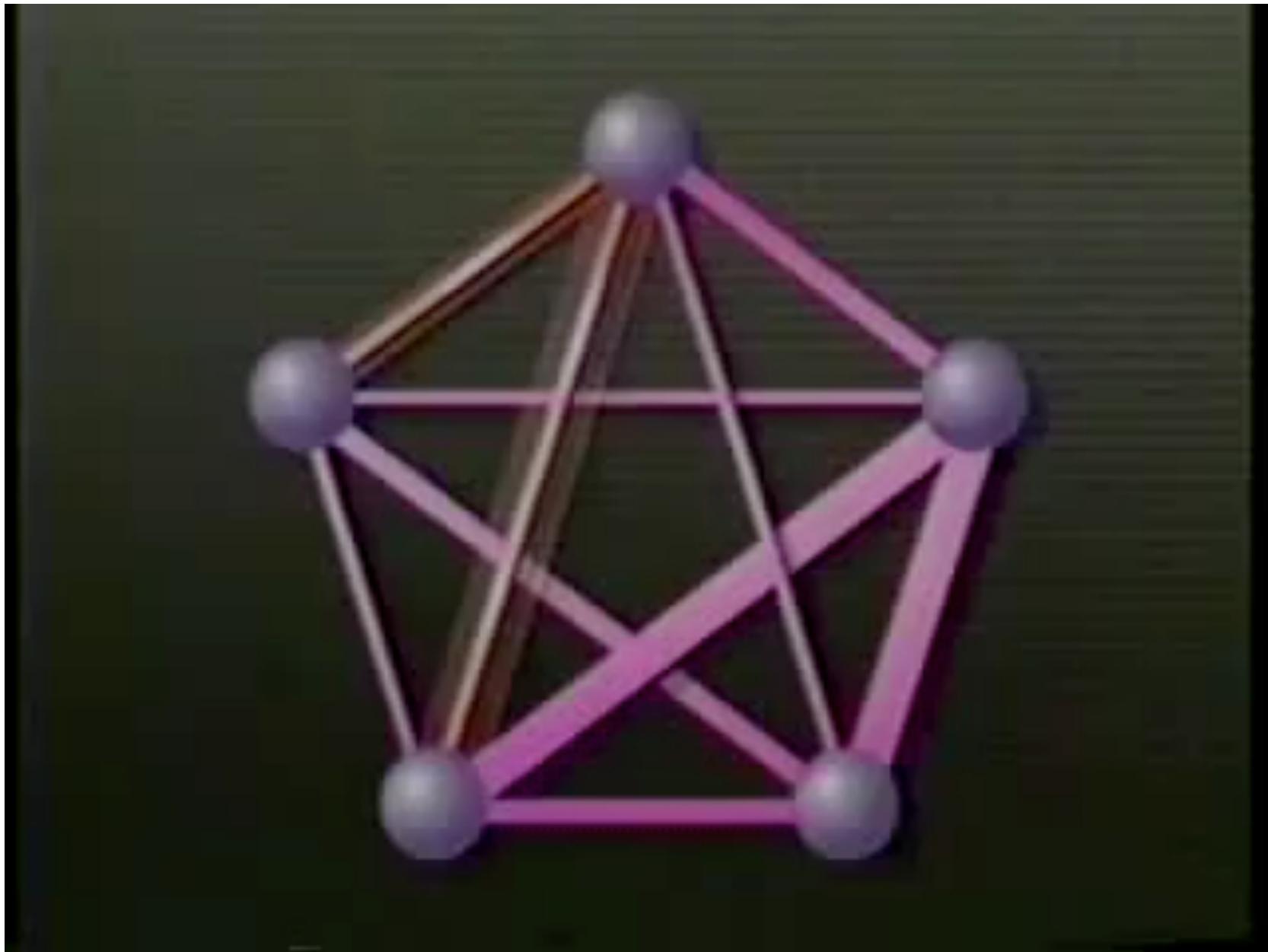


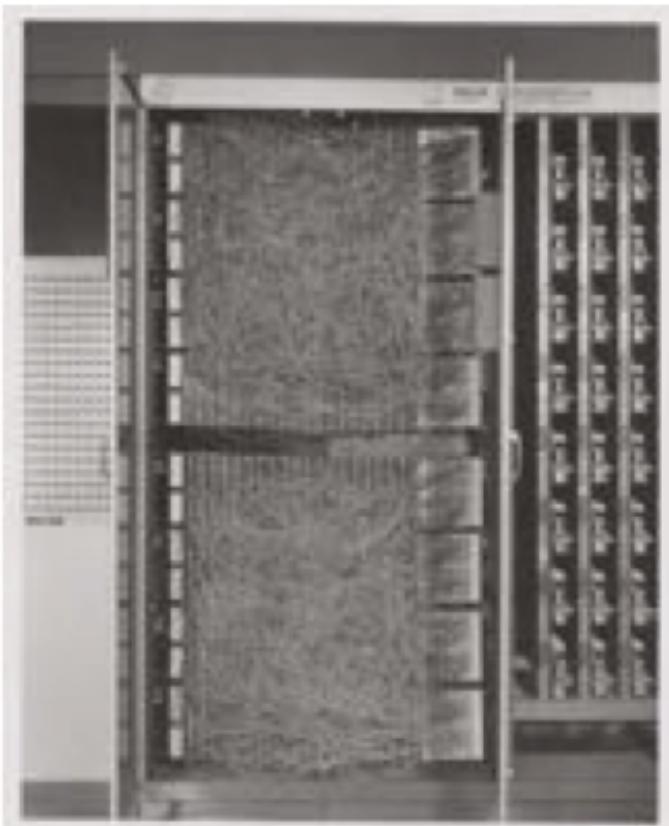
PART II: FOUNDATIONS OF “SHALLOW” NEURAL NETWORKS



Rosenblatt Perceptron

FIRST IMPLEMENTATION OF NEURAL NETWORK [Rosenblatt, 1957!]

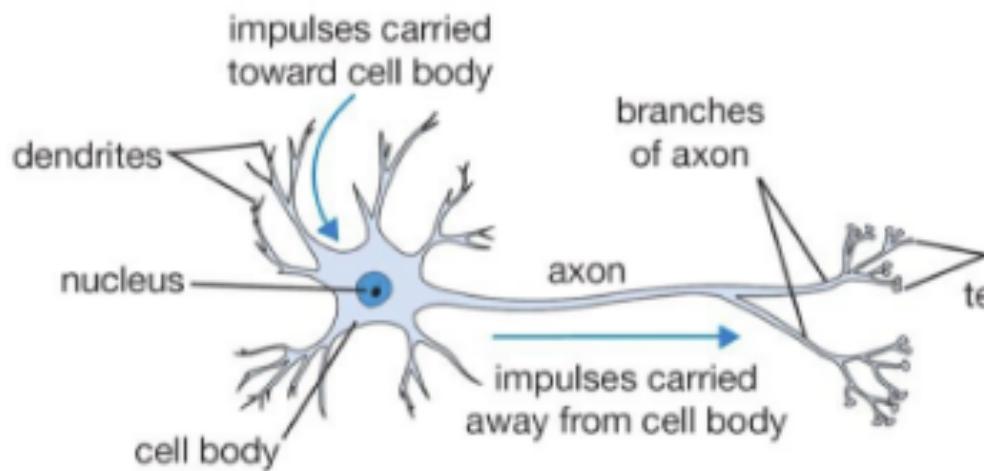
INTENDED TO BE A MACHINE (NOT AN ALGORITHM)



it had an array of 400 photocells, randomly connected to the "neurons".

Weights were encoded in potentiometers, and weight updates during learning were performed by electric motors

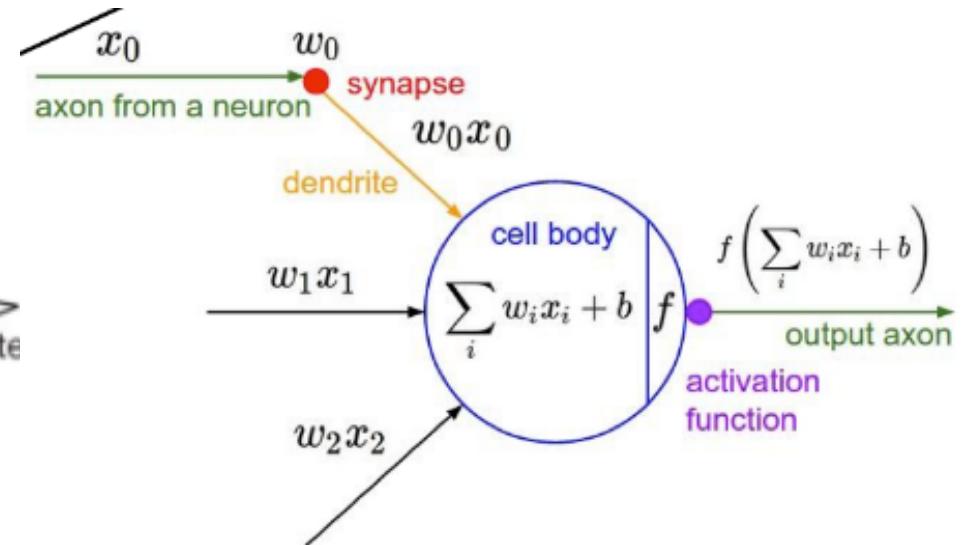
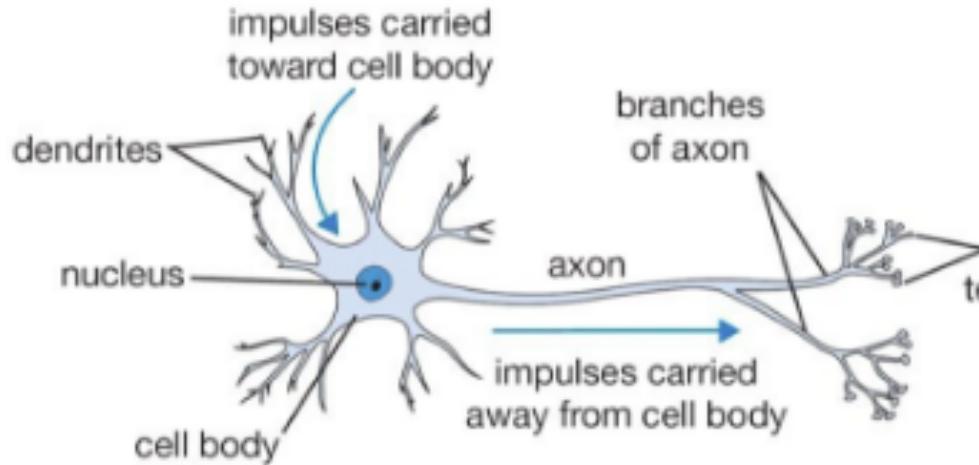
THE NEURON



INSPIRED BY NEURO - SCIENCE?

