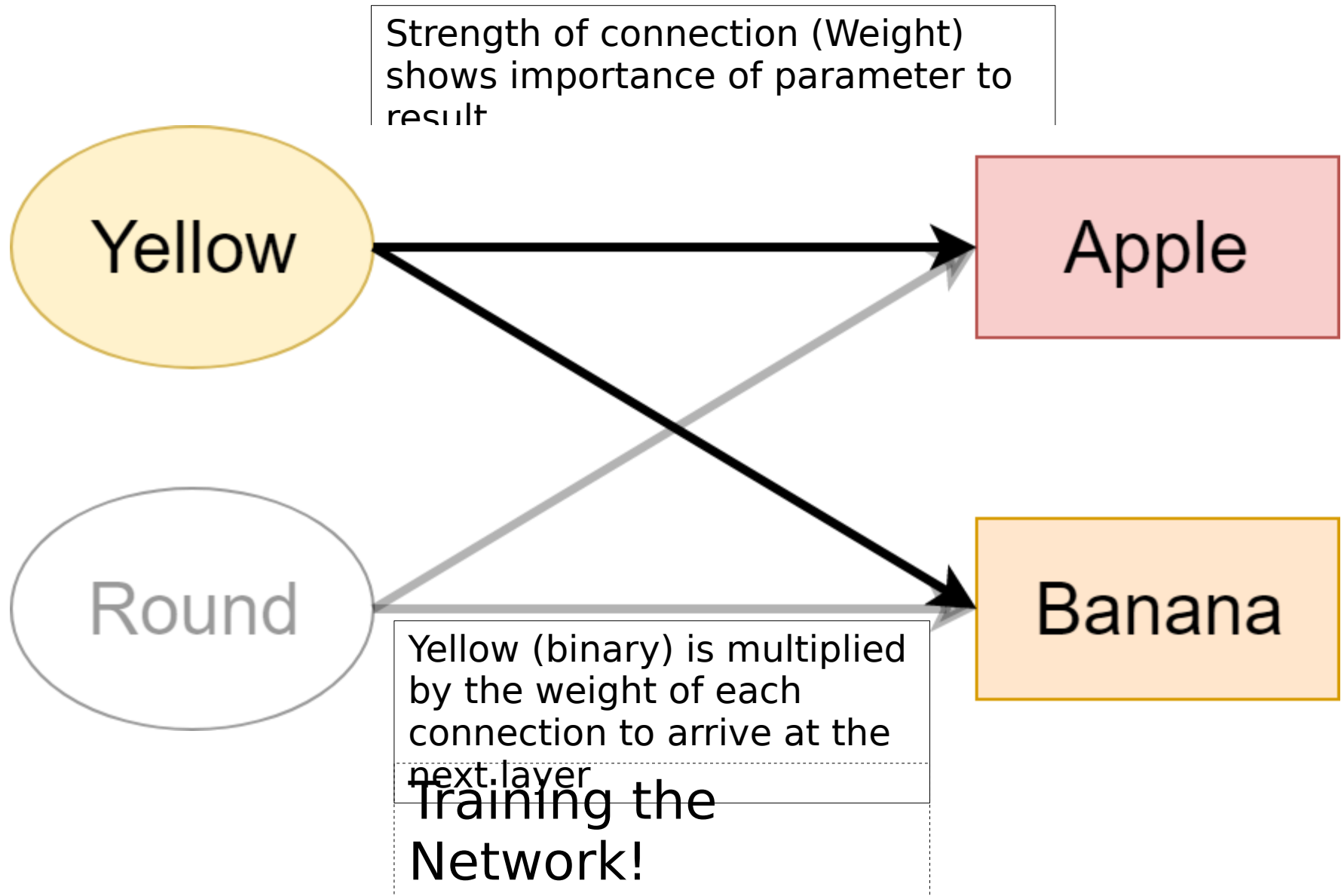


# Banana or Apple?



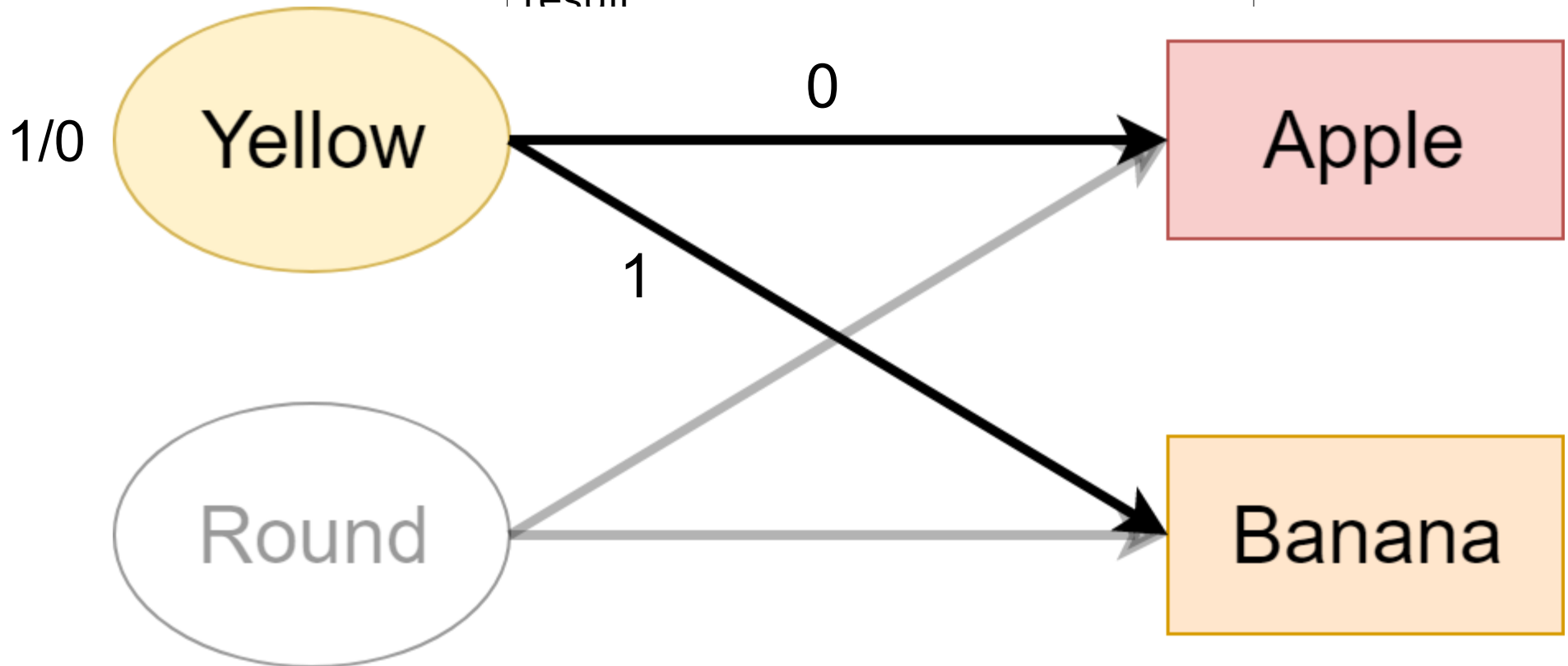
Training the  
Network!

# Banana or Apple?



# Banana or Apple?

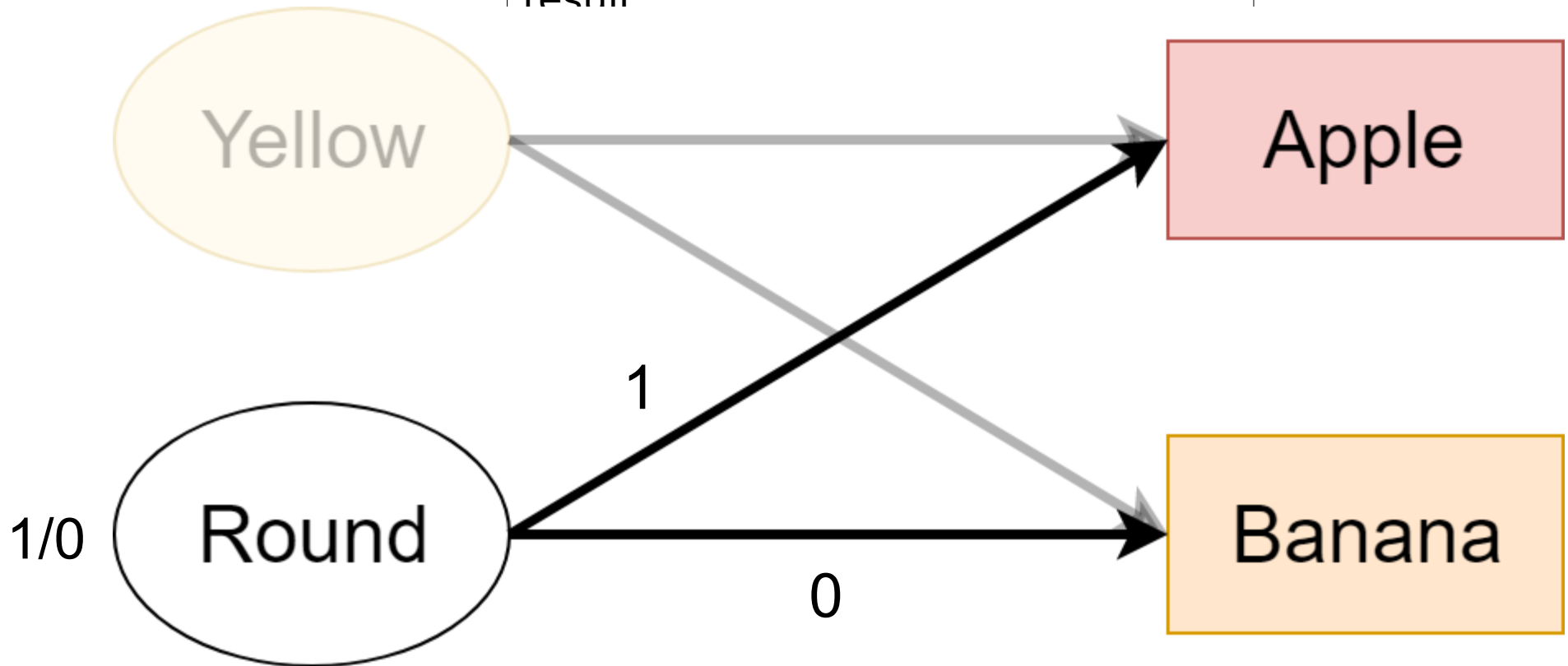
Strength of connection (Weight)  
shows importance of parameter to  
result



Training the  
Network!

# Banana or Apple?

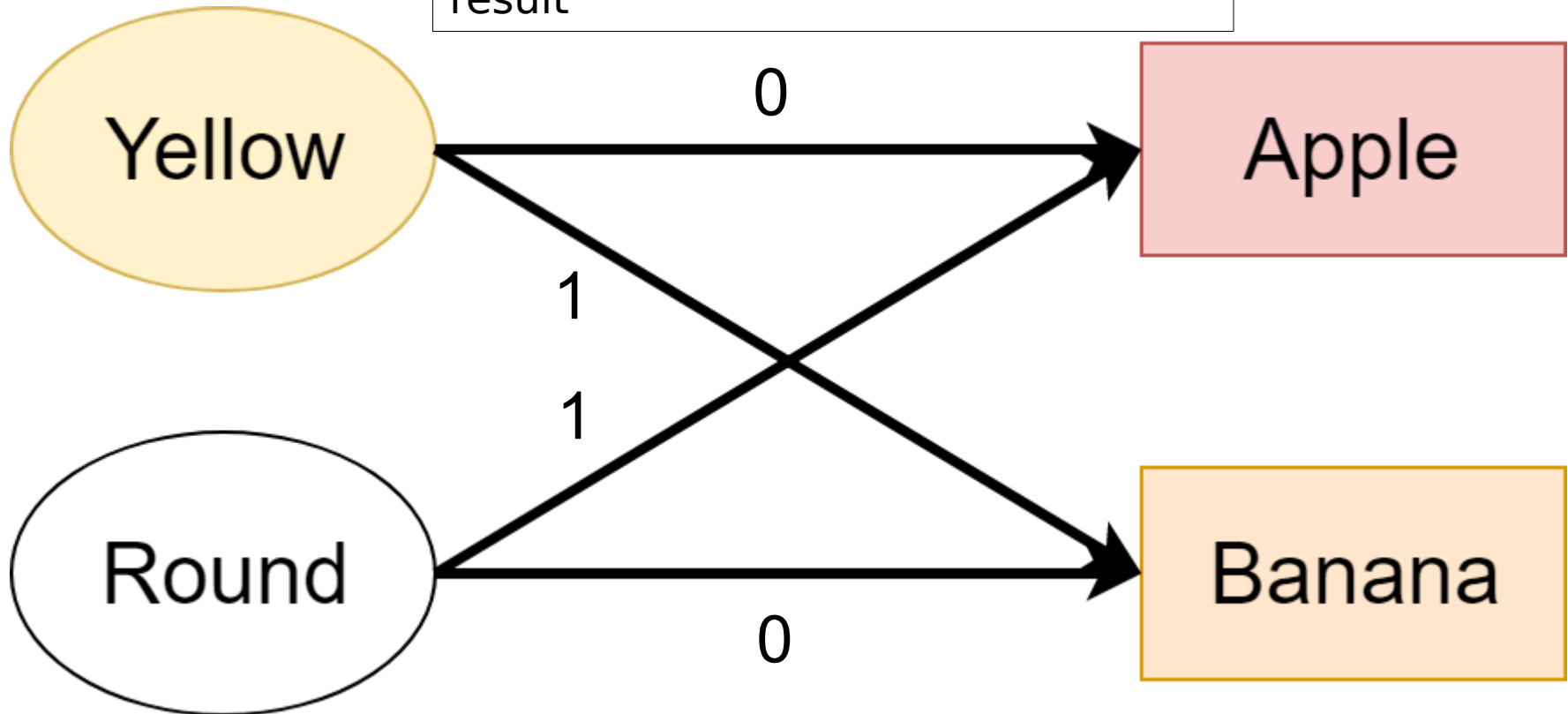
Strength of connection (Weight)  
shows importance of parameter to  
result



Training the  
Network!

