

Foundations of Machine Learning

AI2000 and AI5000

FoML-23

Neural Networks - UAT

Dr. Konda Reddy Mopuri

Department of AI, IIT Hyderabad

July-Nov 2025



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



So far in FoML

- Intro to ML and Probability refresher
- MLE, MAP, and fully Bayesian treatment
- Supervised learning
 - a. Linear Regression with basis functions (regularization, model selection)
 - b. Bias-Variance Decomposition (Bayesian Regression)
 - c. Decision Theory - three broad classification strategies
 - Probabilistic Generative Models - Continuous & discrete data
 - (Linear) Discriminant Functions - least squares solution, Perceptron
 - Probabilistic Discriminative Models - Logistic Regression

Neural Networks - II

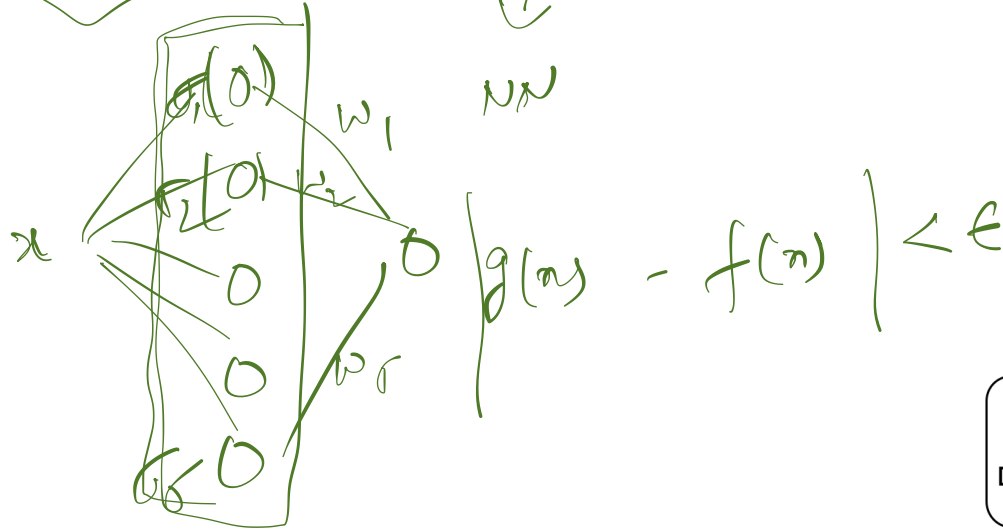


Neural Networks are universal approximators



Universal Approximation Theorem

- Can represent any continuous function ($f: \mathbb{R}^m \rightarrow \mathbb{R}^n$) on a compact area, to any desired approximation ($|g(x) - f(x)| < \epsilon$) with a linear combination of sigmoid neurons



Universal Approximation Theorem

- In other words, NN with a single hidden layer can be used to approximate any continuous function to a desired precision



Universal Approximation Theorem

Math. Control Signals Systems (1989) 2: 303–314

Mathematics of Control,
Signals, and Systems

© 1989 Springer-Verlag New York Inc.

Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†

Neural Networks, Vol. 4, pp. 251–257, 1991
Printed in the USA. All rights reserved.

0893-6080/91 \$3.00 + .00
Copyright © 1991 Pergamon Press plc

ORIGINAL CONTRIBUTION

**Approximation Capabilities of Multilayer
Feedforward Networks**

KURT HORNIK

Technische Universität Wien, Vienna, Austria



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



Universal Approximation Theorem

Theorem 0.1 (UAT, [Cyb89, Hor91]). Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be a *non-constant, bounded, and continuous* function. Let I_m denote the m -dimensional *unit hypercube* $[0, 1]^m$. The space of *real-valued continuous functions on I_m* is denoted by $C(I_m)$. Then, given any $\varepsilon > 0$ and any function $f \in C(I_m)$, *there exist an integer N , real constants $v_i, b_i \in \mathbb{R}$ and real vectors $w_i \in \mathbb{R}^m$ for $i = 1, \dots, N$, such that we may define:*

$$F(\mathbf{x}) = \sum_{i=1}^N v_i \sigma(\mathbf{w}_i^T \mathbf{x} + b_i) = \mathbf{v}^T \sigma(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

as an approximate realization of the function f ; that is,

$$|F(\mathbf{x}) - f(\mathbf{x})| < \varepsilon$$

for all $\mathbf{x} \in I_m$.