Credit: Karpathy

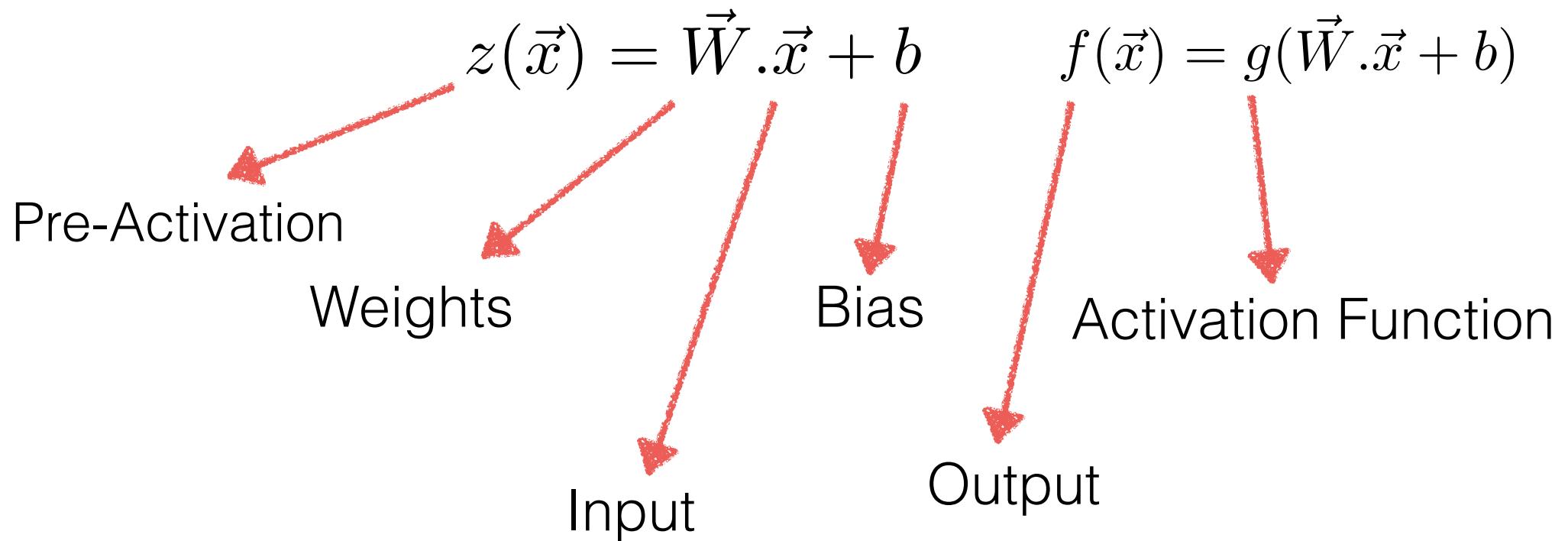
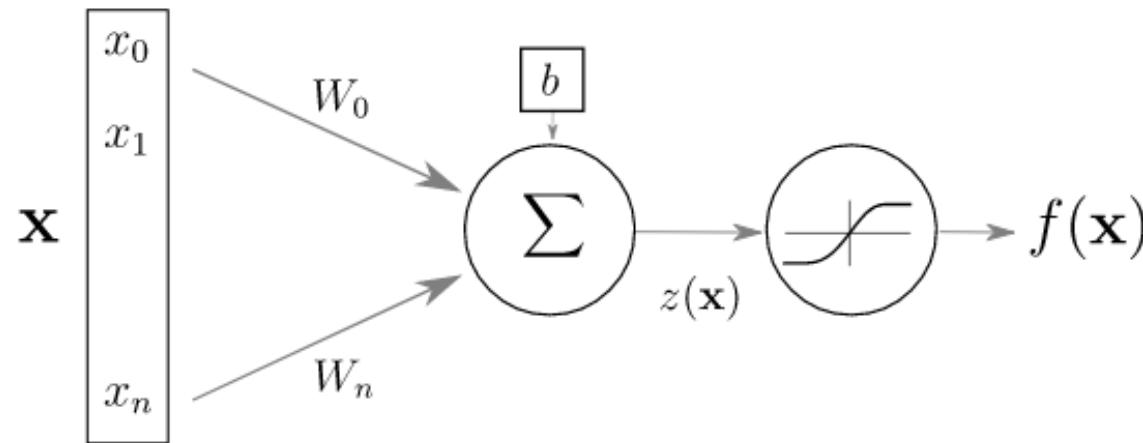
THE NEURON



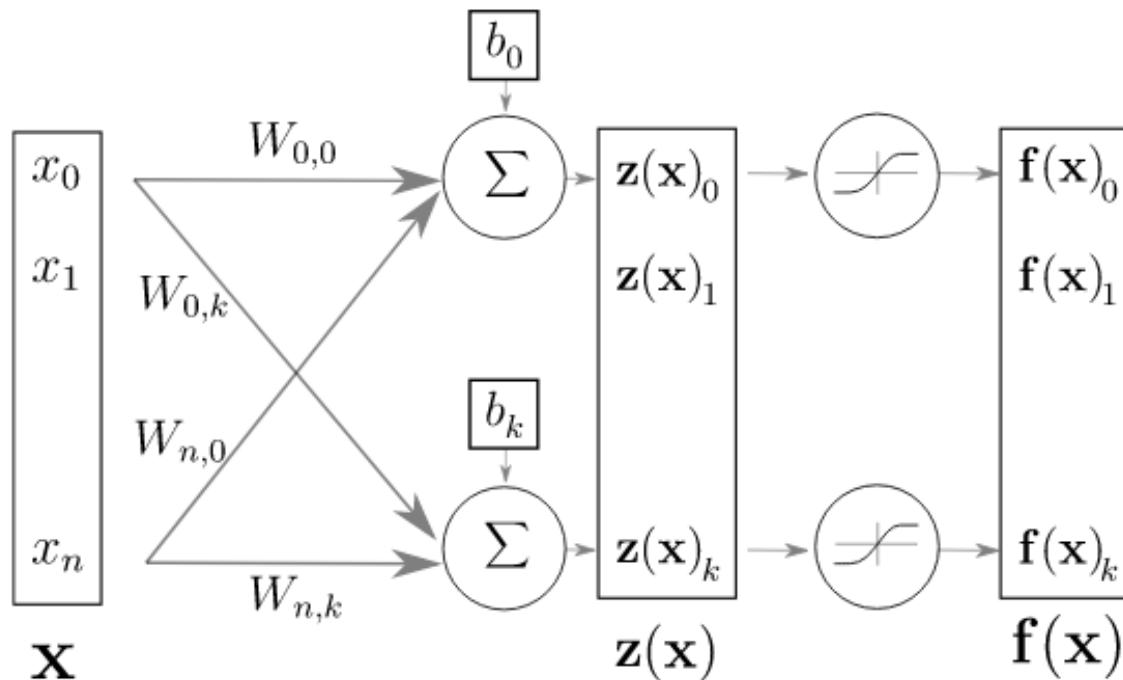
INSPIRED BY NEURO - SCIENCE?

Credit: Karpathy

TODAY'S ARTIFICIAL NEURON



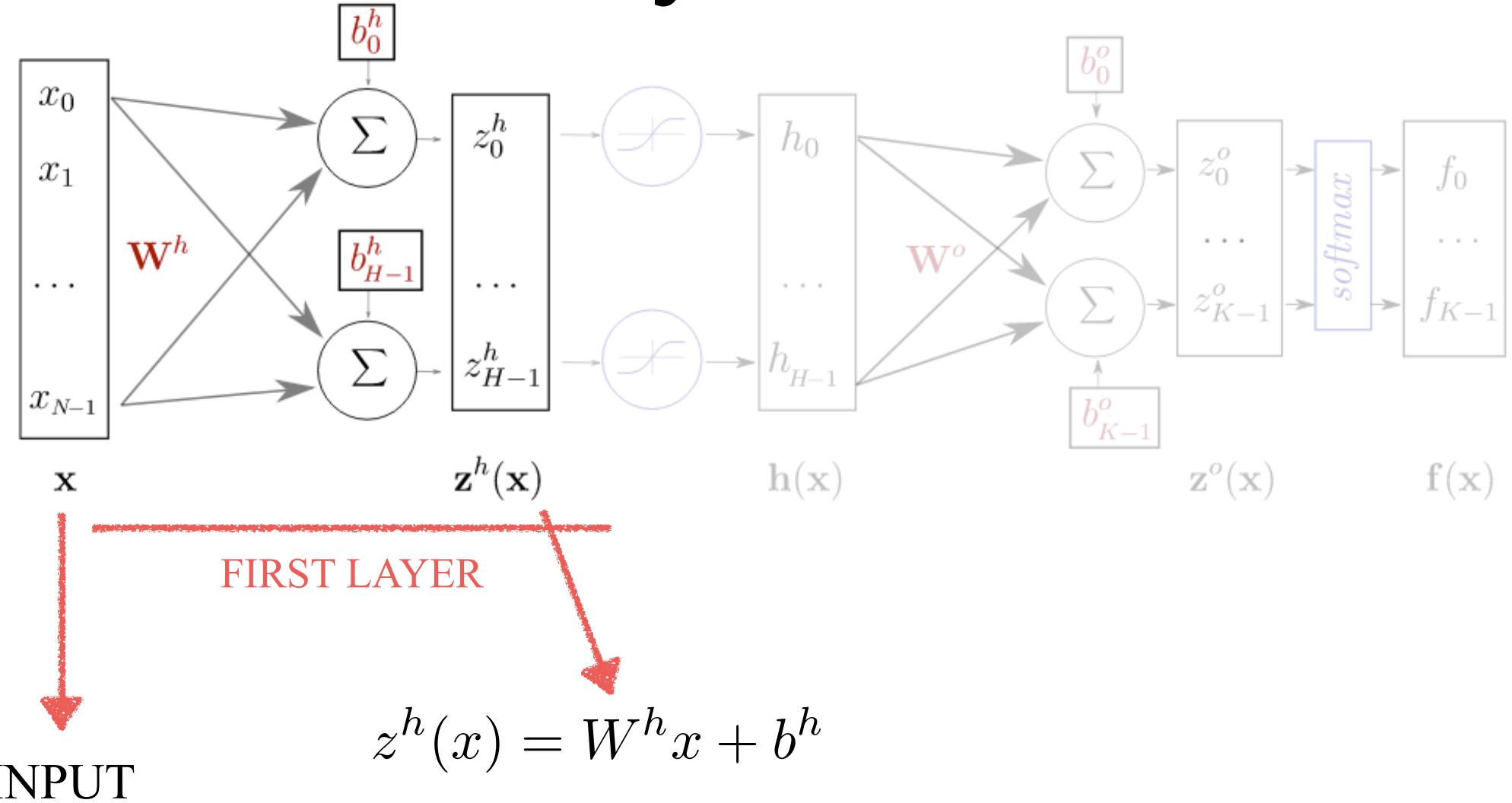
LAYER OF NEURONS



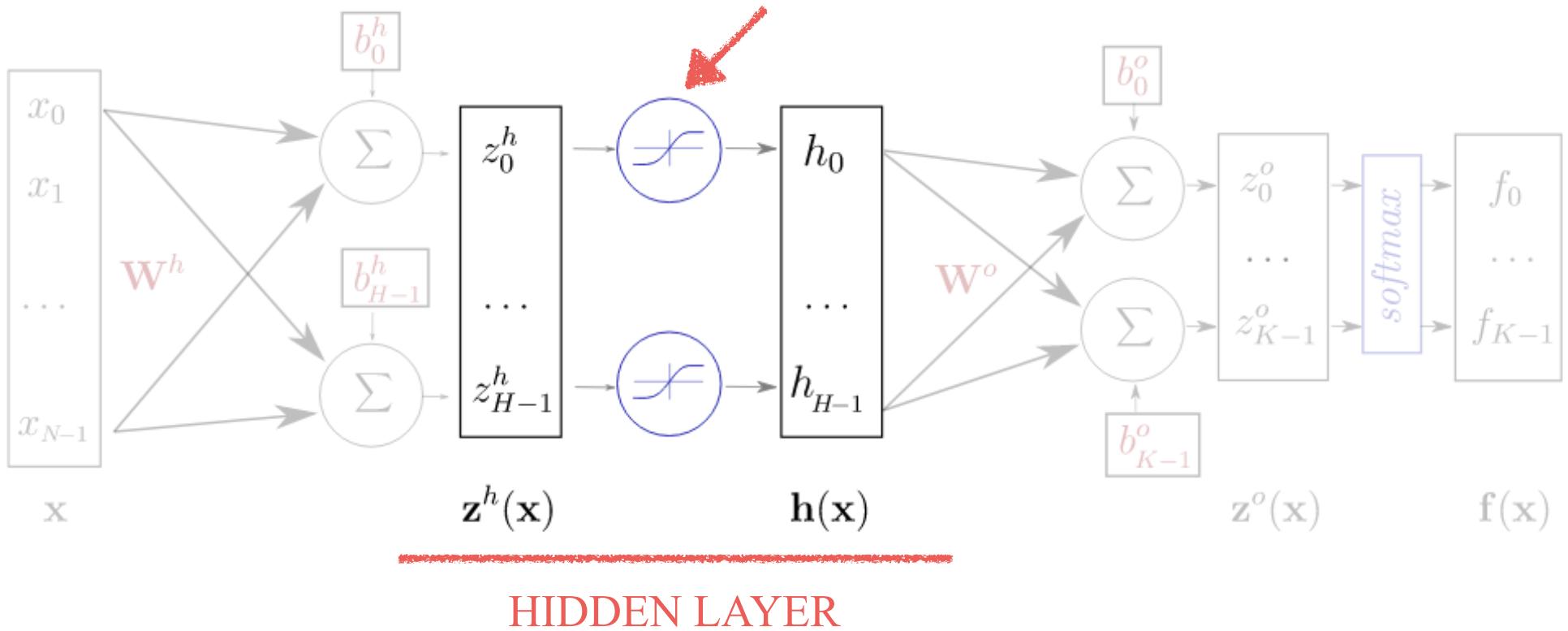
$$f(\vec{x}) = g(\mathbf{W} \cdot \vec{x} + \vec{b})$$