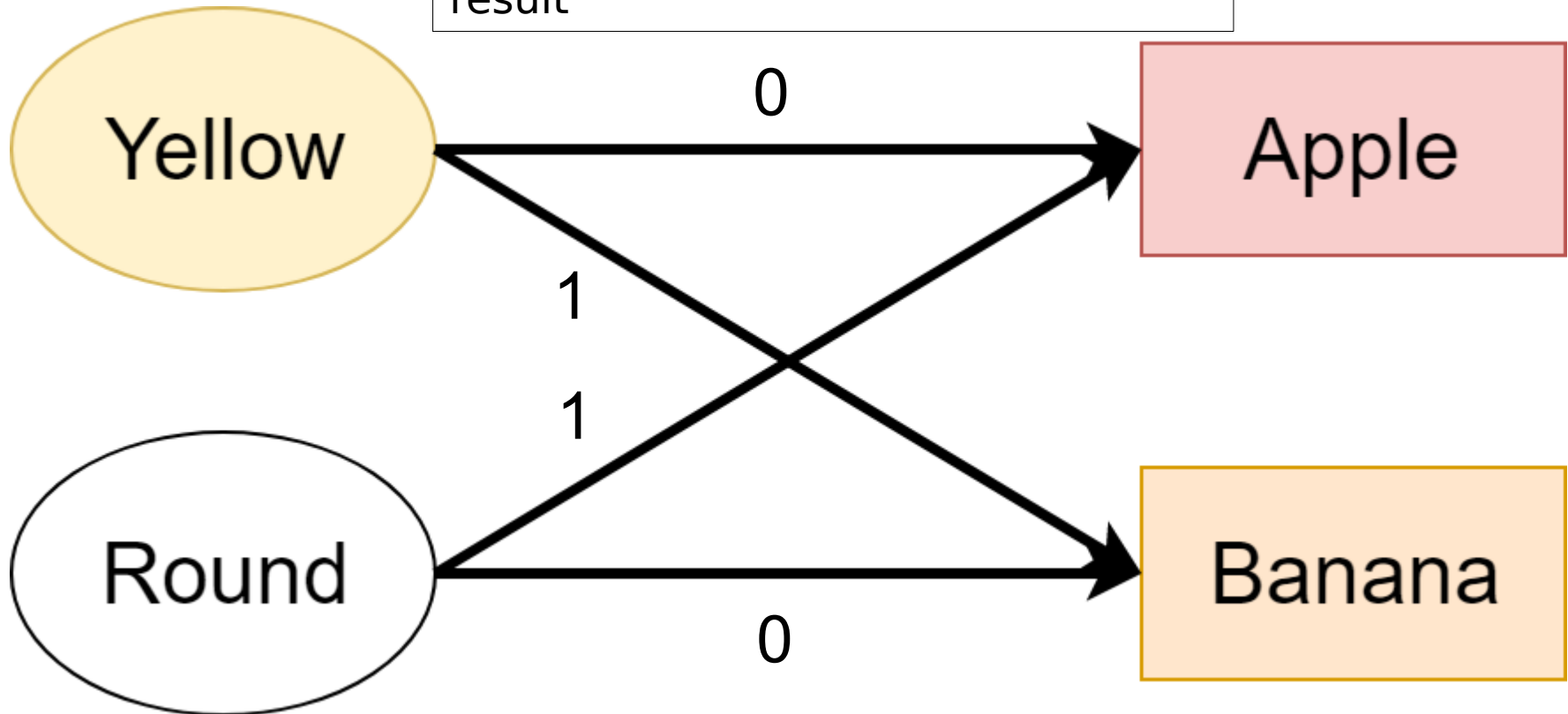
# Banana or Apple?

Strength of connection (Weight)  
shows importance of parameter to  
result



# Banana or Apple?

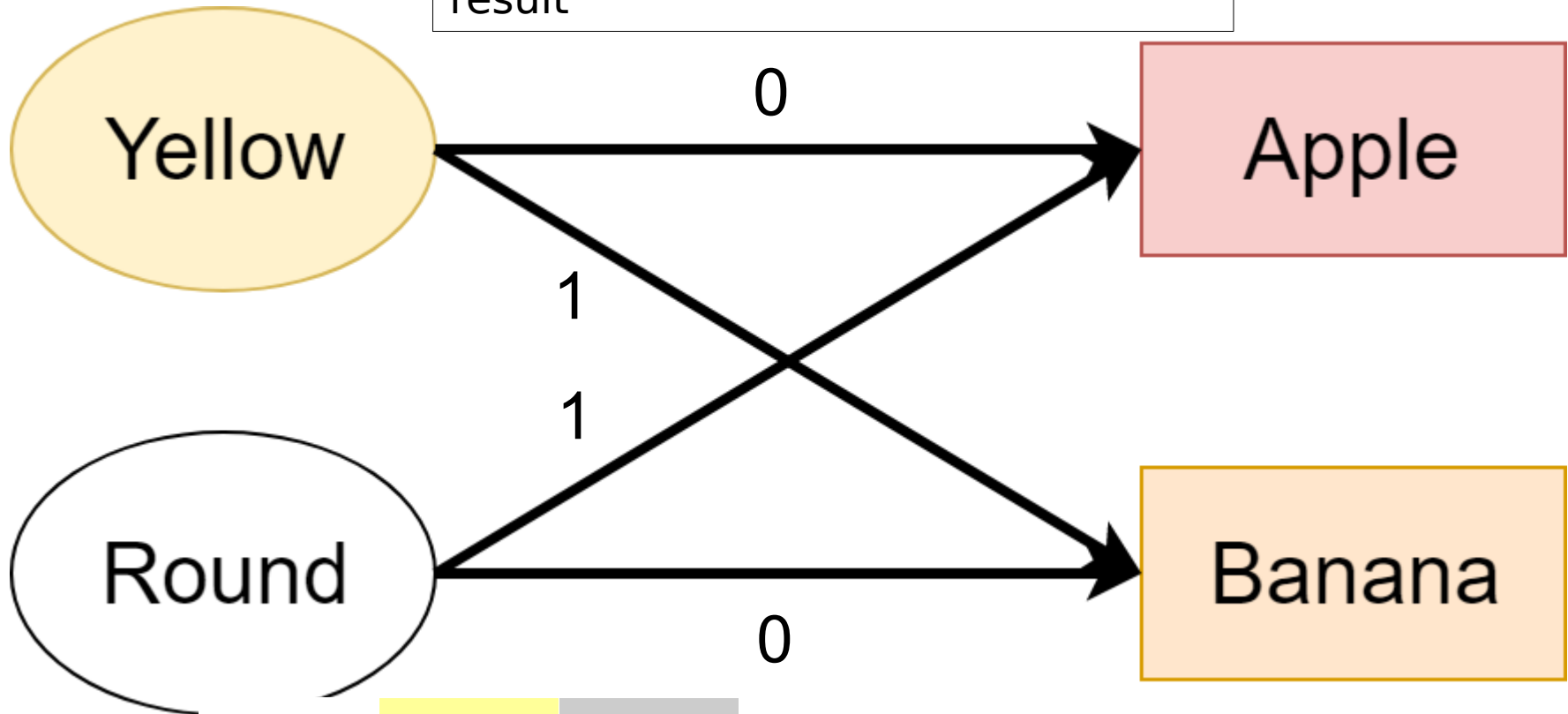
Strength of connection (Weight)  
shows importance of parameter to  
result



	Yellow	Round
Apple	0	1
Banana	1	0

# Banana or Apple?

Strength of connection (Weight)  
shows importance of parameter to  
result

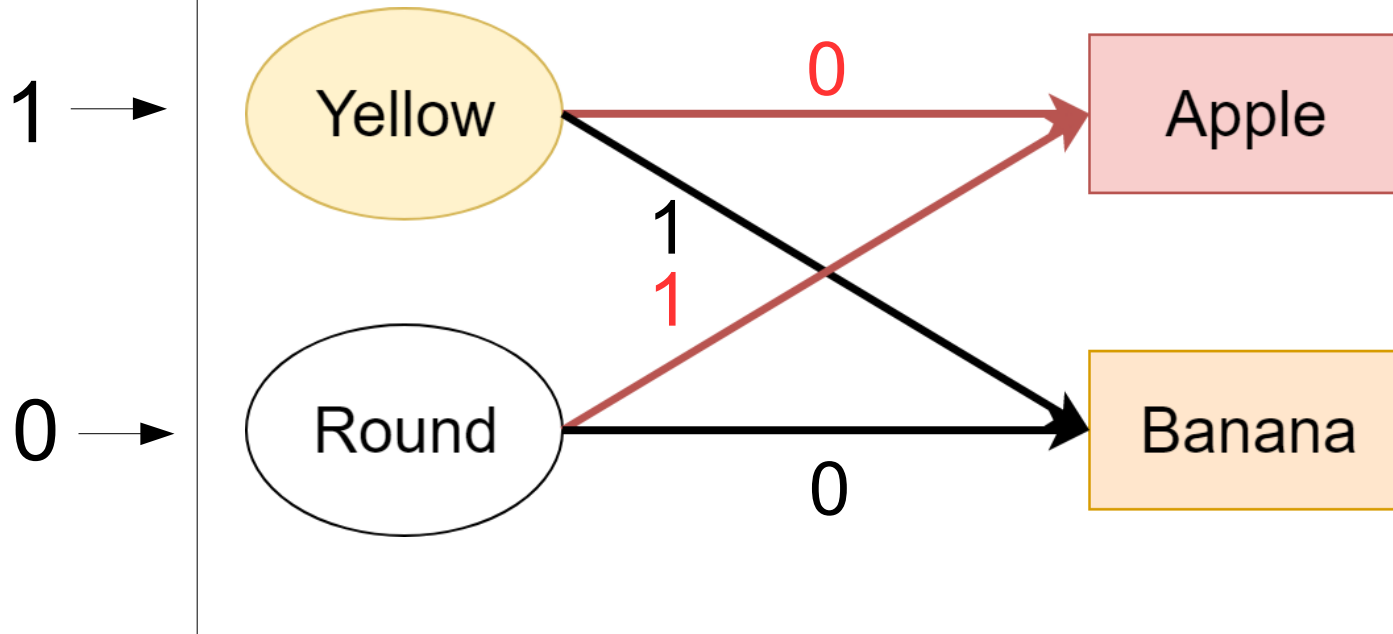


	Yellow	Round
Apple	0	1
Banana	1	0

$$\begin{aligned}\text{Apple} &= \text{Yellow} * 0 + \text{Round} * 1 \\ \text{Banana} &= \text{Yellow} * 1 + \text{Round} * 0\end{aligned}$$

# Banana or Apple?

INPUT



OUTPUT

$$1*0 + 0*1 = 0$$

= No Apple

$$1*1 + 0*0 = 1$$

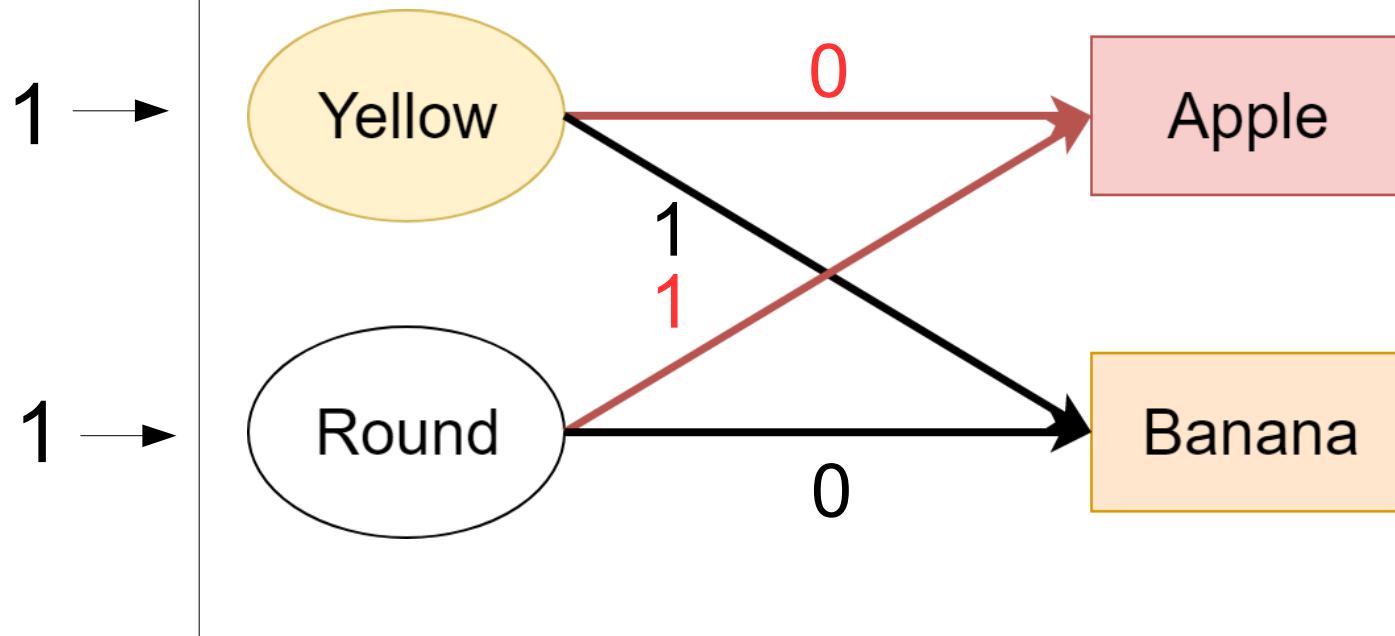
= Is Banana

$$\begin{aligned}\text{Apple} &= \text{Yellow} * 0 + \text{Round} * 1 \\ \text{Banana} &= \text{Yellow} * 1 + \text{Round} * 0\end{aligned}$$



