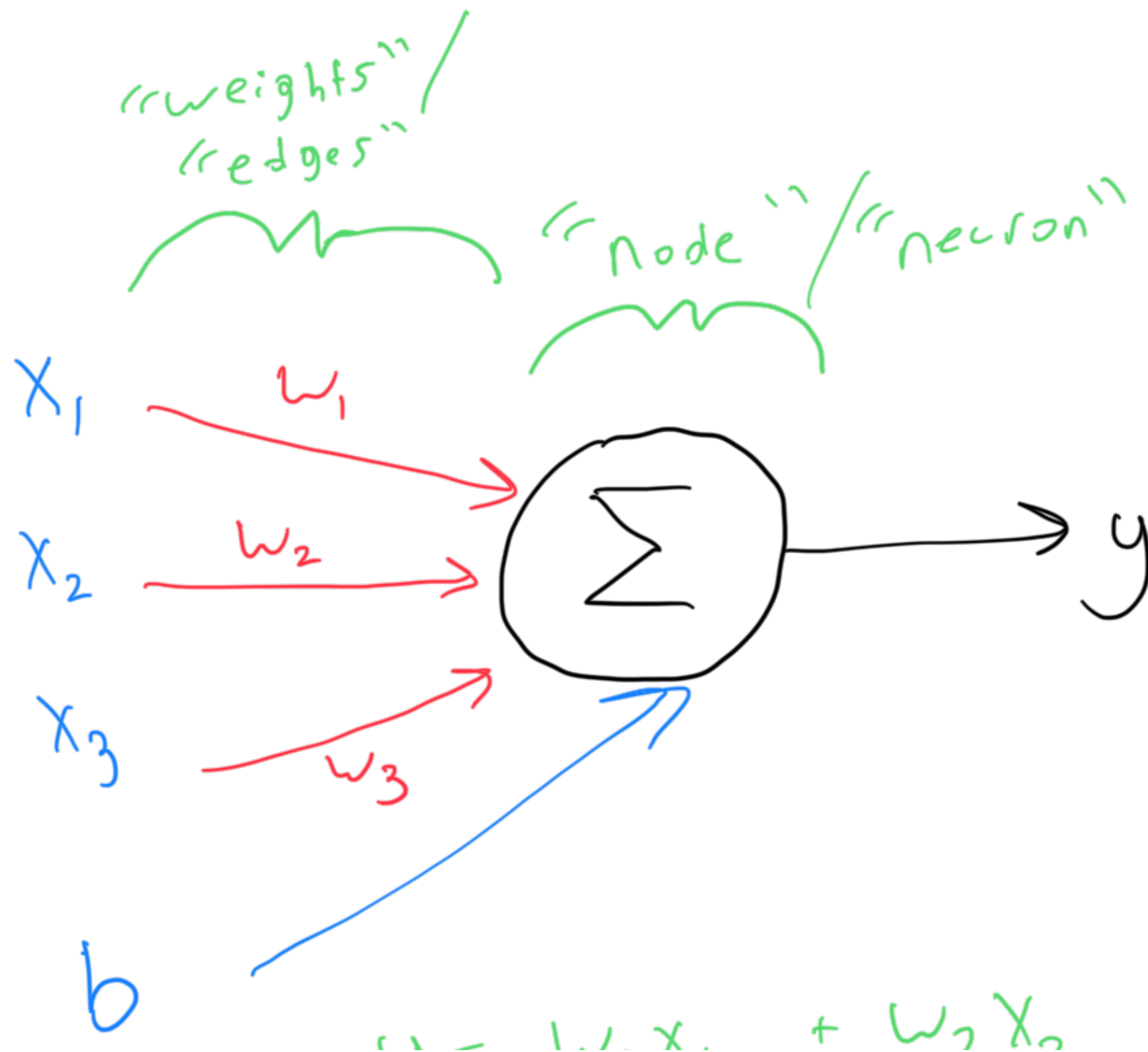


Day 20: Neural Network Fundamentals

Recall: "Linear and Logistic Regression are neural networks"

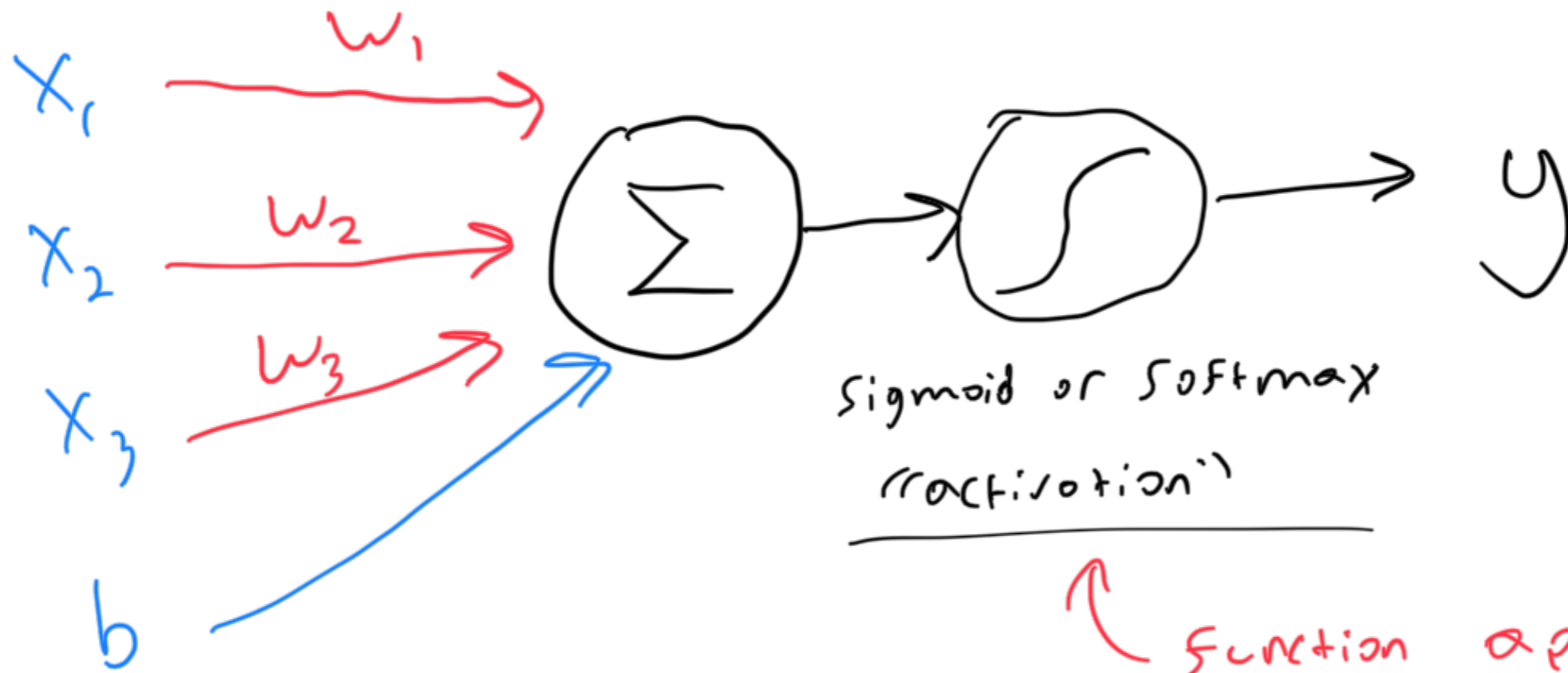
Think of linear regression as:



$$y = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

$$y = w_1 x_1$$

Think of logistic regression as:

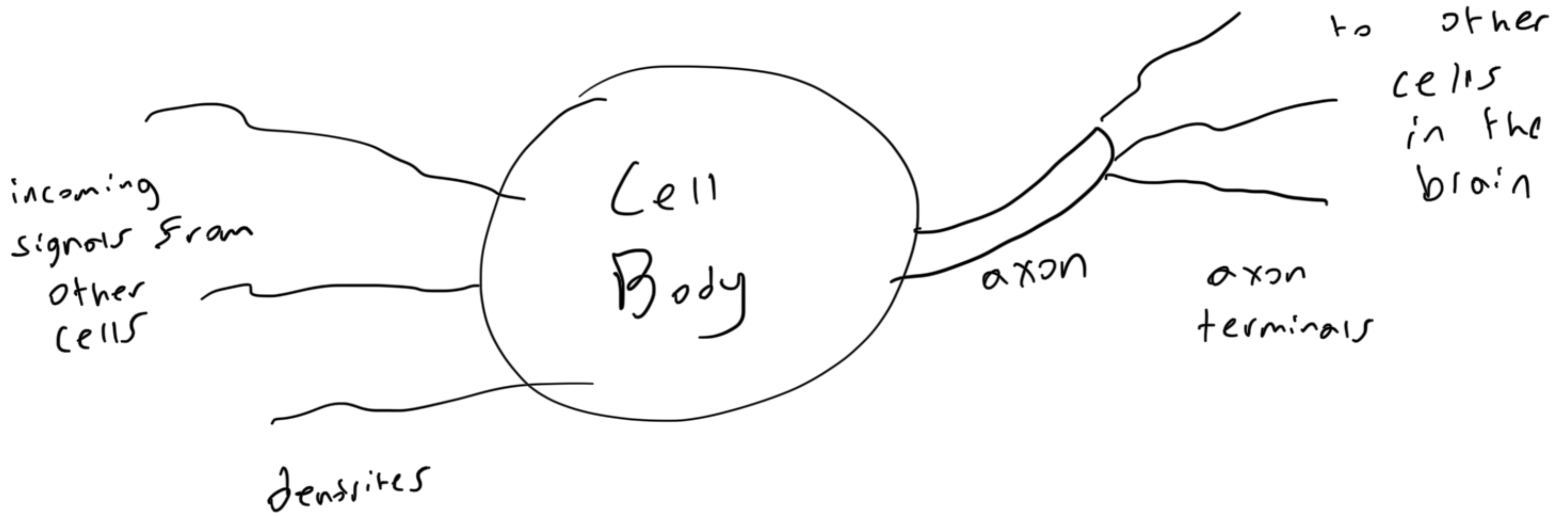


Sigmoid or Softmax
(activation)

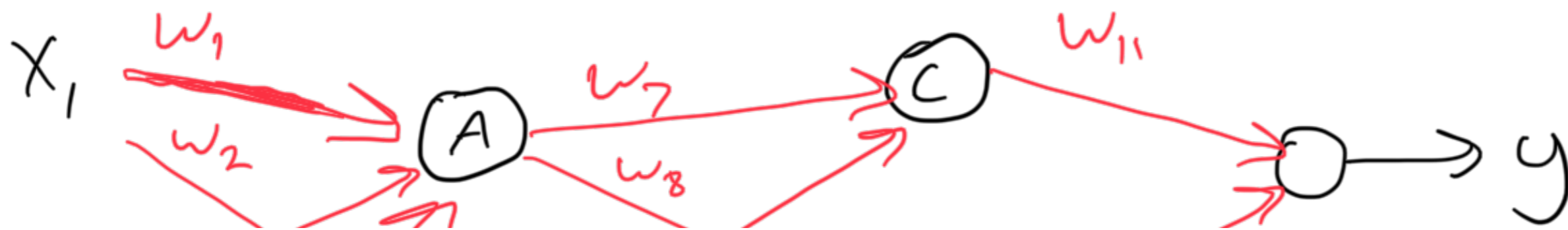
Function applied to the
output of a node