SAME IDEA. NOW **W** becomes a matrix and **b** a vector

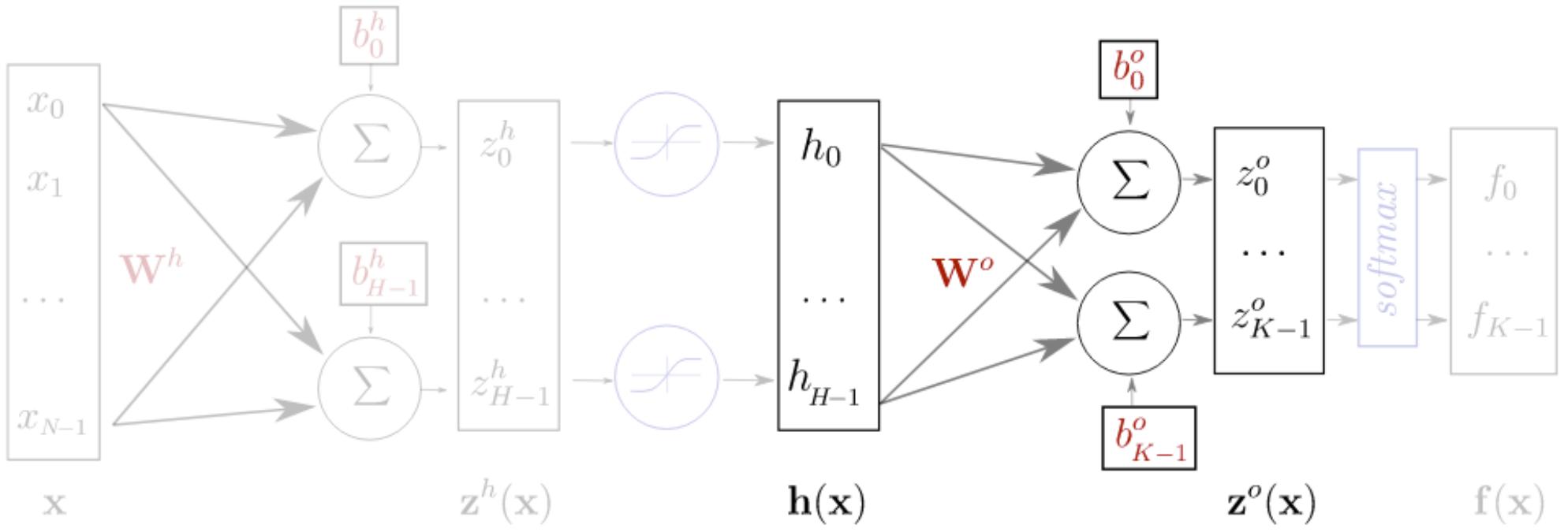
Hidden Layers of Neurons



ACTIVATION FUNCTION

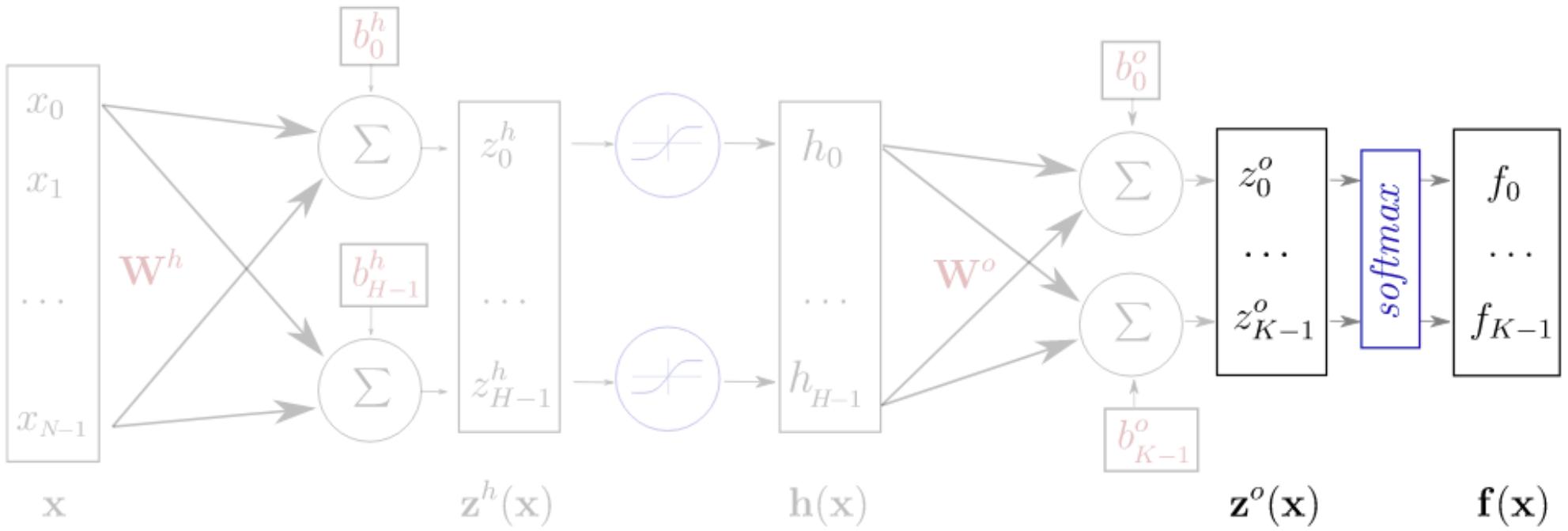


$$h(x) = g(z^h(x)) = g(W^h x + b^h)$$



OUTPUT LAYER

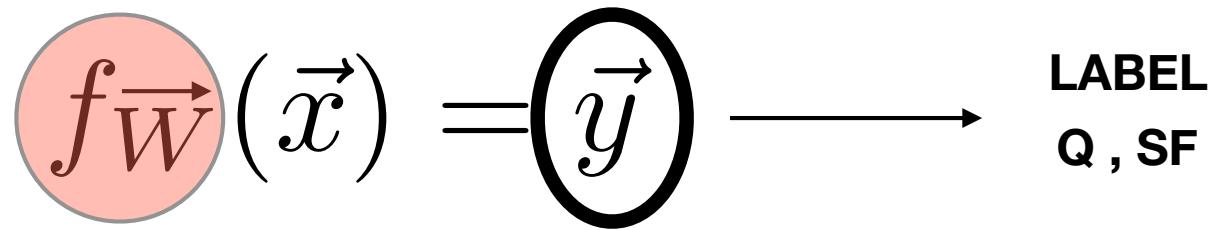
$$z^0(\mathbf{x}) = W^0 h(\mathbf{x}) + b^0$$



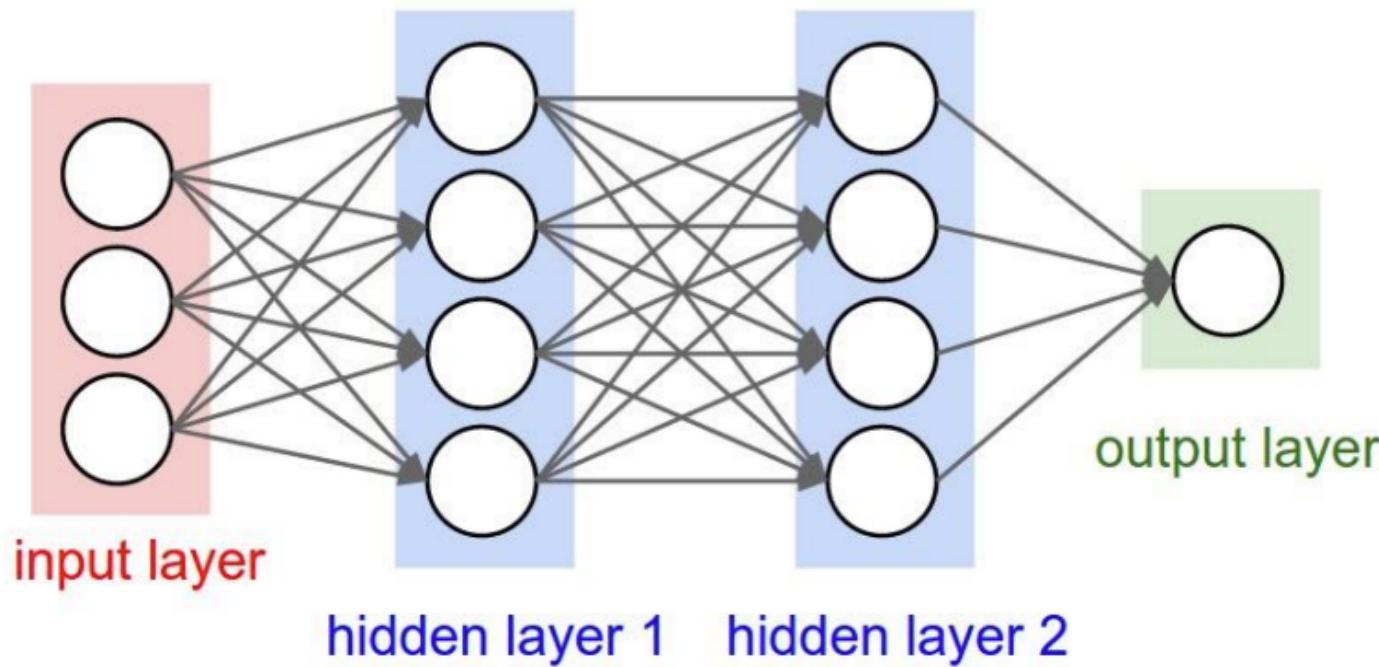
PREDICTION LAYER

$$f(\mathbf{x}) = \text{softmax}(\mathbf{z}^o)$$

**"CLASSICAL"
MACHINE LEARNING**



**REPLACE THIS BY A GENERAL
NON LINEAR FUNCTION WITH SOME PARAMETERS W**

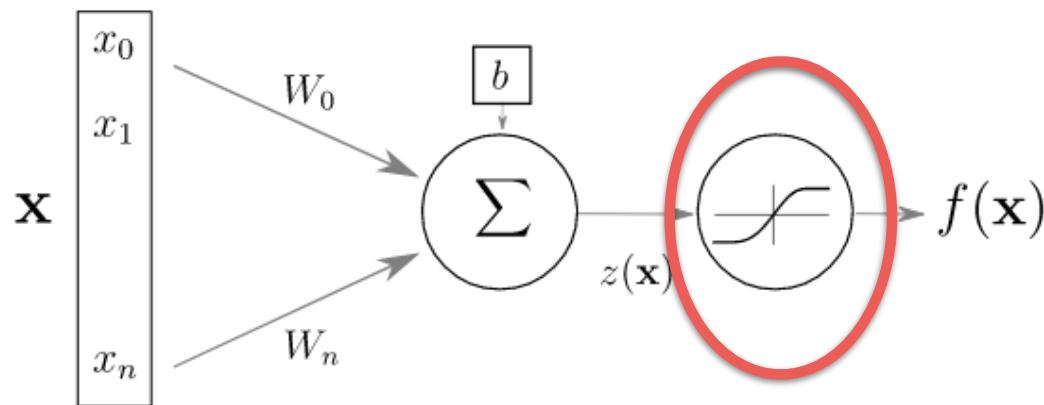


$$p = g_3(W_3g_2(W_2g_1(W_1\vec{x}_0))) \xleftarrow{\text{NETWORK FUNCTION}}$$

SO LET'S GO DEEPER AND DEEPER!