-3
10

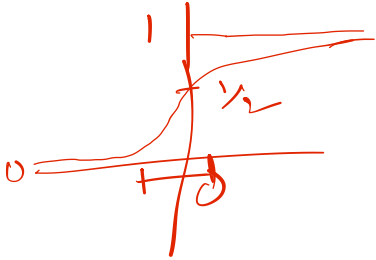
-6
10



Visual proof with one i/p & one o/p and Sigmoid activation

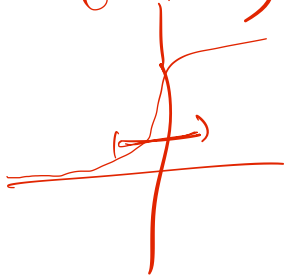
Universality with one i/p & one o/p

$$\sigma(wa + b)$$

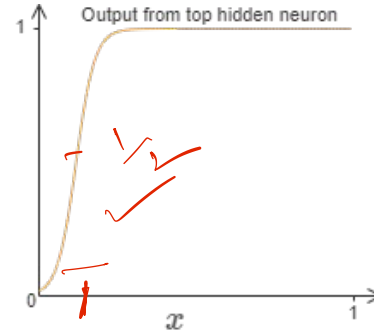
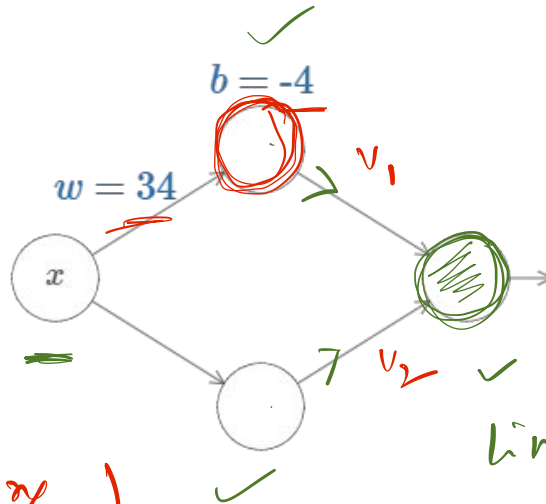


$$\sigma(wa)$$

$$\sigma(100a)$$



$$\sigma(1/100)$$

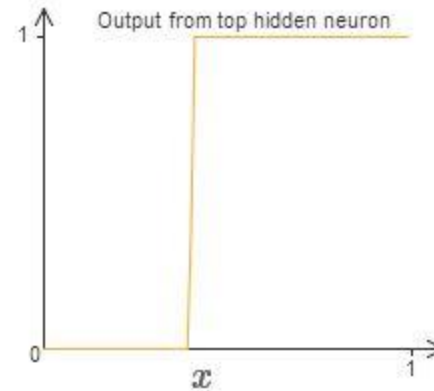
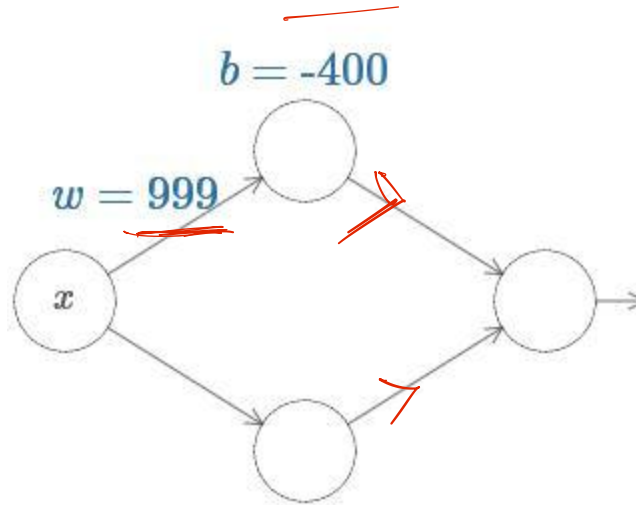


