

Training Neural Networks

SS 2020 – Machine Perception
Otmar Hilliges

12 March 2020

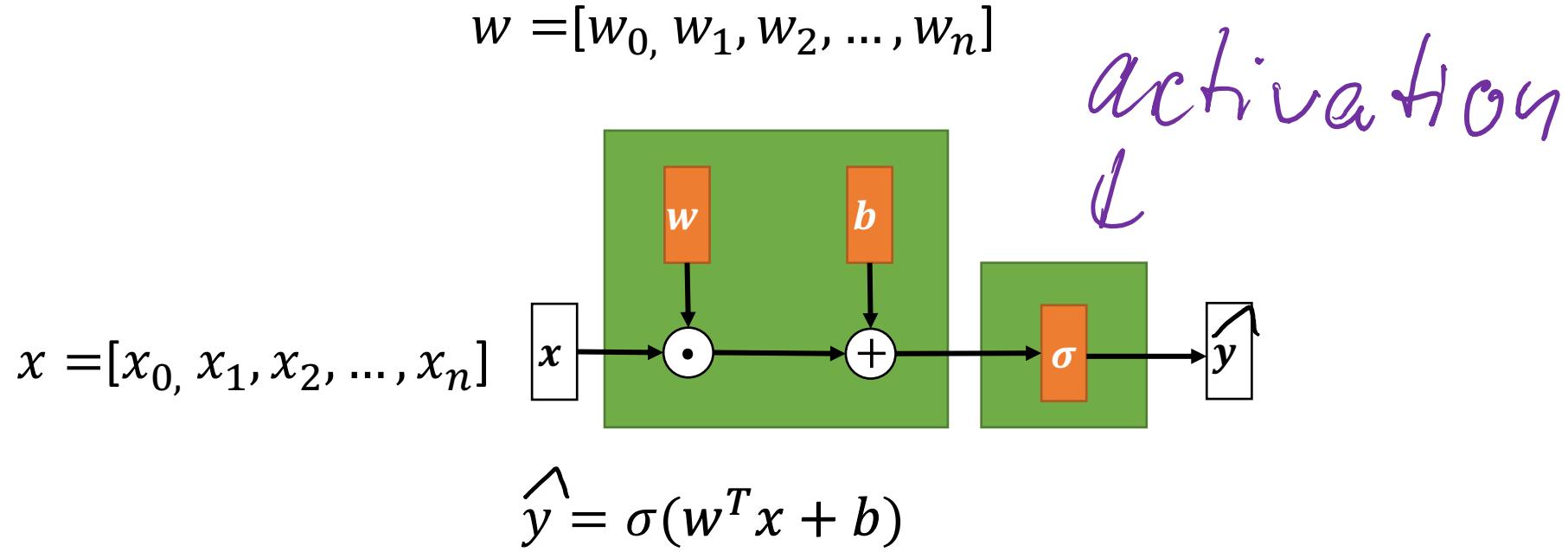
Last lecture

- Perceptron learning algorithm
- MLP as engineering model of a neural network

This lecture

- What types of functions can be approximated by neural networks?
 - Universal approximation theorem
- How do we train neural networks?
 - Backprop algorithm

Perceptron - Block diagram



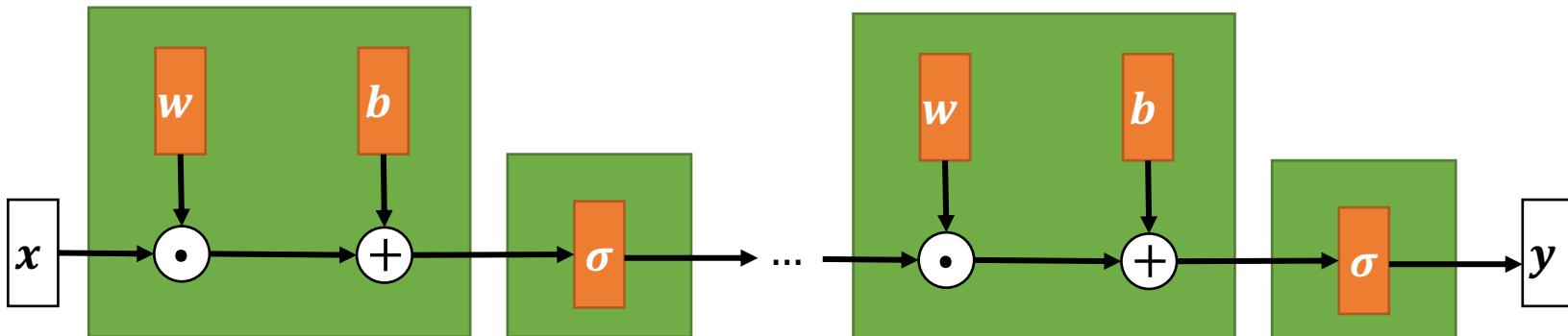
Multi-layered Perceptron - Block diagram

We can combine several layers:

With $x^{(0)} = x$,

$$\forall l = 1, \dots, L, \quad \underline{x}^{(l)} = \sigma \left(\underline{w}^{(l)T} \underline{x}^{(l-1)} + b^{(l)} \right)$$

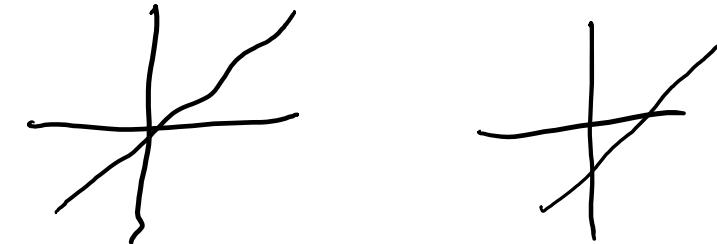
And $f(x; w, b) = \underline{x}^{(L)}$



Can we get rid of activation function?

or use other simple function for that?

Linear activation functions?



Assuming that σ is a *linear* transform,

$$\forall \mathbf{x} \in \mathbb{R}^n, \sigma(\mathbf{x}) = \alpha \mathbf{x} + \beta I$$

with $\alpha, \beta \in \mathbb{R}$, we get:

$$\forall l = 1, \dots, L, \quad \mathbf{x}^{(l)} = \alpha \mathbf{w}^{(l)T} \mathbf{x}^{(l-1)} + \alpha b^{(l)} + \beta I,$$

which results in an *affine* mapping:

$$f(\mathbf{x}; \mathbf{w}, b) = \underline{\mathbf{A}^{(L)} \mathbf{x}} + \underline{\mathbf{B}^{(L)}},$$

where $\mathbf{A}^{(0)} = I, \mathbf{B}^{(0)} = 0$ and

$$\forall l < L \begin{cases} \mathbf{A}^{(l)} = \alpha \mathbf{w}^{(l)T} \mathbf{A}^{(l-1)} \\ \mathbf{B}^{(l)} = \alpha \mathbf{w}^{(l)T} \mathbf{B}^{(l-1)} + \alpha b^{(l)} + \beta I \end{cases}$$

Important message: The activation function should be non-linear, or the resulting MLP is an affine mapping with a peculiar parametrization!