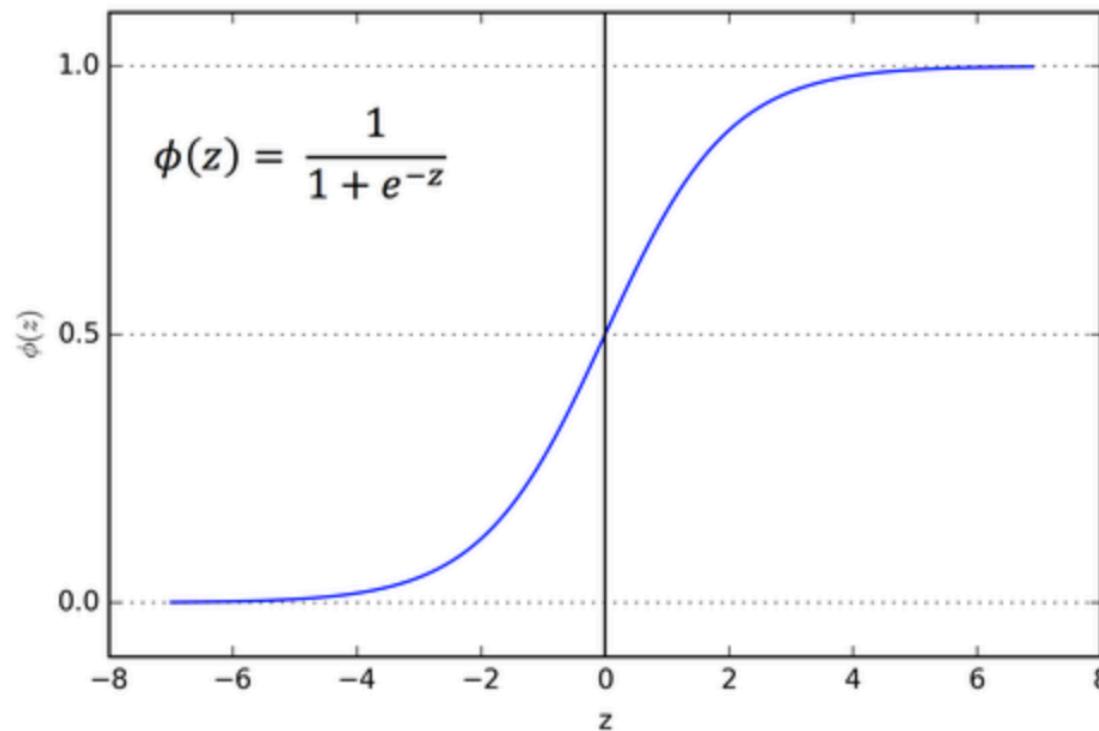
LET'S FIRST EXAMINE IN MORE DETAIL HOW SIMPLE
“SHALLOW” NETWORKS WORK

ACTIVATION FUNCTIONS?

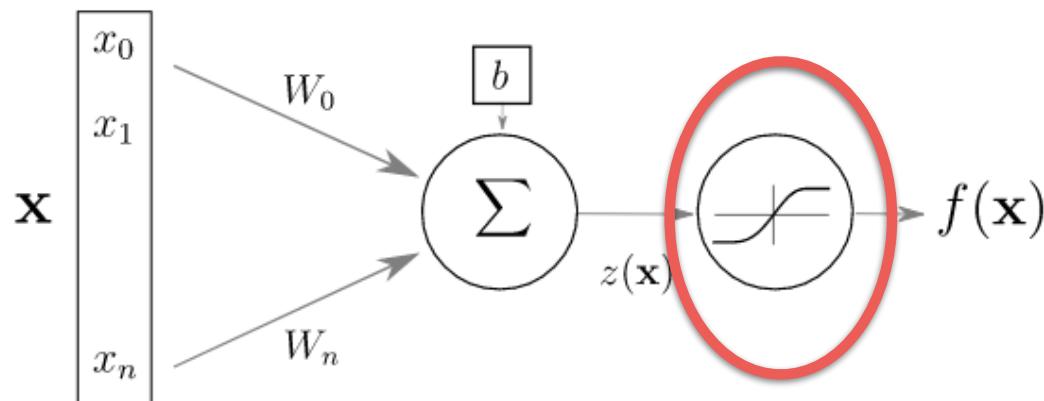


ADD NON LINEARITIES TO THE PROCESS

ACTIVATION FUNCTIONS?

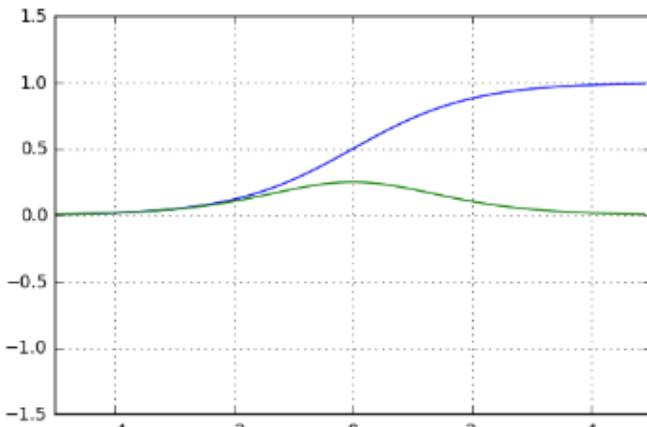


the sigmoid function



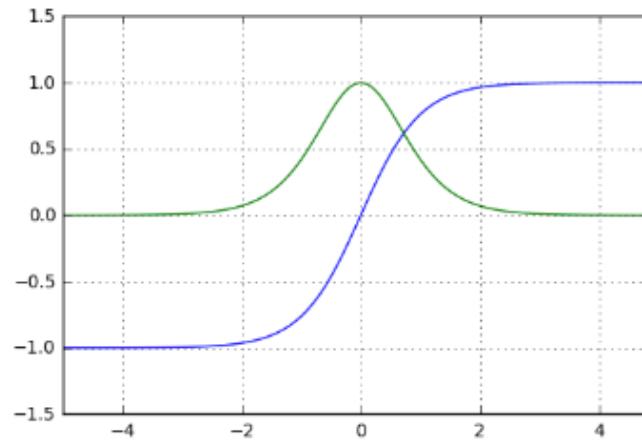
ADD NON LINEARITIES TO THE PROCESS

ACTIVATION FUNCTIONS



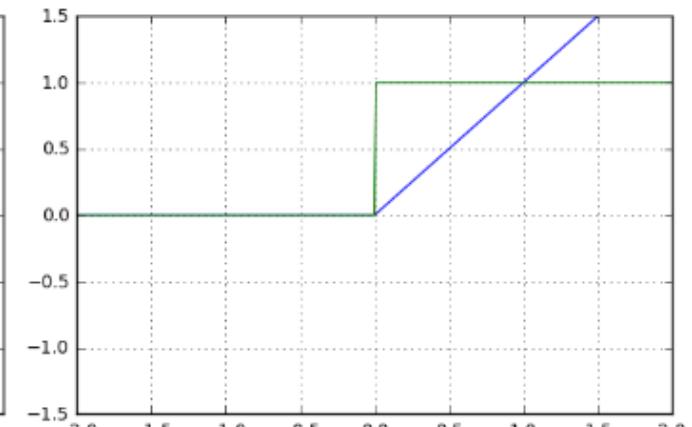
$$\text{sigm}(x) = \frac{1}{1 + e^{-x}}$$

$$\text{sigm}'(x) = \text{sigm}(x)(1 - \text{sigm}(x))$$



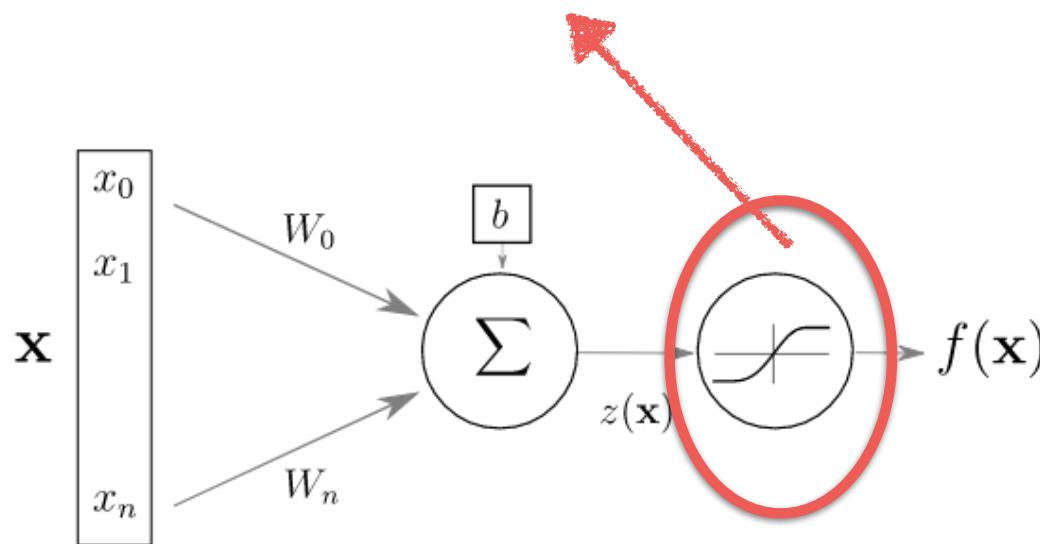
$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

