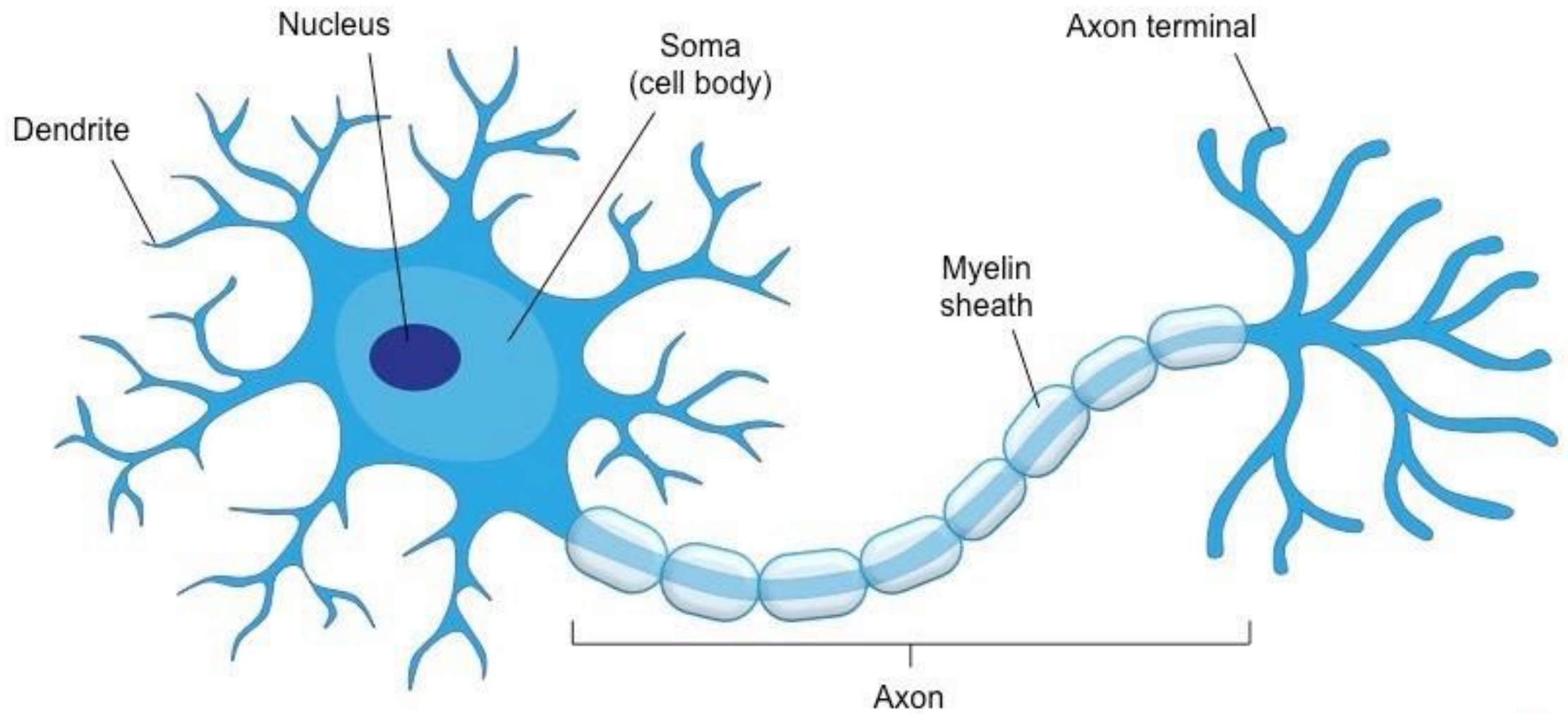


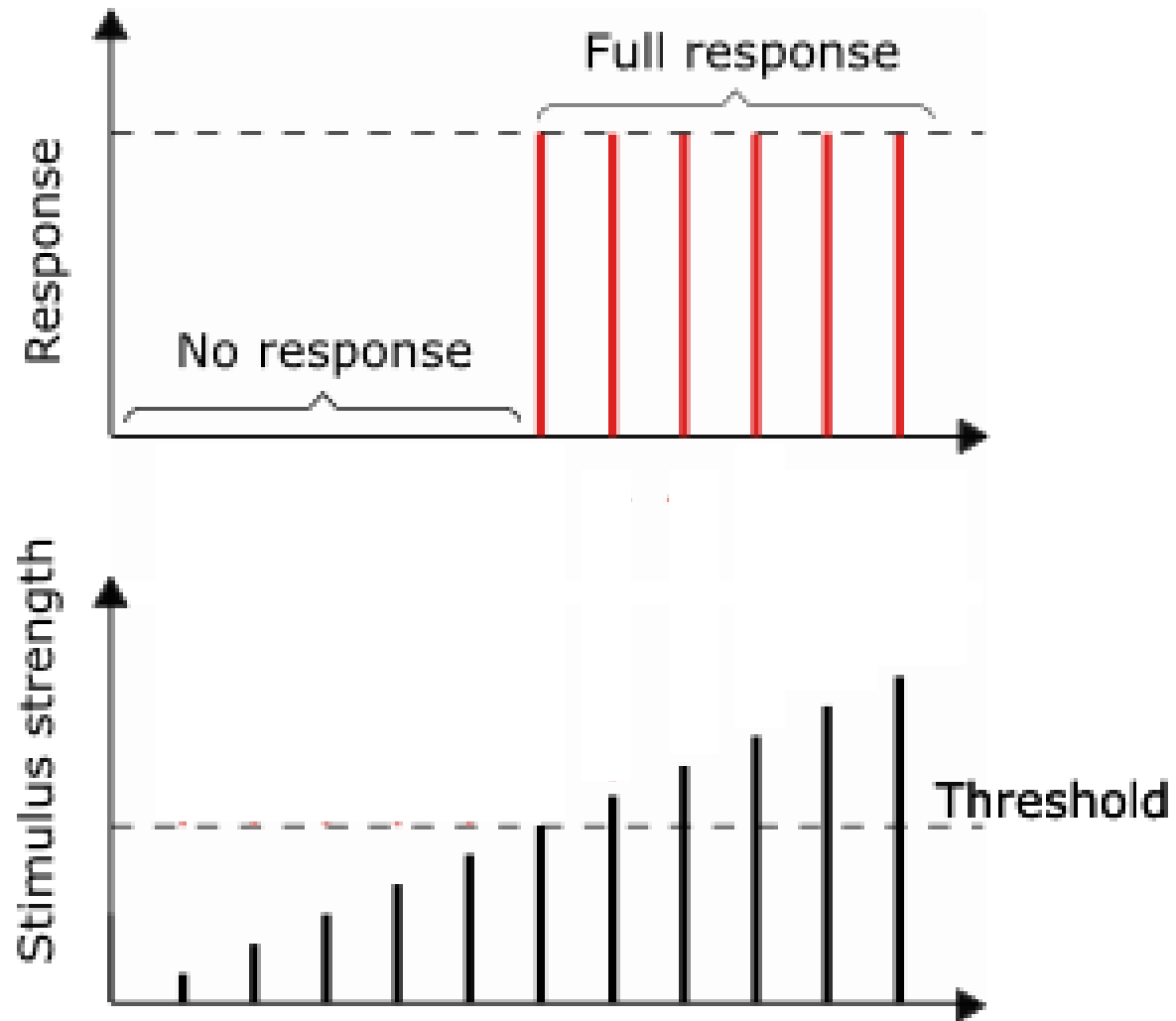


PERCEPTRON

THE DADDY OF NEURAL NETWORKS

FELIPE BUCHBINDER





ALL-OR- NOTHING LAW OF NEURONAL ACTIVATION

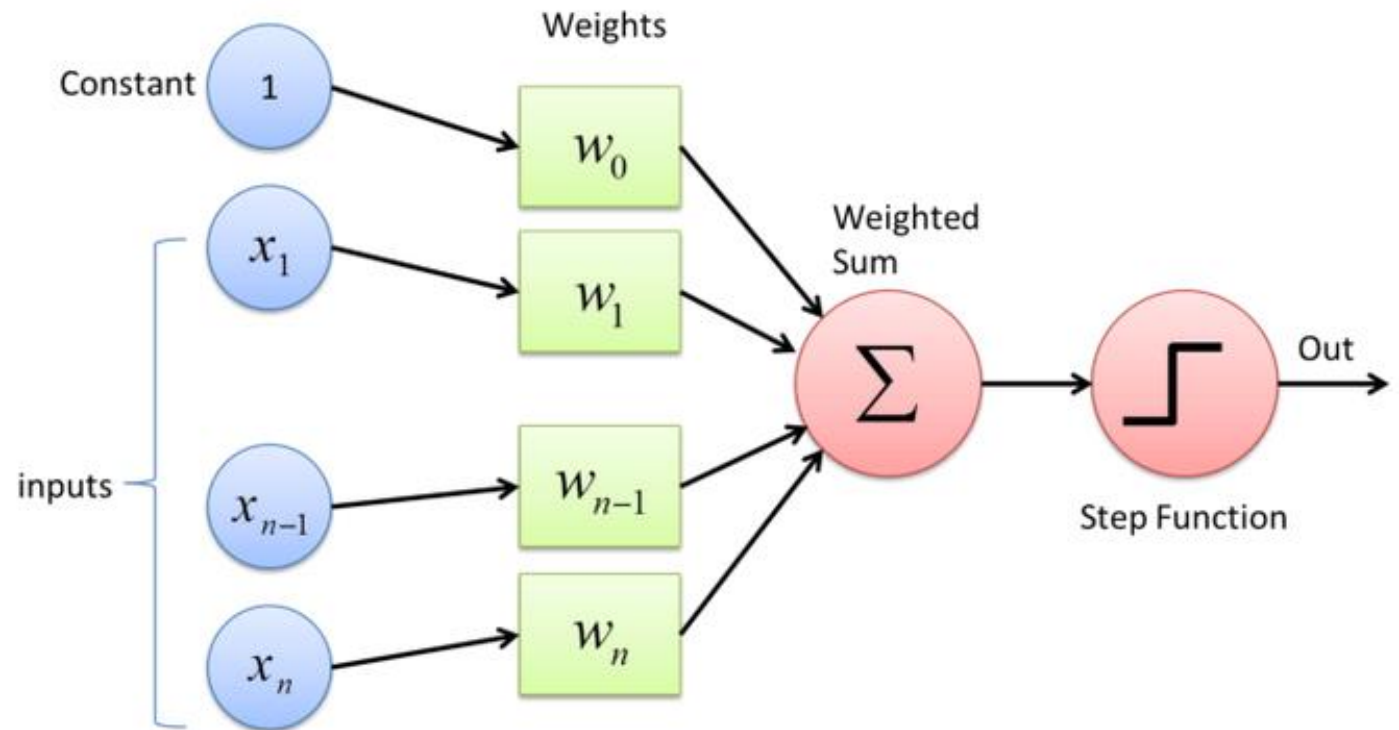