Linear activation functions?

Assuming that σ is a linear transform,

$$\forall x \in \mathbb{R}^n, \sigma(x) = \alpha x + \beta I$$

↓
weight
↓
bias

$$\alpha(w^{(l)\top}x^{(l-1)} + b^{(l)}) + \beta I$$

with $\alpha, \beta \in \mathbb{R}$, we get:

$$\forall l = 1, \dots, L, \quad x^{(l)} = \underline{\alpha w^{(l)\top}} x^{(l-1)} + \underline{\alpha b^{(l)}} + \beta I,$$

which results in an affine mapping:

$$f(x; w, b) = A^{(L)}x + B^{(L)},$$

↓ new weights
← new bias

where $A^{(0)} = I, B^{(0)} = 0$ and

这样 $x^{(0)} = x$

$$\forall l < L \begin{cases} A^{(l)} = \alpha w^{(l)} A^{(l-1)} \\ B^{(l)} = \alpha w^{(l)} B^{(l-1)} + \alpha b^{(l)} + \beta I \end{cases}$$

this does not make sense
to have multi-layers
have one layer is enough.

Important message: The activation function should be non-linear, or the resulting MLP is an affine mapping with a peculiar parametrization!

if the activ. func.
is linear, it does not
matter how many layer
you connected.

It just a linear func.
at last.

$$\text{上页: } \sigma(x) = \alpha x + \beta I$$

$$x^L = \alpha(w^L x^{L-1} + b^L) + \beta I$$

$$= \underline{\alpha w^L x^{L-1}} + \underline{\alpha b^L + \beta I} = \underline{A^L x^{L-1}} + \underline{B^L}$$

$$\begin{aligned}\therefore x^{L+1} &= \alpha(w^{L+1} x^L + b^{L+1}) + \beta I \\&= \alpha(w^{L+1} \boxed{A^L x^{L-1} + B^L} + b^{L+1}) + \beta I \\&= \frac{\alpha w^{L+1} A^L \cdot x^{L-1}}{A^{L+1}} + \frac{\alpha w^{L+1} B^L + \alpha b^{L+1} + \beta I}{B^{L+1}} \\&\therefore A^L = \alpha w^L A^{L-1} \quad B^L = \alpha w^L B^{L-1} + \alpha b^L + \beta I\end{aligned}$$

So what happens with a non-linearity?

Example: we have non-linearity

and

we don't have non-linearity

Example: Solving 'XOR'

Task learn function $y = f^*(x)$ that maps binary variables x_1, x_2 to “true” or “false” if one and only one $x_i == 1$

We will use function $f(x; \Theta)$ and find parameters Θ to make $f^* \approx f$

$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}, \quad Y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

异或门（英語：Exclusive-OR gate，簡稱XOR gate，又稱EOR gate、ExOR gate）是数字逻辑中实现逻辑异或的逻辑门

输入	输出	
A	B	A XOR B
0	0	0
0	1	1
1	0	1
1	1	0

Attempt I: Use a linear model?

Assuming a linear function: *meas square error*

$$J(\Theta) = \frac{1}{(N=4)} \sum_{x \in X} (f^*(x^{(i)}) - f(x^{(i)}; \Theta))^2 \quad \text{with } f(x; w, b) = x^T w + b$$

Solving via normal equation:

$$w = 0 \text{ and } b = \frac{1}{2}$$

Normal Equation

Given a matrix equation

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

the normal equation is that which minimizes the sum of the square differences between the left and right sides:

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}.$$

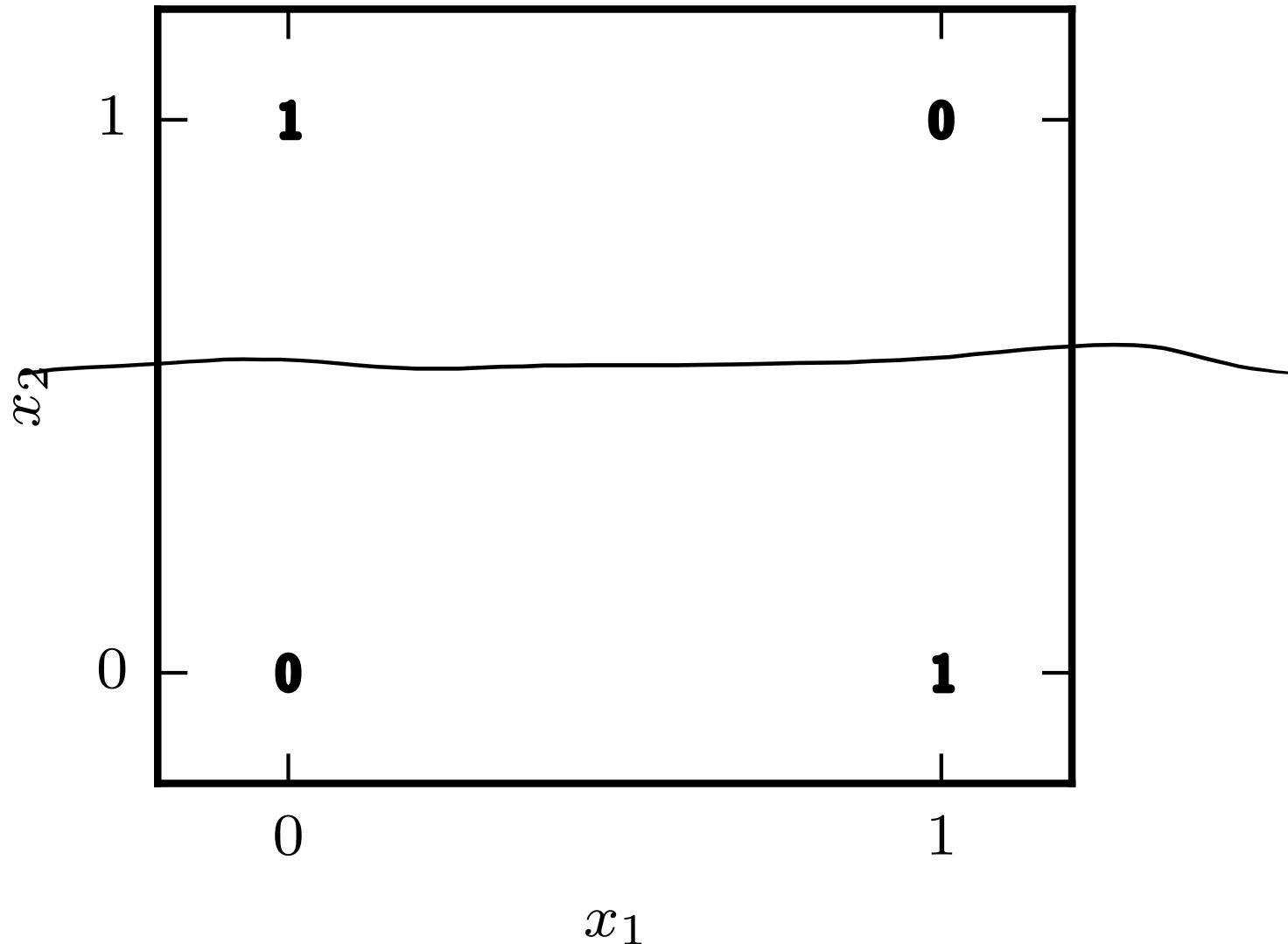
It is called a normal equation because $\mathbf{b} - \mathbf{A}\mathbf{x}$ is normal to the range of \mathbf{A} .