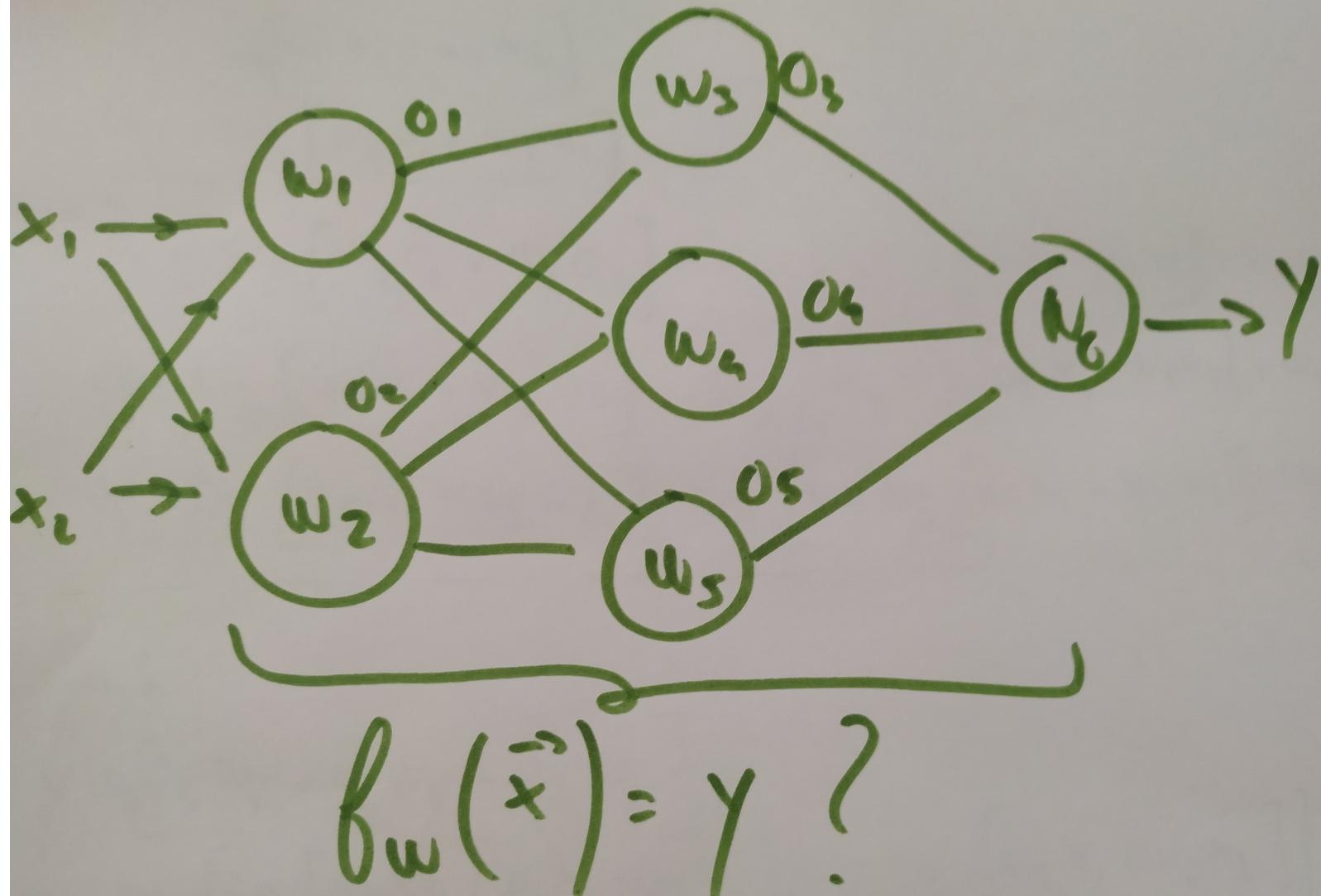
$$\tanh'(x) = 1 - \tanh(x)^2$$



$$\text{relu}(x) = \max(0, x)$$

$$\text{relu}'(x) = 1_{x>0}$$





$$A: \text{sig} [\underbrace{\text{w}_5 \text{sig} [\text{w}_1 x_1 + \text{w}_1 x_2] * w_6}_{+ \text{w}_5 \text{sig} [\text{w}_1 x_1 + \text{w}_1 x_2]}] * w_6$$

$$B: \text{sig} [\text{sig} [\text{w}_5 \text{sig} [\text{w}_1 x_1 + \text{w}_1 x_2] + \text{w}_5 \text{sig} [\text{w}_2 x_1 + \text{w}_2 x_2]] * w_6 + \\ \text{sig} [\text{w}_4 \text{sig} [\text{w}_1 x_1 + \text{w}_1 x_2] + \text{w}_4 \text{sig} [\text{w}_2 x_1 + \text{w}_2 x_2]] * w_6 + \\ \text{sig} [\text{w}_5 \text{sig} [\text{w}_1 x_1 + \text{w}_1 x_2] + \text{w}_5 \text{sig} [\text{w}_2 x_1 + \text{w}_2 x_2]] * w_6]$$

$$C: \text{sig} [\text{sig} [\text{w}_5 \text{sig} [\text{w}_1 x_1 + \text{w}_1 x_2] * w_6 + \text{w}_4 \text{sig} [\text{w}_2 x_1 + \text{w}_2 x_2] * w_6 \\ + \text{w}_5 \text{sig} [\text{w}_1 x_1 + \text{w}_1 x_2] * w_6]$$

When poll is active, respond at **pollev.com/marchuertasc257**

Text **MARCHUERTASC257** to **22333** once to join



Which expression corresponds to the Neural Network shown in the figure?

A
B
C

AND REMEMBER WE NEED A LOSS FUNCTION ...

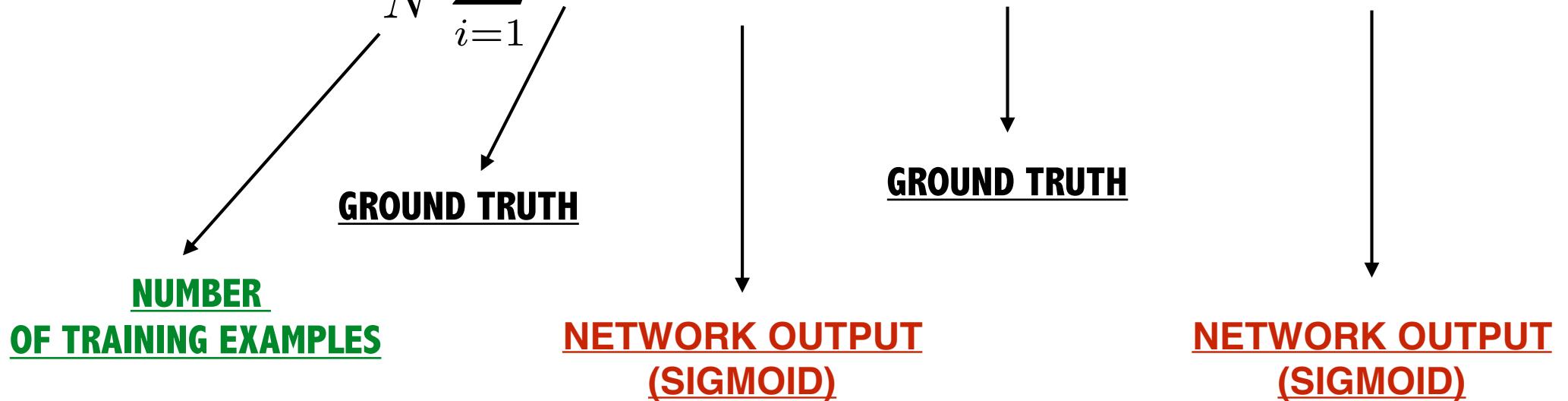
AND REMEMBER WE NEED A LOSS FUNCTION ...

LET'S SART WITH A SIMPLE CLASSIFICATION PROBLEM

BINARY CLASSIFICATION LOSS

WE TYPICALLY USE THE BINARY CROSS-ENTROPY LOSS:

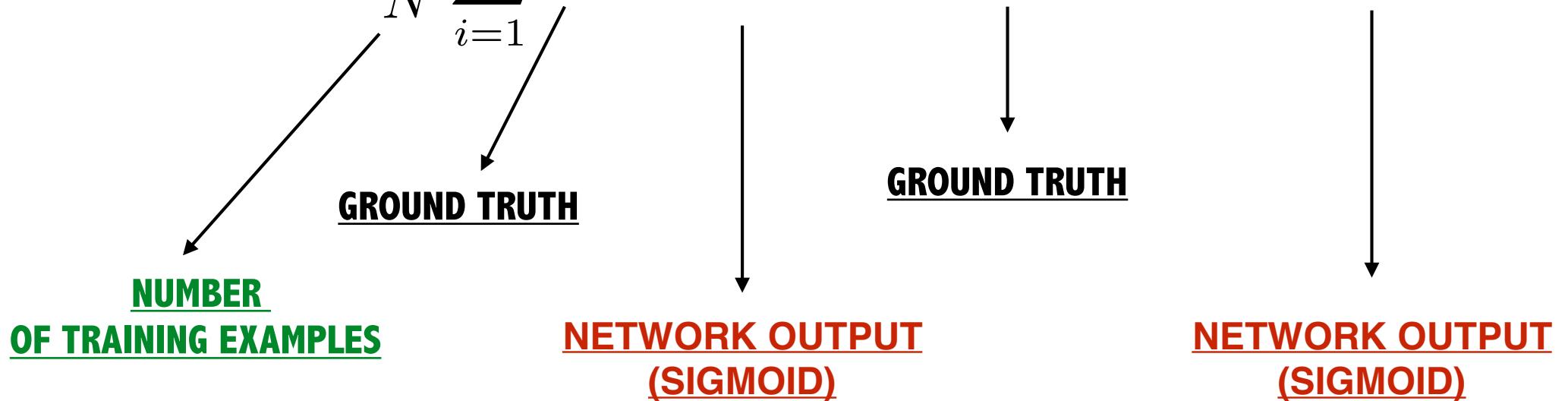
$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \log(p(y_i)) + (1 - y_i) \log(1 - p(y_i))$$



BINARY CLASSIFICATION LOSS

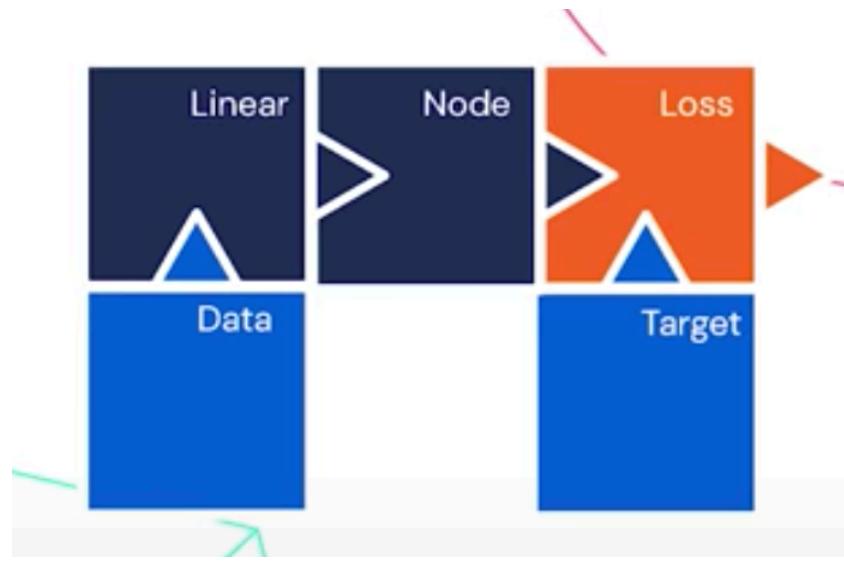
WE MINIMIZE THE CROSS ENTROPY BETWEEN THE TRUE $[q(y)]$ AND PREDICTED $[p(y)]$ DISTRIBUTIONS

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \log(p(y_i)) + (1 - y_i) \log(1 - p(y_i))$$



OK, SO NOW WE HAVE THE 3 MAIN INGREDIENTS THAT COMPOSE AN ANN:

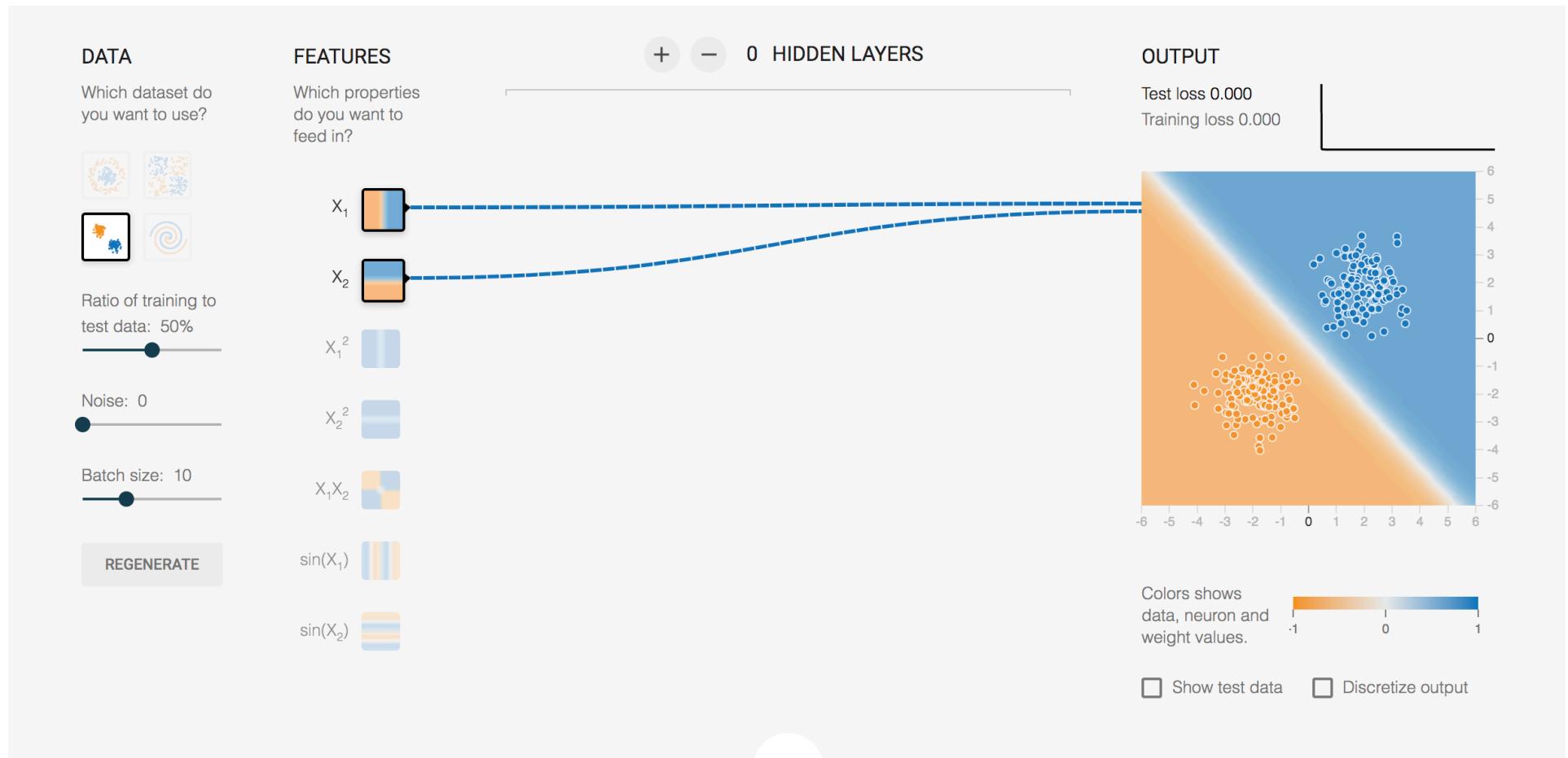
1. LINEAR NEURON OR UNIT
2. NON LINEARITIES (ACTIVATION FUNCTION)
3. LOSS FUNCTION



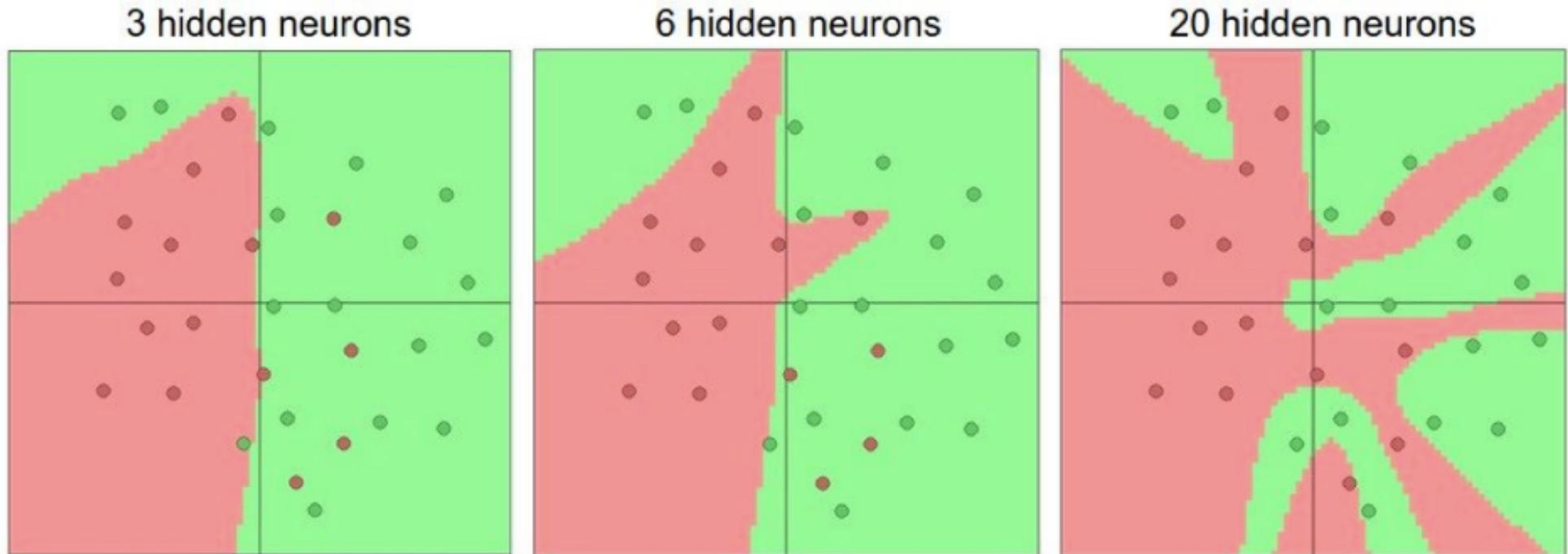
(deepmind
@Czarnecki)

LET'S SEE WHAT WE CAN DO WITH THIS SIMPLE MODEL!

<https://playground.tensorflow.org/>



HIDDEN LAYERS ALLOW INCREASING COMPLEXITY



More complex functions allow increasing complexity

Credit: Karpathy

GO HERE:

[https://github.com/mhuertascompany/SaaS-Fee/blob/main/
hands-on/session1/hello_ANN.ipynb](https://github.com/mhuertascompany/SaaS-Fee/blob/main/hands-on/session1/hello_ANN.ipynb)

UNIVERSAL APPROXIMATION THEOREM

FOR ANY CONTINUOS FUNCTION FOR A HYPERCUBE $[0,1]^d$ TO REAL NUMBERS, AND EVERY POSITIVE EPSILON, THERE EXISTS A SIGMOID BASED 1-HIDDEN LAYER NEURAL NETWORK THAT OBTAINS AT MOST EPSILON ERROR IN FUNCTIONAL SPACE

Cybenko+89

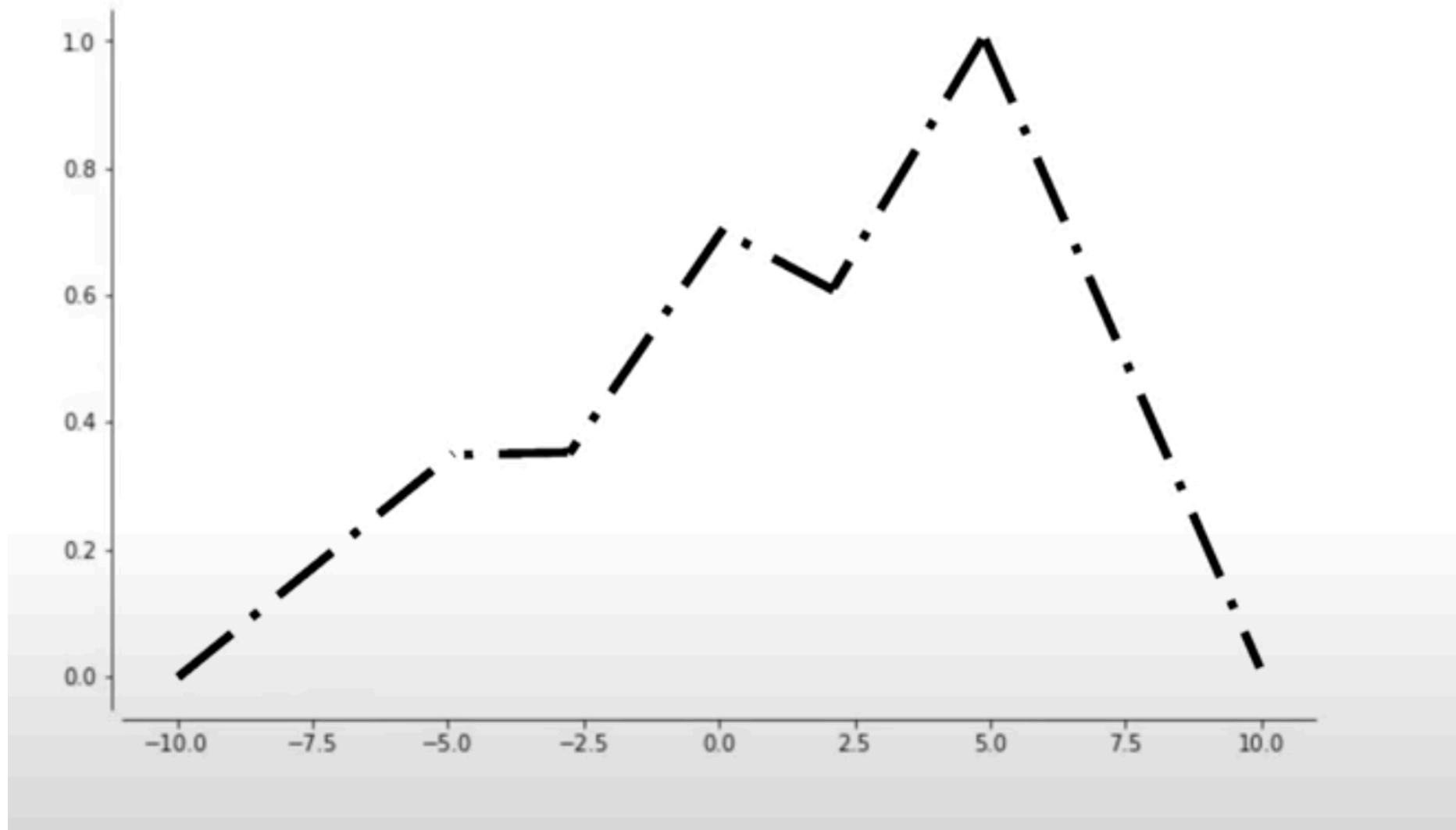
“BIG ENOUGH NETWORK CAN APPROXIMATE, BUT NOT REPRESENT ANY SMOOTH FUNCTION. THE MATH DEMONSTRATION IMPLIES SHOWING THAT NETWORS ARE DENSE IN THE SPACE OF TARGET FUNCTIONS”