linear comb. of sigmoid neurons

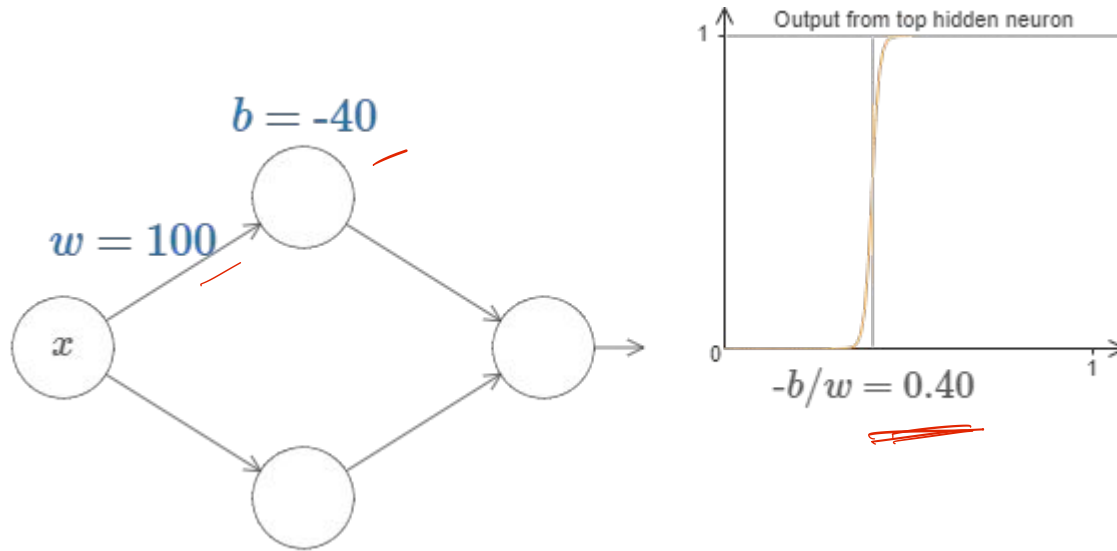
$$\sigma(a+b)$$

$$\begin{cases} a+b=0 \\ a=-b \end{cases}$$

Universality with one i/p & one o/p



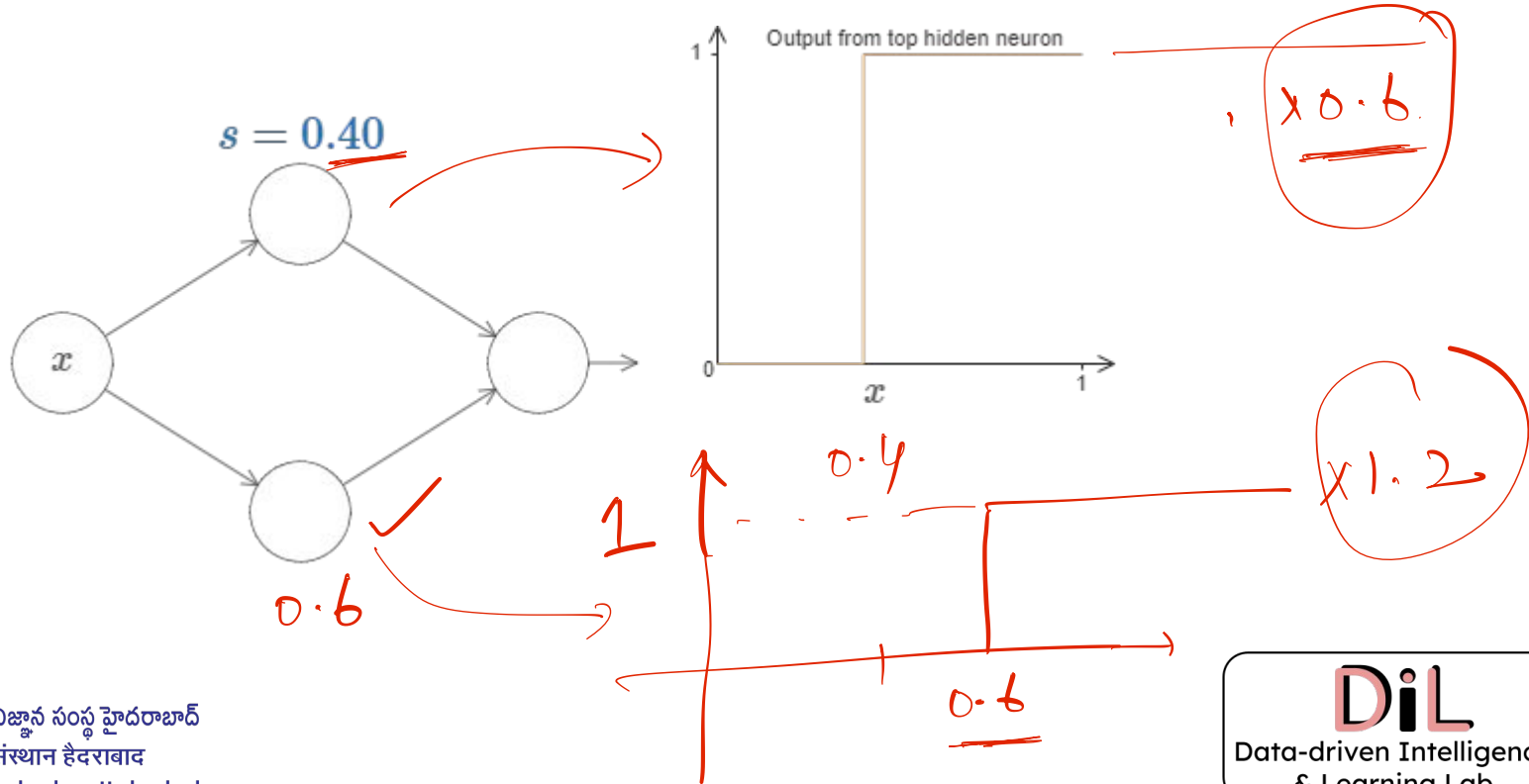
Universality with one i/p & one o/p