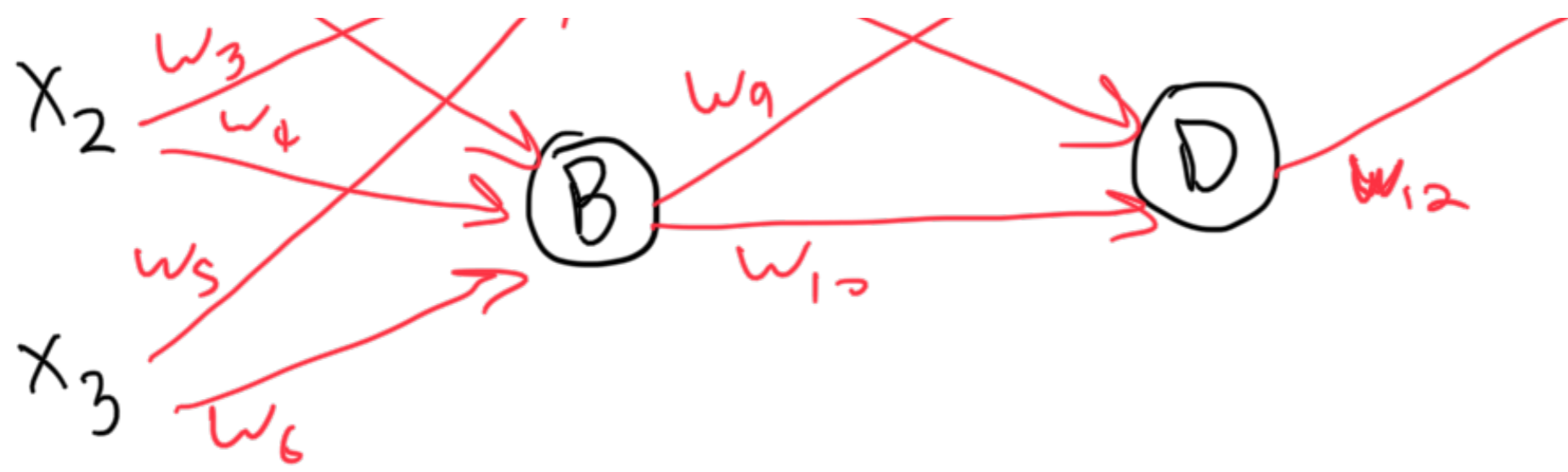
$$y = \text{sigmoid}(w_1 x_1 + w_2 x_2 + w_3 x_3 + b)$$

Above formulation is called a "Perceptron"
output signals



An artificial neural network is a set of connected layers of perceptrons, where the output of one perceptron is the input to the next perceptron





Input
Layer

Hidden
Layer

Hidden
Layer

Output
Layer

$$y = w_{11}C + w_{12}D$$

$$= w_{11}(w_7A + w_9B) + w_{12}(w_8A + w_{10}B)$$

$$w_{11} \left(w_7 (w_1 X_1 + w_3 X_2 + w_5 X_3) + w_9 (w_2 X_1 + w_4 X_2 + w_6 X_3) \right) +$$

$$+ w_{12} (w_8 (w_1 X_1 + w_3 X_2 + w_5 X_3) + w_{10} (w_2 X_1 + w_4 X_2 + w_6 X_3))$$

$$\begin{aligned}
 & \cup 12 \left(w_8 (w_1 x_1 + w_3 x_2 + w_5 x_3) \right. \\
 & \quad \left. + w_{10} (w_2 x_1 + w_4 x_2 + w_6 x_3) \right) \\
 & = w_{11} w_7 w_1 x_1 + w_{11} w_7 w_3 x_2 \\
 & \quad + w_{11} w_7 w_5 x_3 + w_{11} w_9 w_2 x_1 + \dots
 \end{aligned}$$

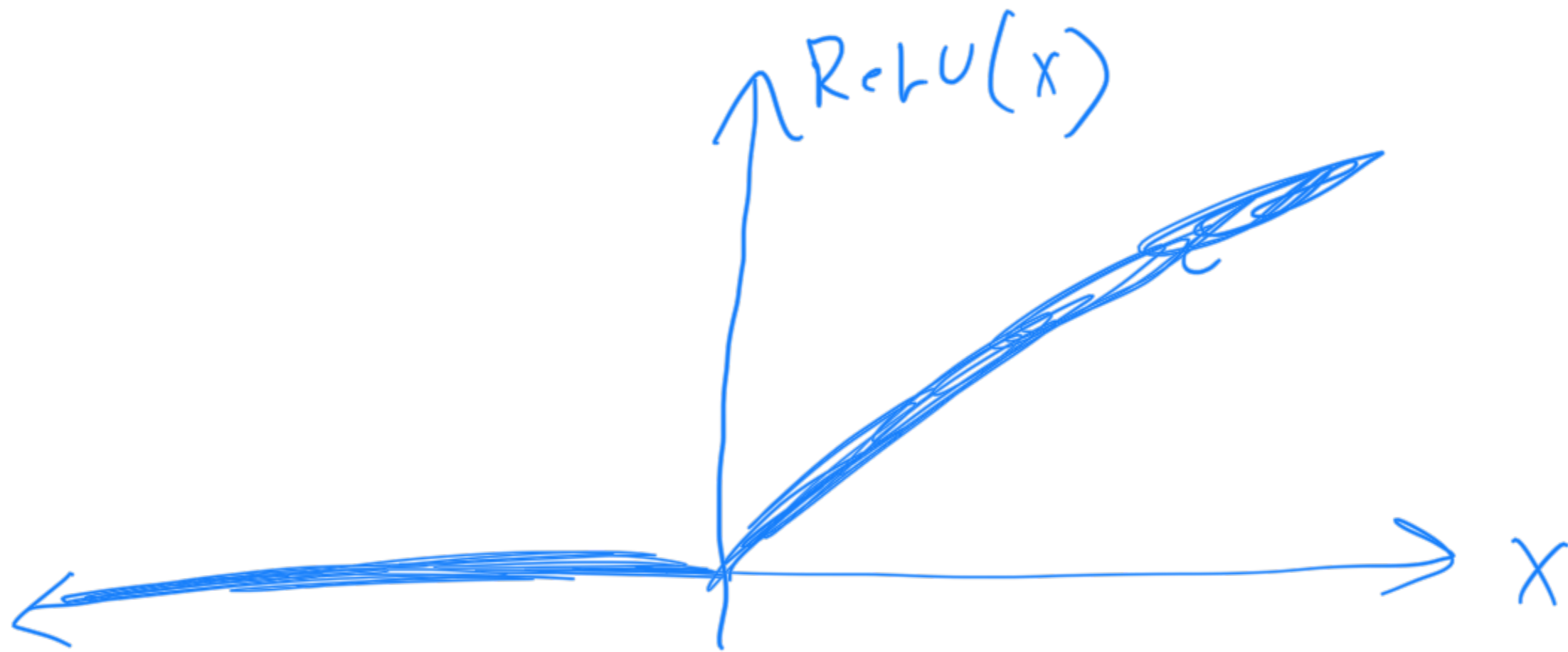
Isn't this a linear function?