Here, $\mathbf{A}^T \mathbf{A}$ is a **normal matrix**.

Visualization

bold output one the plane is the value of y .

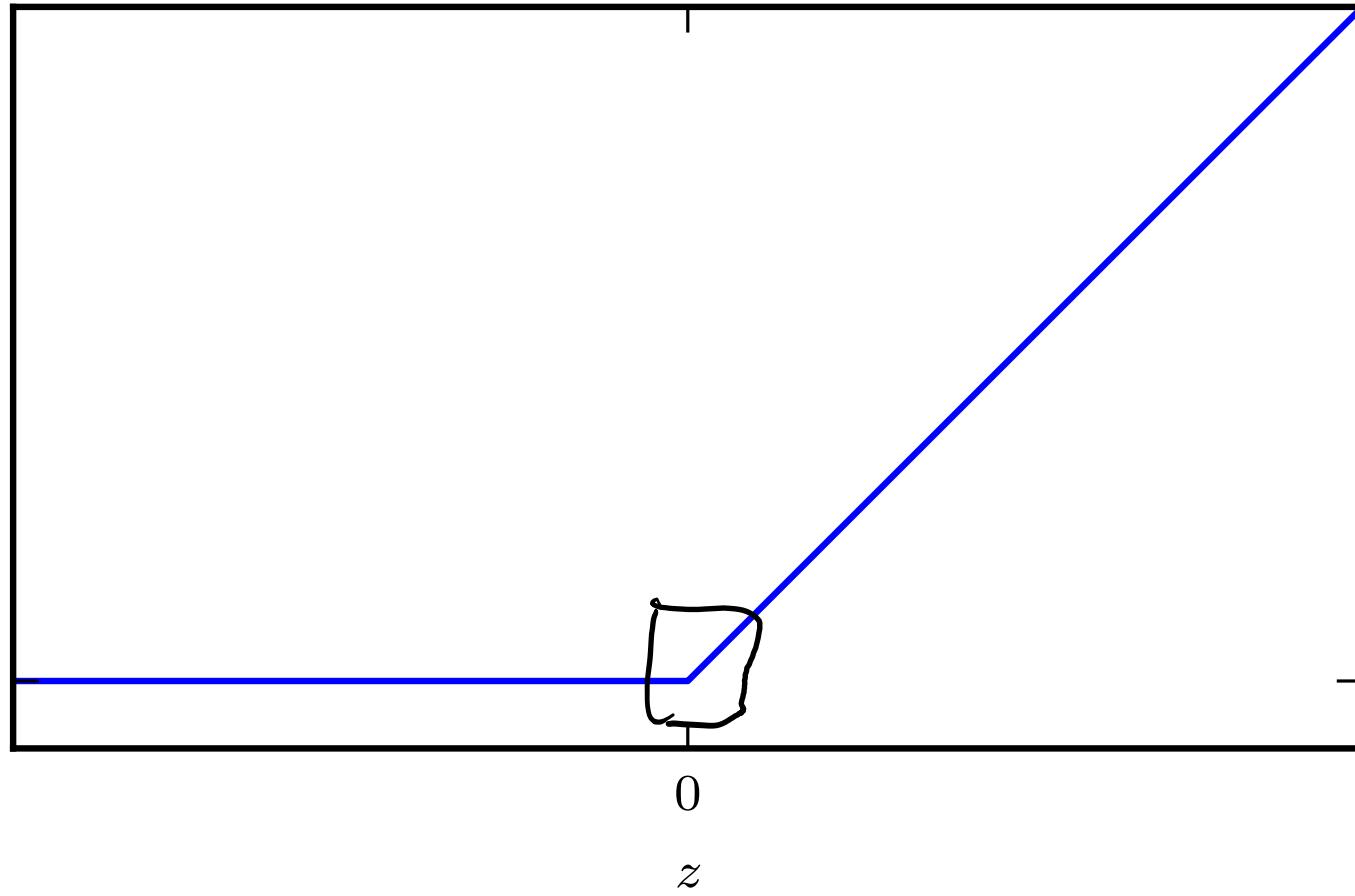
Original \boldsymbol{x} space



ReLU

(non-linear)