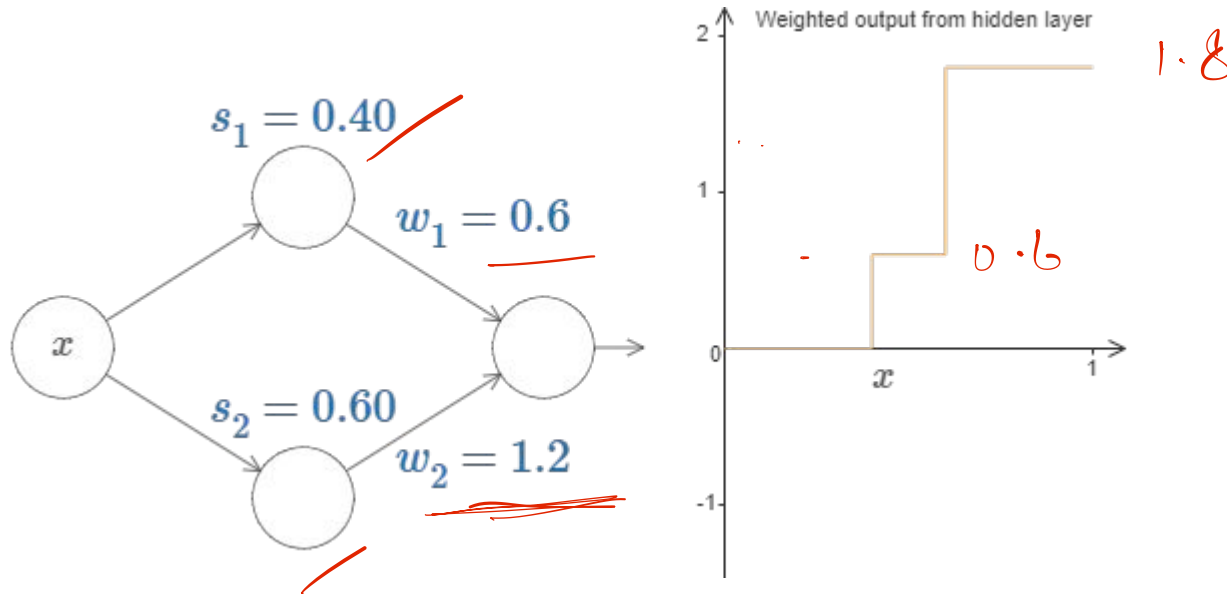
$$-b/w =$$



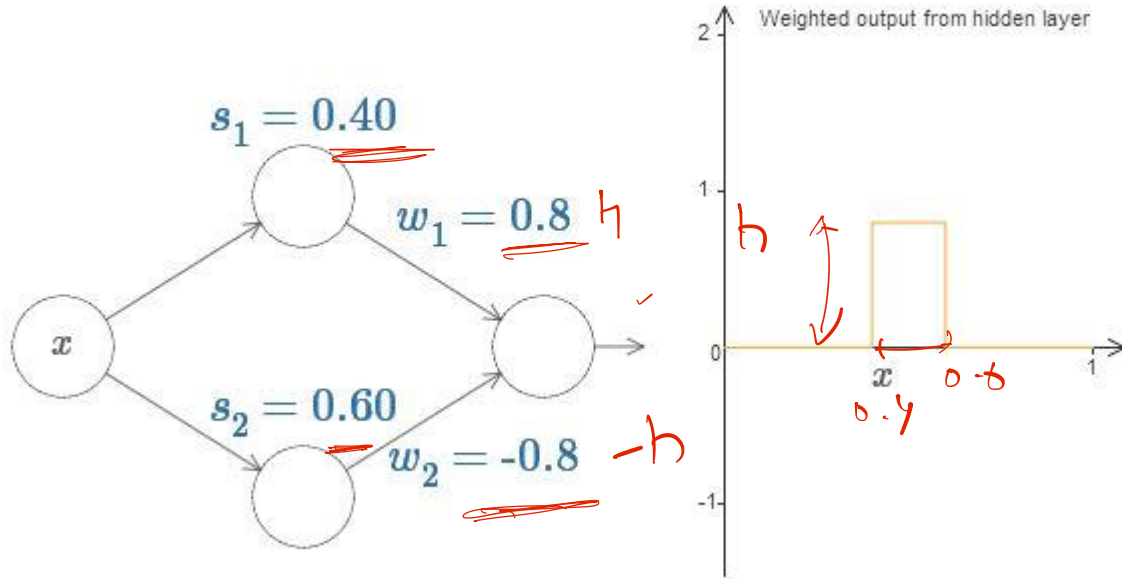
Universality with one i/p & one o/p



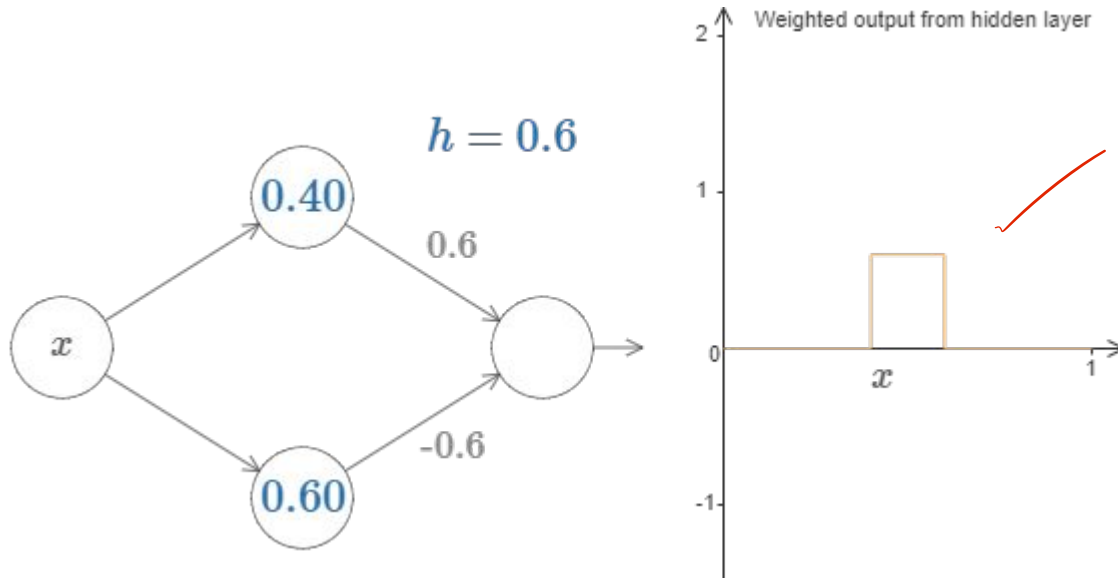
Universality with one i/p & one o/p