(WHY IS IT, THEN, THAT WE SOMETIMES FEEL
MORE PAIN OR LESS PAIN?)



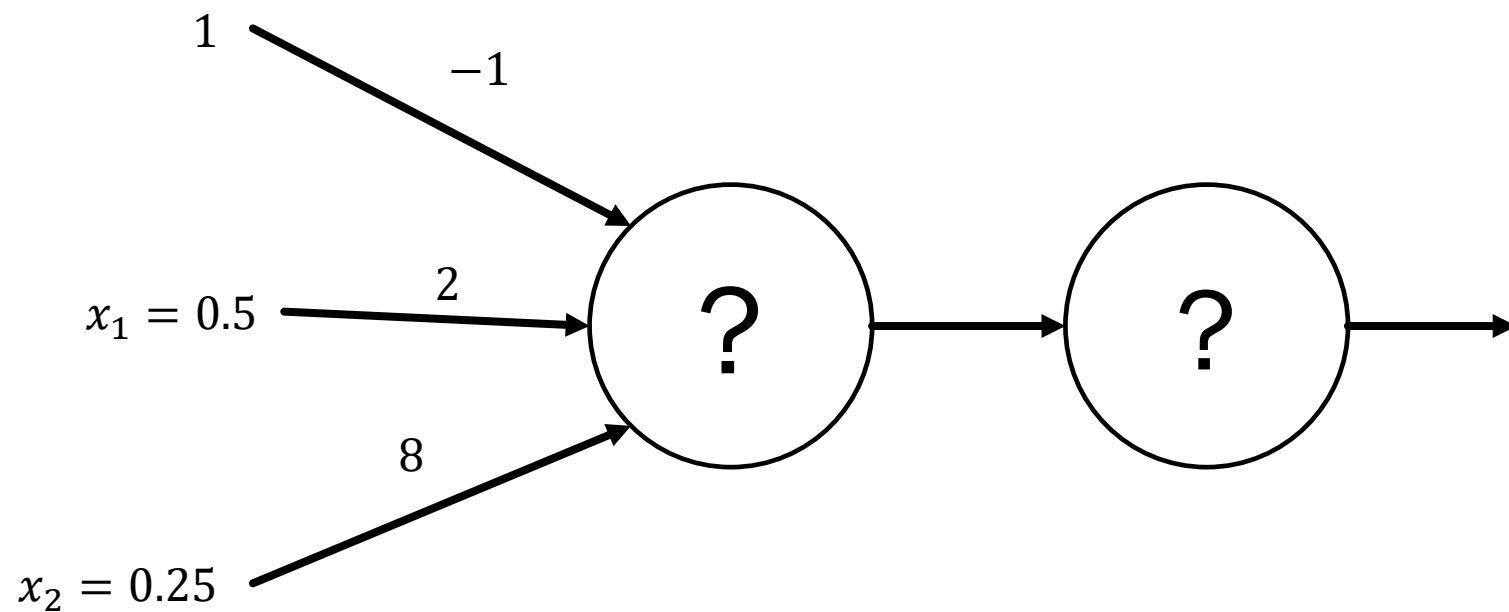
ALL-OR- NOTHING LAW OF NEURONAL ACTIVATION

(WHY IS IT, THEN, THAT WE SOMETIMES FEEL
MORE PAIN OR LESS PAIN?)