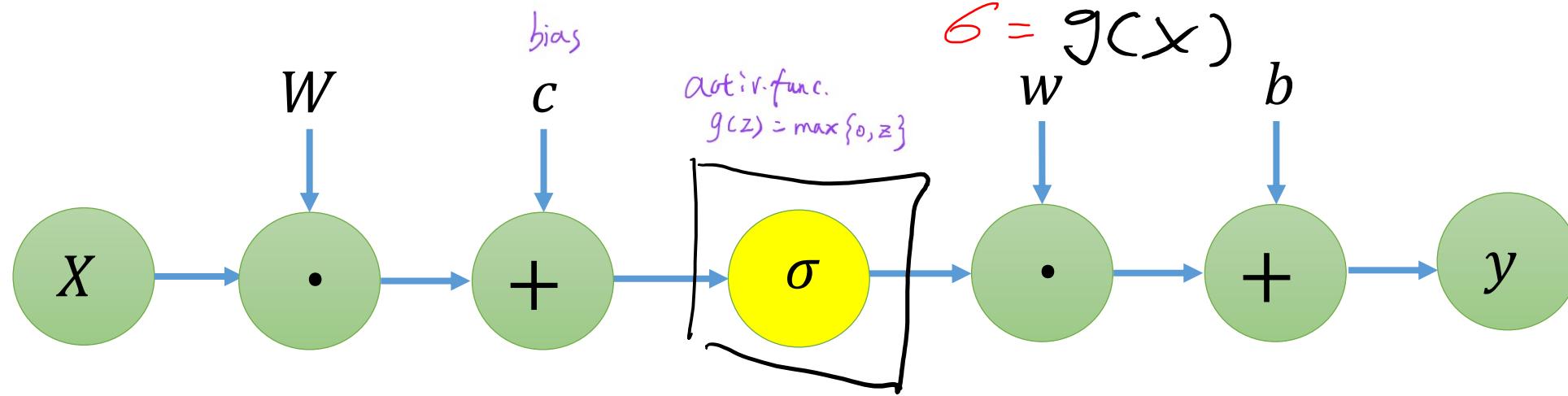
$$g(z) = \max\{0, z\}$$



$$g(z) = \max\{0, z\}$$

Feedforward Neural Network

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \underbrace{\mathbf{w}^T \max\{0, \mathbf{X}\mathbf{W} + \mathbf{c}\}}_{f(x)} + b$$



The XOR Multi-Layered Perceptron

$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}, \quad Y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

Function we want to learn:

$$f(x; W, c, w, b) = w^T \max\{0, XW + c\} + b$$

Oracle solution:

$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \quad b = 0$$

$$XW = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

The XOR Multi-Layered Perceptron

$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}, \quad Y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

Function we want to learn:

$$f(x; W, c, w, b) = w^T \max\{0, XW + c\} + b$$

Oracle solution:

$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \quad b = 0$$

$$XW + c = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

The XOR Multi-Layered Perceptron

$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}, \quad Y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

Function we want to learn:

$$f(x; W, c, w, b) = w^T \max\{0, XW + c\} + b$$

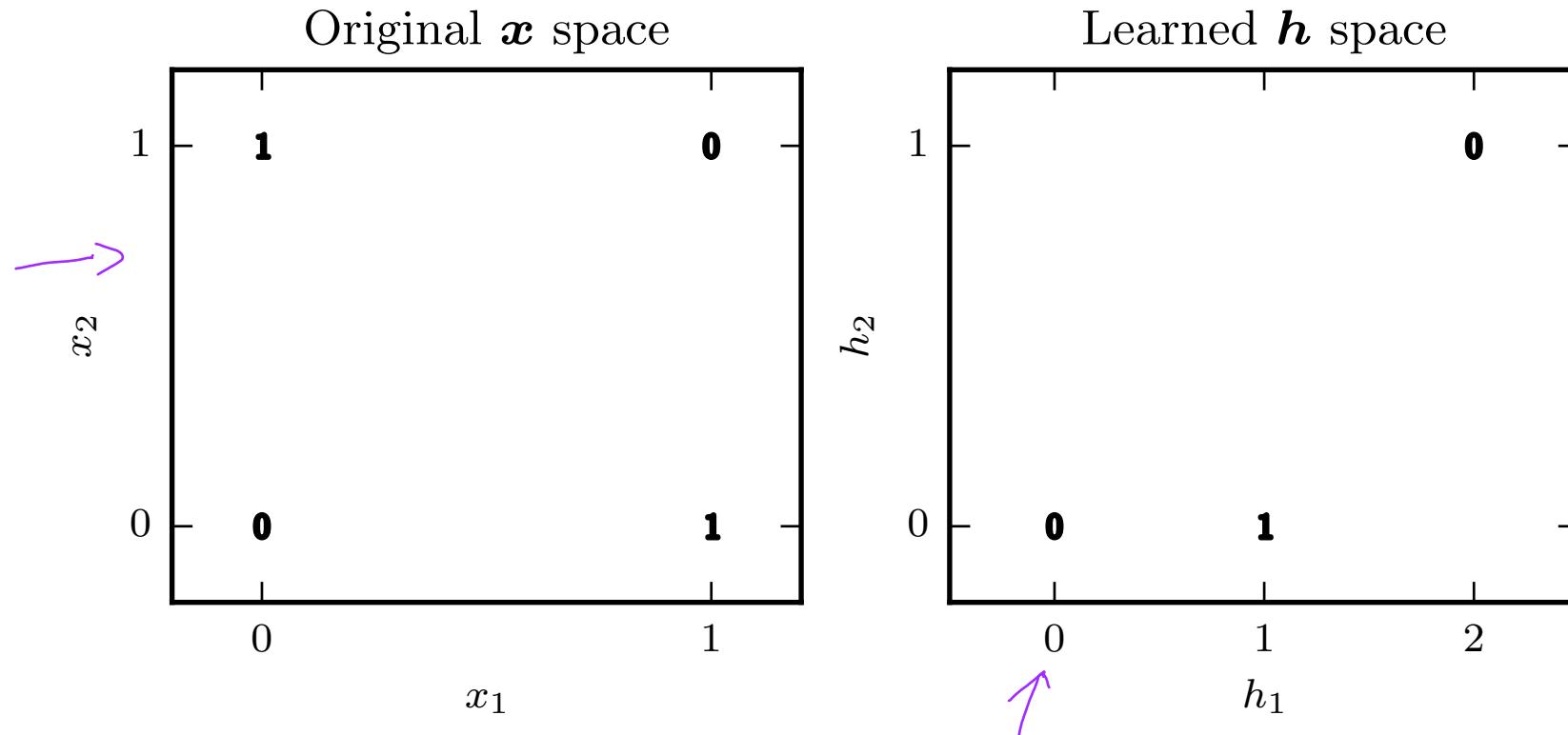
Oracle solution:

$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \quad b = 0$$

$$\max\{0, XW + c\} = \max\{0, \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}\} \xrightarrow{\text{max } \{0, -1\}} \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

$h_1 \quad h_2$

$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$



$$\max\{0, XW + c\} = \max\{0, \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}\} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ -2 & 1 \end{bmatrix}$$

h_1 h_2

The XOR Multi-Layered Perceptron

$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}, \quad Y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

Function we want to learn:

$$f(x; W, c, w, b) = w^T \max\{0, XW + c\} + b$$

Oracle solution:

$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \quad b = 0$$

$$w^T \max\{0, WX + c\} + b = [1 \quad -2] \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} + 0 = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = Y$$

identical

we solve the problem
by introduce the non-
linearity into the
network.