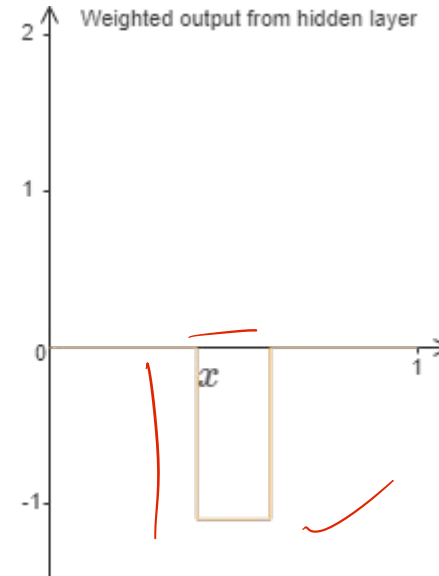
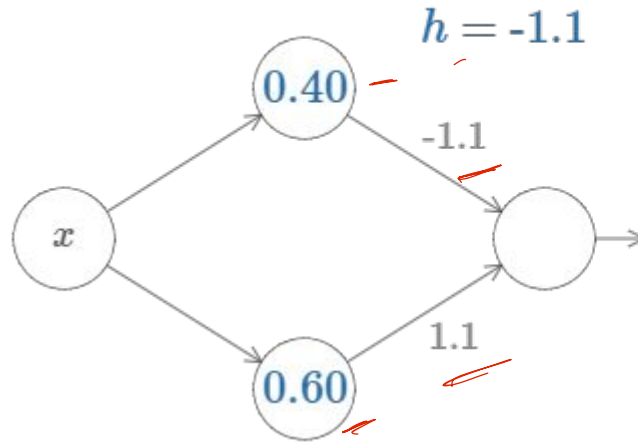
Universality with one i/p & one o/p



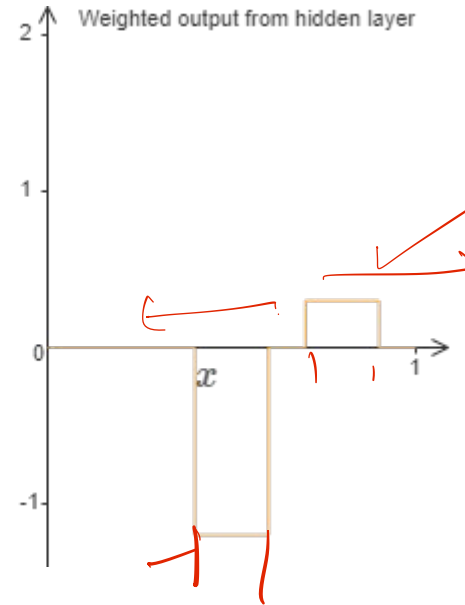
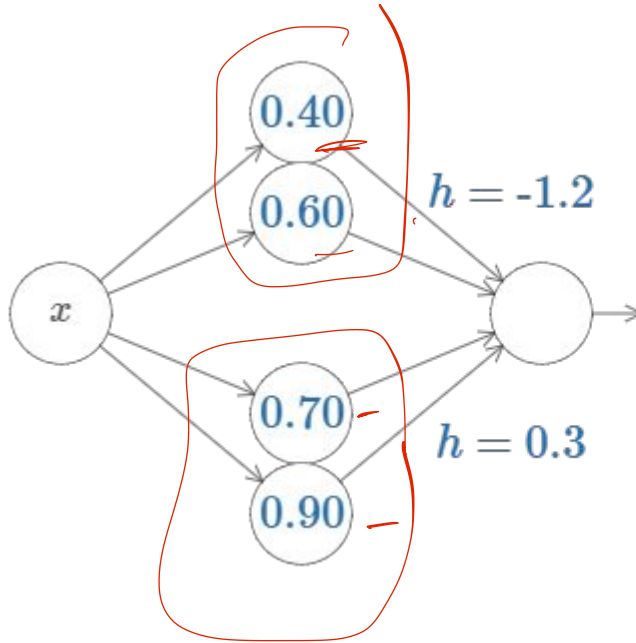
Universality with one i/p & one o/p