Why not use linear/logistic regression?

Answer; add non-linearity to the network!!!

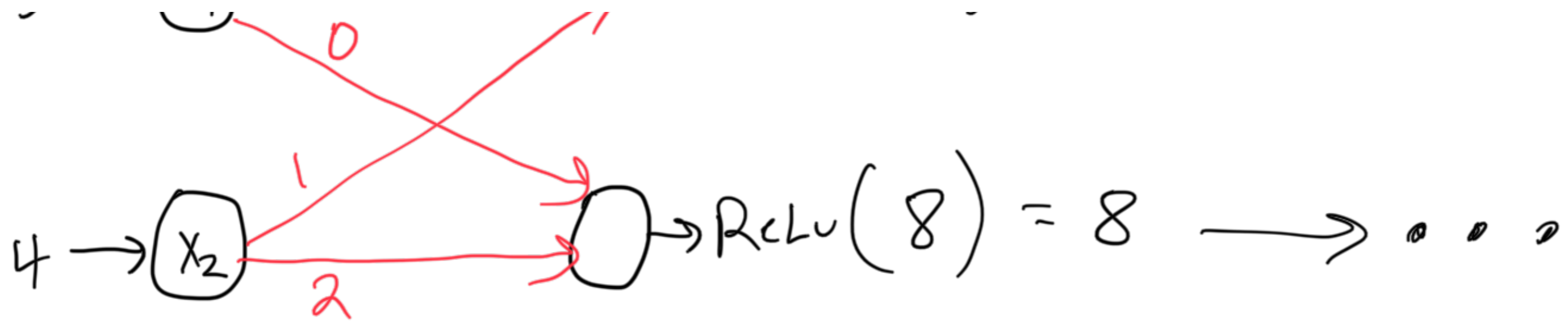
Take each node's output, and apply a non-linear activation function to it.

Popular activation function: ReLU
"Rectified Linear Unit"



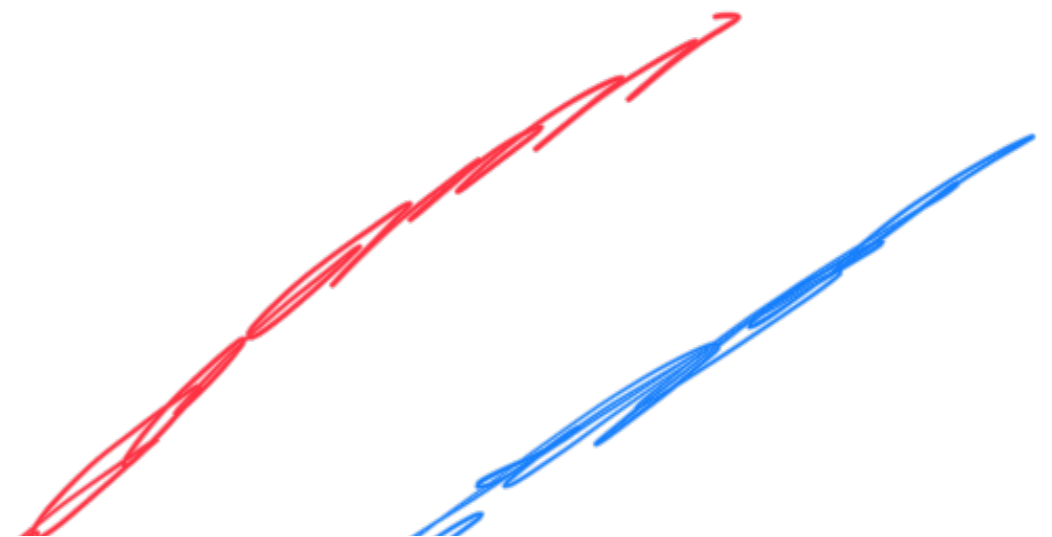
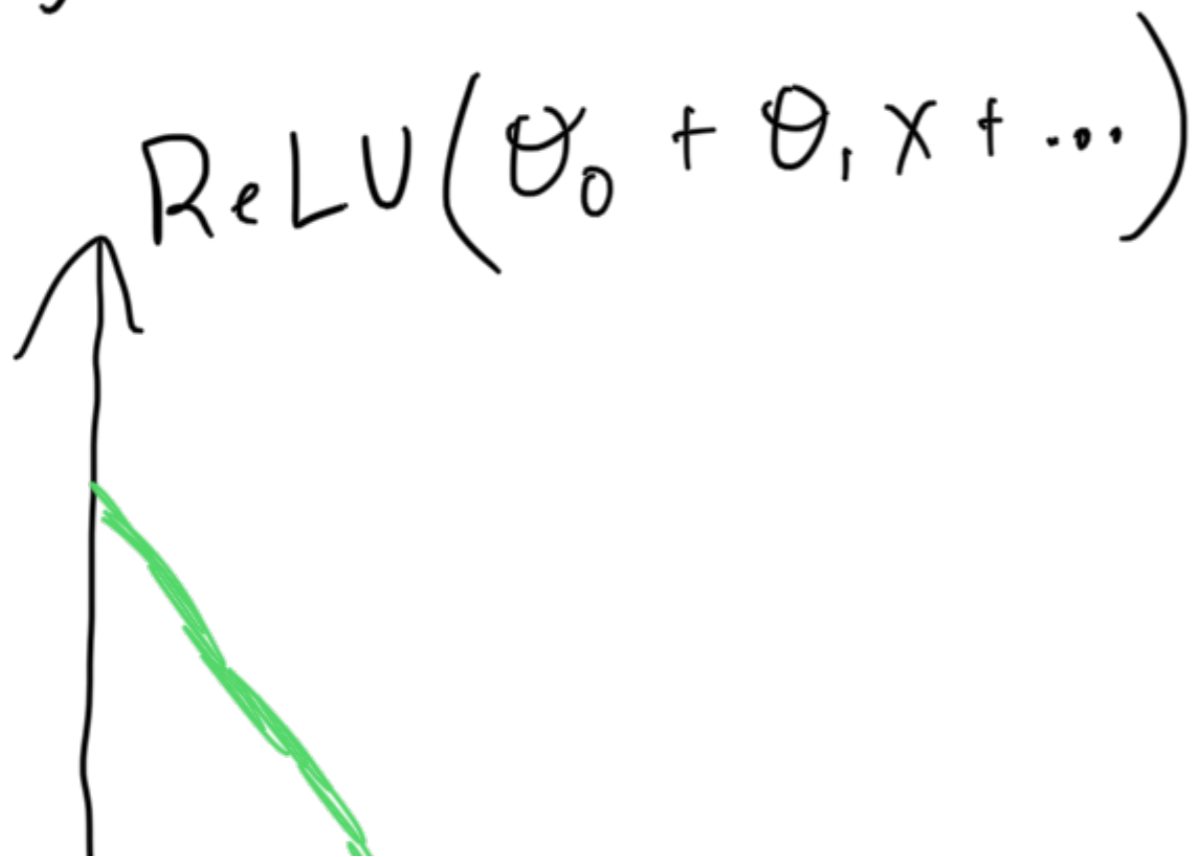
$$\begin{aligned} \text{ReLU}(x) &= \max(0, x) \\ &= \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases} \end{aligned}$$