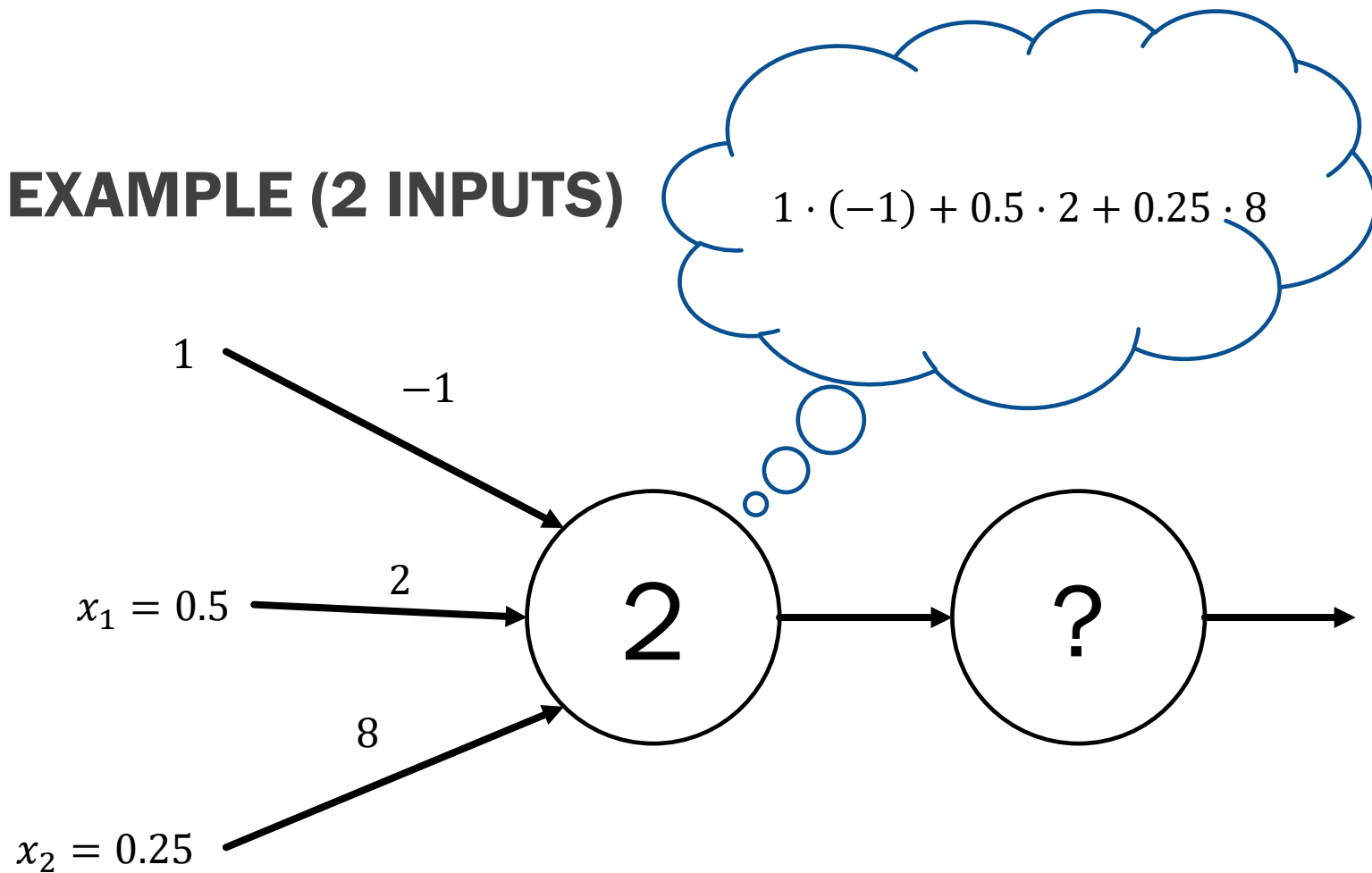
THE PERCEPTRON



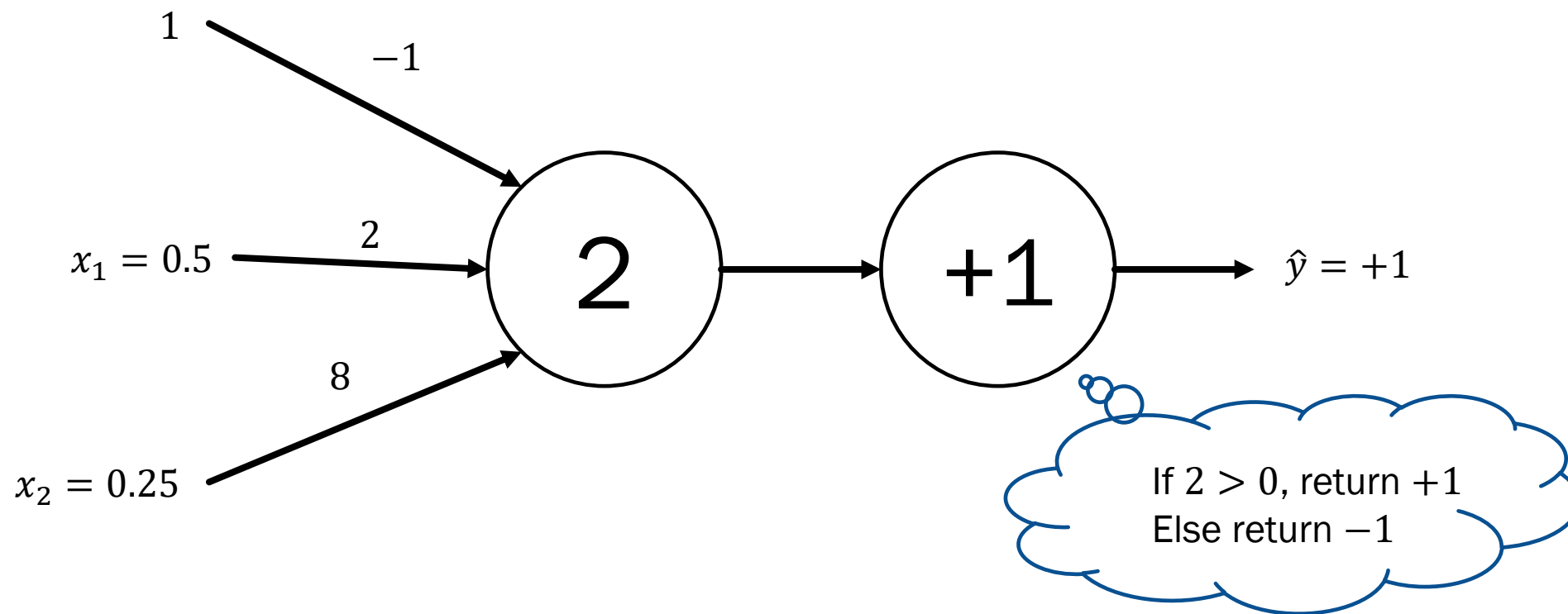
SIMPLE EXAMPLE (2 INPUTS)