UNIVERSAL APPROXIMATION THEOREM

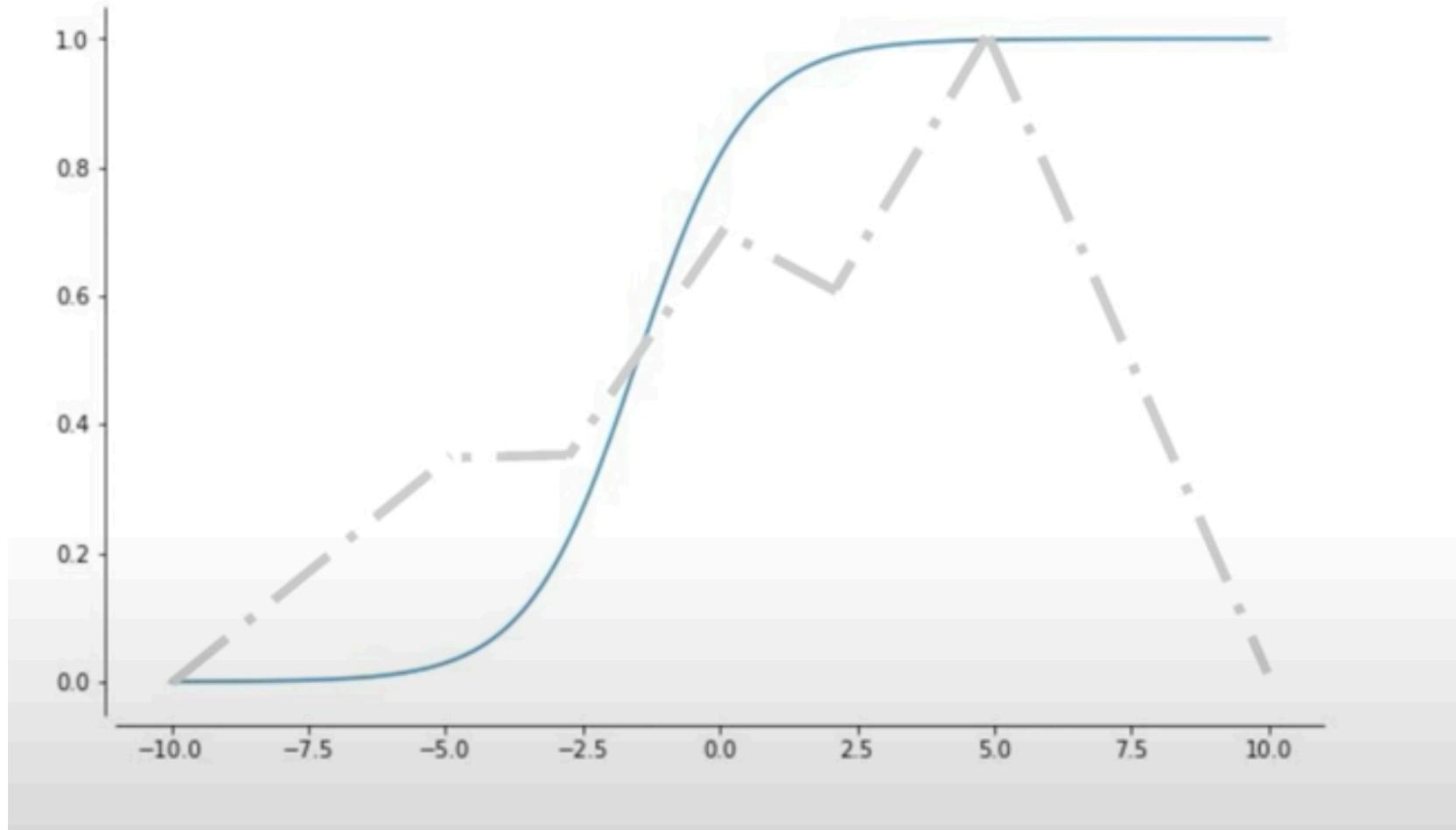
FOR ANY CONTINUOS FUNCTION FOR A HYPERCUBE $[0,1]^d$ TO REAL NUMBERS, NON-CONSTANT, BOUNDED AND CONTINUOUS ACTIVATION FUNCTION f , AND EVERY POSITIVE EPSILON, THERE EXISTS A 1-HIDDEN LAYER NEURAL NETWORK USING f THAT OBTAINS AT MOST EPSILON ERROR IN FUNCTIONAL SPACE

Horvik+91

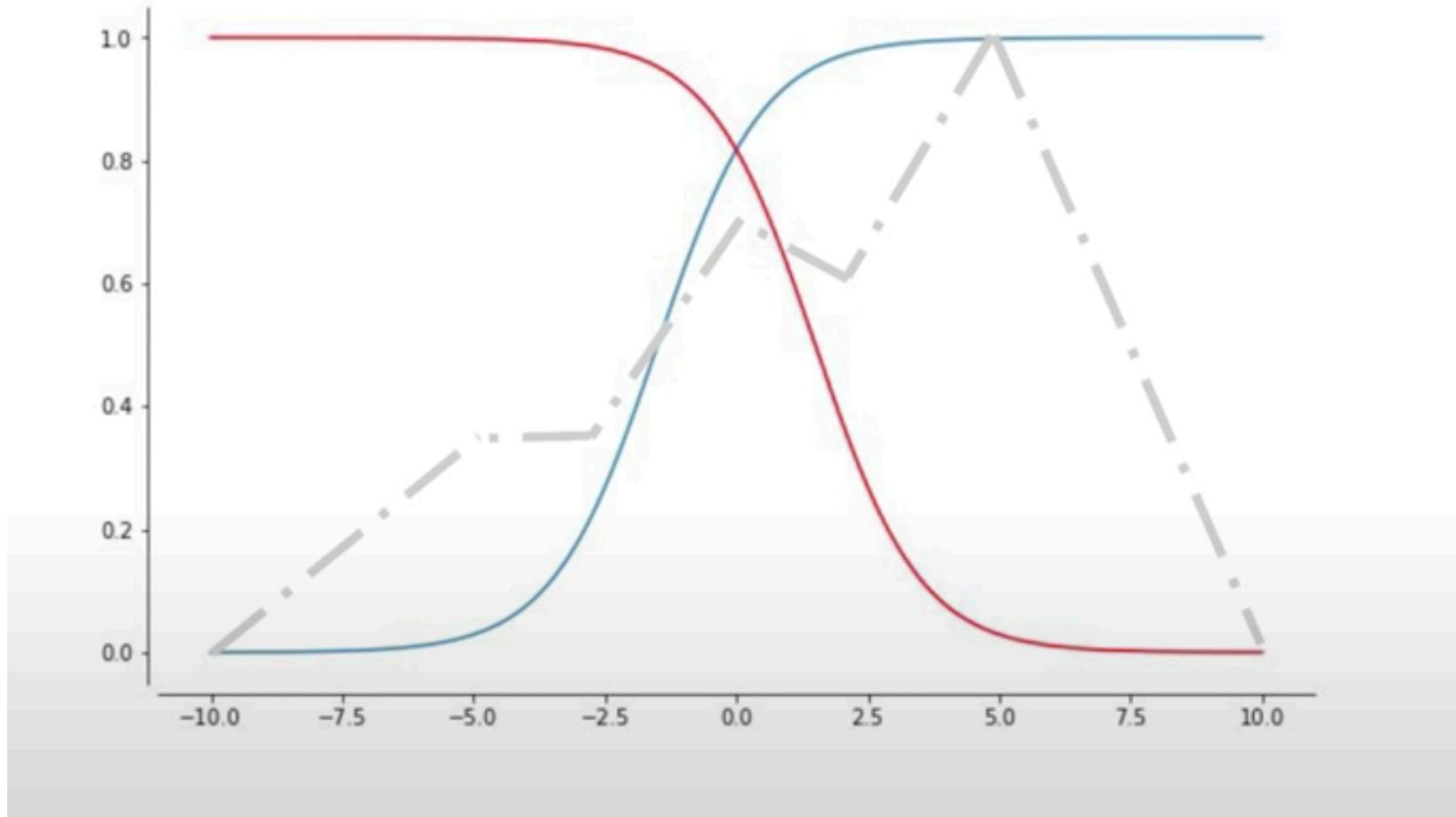
“BIG ENOUGH NETWORK CAN APPROXIMATE, BUT NOT REPRESENT ANY SMOOTH FUNCTION. THE MATH DEMONSTRATION IMPLIES SHOWING THAT NETWORS ARE DENSE IN THE SPACE OF TARGET FUNCTIONS”



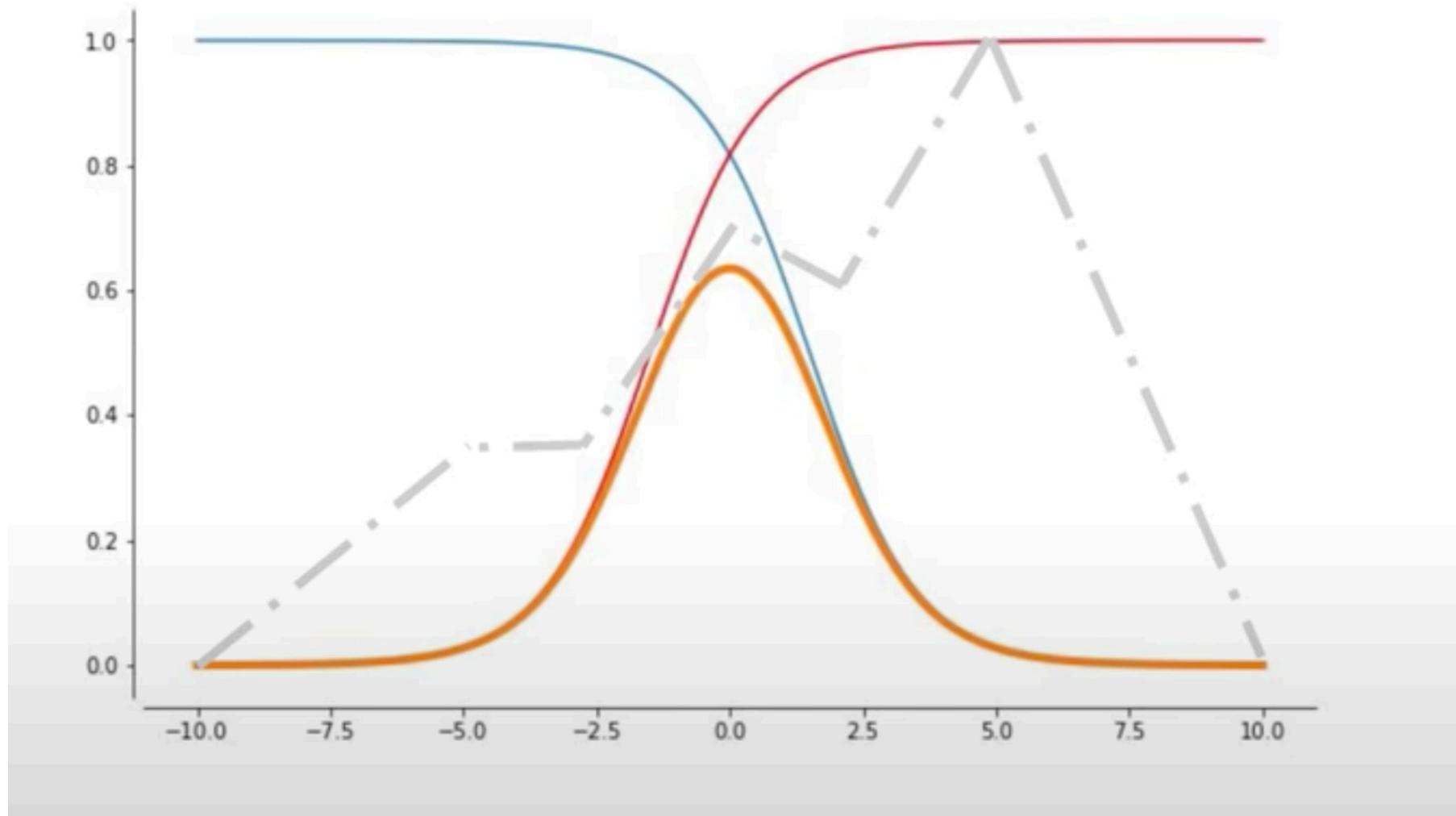
(deepmind
@Czarnecki)