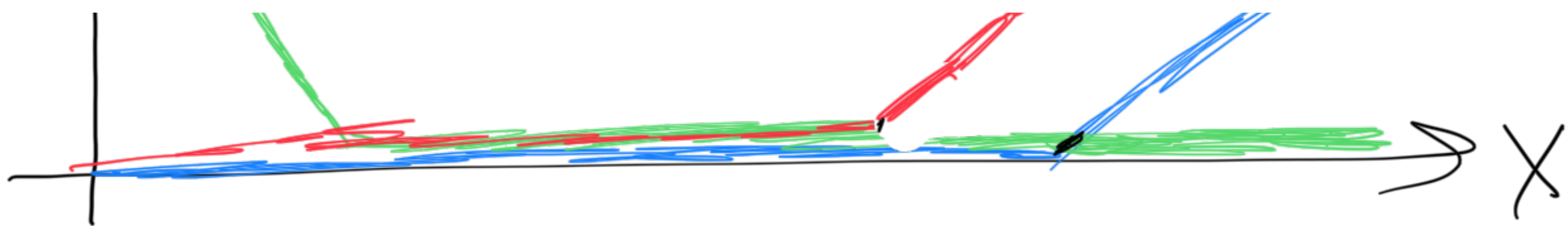
$5 \rightarrow (x_1) \xrightarrow{-1} 0 \rightarrow \text{ReLU}(-5) = 0 \rightarrow \dots$





Activation functions like ReLU allow networks
 to learn ANY function (given enough
 nodes and edges)

Why?





-  = Θ set 1
-  = Θ set 2
-  = Θ set 3

When adding these up, we get:

Sum of the 3 functions above



... and more functions,

If you can approximate anything!

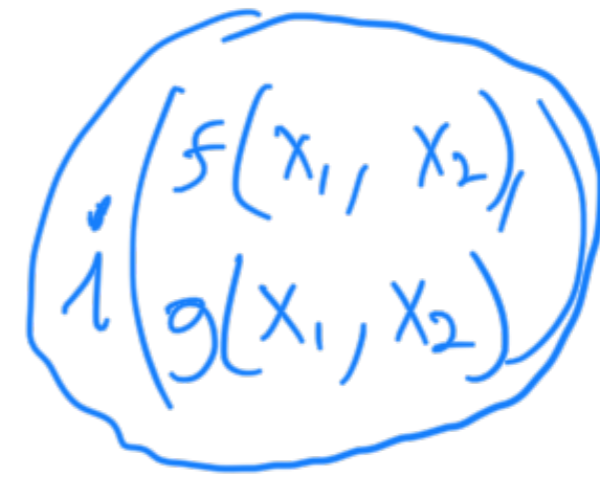
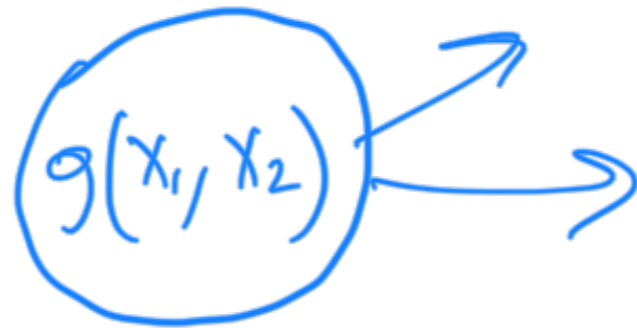
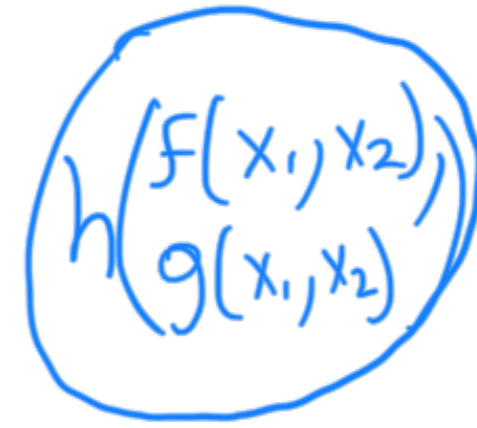
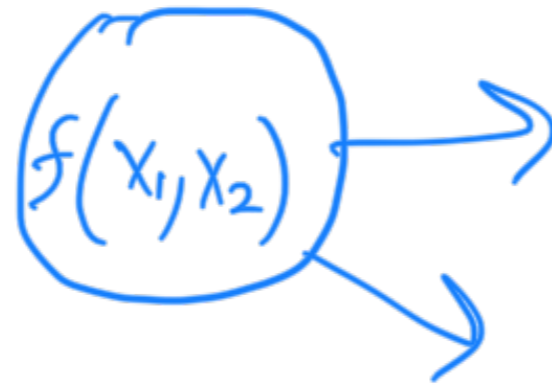
Universal Approximation Theorem

"Deep Learning" = Neural net (NN)
w/ many layers

Why use deep learning (DL) when we
can approximate anything with
1 hidden layer?