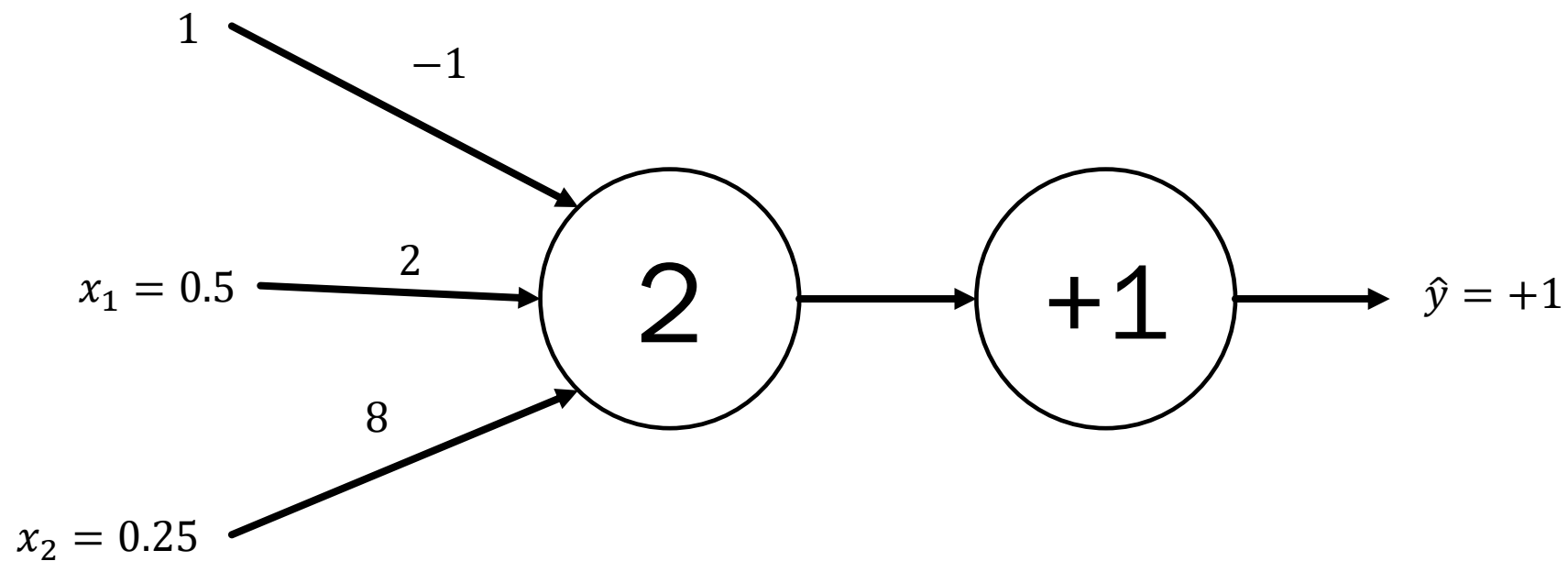
SIMPLE EXAMPLE (2 INPUTS)



SIMPLE EXAMPLE (2 INPUTS)



SIMPLE EXAMPLE (2 INPUTS)

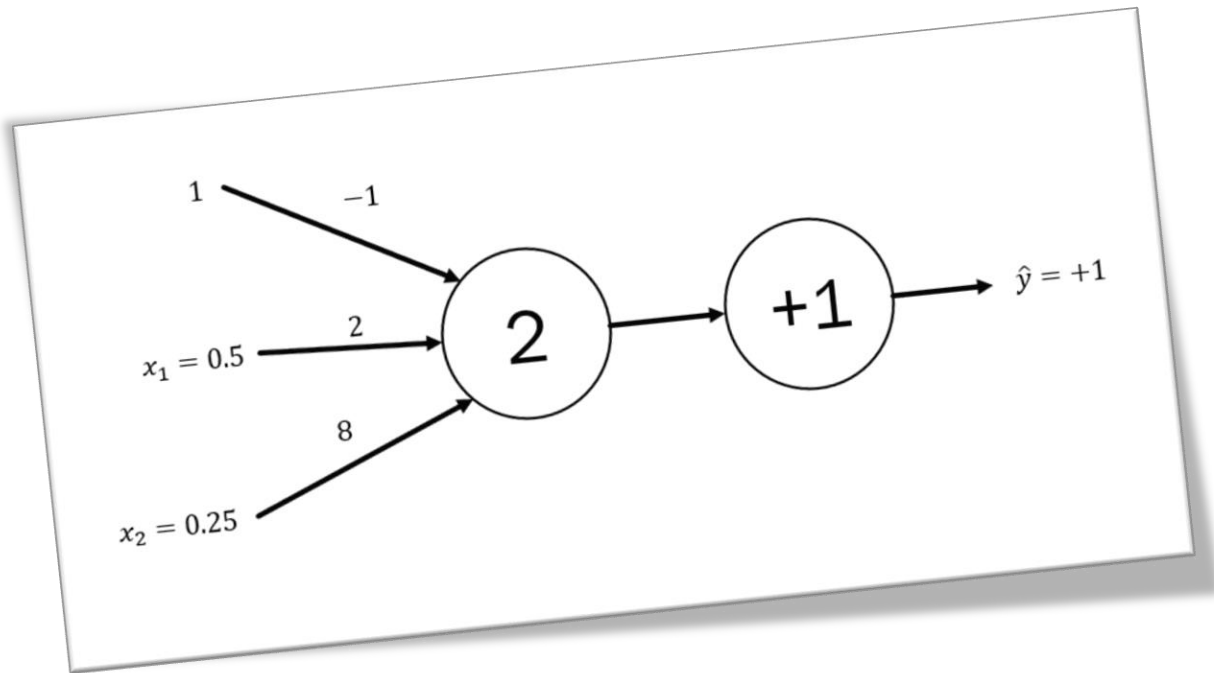


OUR PERCEPTRON'S PREDICTION CAN BE WRITTEN IN A SINGLE LINE

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$

OUR PERCEPTRON'S PREDICTION CAN BE WRITTEN IN A SINGLE LINE

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$



$$\begin{aligned}\hat{y} &= \text{sign}((-1) \cdot 1 + 2 \cdot 0.5 + 8 \cdot 0.25) \\ &= \text{sign}(2) \\ &= +1\end{aligned}$$

PERCEPTRON'S CAN HAVE DIFFERENT **ACTIVATION FUNCTIONS**

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$

Heaviside (step) function

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{x})$$

Sigmoid function

$$\hat{y} = \tanh(\mathbf{w}^T \mathbf{x})$$

Hyperbolic tangent function

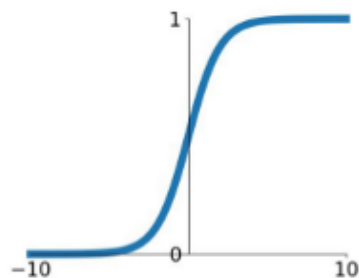
$$\hat{y} = \text{ReLU}(\mathbf{w}^T \mathbf{x})$$