# Banana or Apple?

INPUT



OUTPUT

$$1*0 + 1*1 = 1$$

= Is Apple

$$1*1 + 1*0 = 1$$

= Is Banana

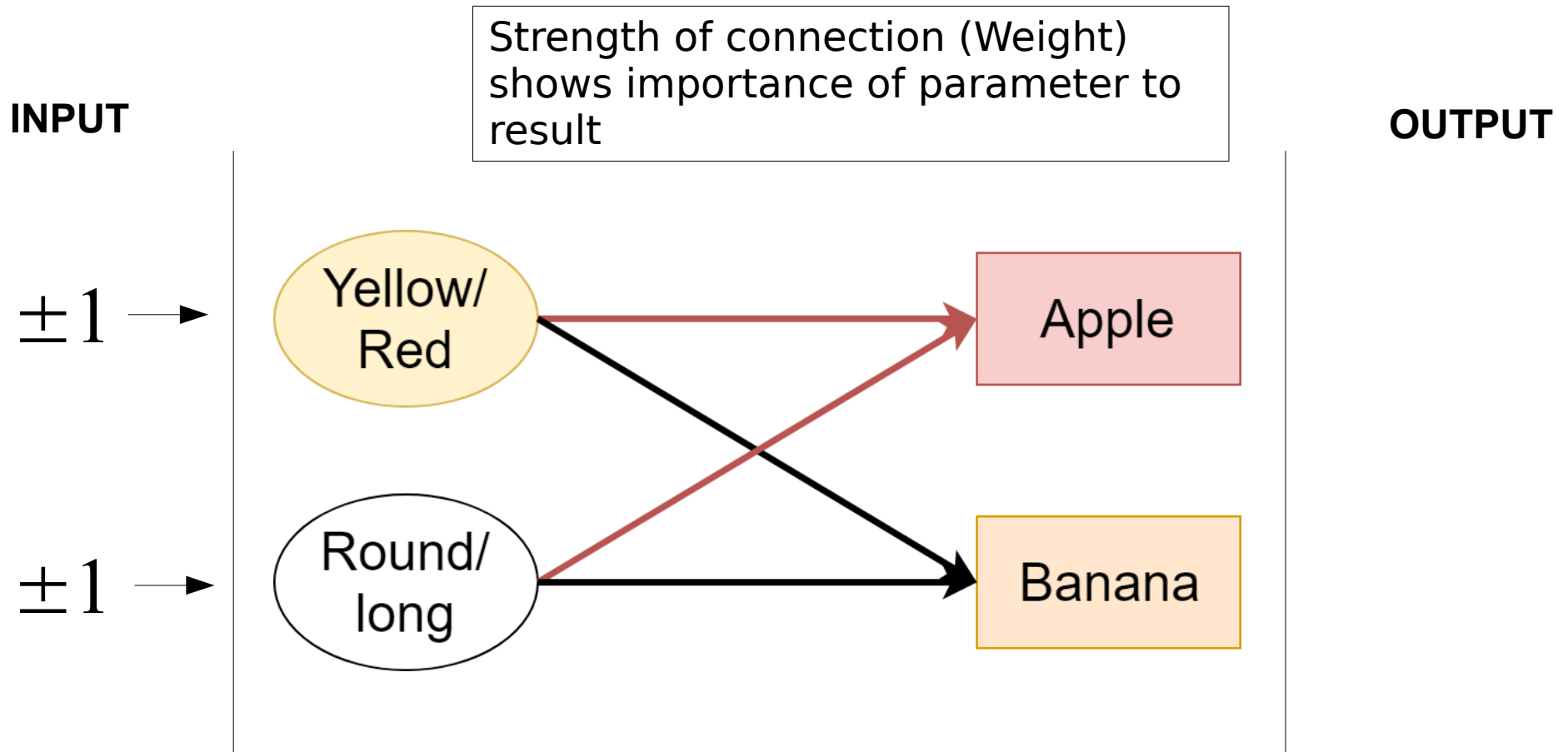
Our model is too simple

$$\begin{aligned}\text{Apple} &= \text{Yellow} * 0 + \text{Round} * 1 \\ \text{Banana} &= \text{Yellow} * 1 + \text{Round} * 0\end{aligned}$$

# Solution:

- Expand inputs and weight range
- Hopefully this will add the needed Complexity to the model

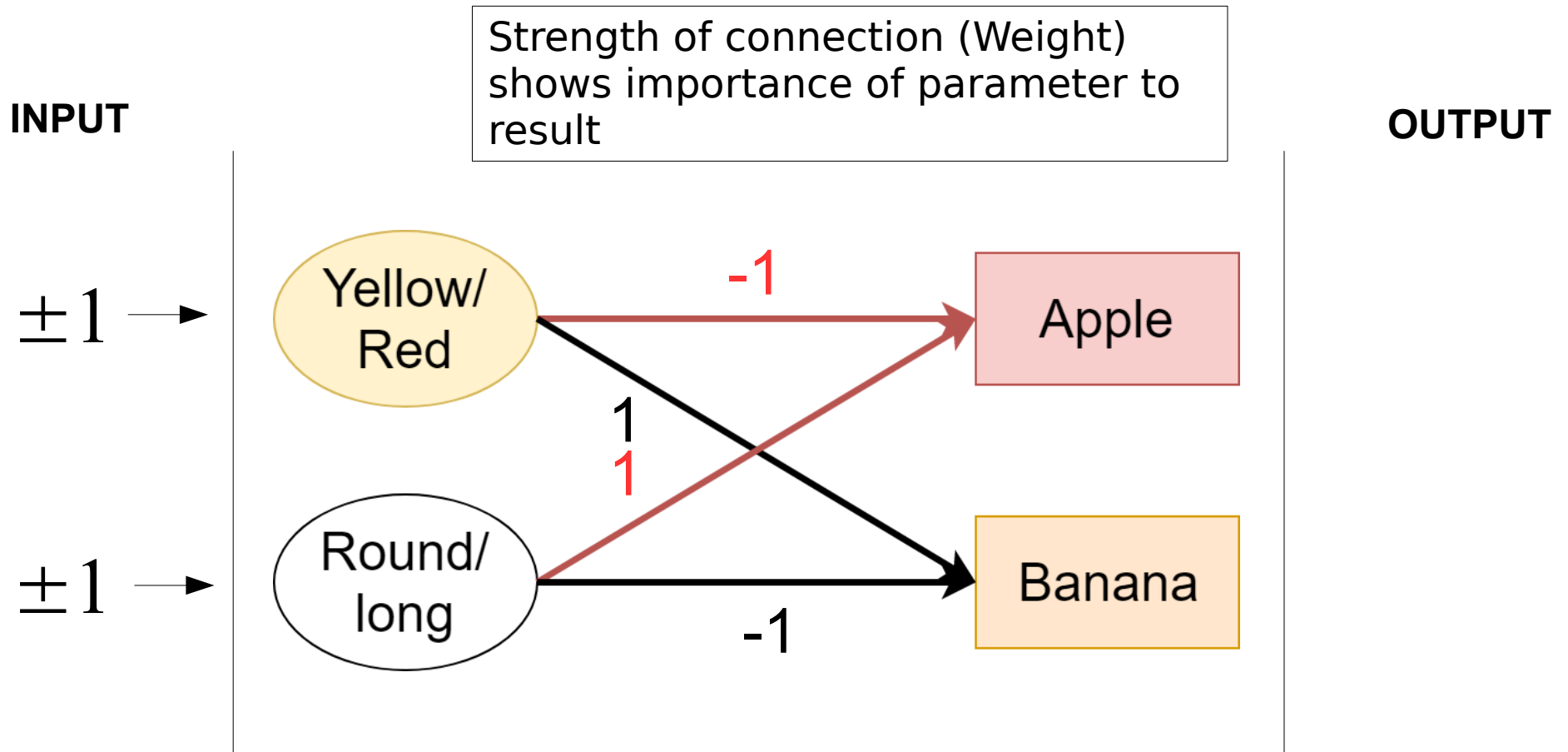
# Banana or Apple?



+1 is Yellow or Round like before,  
-1 is Green or Long,  
0 is lack of data.

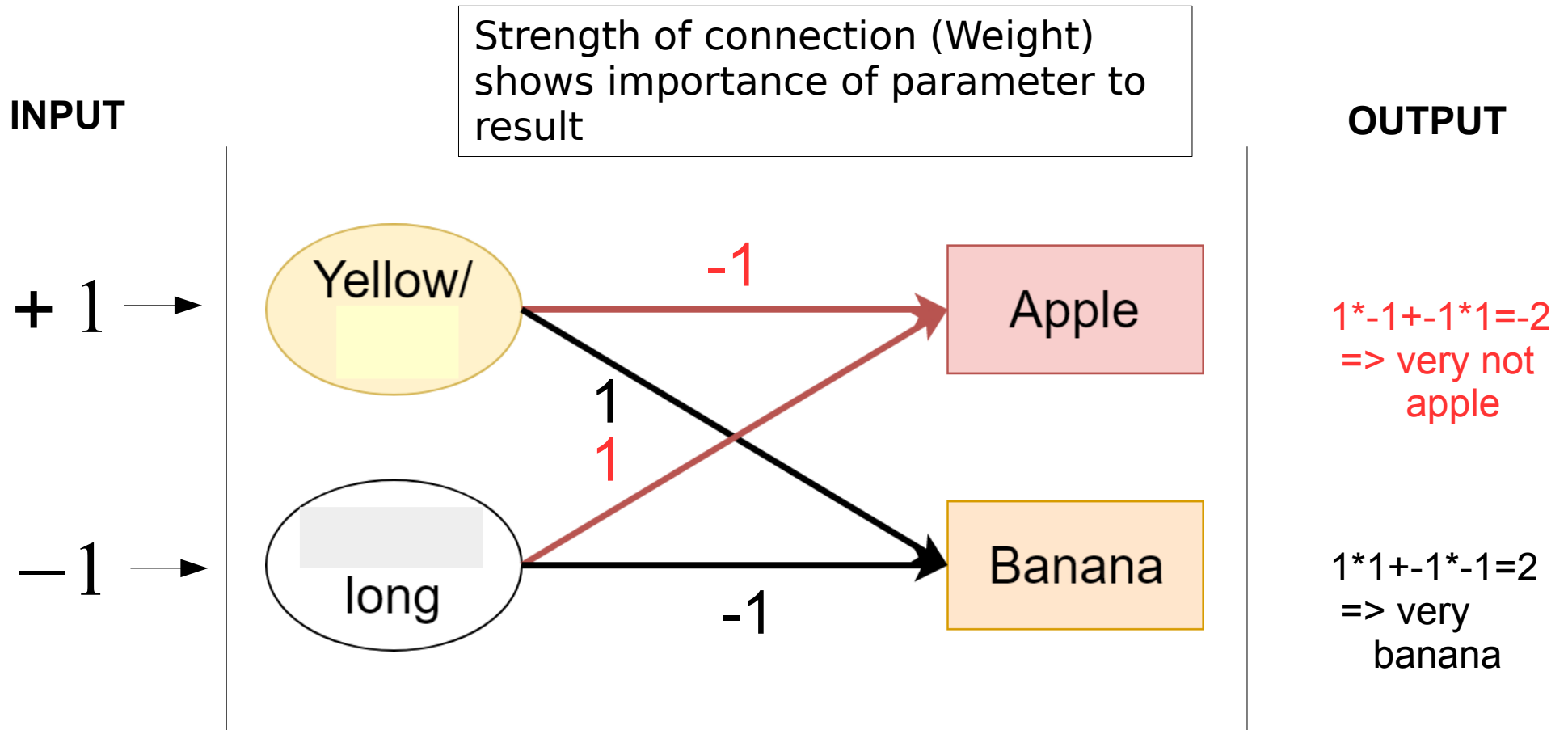
Value	Input 1	Input 2
+1	Yellow	Round
-1	Red	Long
0	No Data	No Data

# Banana or Apple?



+1 is Yellow or Round like before,  
-1 is Green or Long,  
0 is lack of data.

# Banana or Apple?



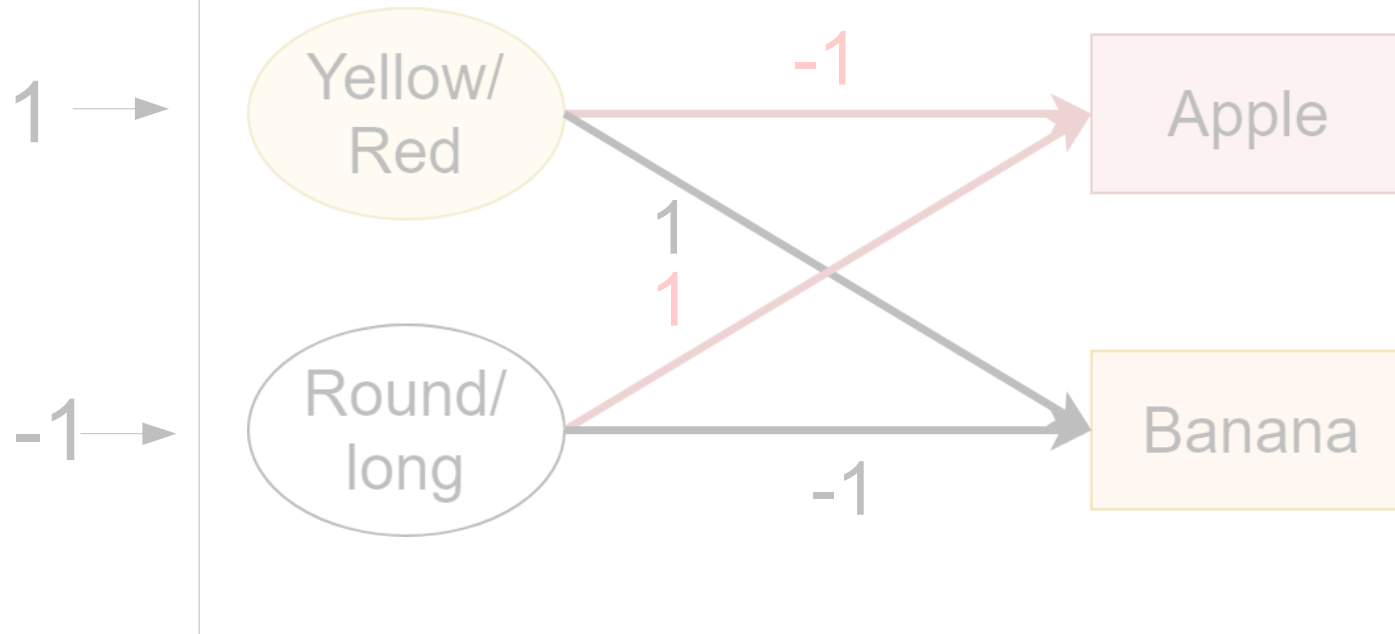
+1 is Yellow or Round like before,  
-1 is Green or Long,  
0 is lack of data.

Apple = Yellow \* 0 + Round \* 1  
Banana = Yellow \* 1 + Round \* 0

# Banana or Apple?

Strength of connection (Weight) shows importance of parameter to result

INPUT



OUTPUT

$1 * -1 + -1 * 1 = -2$   
 $\Rightarrow$  very not apple

$1 * 1 + -1 * -1 = 2$   
 $\Rightarrow$  very banana

+1 is Yellow or Round like before,  
-1 is Green or Long,  
0 is lack of data.

Normalize to  $\{0,1\}$  range  
to ease comparison

# Banana or Apple?

Strength of connection (Weight) shows importance of parameter to result

INPUT

How about we cut and squeeze the output into a  $\{0,1\}$  range?