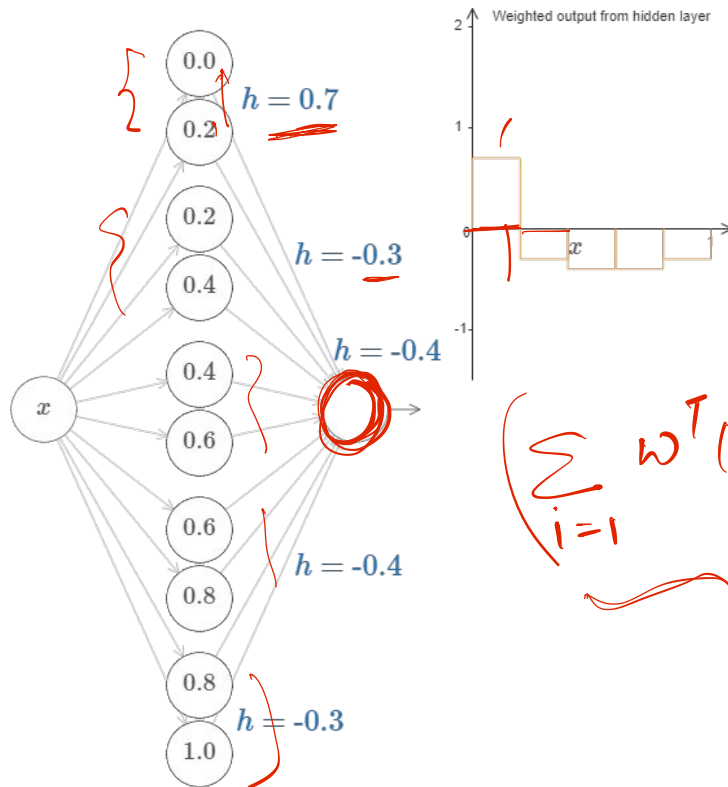
Universality with one i/p & one o/p



Universality with one i/p & one o/p



Universality with one i/p & one o/p

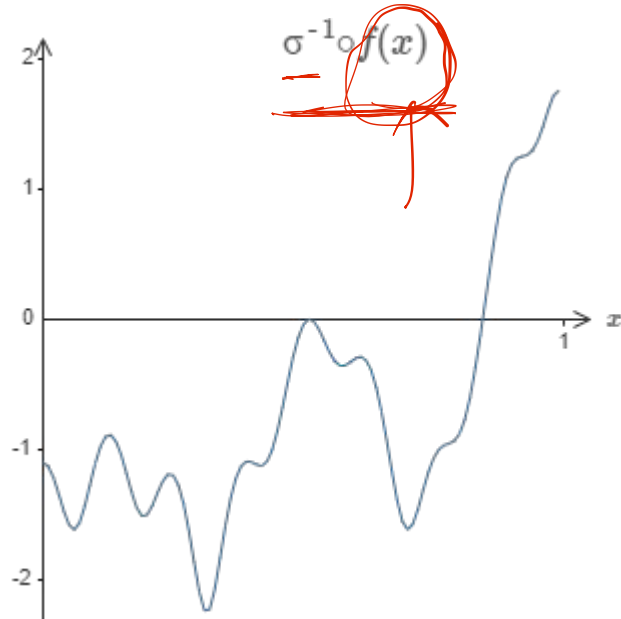


$$\left(\sum_{i=1} w^T(p_i) \right)$$

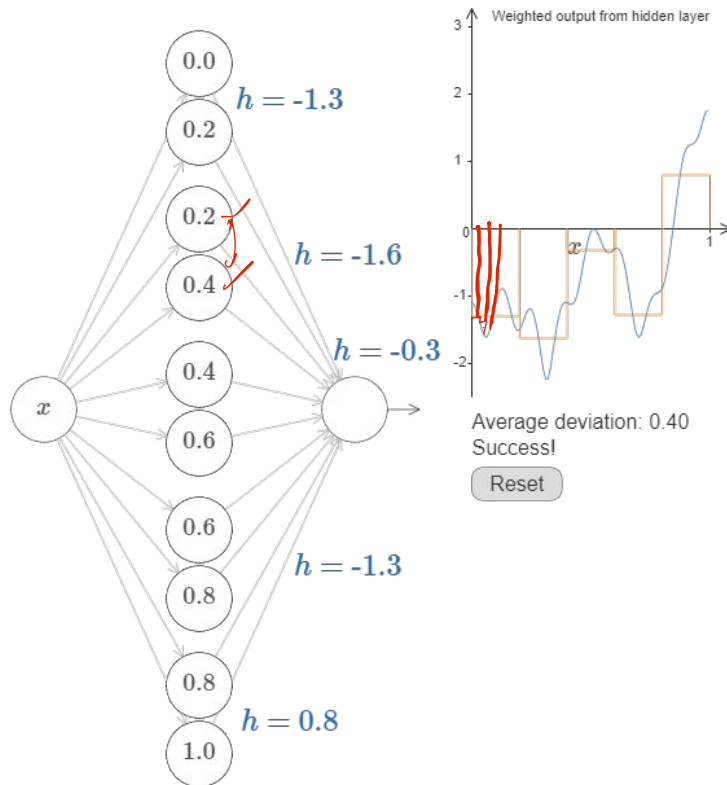
$$= \sigma \left(\frac{-1}{4} \right)$$



Universality with one i/p & one o/p



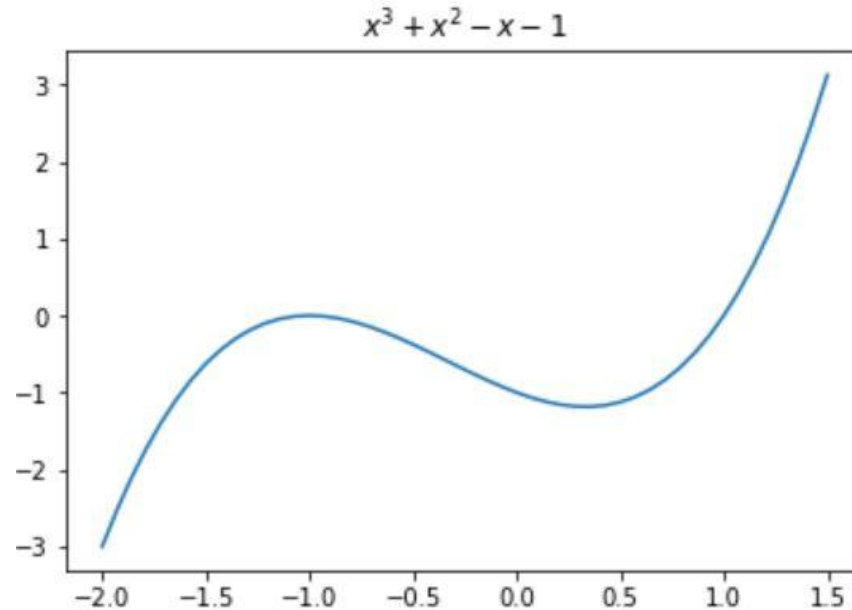
Universality with one i/p & one o/p



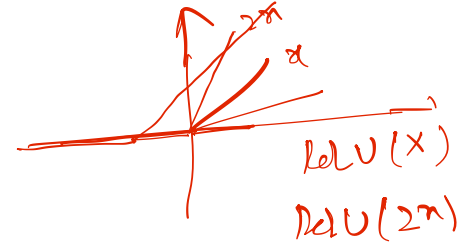
With ReLU activation



Universality with one i/p & one o/p



Universality with one i/p & one o/p



$$n_1 = \text{ReLU}(-5x - 7.7)$$