Rectified Linear Unit

Activation Functions

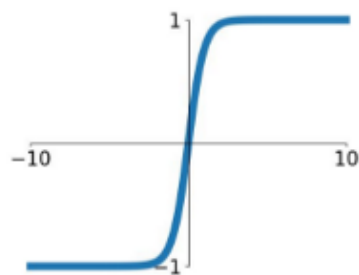
Sigmoid

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



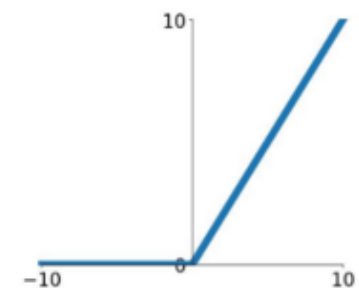
tanh

$$\tanh(x)$$



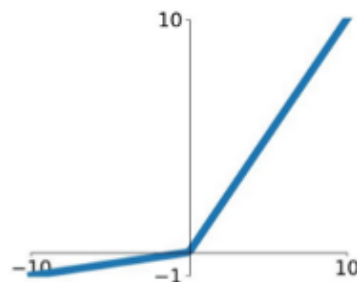
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

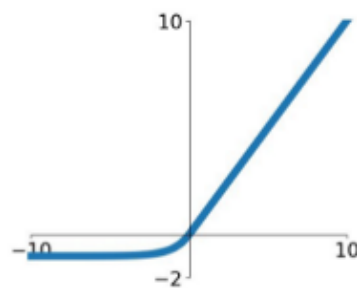


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

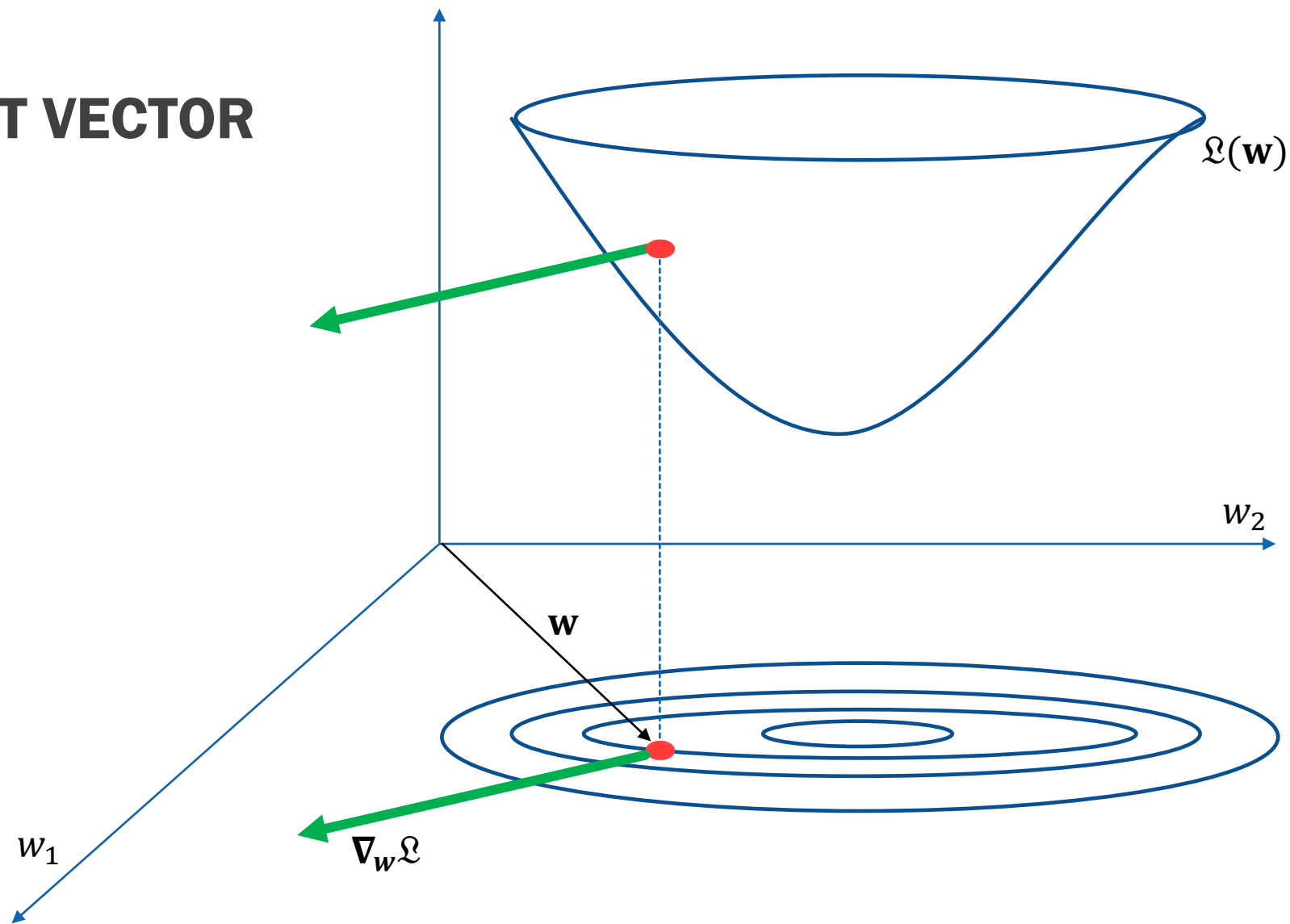




HOW CAN THE PERCEPTRON LEARN WHICH WEIGHTS TO USE?