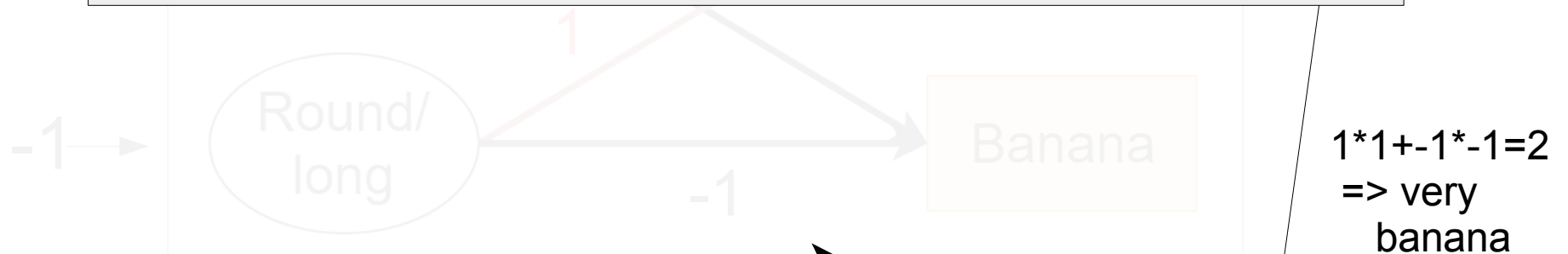
For  $x > 0$  :  $f(x) = x/2$

For  $x < 0$  :  $f(x) = 0$

Apple = -2  $\rightarrow$  =0

Banana = 2  $\rightarrow$  =1

$+ -1 * 1 = -2$   
very not  
apple

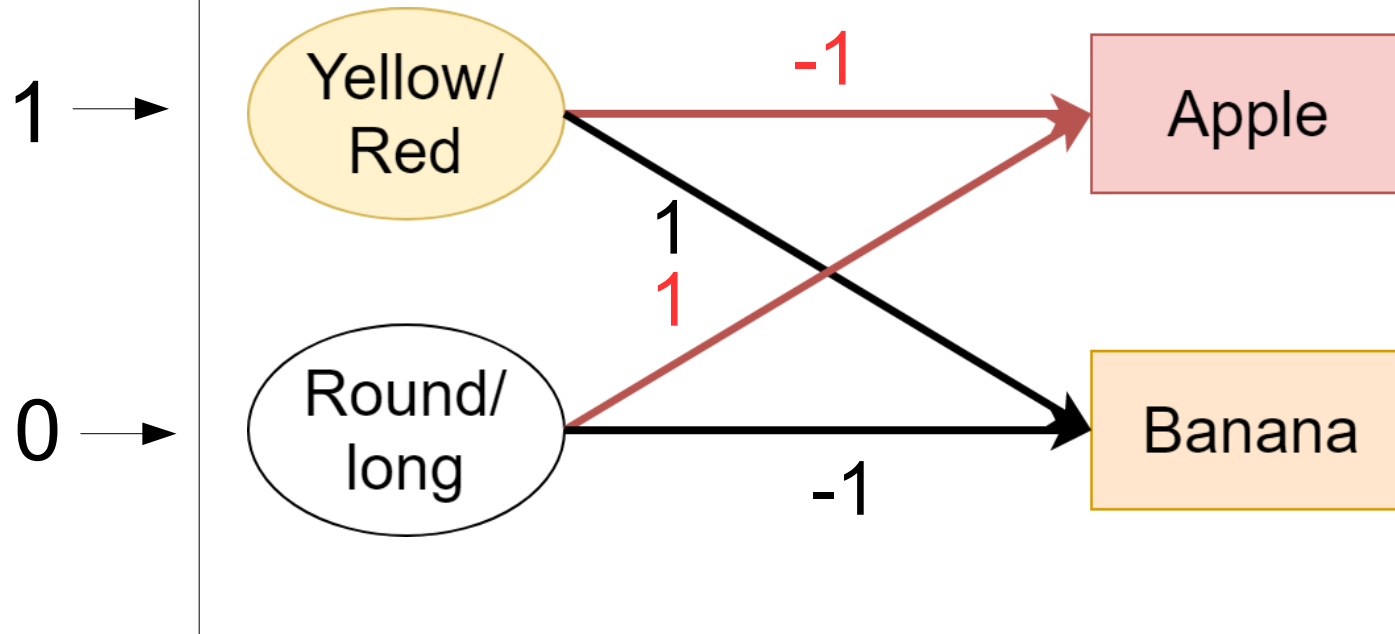


+1 is Yellow or Round like before,  
-1 is Green or Long,  
0 is lack of data.

Normalize to  $\{0,1\}$  range  
to ease comparison

# Banana or Apple?

INPUT



OUTPUT  
 $f(x) = \text{abs}(x/2)$

$1 * -1 + 0 * 1 = -1$   
 $f(-1) = 0$   
 $\Rightarrow$  not apple

$1 * 1 + 0 * -1 = 1$   
 $F(1) = 0.5$   
 $\Rightarrow$  50% banana

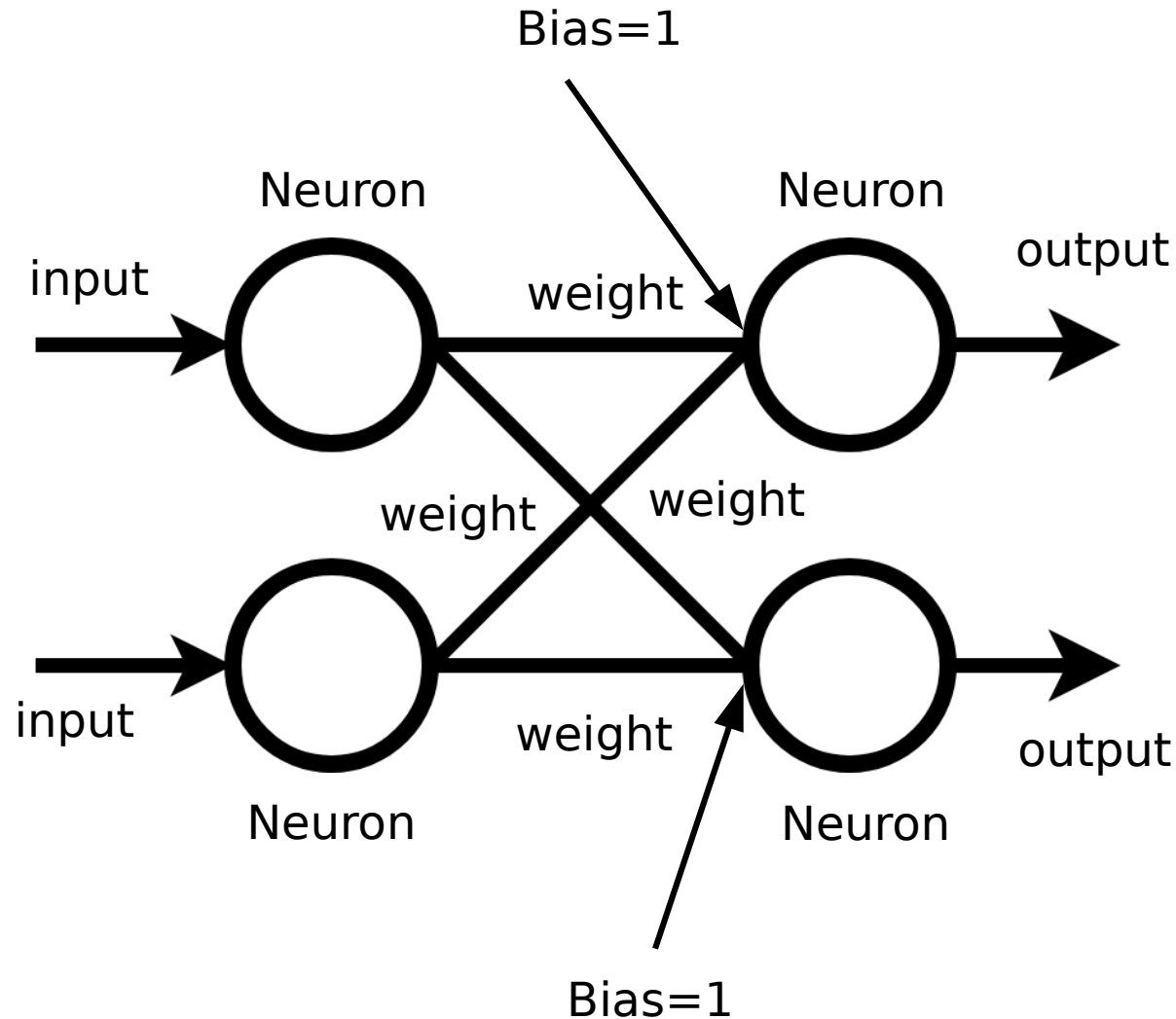
Apple = Yellow \* 0 + Round \* 1  
Banana = Yellow \* 1 + Round \* 0

$\Rightarrow$

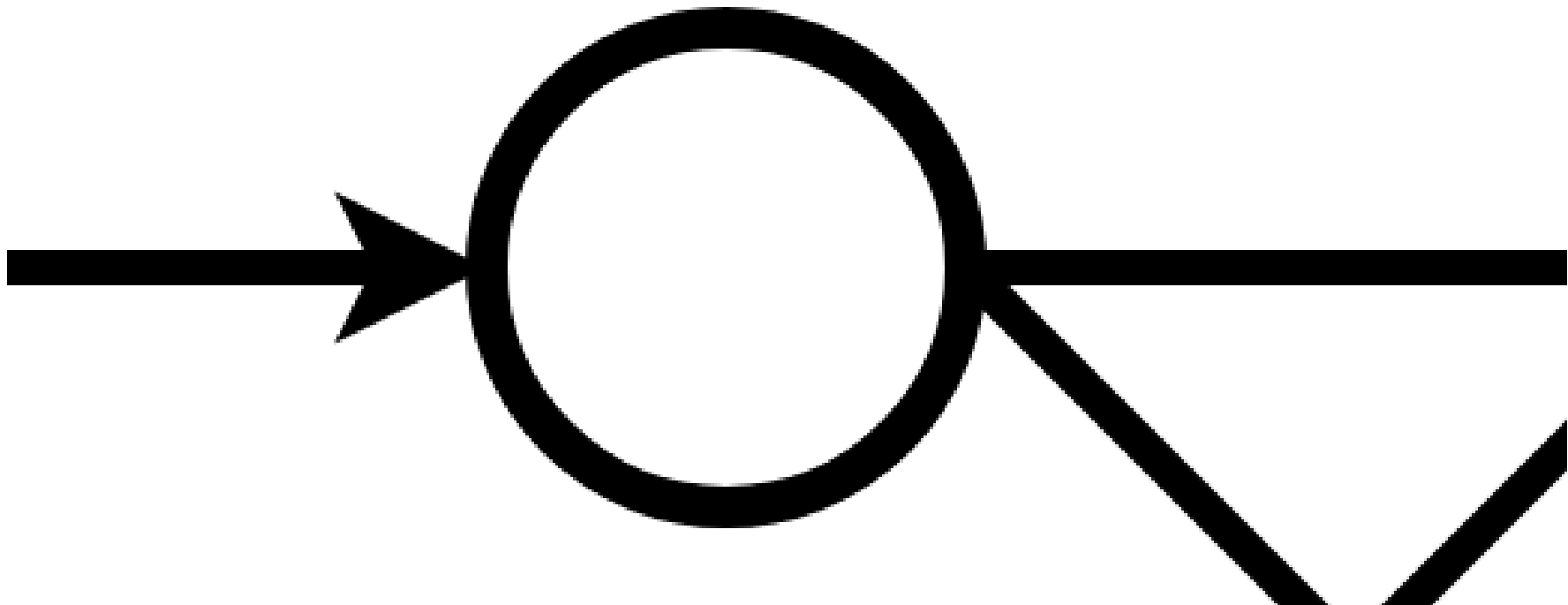
For  $x > 0$  :  $f(x) = x/2$   
For  $x < 0$  :  $f(x) = 0$