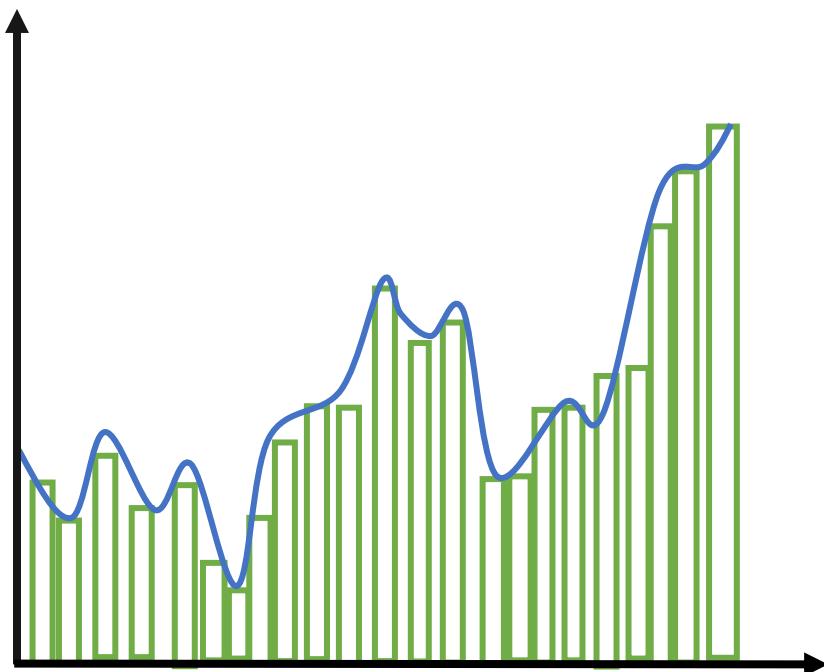
So far: The non-linearity is critically important

Next: What *type* of functions can we approximate?

Universal approximation theorem

Given that $\sigma \in C^\infty(\mathbb{R})$ is non-linear (e.g., sigmoid), what type of functions can we learn?

as long as $f(x)$ continuous function



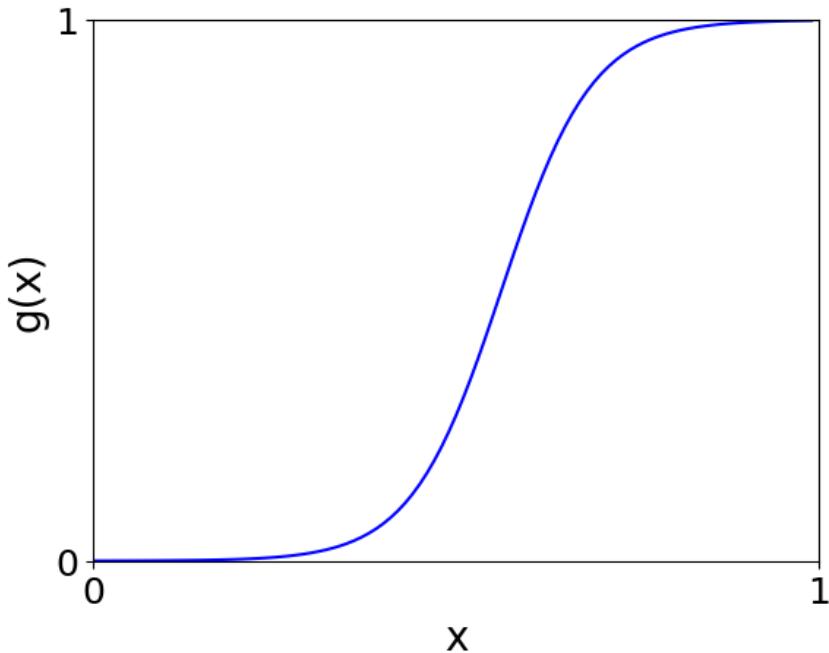
Given enough hidden units.
One layer is enough in theory.
In practice deeper is better.

$\exists g(x)$ as NN, $g(x) \approx f(x)$, and $|g(x) - f(x)| < \epsilon$

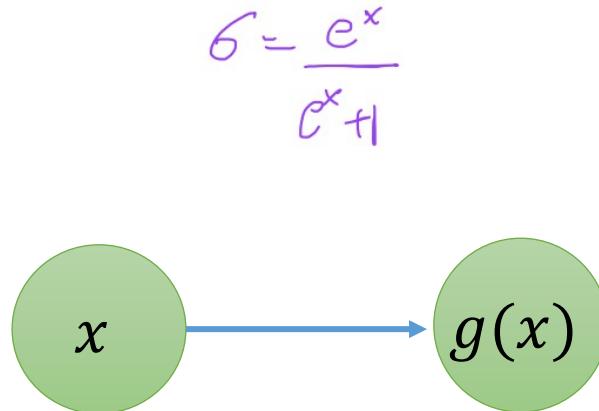
It is an approximation

Universal approximation theorem

$$g(x) = \sigma(w^T x + b) = \frac{\exp(w^T x + b)}{\exp(w^T x + b) + 1} \quad \text{Where } \sigma \text{ is the sigmoid function}$$