CHOOSE WEIGHTS TO MINIMIZE SOME LOSS FUNCTION

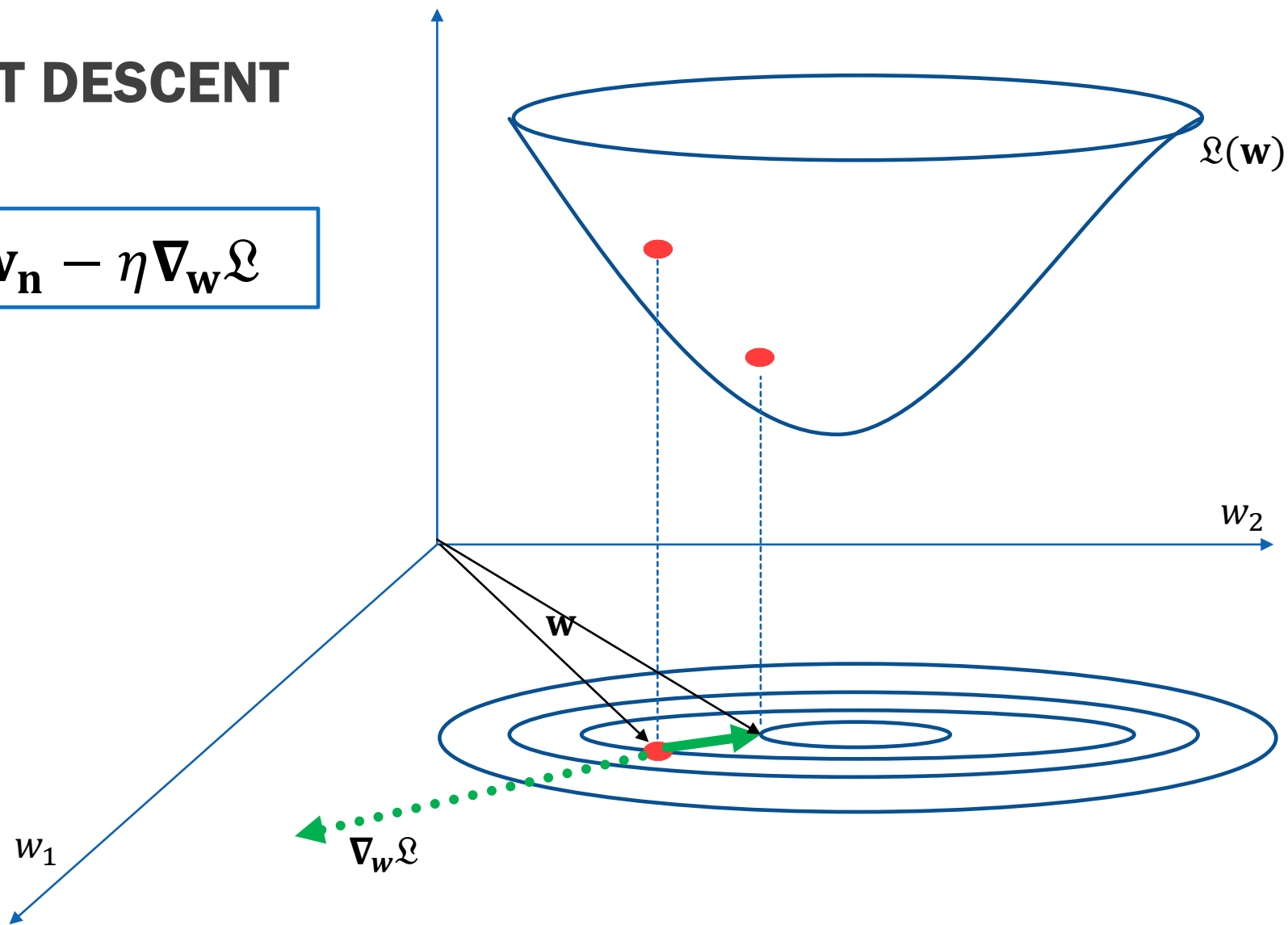
GRADIENT VECTOR



GRADIENT DESCENT

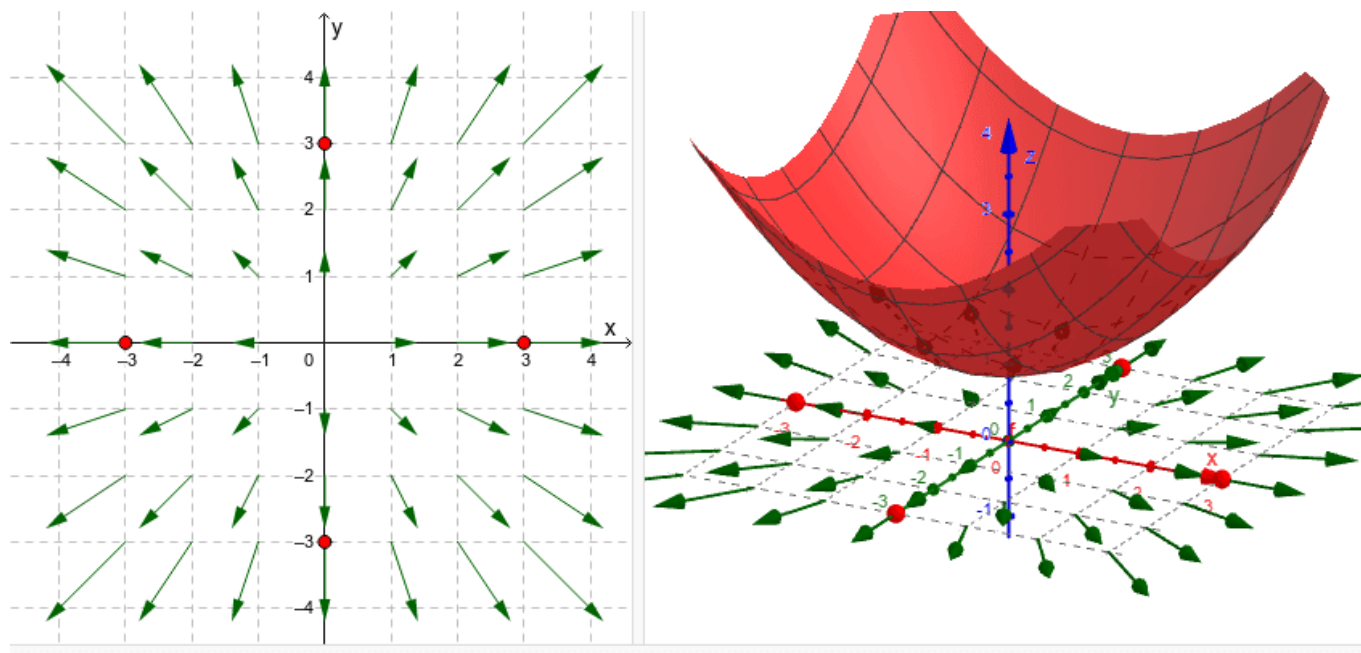
(WITH 1 POINT)

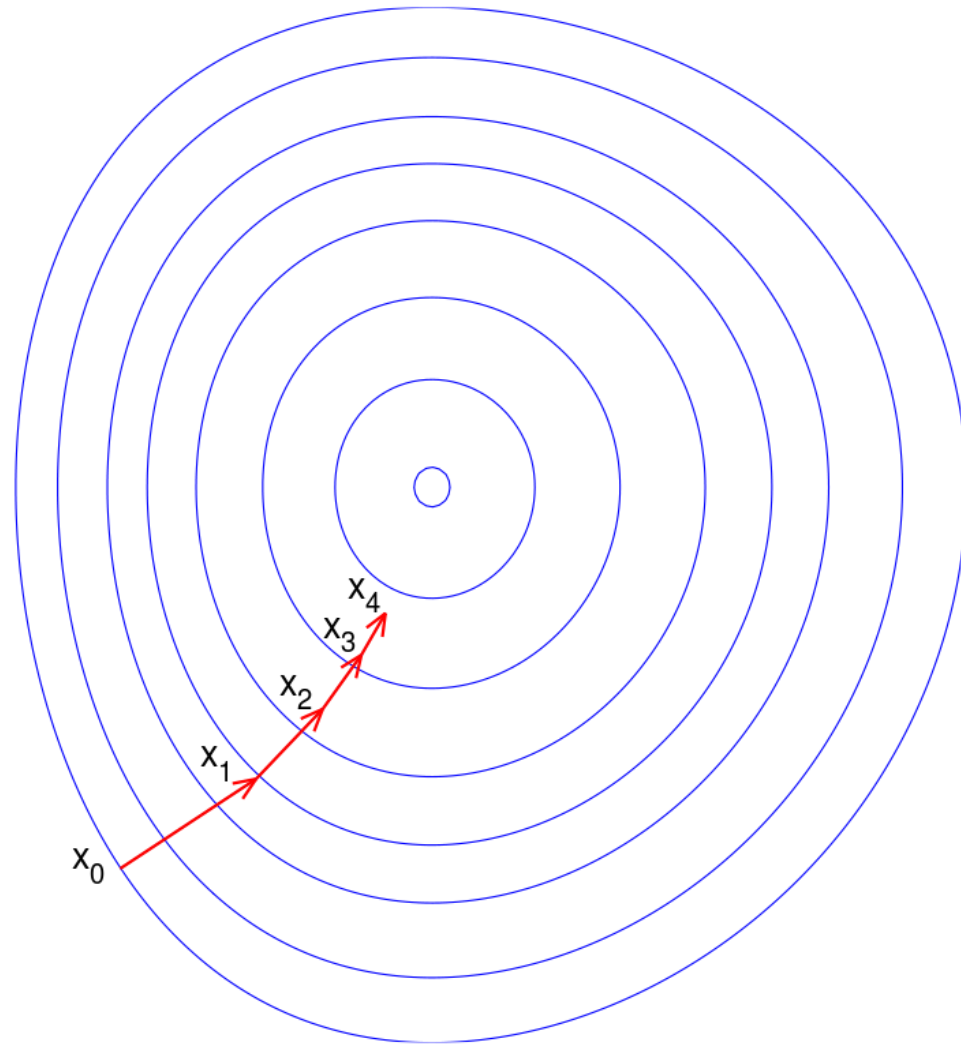
$$\mathbf{w}_{n+1} \leftarrow \mathbf{w}_n - \eta \nabla_{\mathbf{w}} \mathcal{L}$$