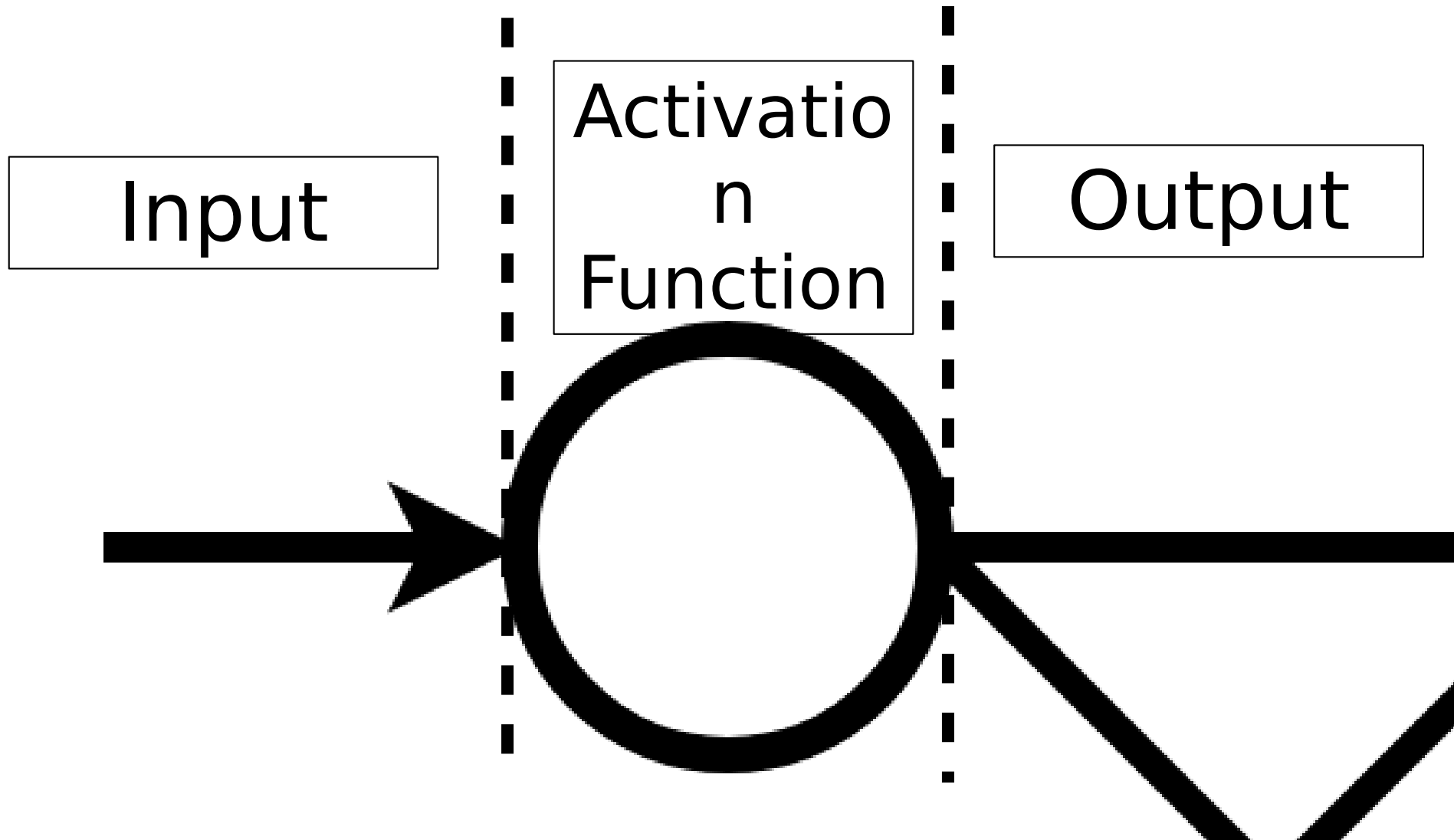
# Multilayer Perceptrons



# Multilayer Perceptrons



# Multilayer Perceptrons

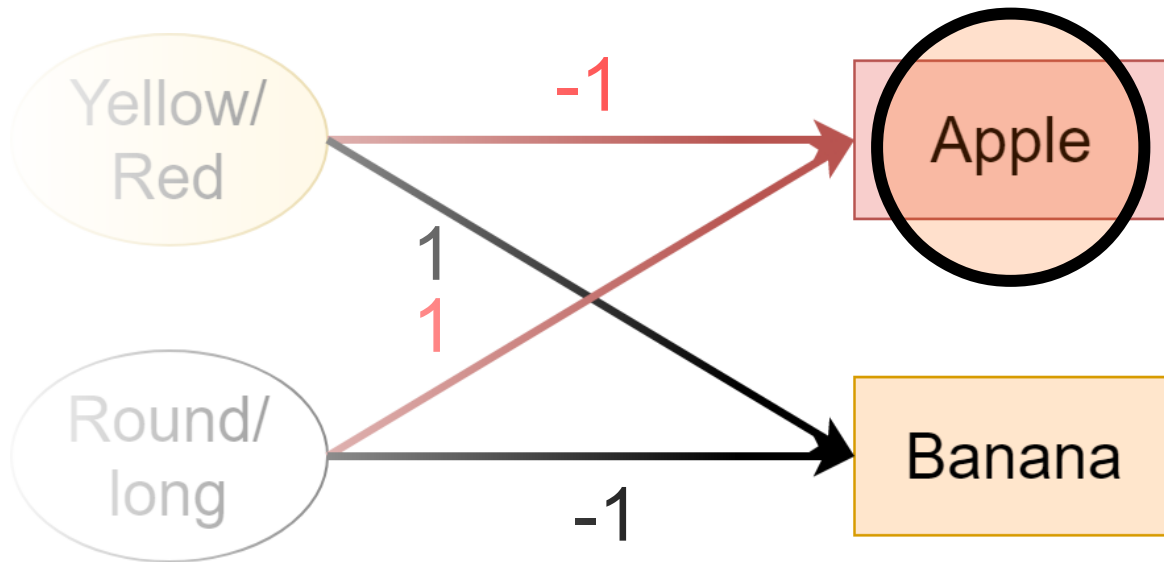


# Multilayer Perceptrons

Input

Activation  
Function

Output



$$\sum (Weights * Inputs) + Bias$$

=

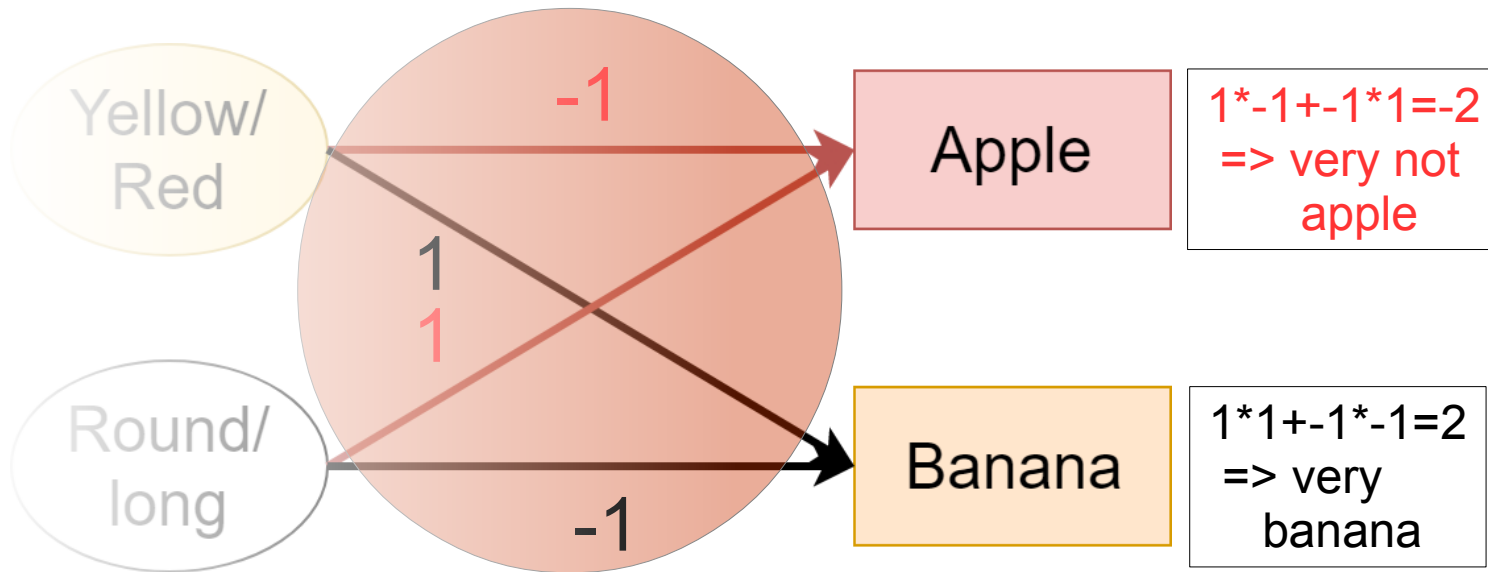
$$\sum_{i=0}^m (w_i * x_i) + Bias$$

# Multilayer Perceptrons

Input

Activation  
Function

Output



$$\sum (Weights * Inputs)$$

=

$$\sum_{i=0}^m (w_i * x_i)$$

# Multilayer Perceptrons

Input

Activation  
Function

Output

How about we cut and squeeze the output into a  $\{0,1\}$  range?

For  $x > 0$  :  $f(x) = x/2$

For  $x < 0$  :  $f(x) = 0$

Apple = -2  $\rightarrow$  =0

Banana = 2  $\rightarrow$  =1

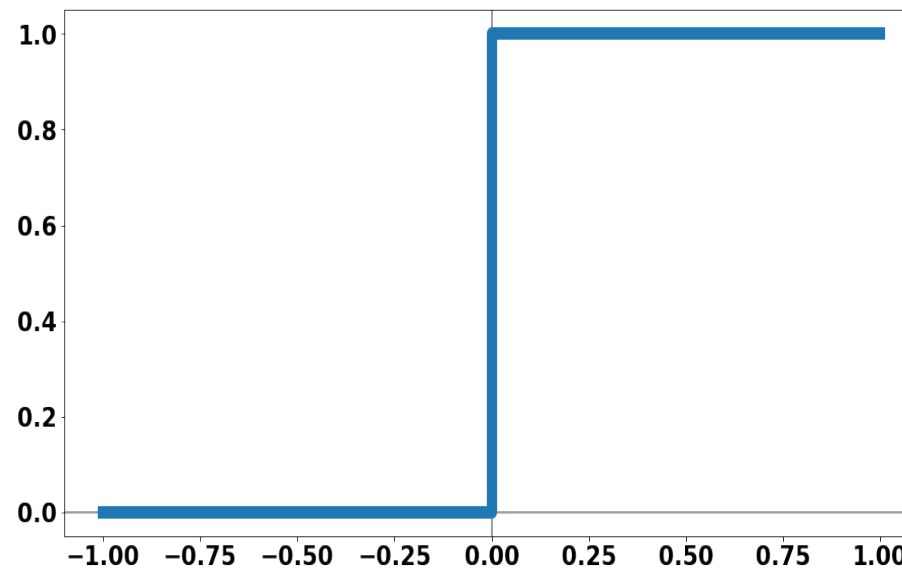
$f(Inputs)$

# Multilayer Perceptrons

Input

Activation  
Function

Output



Step function (Heaviside)

$$f(x) = \begin{cases} 0 & \text{for } x \leq 0 \\ 1 & \text{for } x > 0 \end{cases}$$

$1 * -1 + -1 * 1 = -2$   
 $\Rightarrow$  very not  
Apple  
 $\Rightarrow f(x) = 0$

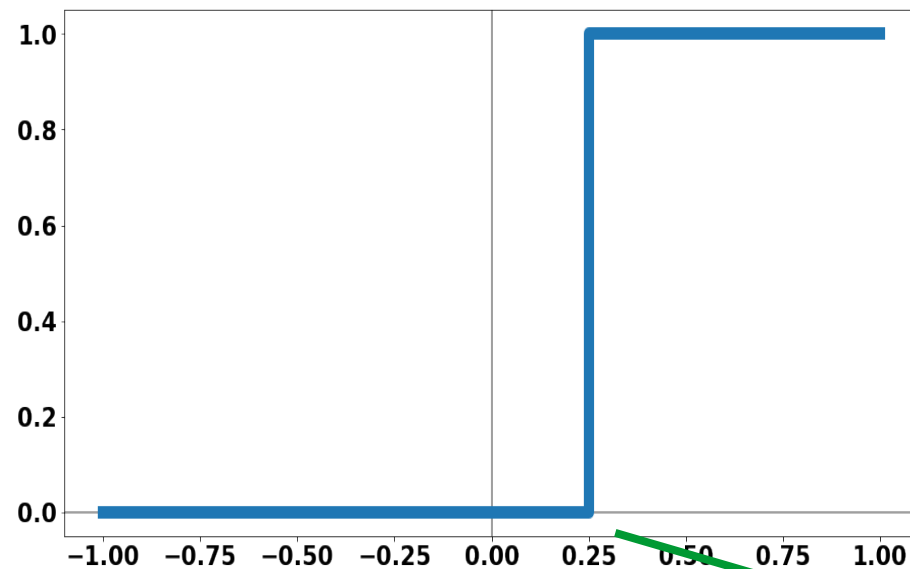
$1 * 1 + -1 * -1 = 2$   
 $\Rightarrow$  very  
Banana  
 $\Rightarrow f(x) = 1$

# Multilayer Perceptrons

Input

Activation  
Function

Output



Step function (Heaviside)

$$f(x) = \begin{cases} 0 & \text{for } x \leq 0 + 0.25 \\ 1 & \text{for } x > 0 + 0.25 \end{cases}$$

$$f(x) = \text{step function}(x - 0.25)$$

Bias

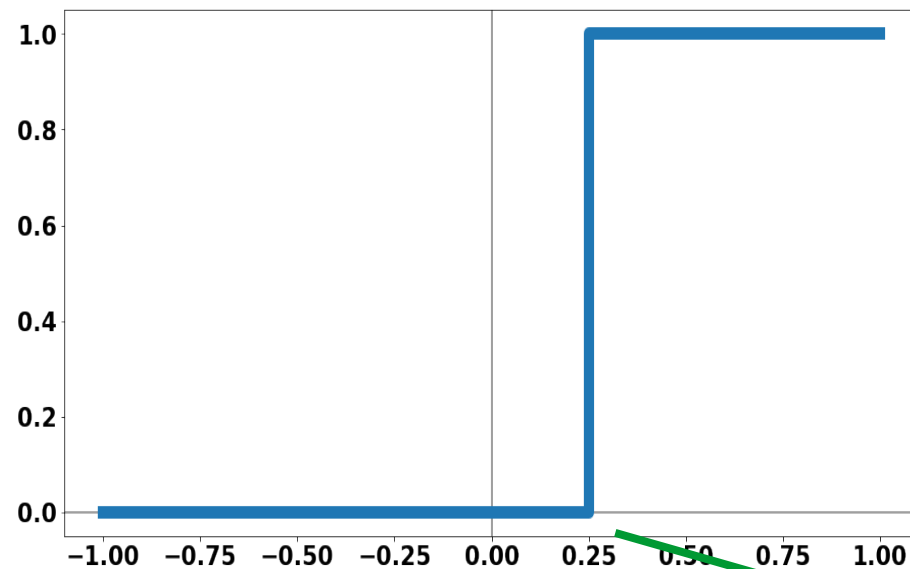


# Multilayer Perceptrons

Input

Activation  
Function

Output



Step function (Heaviside)

$$f(x) = \begin{cases} 0 & \text{for } x \leq 0 \\ 1 & \text{for } x > 0 \end{cases} + 0.25$$

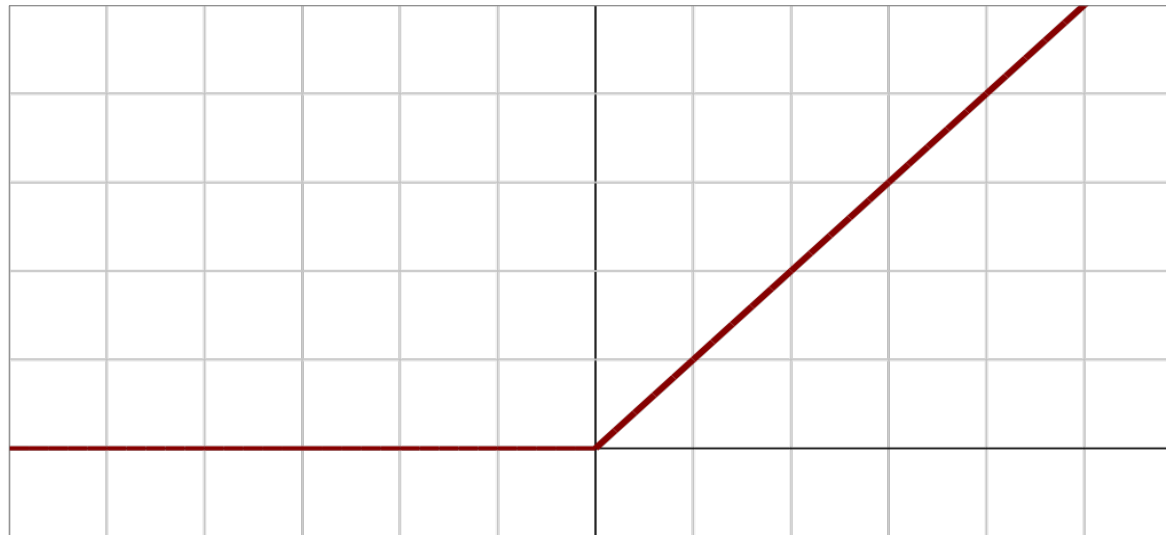
$$\text{Bias} = 1 * w_b$$

# Multilayer Perceptrons

Input

Activation  
Function

Output



$1 * -1 + -1 * 1 = -2$   
 $\Rightarrow$  very not  
Apple  
 $\Rightarrow f(x) = 0$