Deep networks need exponentially fewer θ

to achieve similar performance



Stacking layers allows using prior layers as modular functions

→ Can create more complex functions w/ same # of weights

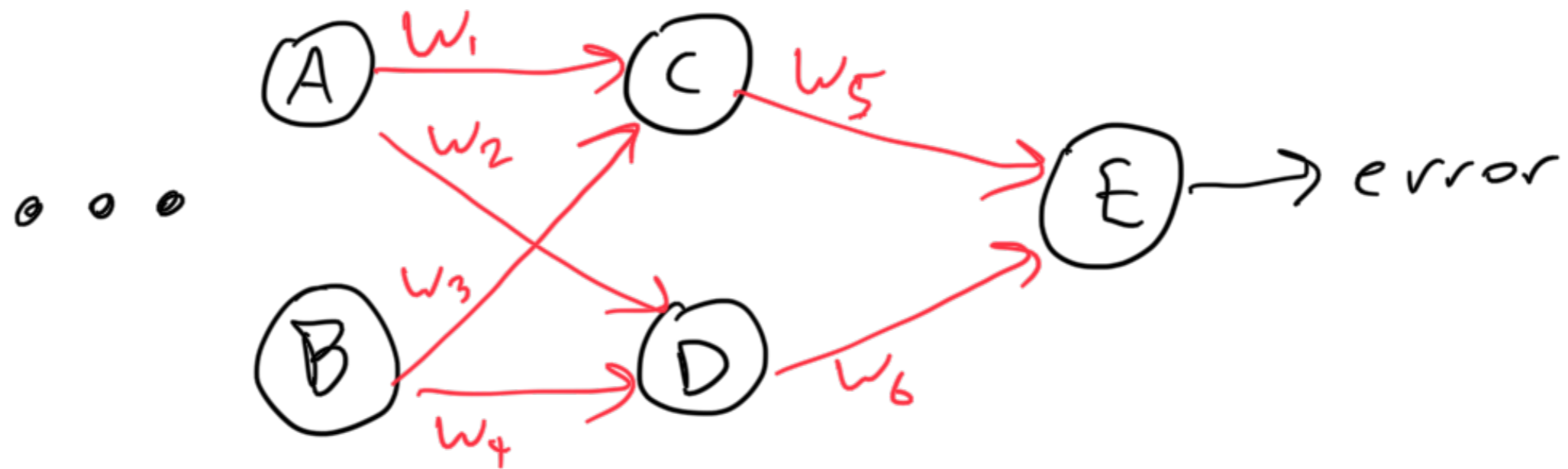
Progressive layers learn more high-level/abstract concepts

How are NN trained?

Gradient Descent, like before

$$W_i = W_i - \alpha \frac{\partial \text{Loss}}{\partial W_i}$$

But weights are layered;



$\frac{\partial \text{Loss}}{\partial W_1}$ depends on $W_5, C,$ and E

chain rule

Use the chain rule

If $y = f(u)$ where $u = g(x)$

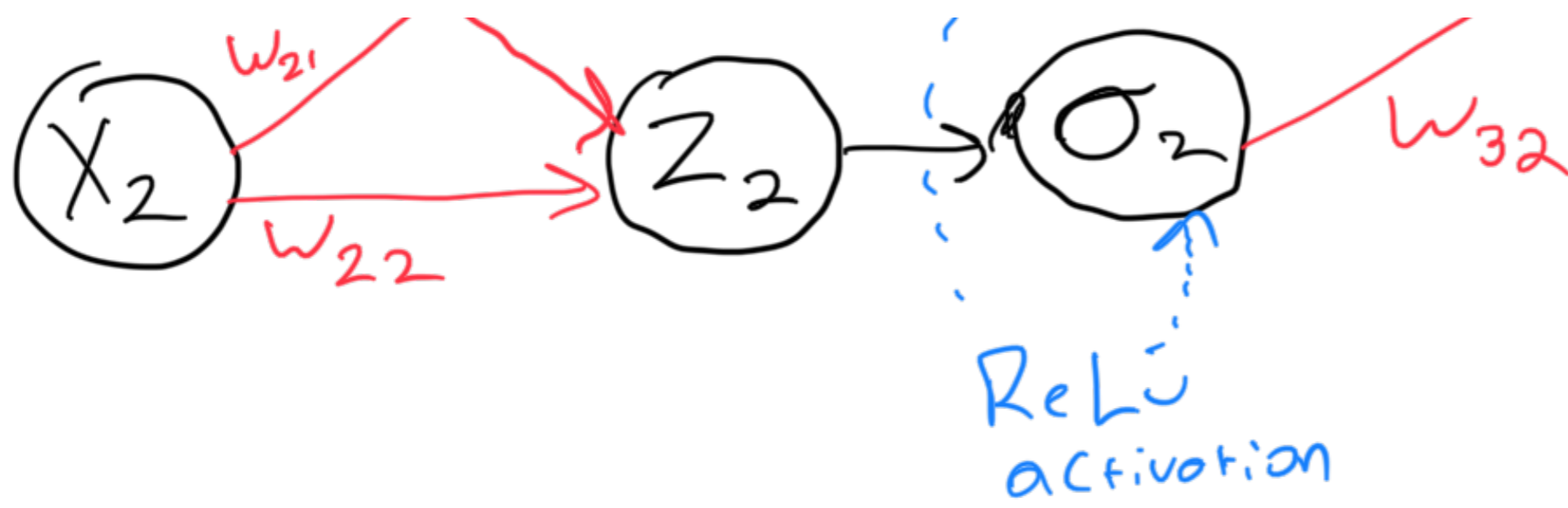
Then:

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

Backpropagation

Backprop Example



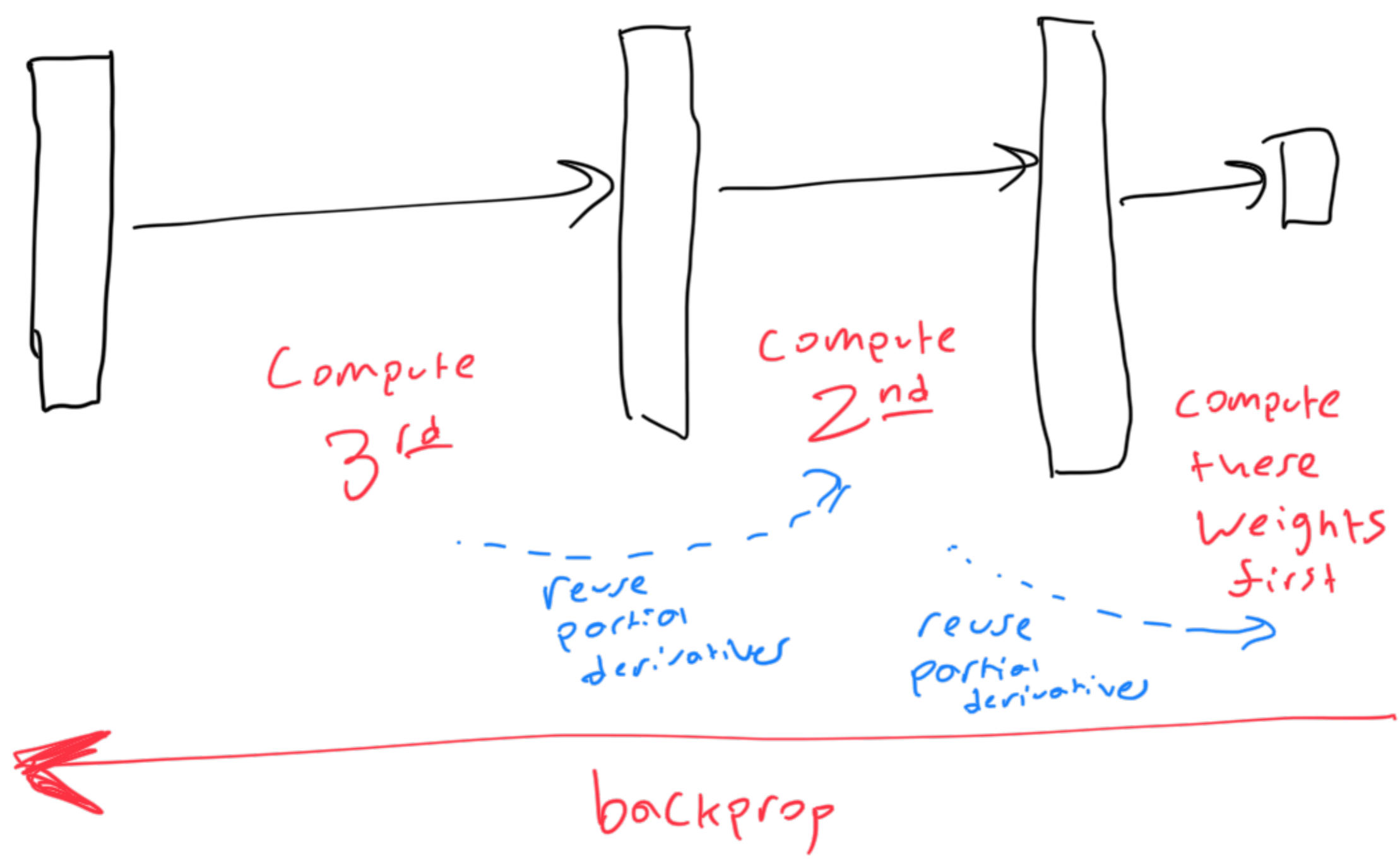


$$\frac{\partial \text{Loss}}{\partial w_{31}} = \underbrace{\frac{\partial \text{Loss}}{\partial \sigma_3} \cdot \frac{\partial \sigma_3}{\partial z_3}}_{\text{already computed in}} \cdot \frac{\partial z_3}{\partial w_{31}}$$

$$\frac{\partial \text{Loss}}{\partial w_{21}} = \underbrace{\frac{\partial \text{Loss}}{\partial \sigma_3} \cdot \frac{\partial \sigma_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial \sigma_1}}_{\text{already computed in}} \cdot \frac{\partial \sigma_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_{21}}$$

the previous step

We can use dynamic programming to cache results of prior computations



"Multi-Layer Perceptron (MLP)"

or
"Dense Neural Network"

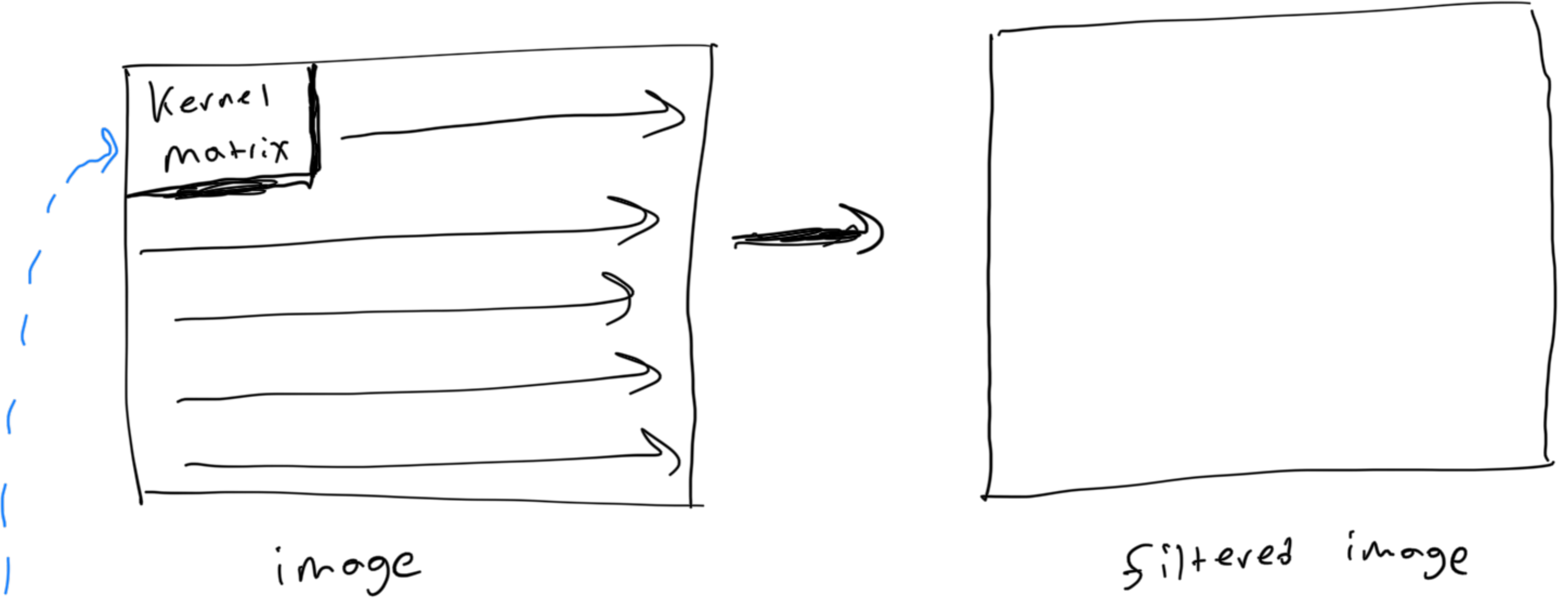
are the "basic" type of NN
we have learned about so far