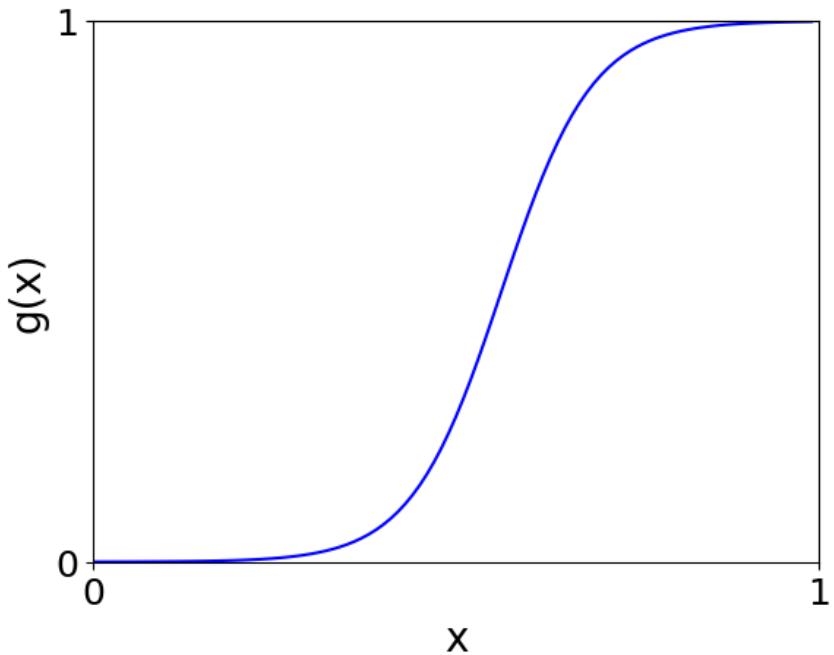
$$w = 16, b = -9$$



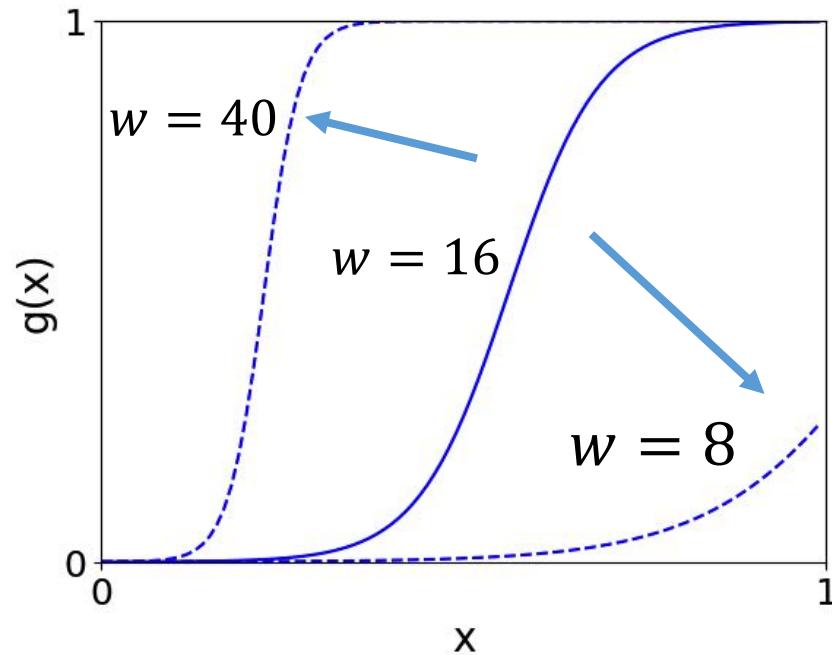
$$\sigma = \frac{e^x}{e^x + 1}$$

Universal approximation theorem

$$g(x) = \sigma(w^T x + b) = \frac{\exp(w^T x + b)}{\exp(w^T x + b) + 1}$$

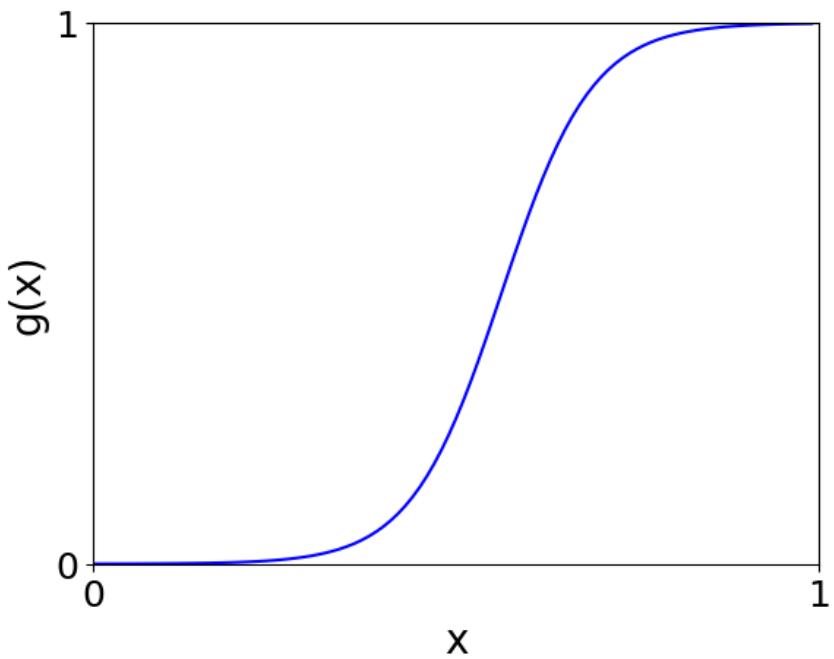


$$w = 16, b = -9$$

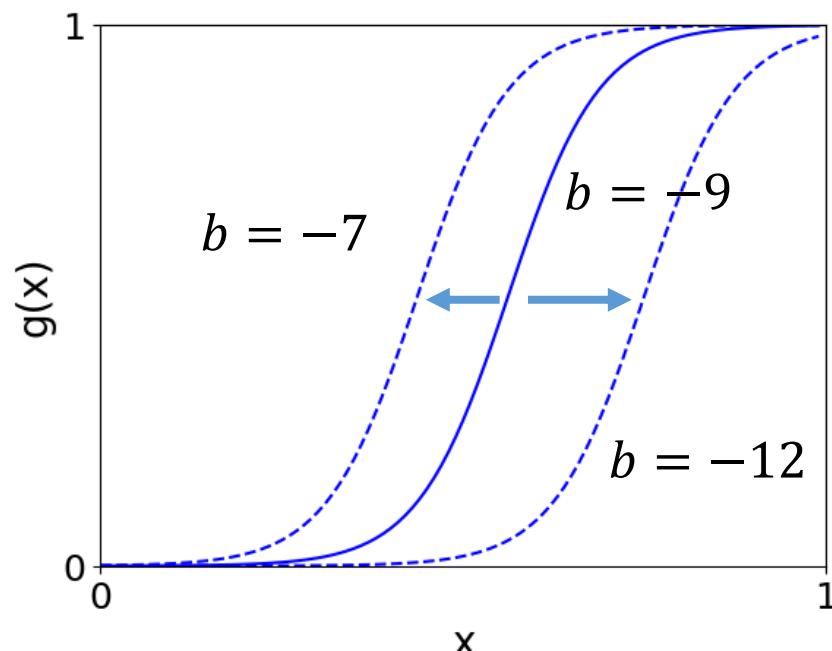


Universal approximation theorem

$$g(x) = \sigma(w^T x + b) = \frac{\exp(w^T x + b)}{\exp(w^T x + b) + 1} \quad \text{Where } \sigma \text{ is the sigmoid function}$$