$1 * 1 + -1 * -1 = 2$   
 $\Rightarrow$  very  
Banana  
 $\Rightarrow f(x) = 2$

RELU (Rectified Linear Units)

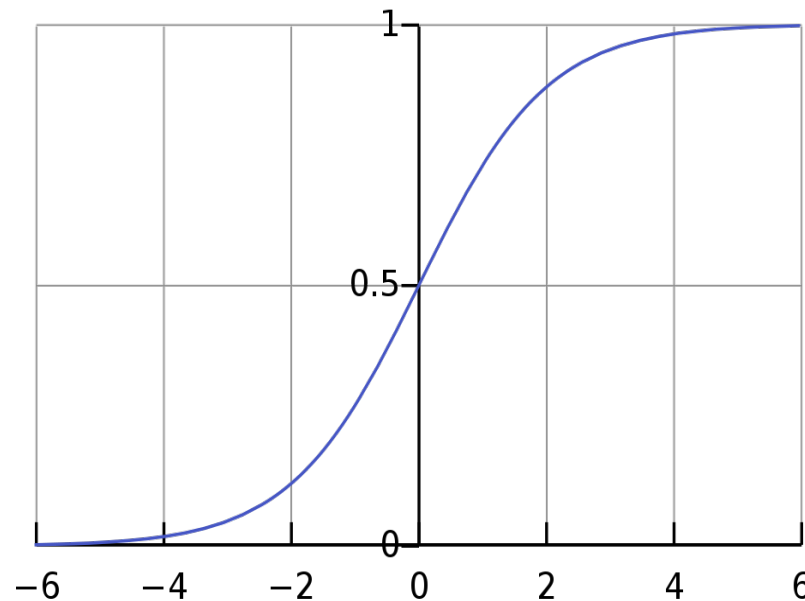
$$f(x) = \begin{cases} 0 & \text{for } x \leq 0 \\ x & \text{for } x > 0 \end{cases}$$

# Multilayer Perceptrons

Input

Activation  
Function

Output



Sigmoid function (logistic  
function)

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

$1 * -1 + -1 * 1 = -2$   
 $\Rightarrow$  very not  
Apple  
 $\Rightarrow f(x) = 0.1192$

$1 * 1 + -1 * -1 = 2$   
 $\Rightarrow$  very  
Banana  
 $\Rightarrow f(x) = 0.8807$

...aka Softstep..

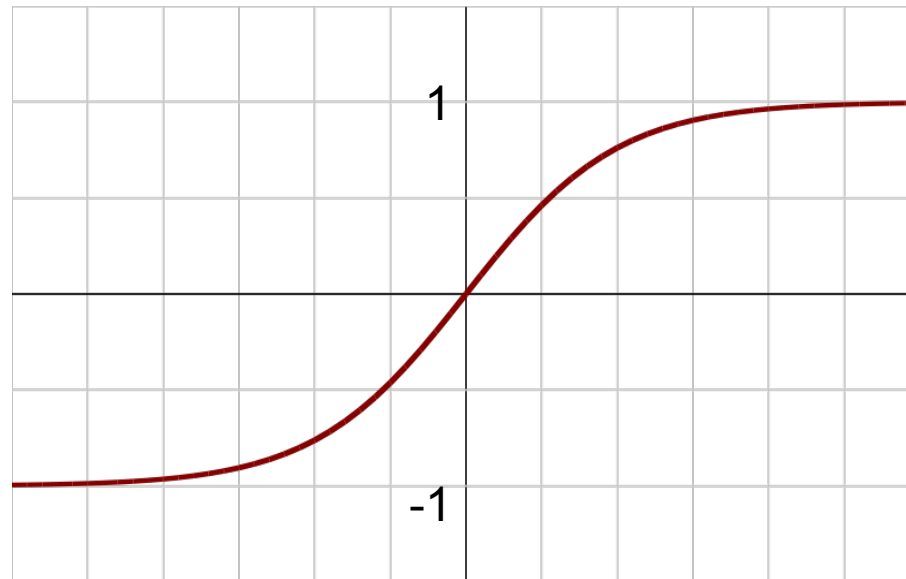
...aka Logistic Curve...

# Multilayer Perceptrons

Input

Activation  
Function

Output



Tanh function

$$f(x) = \frac{2}{1 + e^{-2x}} - 1 = 2 \operatorname{sigmoid}(2x) - 1$$

$1 * -1 + -1 * 1 = -2$   
 $\Rightarrow$  very not  
Apple  
 $\Rightarrow f(x) = -0.964$

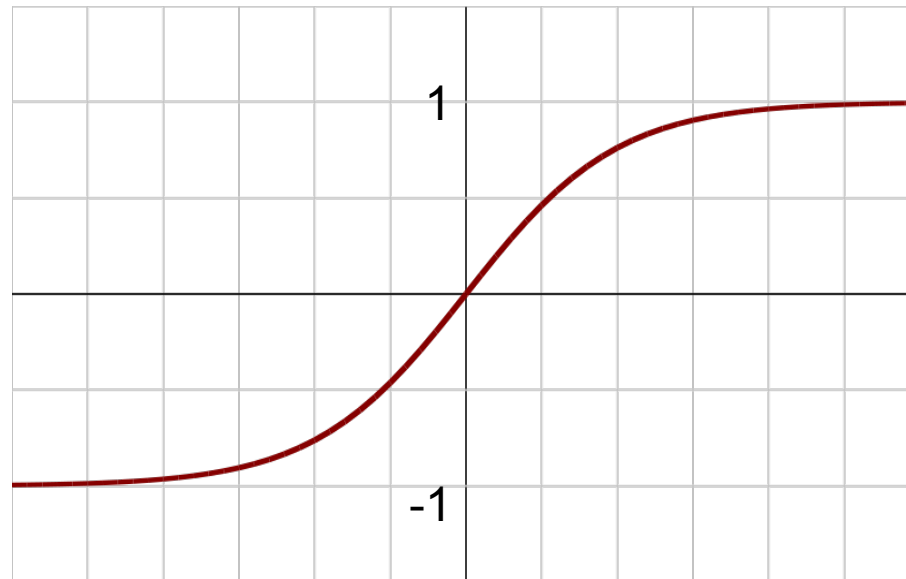
$1 * 1 + -1 * -1 = 2$   
 $\Rightarrow$  very  
Banana  
 $\Rightarrow f(x) = 0.964$

# Multilayer Perceptrons

Input

Activation  
Function

Output



Tanh function

$1 * -1 + -1 * 1 = -2$   
 $\Rightarrow$  very not  
apple

$1 * 1 + -1 * -1 = 2$   
 $\Rightarrow$  very  
banana

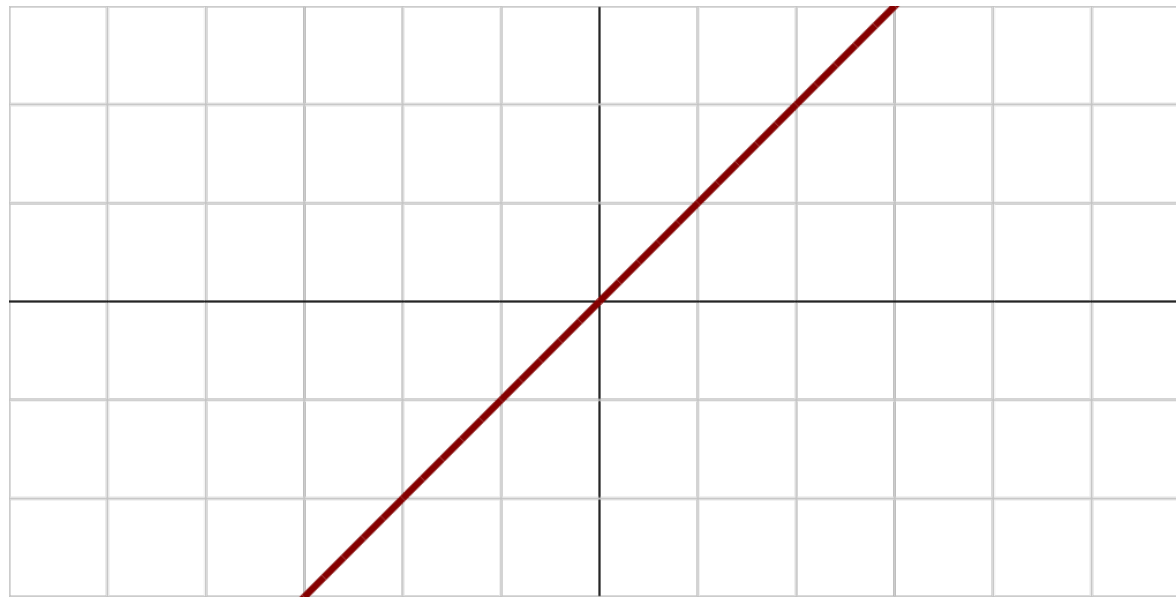
Range  $\{-1, 1\}$  might be used for when you want to make negative predictions, such as Yellow + Long has not just a zero chance of being apple, but a negative chance.

# Multilayer Perceptrons

Input

Activation  
Function

Output



Identity function

$$f(x) = x$$

$1 * -1 + -1 * 1 = -2$   
 $\Rightarrow$  very not  
apple

$1 * 1 + -1 * -1 = 2$   
 $\Rightarrow$  very  
banana



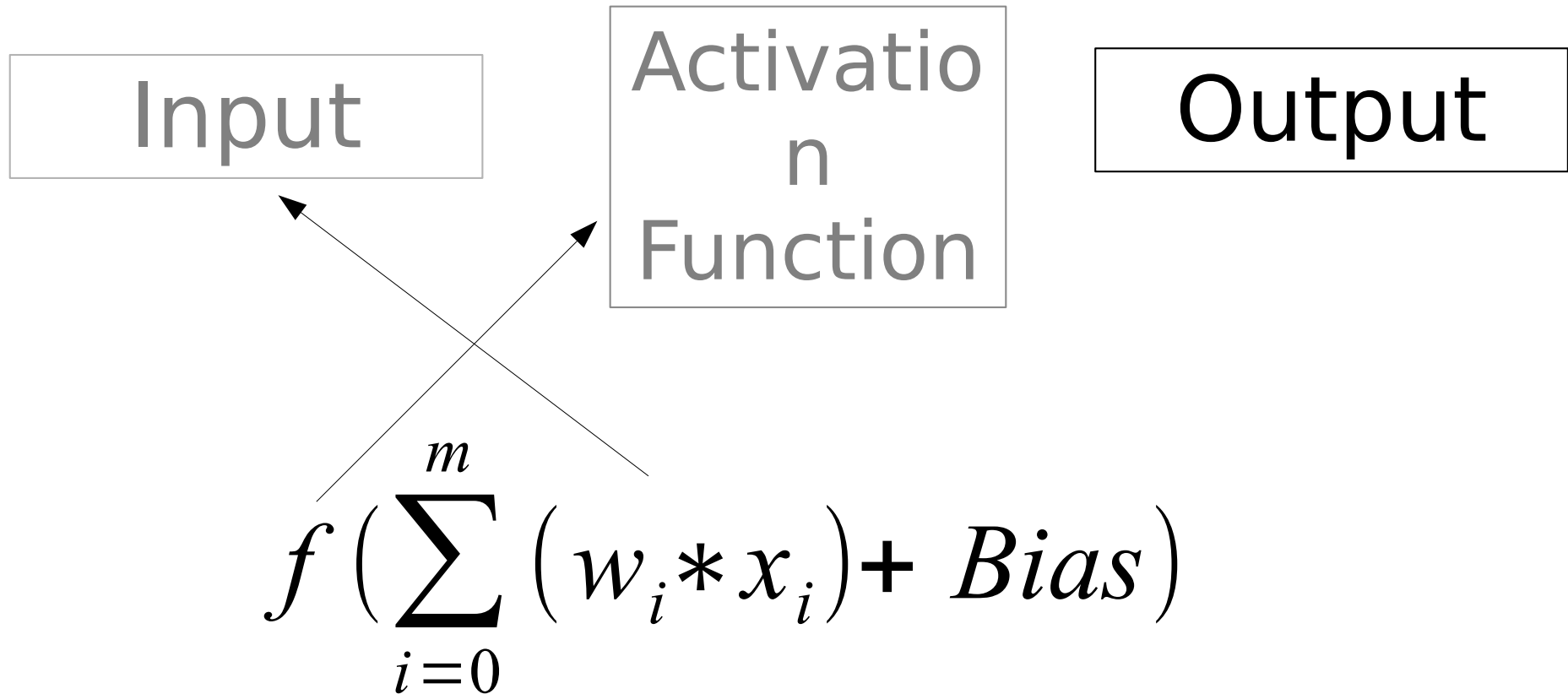
- Monotonic – When the activation function is monotonic, the error surface associated with a single-layer model is guaranteed to be convex.<sup>[4]</sup>
- Smooth Functions with a Monotonic derivative – These have been shown to generalize better in some cases. The argument for these properties suggests that such activation functions are more consistent with Occam's razor.<sup>[5]</sup>
- Approximates identity near the origin – When activation functions have this property, the neural network will learn efficiently when its weights are initialized with small random values. When the activation function does not approximate identity near the origin, special care must be used when initializing the weights.