Convolutional ^{Neural} Networks (CNNs)

Short Review

extract visual properties of

Recall: We can filter an image through convolutions:



Recall: Kernels are specialized to filter images differently.

-1	-1	-1
-1	8	-1

outline kernel

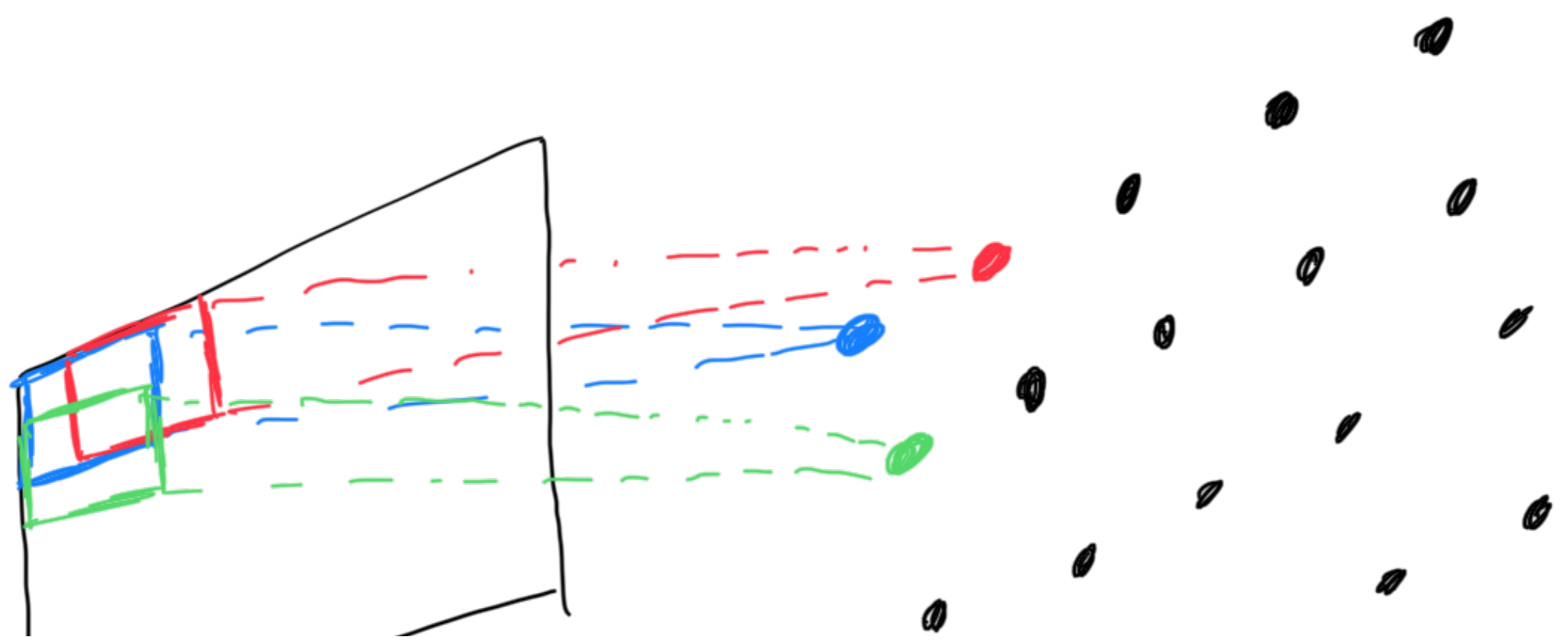
-1	-1	-1
----	----	----


1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

blur kernel

0	-1	0
-1	5	-1
0	-1	0

sharper kernel






each pixel in new image
is dot product between
the kernel and
corresponding image portion

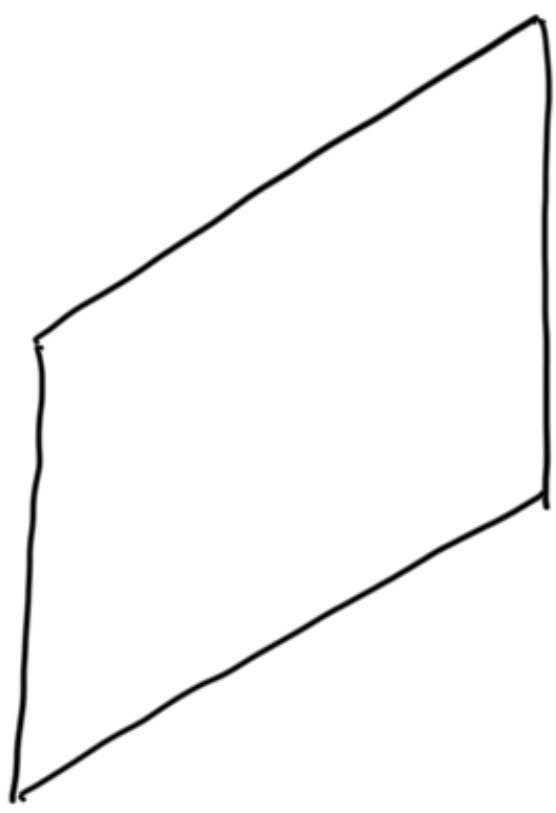
Key ideas of CNNs:




- ① use several kernels
- ② learn the best kernels to use
(i.e., best features to extract)

Convolutional Layers



11



-  Kernel 1 →
-  Kernel 2 →
- ...
-  Kernel K →