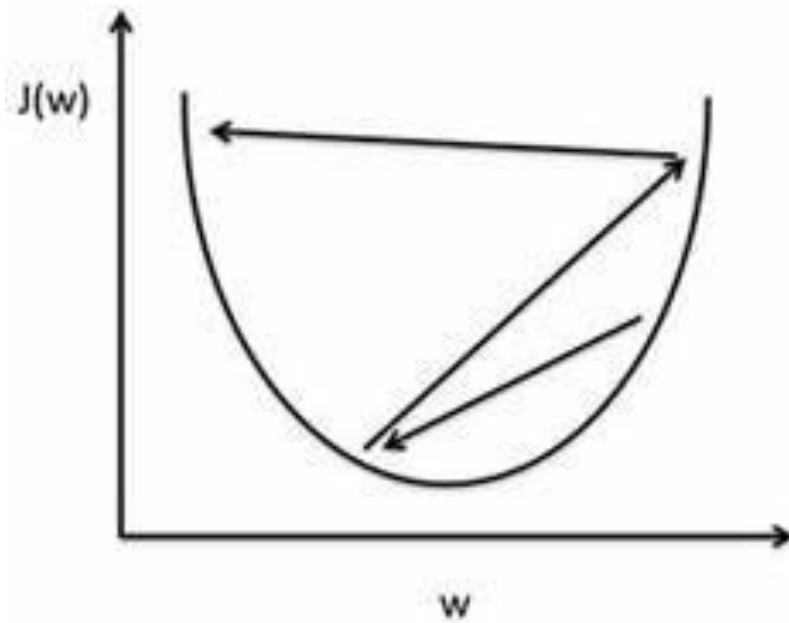
GRADIENT DESCENT

$$\mathbf{w}_{n+1} \leftarrow \mathbf{w}_n - \eta \sum_i \nabla_{\mathbf{w}} \mathcal{L}$$

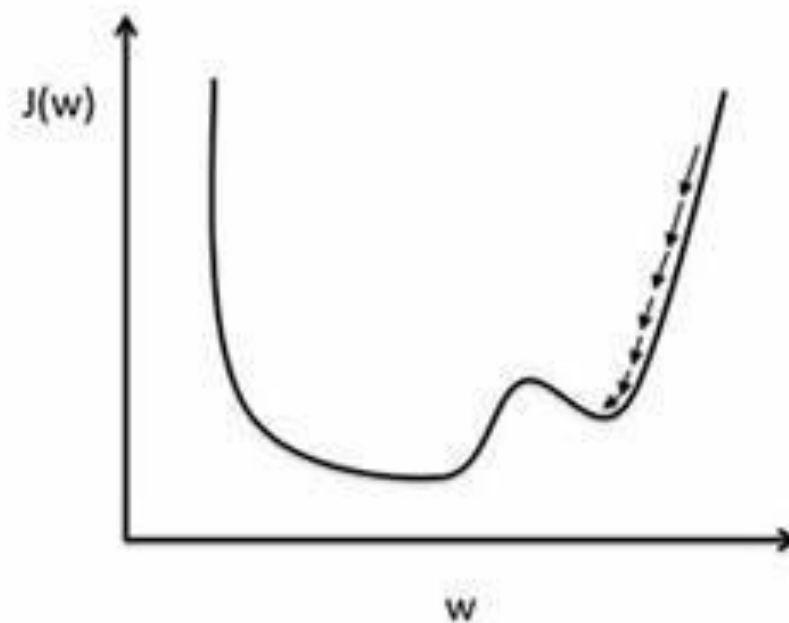




EFFECT OF THE LEARNING RATE (η)



Large Learning Rate



Small Learning Rate

EXAMPLE: STEP ACTIVATION FUNCTION WITH HINGE LOSS

$$\mathcal{L} = \max(0; 1 - y\hat{y})$$

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta \nabla_{\mathbf{w}} \mathcal{L}$$

$$\nabla_{\mathbf{w}} \mathcal{L} = y\mathbf{x}$$

\therefore

$$\mathbf{w}_{n+1} \leftarrow \mathbf{w}_n - \eta y \mathbf{x}$$

EXAMPLE: LOGISTIC ACTIVATION FUNCTION WITH HINGE LOSS

$$\mathcal{L} = \max(0; 1 - y\hat{y})$$

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{x})$$

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta \nabla_{\mathbf{w}} \mathcal{L}$$

$$\nabla_{\mathbf{w}} \mathcal{L} = y\sigma(\mathbf{w}^T \mathbf{x})[1 - \sigma(\mathbf{w}^T \mathbf{x})]\mathbf{x}$$

\therefore

$$\mathbf{w}_{n+1} \leftarrow \mathbf{w}_n - \eta y \sigma(\mathbf{w}^T \mathbf{x})[1 - \sigma(\mathbf{w}^T \mathbf{x})]\mathbf{x}$$

EXAMPLE: LOGISTIC ACTIVATION FUNCTION WITH HINGE LOSS

This tells us about the importance
of normalizing the inputs.
Why?

$$\mathcal{L} = \max(0; 1 - y\hat{y})$$

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{x})$$

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta \nabla_{\mathbf{w}} \mathcal{L}$$

$$\nabla_{\mathbf{w}} \mathcal{L} = y \sigma(\mathbf{w}^T \mathbf{x}) [1 - \sigma(\mathbf{w}^T \mathbf{x})] \mathbf{x}$$

\therefore

$$\mathbf{w}_{n+1} \leftarrow \mathbf{w}_n - \eta y \sigma(\mathbf{w}^T \mathbf{x}) [1 - \sigma(\mathbf{w}^T \mathbf{x})] \mathbf{x}$$

EXAMPLE: LOGISTIC ACTIVATION FUNCTION WITH HINGE LOSS

This tells us about the importance
of normalizing the inputs.
Why?

$$\mathcal{L} = \max(0; 1 - y\hat{y})$$

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{x})$$

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \eta \nabla_{\mathbf{w}} \mathcal{L}$$

$$\nabla_{\mathbf{w}} \mathcal{L} = y \sigma(\mathbf{w}^T \mathbf{x}) [1 - \sigma(\mathbf{w}^T \mathbf{x})] \mathbf{x}$$