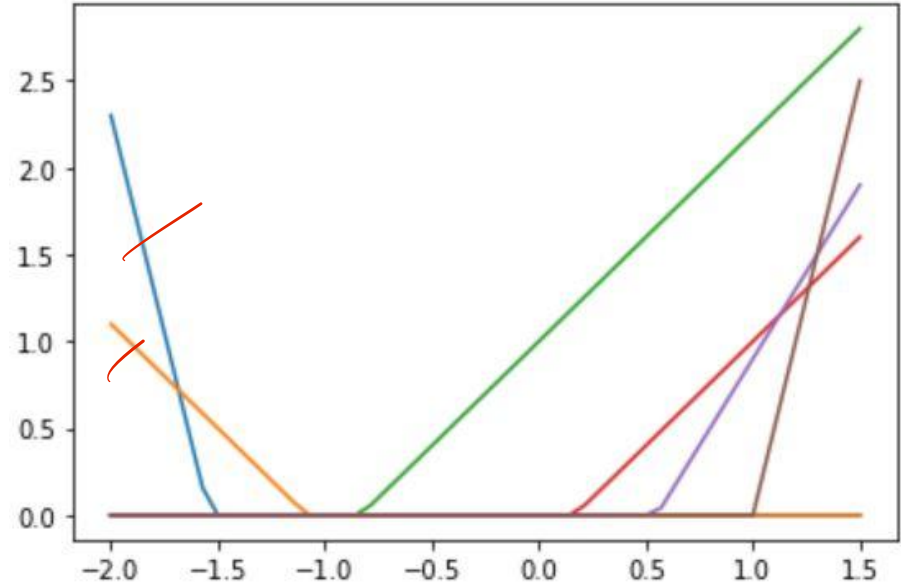
$$n_2 = \text{ReLU}(-1.2x - 1.3)$$

$$n_3 = \text{ReLU}(1.2x + 1)$$

$$n_4 = \text{ReLU}(1.2x - 0.2)$$

$$n_5 = \text{ReLU}(2x - 1.1)$$

$$n_6 = \text{ReLU}(5x - 5)$$



Universality with one i/p & one o/p

width ✓
depth ✓

complexity } sophistication
capacity }

(wid) ^{depth} approximation using ReLUs

$$n_1 = \text{ReLU}(-5x - 7.7)$$

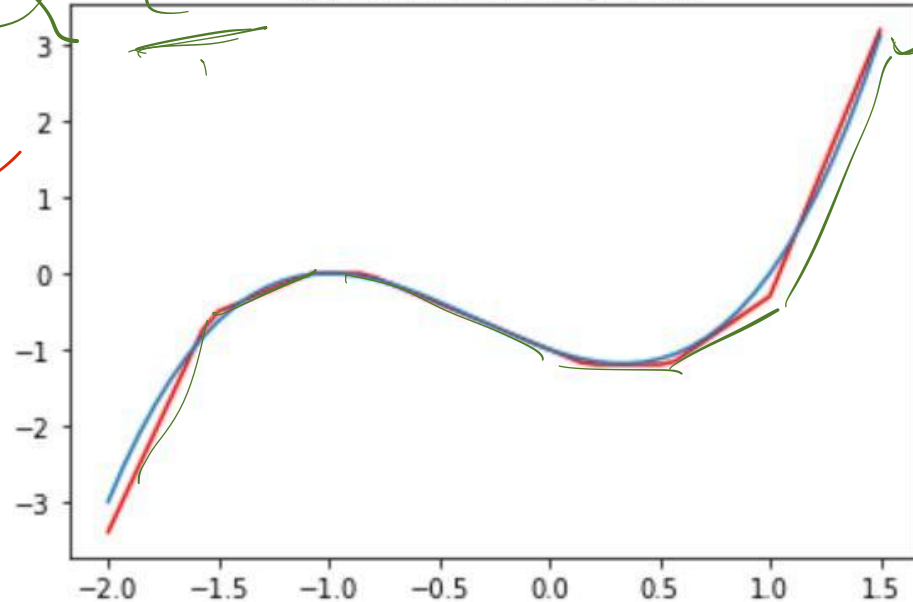
$$n_2 = \text{ReLU}(-1.2x - 1.3)$$

$$n_3 = \text{ReLU}(1.2x + 1)$$

$$n_4 = \text{ReLU}(1.2x - 0.2)$$

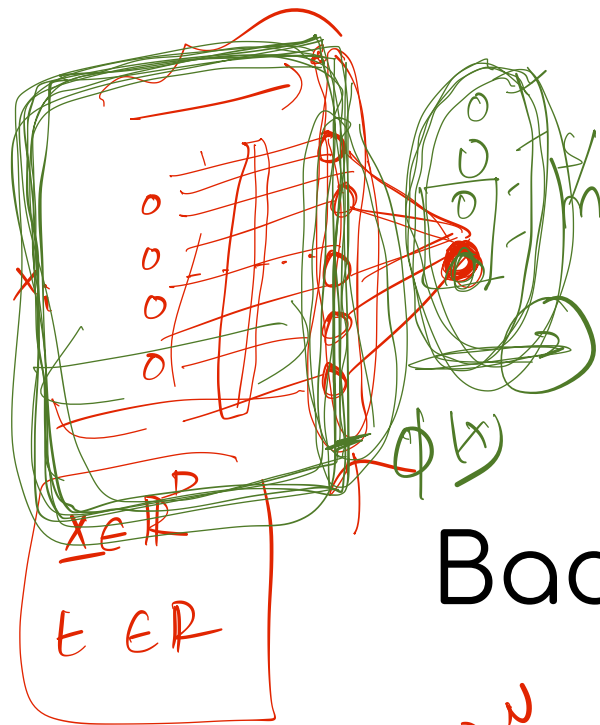
$$n_5 = \text{ReLU}(2x - 1.1)$$

$$n_6 = \text{ReLU}(5x - 5)$$



NAS
D
W
A





$$\begin{aligned}
 &P(y/x) \\
 &P(\omega/x) \\
 &\vdots \\
 &P(L_y/x)
 \end{aligned}$$

$$\begin{aligned}
 x &\in \mathbb{R}^D \quad t = \{0, 1\} \\
 D &= \{ (x_i, y_i) \}_{i=1}^N
 \end{aligned}$$

$$\begin{aligned}
 h_i (1-h_i) \sigma(\omega^T x) &= P(y/x) \\
 1 - \sigma(\omega^T x) &= P(\omega/x)
 \end{aligned}$$

Next Backpropagation

k -class classification

$$D = \{ (x_i, h_i) \}_{i=1}^N$$

$$\begin{aligned}
 P(\underline{y}_k/x, w) &= \text{softmax}(\omega^T \phi(x)) \\
 &= \frac{e^{\omega_k^T \phi(x)}}{\sum_{i=1}^K e^{\omega_i^T \phi(x)}}
 \end{aligned}$$