The following table compares the properties of several activation functions that are functions of one fold  $x$  from the previous layer or layers:

Name	Plot	Equation	Derivative (with respect to $x$ )	Range	Order of continuity	Monotonic	Derivative Monotonic	Approximates identity near the origin
Identity		$f(x) = x$	$f'(x) = 1$	$(-\infty, \infty)$	$C^\infty$	Yes	Yes	Yes
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$	$\{0, 1\}$	$C^{-1}$	Yes	No	No
Logistic (a.k.a. Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$	$(0, 1)$	$C^\infty$	Yes	No	No
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$	$(-1, 1)$	$C^\infty$	Yes	No	Yes
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$	$(-\frac{\pi}{2}, \frac{\pi}{2})$	$C^\infty$	Yes	No	Yes
Softsign <sup>[7][8]</sup>		$f(x) = \frac{x}{1 +  x }$	$f'(x) = \frac{1}{(1 +  x )^2}$	$(-1, 1)$	$C^1$	Yes	No	Yes
Inverse square root unit (ISRU) <sup>[9]</sup>		$f(x) = \frac{x}{\sqrt{1 + \alpha x^2}}$	$f'(x) = \left( \frac{1}{\sqrt{1 + \alpha x^2}} \right)^3$	$(-\frac{1}{\sqrt{\alpha}}, \frac{1}{\sqrt{\alpha}})$	$C^\infty$	Yes	No	Yes
Rectified linear unit (ReLU) <sup>[10]</sup>		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$[0, \infty)$	$C^0$	Yes	Yes	No
Leaky rectified linear unit (Leaky ReLU) <sup>[11]</sup>		$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	$C^0$	Yes	Yes	No
Parametric rectified linear unit (PReLU) <sup>[12]</sup>		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	$C^0$	Yes if $\alpha \geq 0$	Yes	Yes if $\alpha = 1$
Randomized leaky rectified linear unit (RRReLU) <sup>[13]</sup>		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	$C^0$	Yes	Yes	No
Exponential linear unit (ELU) <sup>[14]</sup>		$f(\alpha, x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha e^x & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	$C^0$ when $\alpha = 1$ otherwise $C^1$	Yes if $\alpha \geq 0$	Yes if $0 \leq \alpha \leq 1$	Yes if $\alpha = 1$
Scaled exponential linear unit (SELU) <sup>[15]</sup>		$f(\alpha, x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ with $\lambda = 1.67326$ and $\alpha = 1.67326$	$f'(\alpha, x) = \begin{cases} \alpha e^x & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	$C^0$	Yes	No	No
S-shaped rectified linear activation unit (SRReLU) <sup>[16]</sup>		$f_{t_l, a_l, t_r, a_r}(x) = \begin{cases} 0 & \text{for } x \leq t_l \\ a_l(x - t_l) & \text{for } t_l < x < t_r \\ t_r + a_r(x - t_r) & \text{for } x \geq t_r \end{cases}$ $t_l, t_r, a_l, a_r$ are parameters.	$f'(x) = \begin{cases} 0 & \text{for } x \leq t_l \\ a_l & \text{for } t_l < x < t_r \\ a_r & \text{for } x \geq t_r \end{cases}$	$(-\infty, \infty)$	$C^0$	No	No	No
Inverse square root linear unit (ISRLU) <sup>[2]</sup>		$f(x) = \begin{cases} \frac{x}{\sqrt{1 + \alpha x^2}} & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \frac{1}{(1 + \alpha x^2)^{3/2}} & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\frac{1}{\sqrt{\alpha}}, \infty)$	$C^2$	Yes	Yes	Yes
Adaptive piecewise linear (APL) <sup>[17]</sup>		$f(x) = \sum_{i=1}^n a_i \max(0, x + b_i)$	$f'(x) = H(x) - \sum_{i=1}^n a_i H(-x + b_i)$ <sup>[2]</sup>	$(-\infty, \infty)$	$C^0$	No	No	No
SoftPlus <sup>[18]</sup>		$f(x) = \ln(e^x + 1)$	$f'(x) = \frac{1}{1 + e^{-x}}$	$(0, \infty)$	$C^\infty$	Yes	Yes	No
Bent identity		$f(x) = \frac{\sqrt{x^2 + 1} - 1}{2} + x$	$f'(x) = \frac{x}{2\sqrt{x^2 + 1}} + 1$	$(-\infty, \infty)$	$C^\infty$	Yes	Yes	Yes
SoftExponential <sup>[19]</sup>		$f(\alpha, x) = \begin{cases} -\frac{\ln(1 - \alpha(x + \alpha))}{\alpha} & \text{for } x < 0 \\ x & \text{for } x = 0 \\ \frac{e^{\alpha x} - 1}{\alpha} + \alpha & \text{for } x > 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \frac{1}{1 - \alpha(x + \alpha)} & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	$C^\infty$	Yes	Yes	Yes if $\alpha = 0$
Sinoid <sup>[20]</sup>		$f(x) = \sin(x)$	$f'(x) = \cos(x)$	$[-1, 1]$	$C^\infty$	No	No	Yes
Sinc		$f(x) = \begin{cases} 1 & \text{for } x = 0 \\ \frac{\sin(x)}{x} & \text{for } x \neq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x = 0 \\ \frac{\cos(x)}{x} - \frac{\sin(x)}{x^2} & \text{for } x \neq 0 \end{cases}$	$[\approx -0.217234, 1]$	$C^\infty$	No	No	No
Gaussian		$f(x) = e^{-x^2}$	$f'(x) = -2xe^{-x^2}$	$(0, 1]$	$C^\infty$	No	No	No

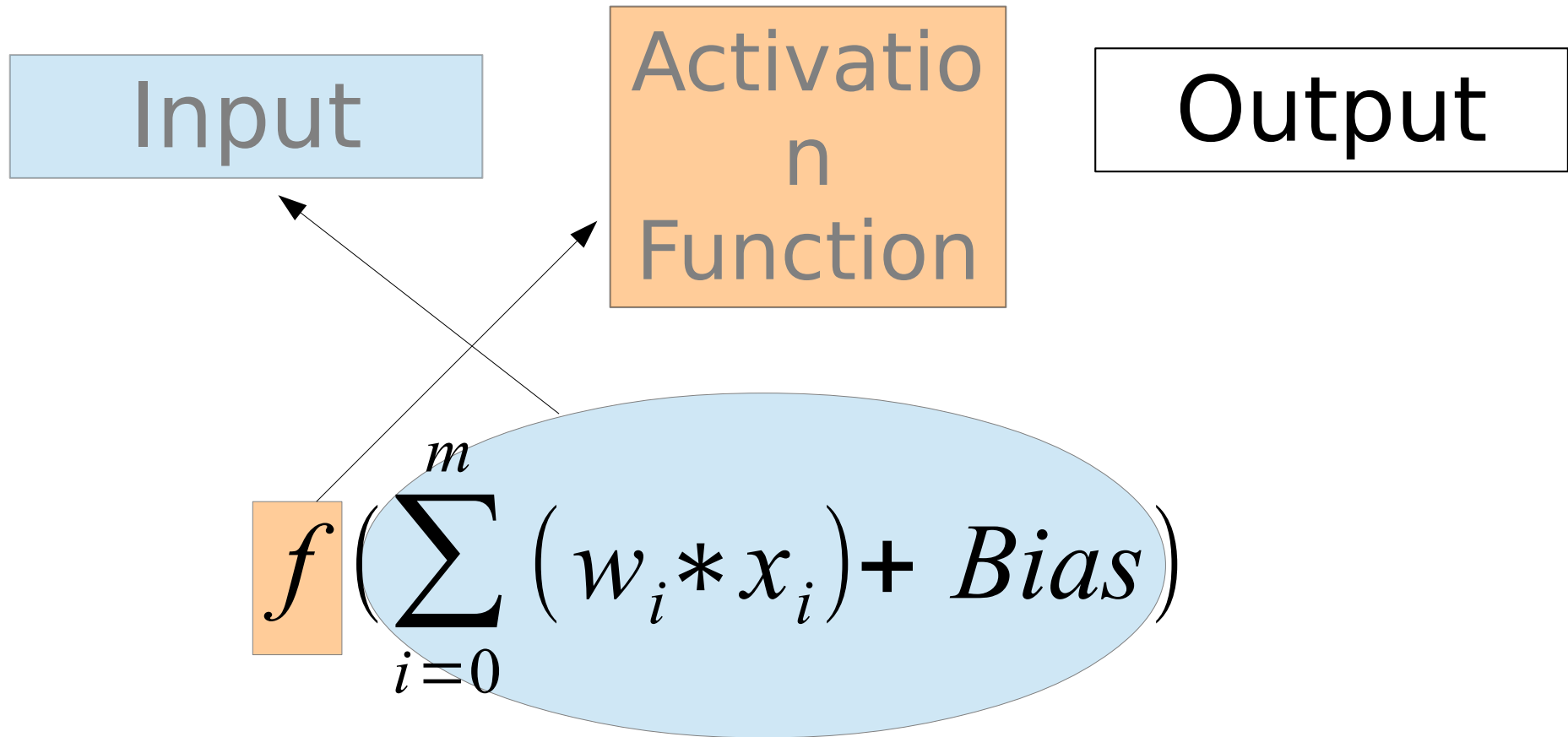
<sup>^</sup> Here,  $H$  is the Heaviside step function.

# Multilayer Perceptrons



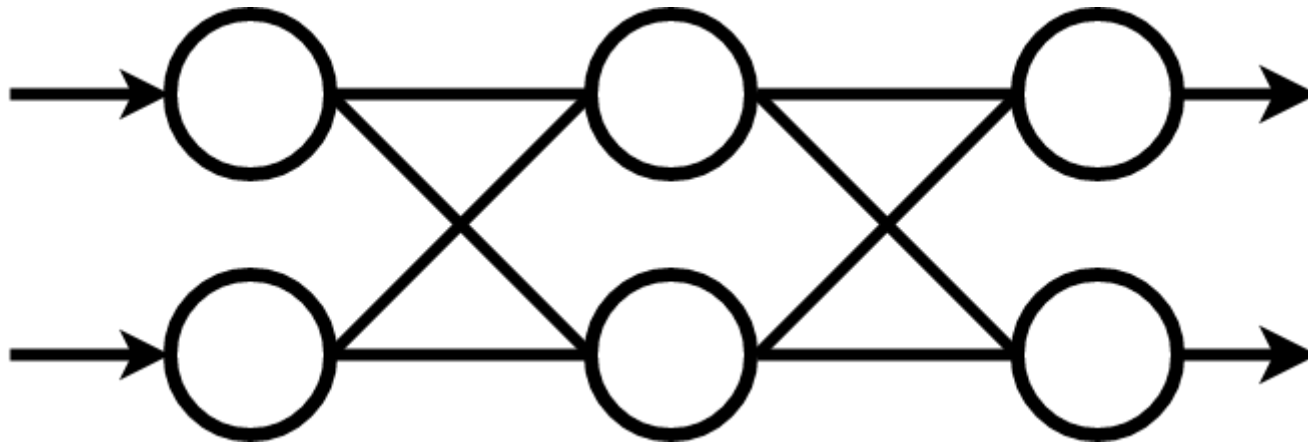


# Multilayer Perceptrons



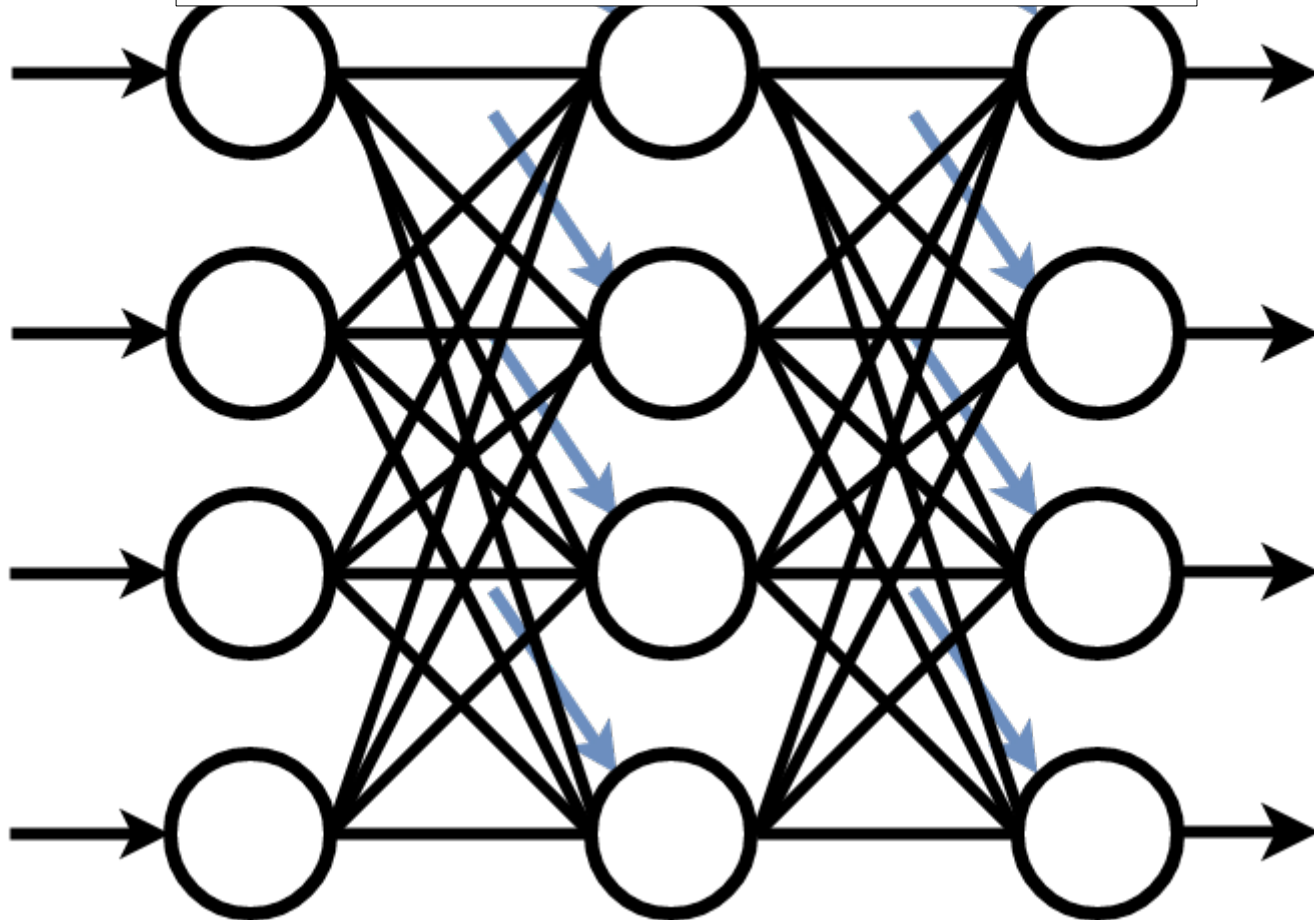
# Multilayer Perceptrons

Where's the  
Advantage?