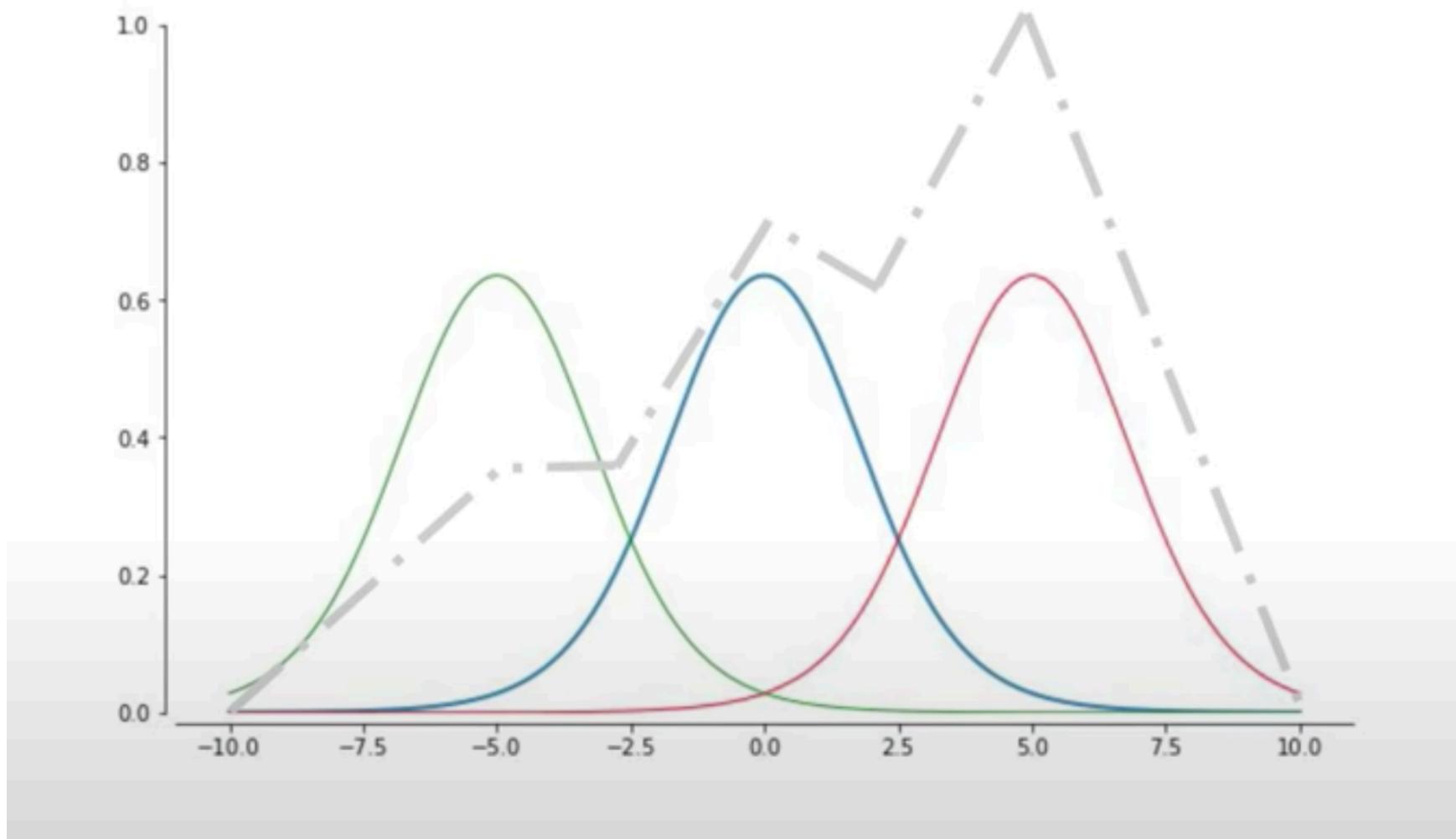
(deepmind
@Czarnecki)



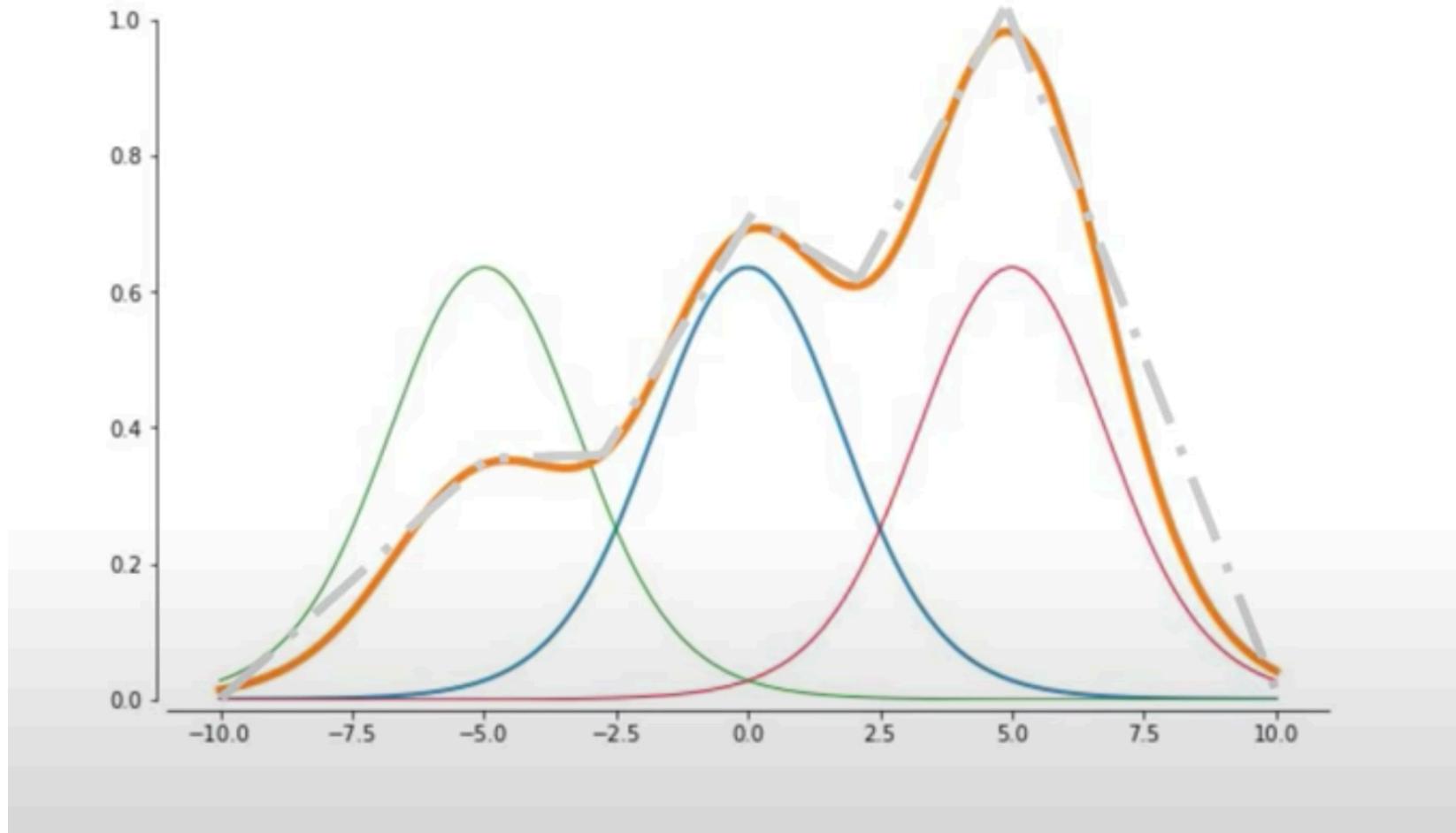
(deepmind
@Czarnecki)



(deepmind
@Czarnecki)

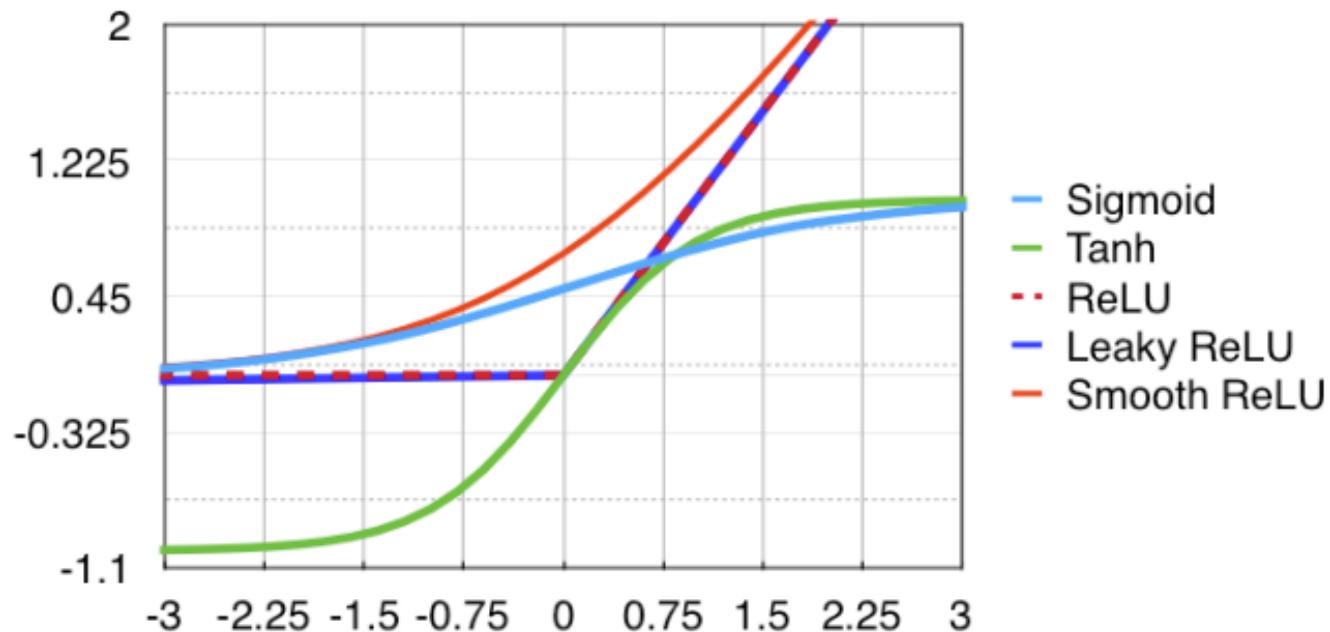


(deepmind
@Czarnecki)



(deepmind
@Czarnecki)

ACTIVATION FUNCTIONS



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

Tanh: $f(x) = \tanh(x)$

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

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

Leaky ReLU: $f(x) = \epsilon x + (1 - \epsilon) \max(0, x)$

SOFTMAX

A generalization of the SIGMOID ACTIVATION

$$\text{softmax}(\mathbf{x}) = \frac{e^{\mathbf{x}}}{\sum_{i=1}^n e^{x_i}}$$

THE OUTPUT IS NORMALIZED BETWEEN 0 AND 1

THE COMPONENTS ADD TO 1

CAN BE INTERPRETED AS A PROBABILITY

$$p(Y = c | X = \mathbf{x}) = \text{softmax}(z(\mathbf{x}))_c$$

SO LET'S GO DEEPER AND DEEPER!

YES BUT...

NOT SO STRAIGHTFORWARD, DEEPER MEANS MORE
WEIGHTS, MORE DIFFICULT OPTIMIZATION, RISK OF
OVERFITTING...