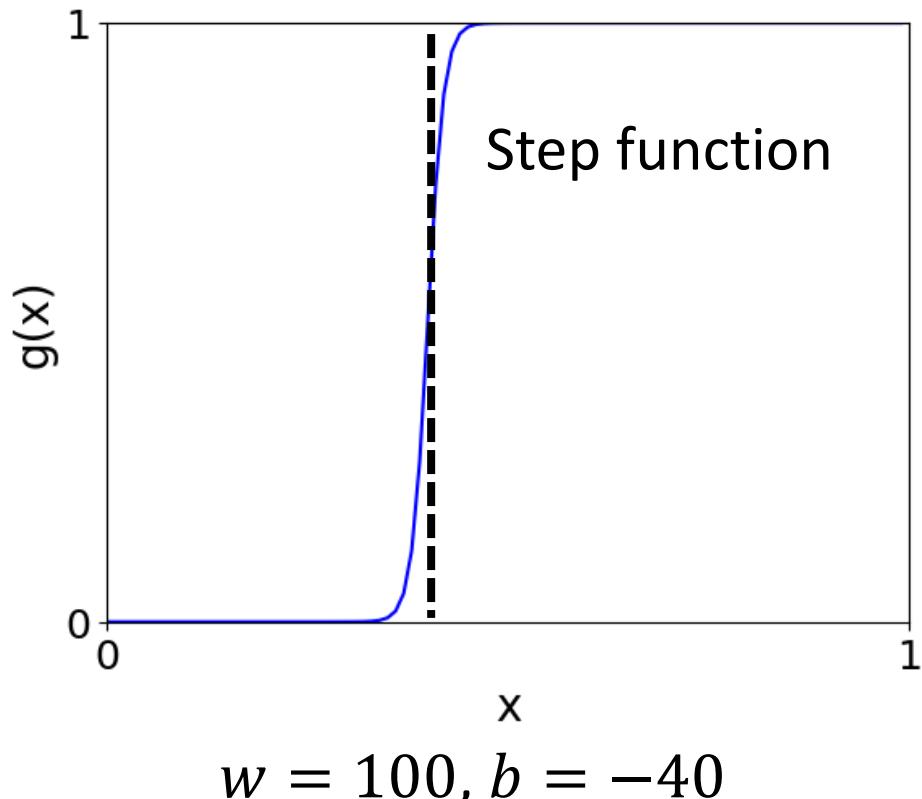


$$w = 16, b = -9$$



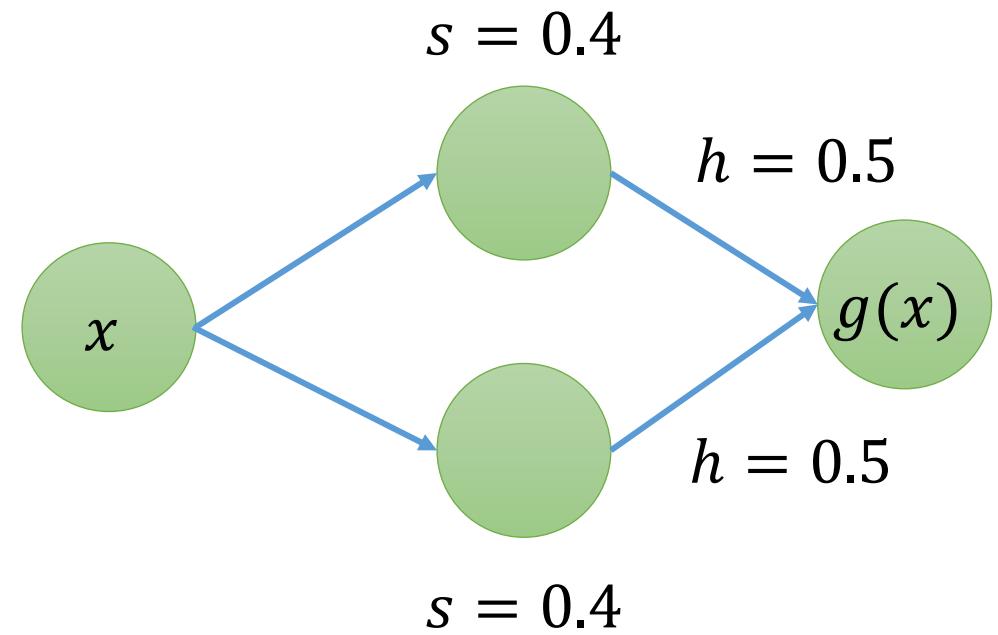
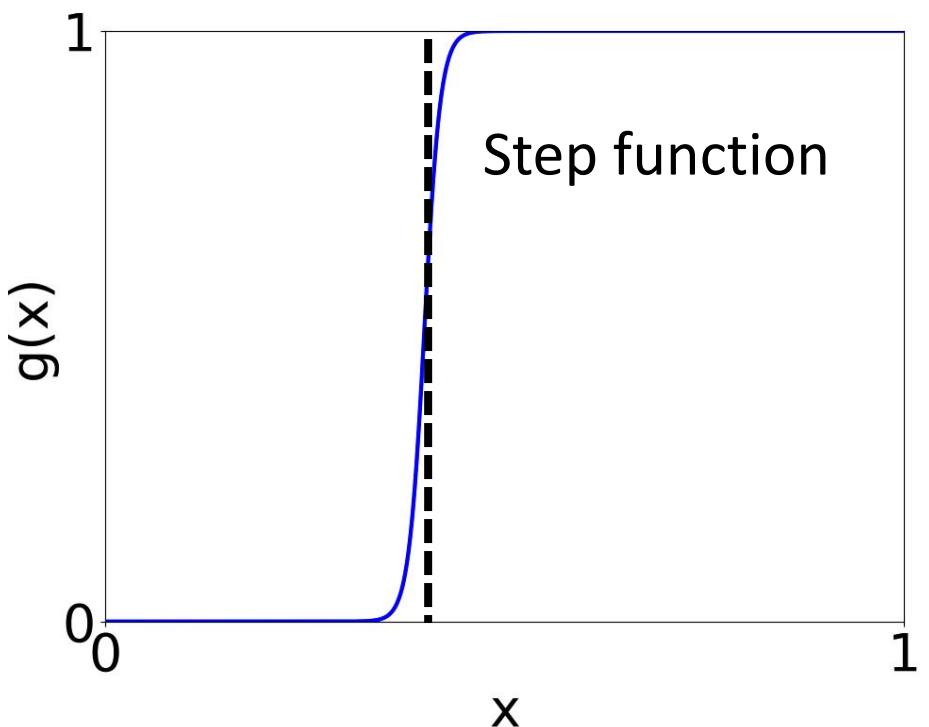
Universal approximation theorem