# Multilayer Perceptrons

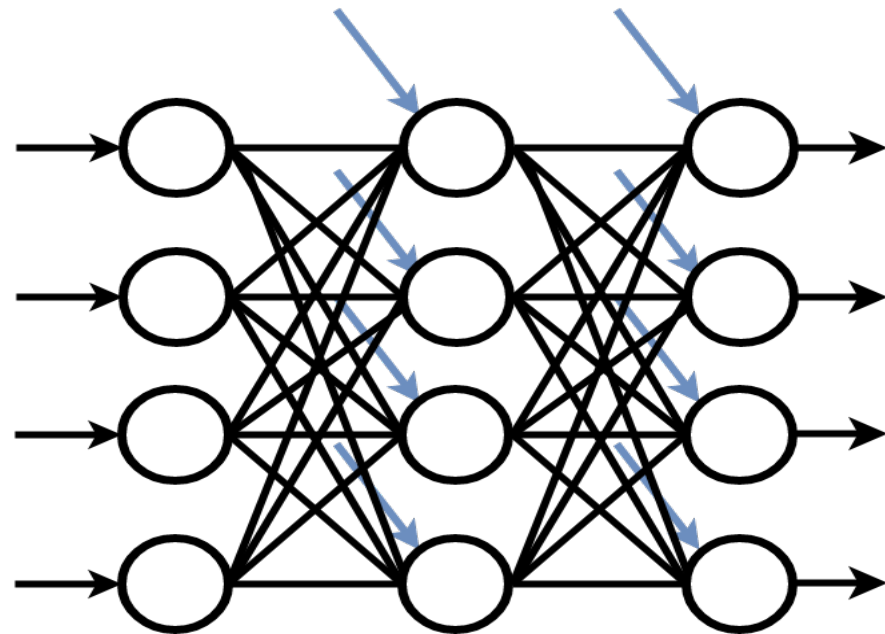
Where's the  
Advantage?



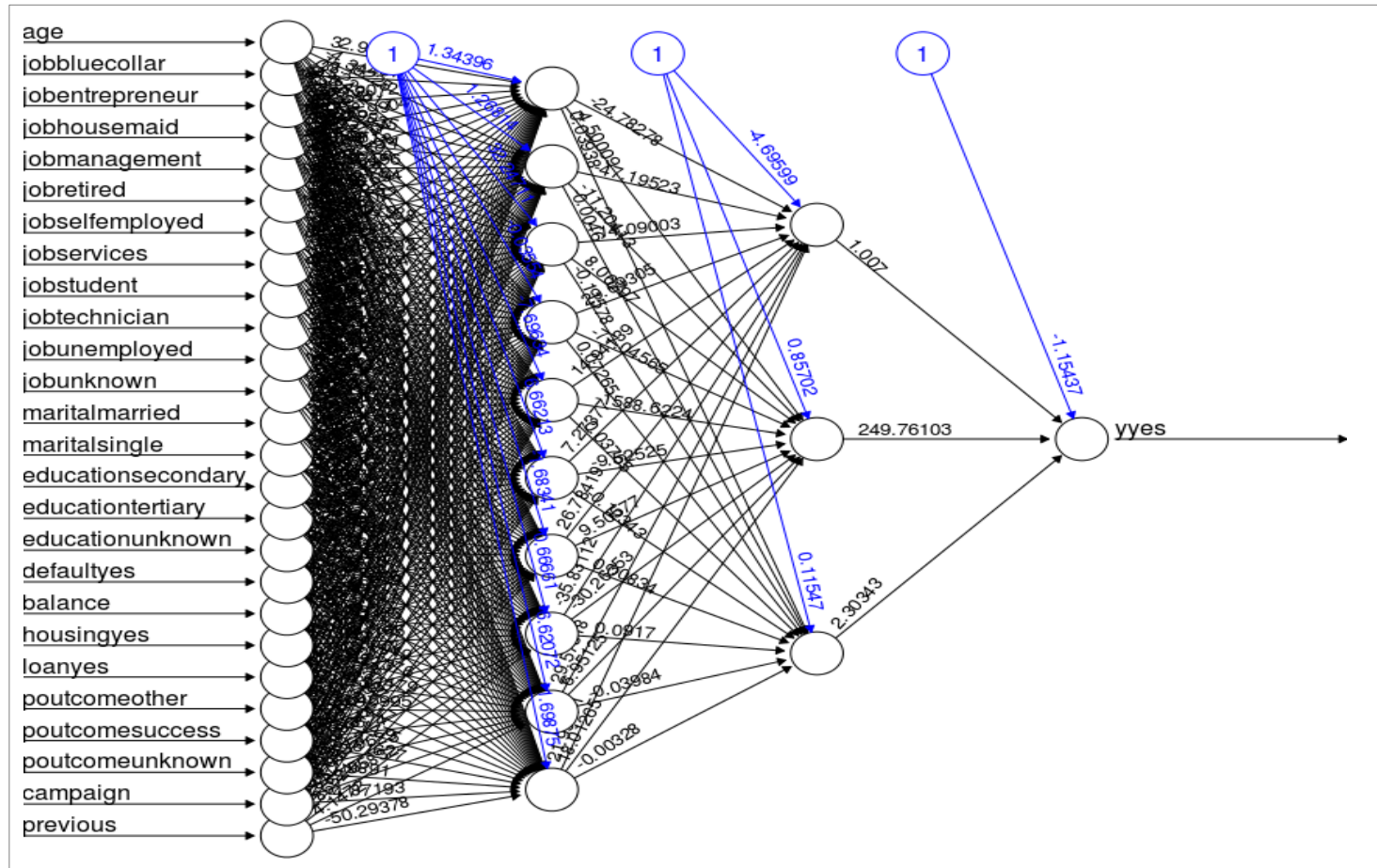
# Multilayer Perceptrons

Where's the  
Advantage?

$$f\left(\sum_{j=0}^m f\left(\sum_{i=0}^m w_{i,j} * x_{i,j}\right) + \textit{Bias}\right)$$



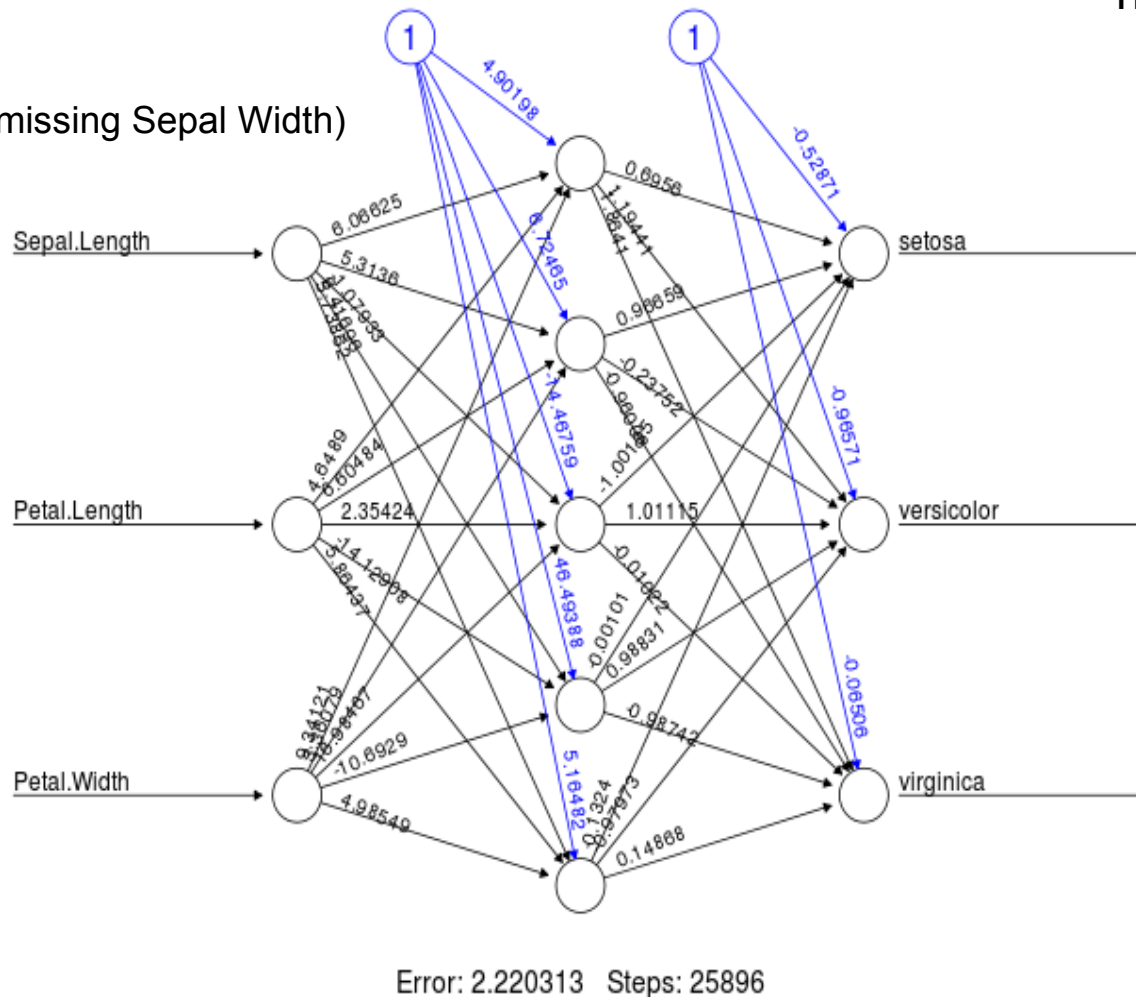
# Multilayer Perceptrons



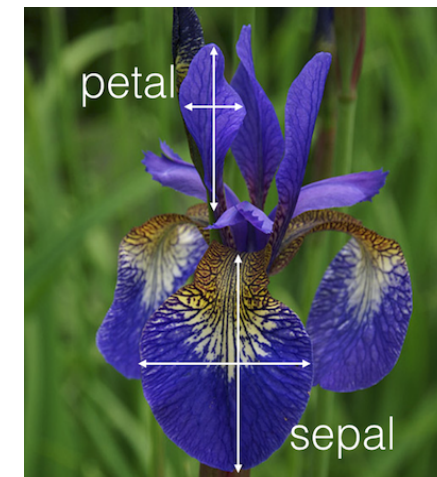
# Multilayer Perceptrons

<https://archive.ics.uci.edu/ml/datasets/Iris/>

(missing Sepal Width)



Famous Iris dataset which measured the Different Length of Iris petals and classifies them as one of 3 Species



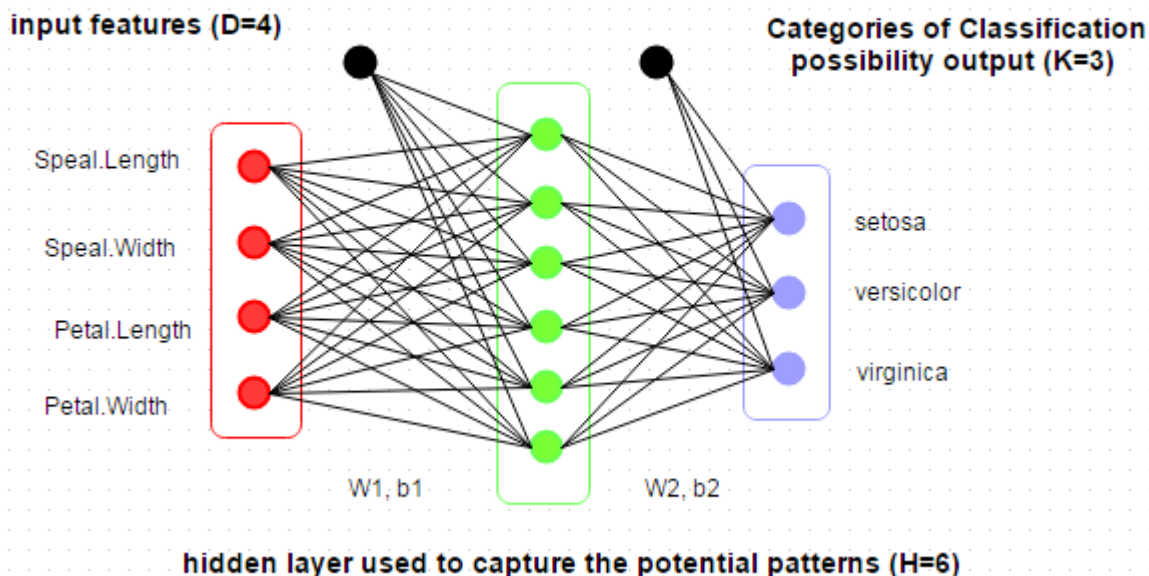
<http://www.learnbymarketing.com/tutorials/neural-networks-in-r-tutorial/>

<http://blog.kaggle.com/2015/04/22/scikit-learn-video-3-machine-learning-first-steps-with-the-iris-dataset/>

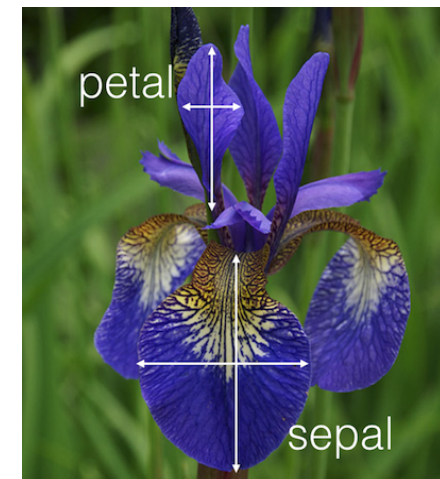
# Multilayer Perceptrons

<https://archive.ics.uci.edu/ml/datasets/Iris/>

## Classification Example for IRIS data by DNN



Famous Iris dataset which measured the Different Length of Iris petals and classifies them as one of 3 Species



<https://stats.stackexchange.com/questions/268202/backpropagation-algorithm-nn-with-rectified-linear-unit-relu-activation>

<http://blog.kaggle.com/2015/04/22/scikit-learn-video-3-machine-learning-first-steps-with-the-iris-dataset/>

# Machine Learning, An introduction 1: Conclusion

Machine Learning is the training of a model which learns based on data  $X$  to predict  $Y$

It imitates the system which generated data  $X$  and  $Y$ , hopefully well...



# See you next time for:

- Math!
- Optimization methods!
- More Math!!
- Eating habits of Neural Networks!
- Descending Gradients!

Written by Nicolas Symeou  
for the Physical-Digital Affordances Group of  
the University of Regensburg