\therefore

$$\mathbf{w}_{n+1} \leftarrow \mathbf{w}_n - \eta y \sigma(\mathbf{w}^T \mathbf{x}) [1 - \sigma(\mathbf{w}^T \mathbf{x})] \mathbf{x}$$

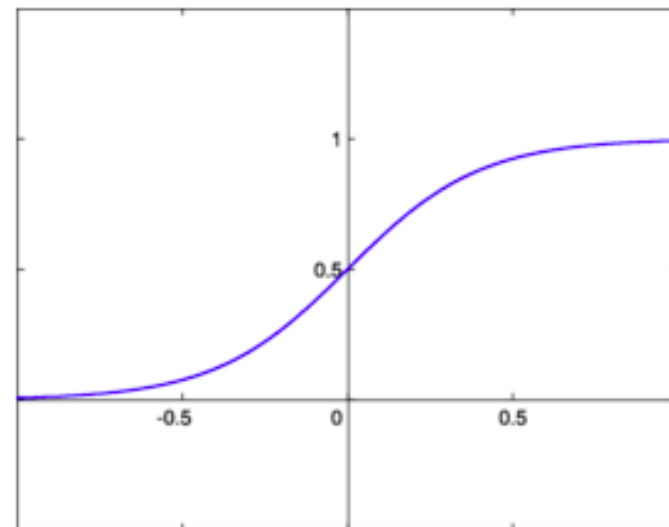
This also offers a nice illustration of
the vanishing gradient problem.
Why? – and what to do about it?

THE VANISHING GRADIENT PROBLEM

$$\mathbf{w}_{n+1} \leftarrow \mathbf{w}_n - \eta \sum_i \nabla_{\mathbf{w}} \mathcal{L}$$

When $\nabla_{\mathbf{w}} \mathcal{L} \rightarrow 0$, learning stops ($\mathbf{w}_{n+1} \approx \mathbf{w}_n$)

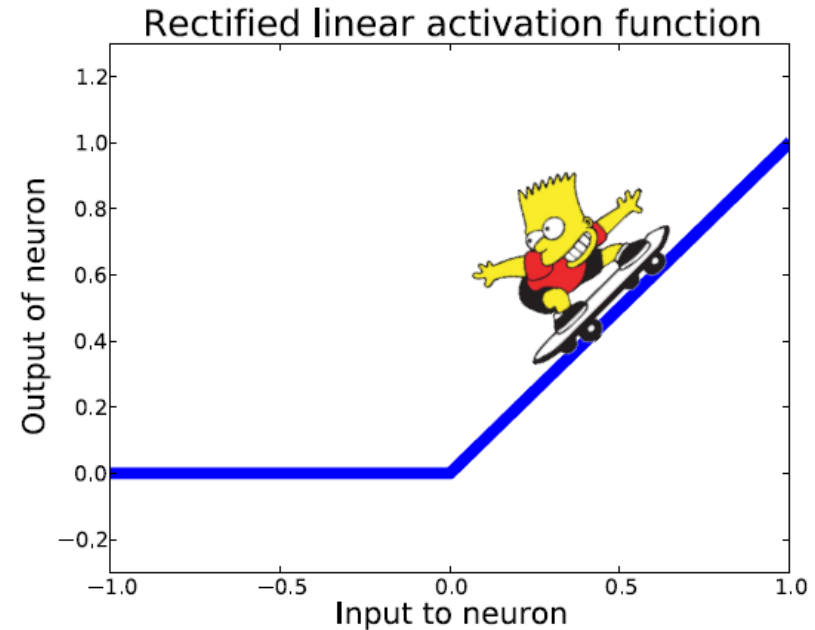
What can we do about it?



SOLUTIONS TO THE VANISHING GRADIENT PROBLEM

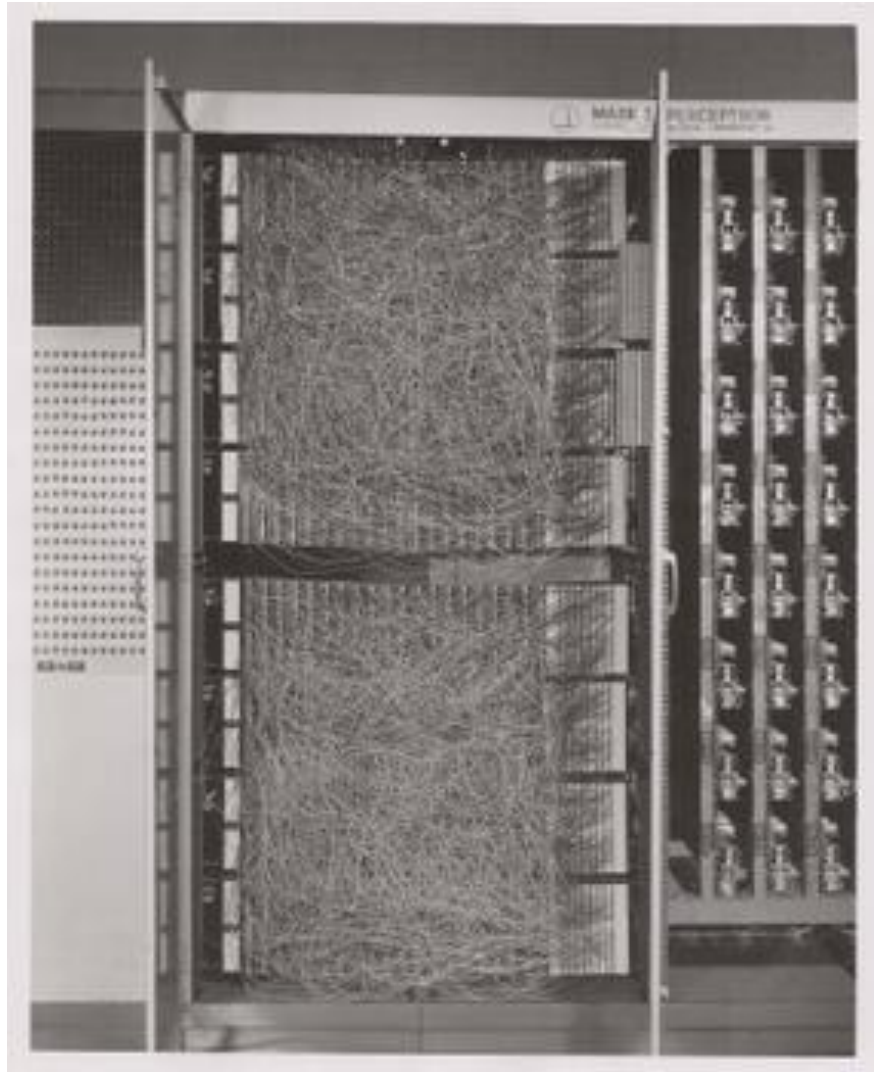
$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

Normalization



Changing the activation function

and other more advanced (or pricy) options...



THE FIRST PERCEPTRON (IBM, 1958)

**COMING UP
NEXT:**

**NEURAL
NETWORKS**