Leading edge: $s = -\frac{b}{w} = 0.4$

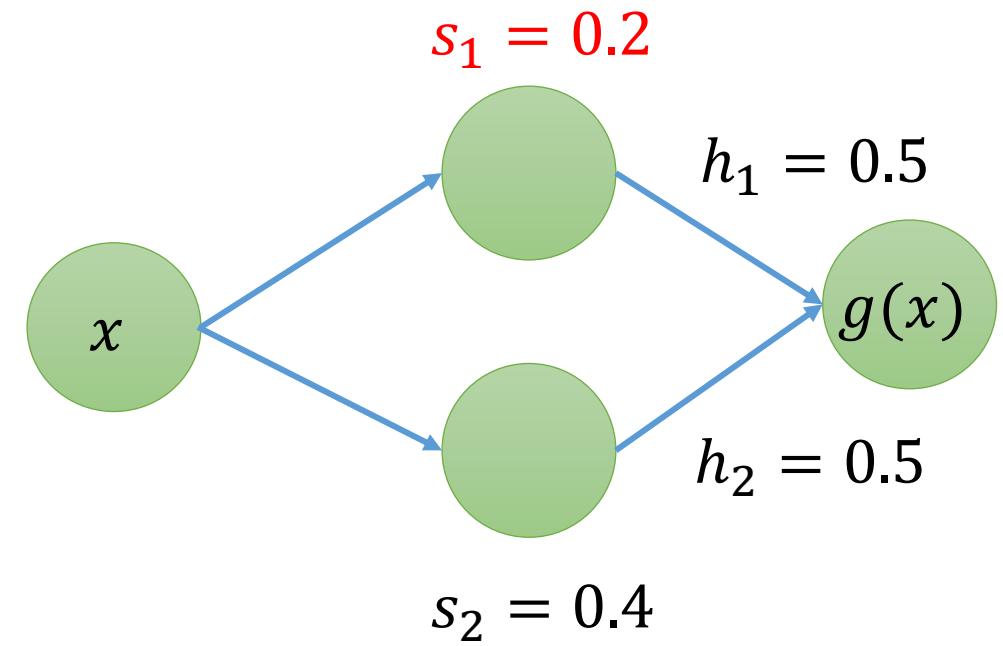
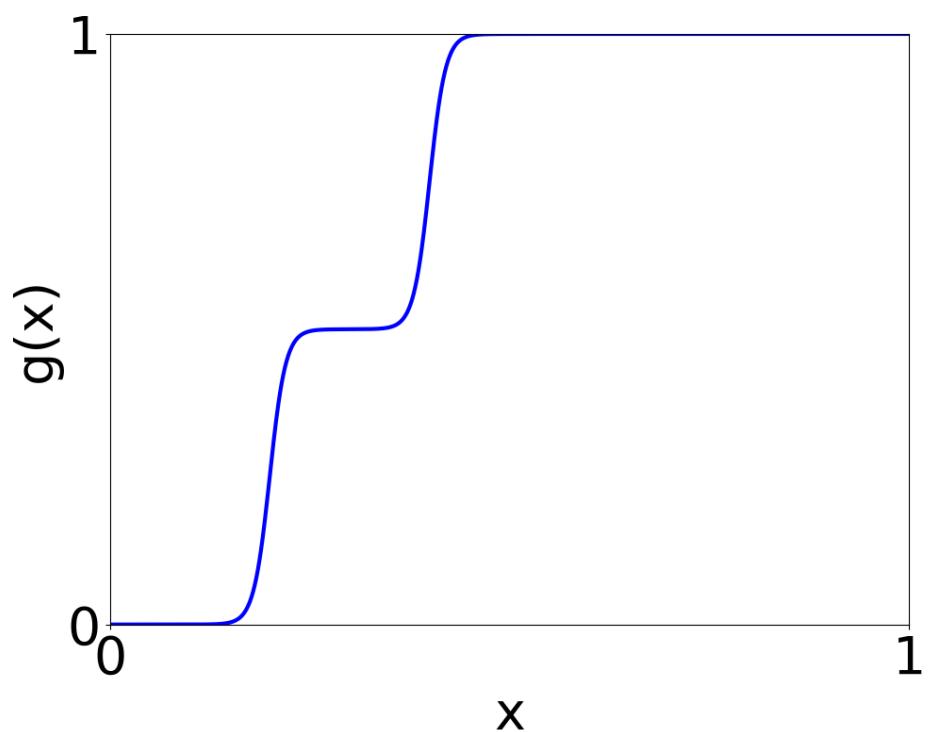


Universal approximation theorem

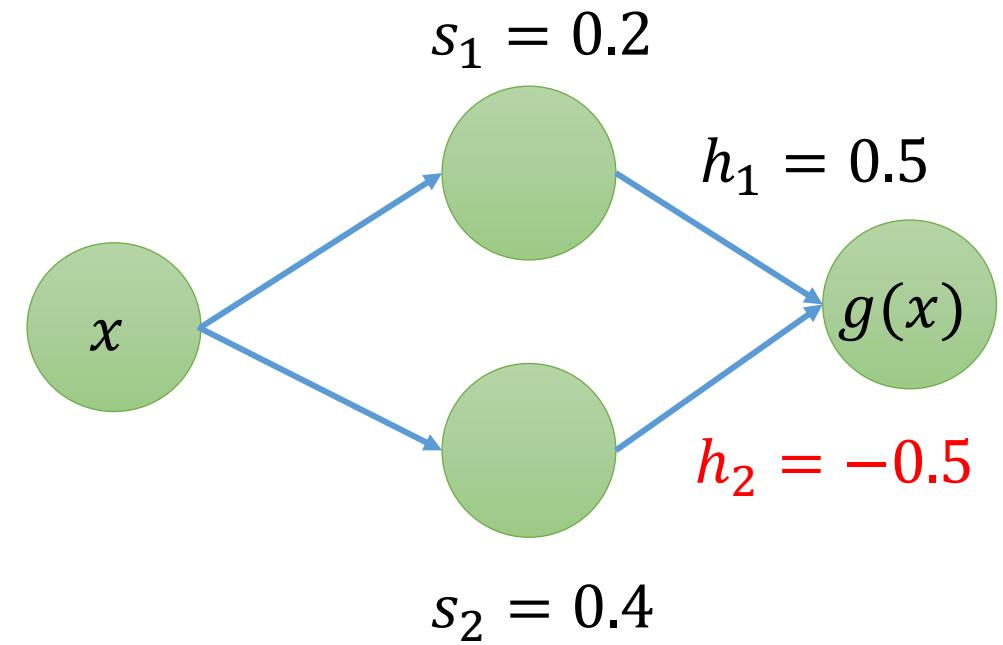
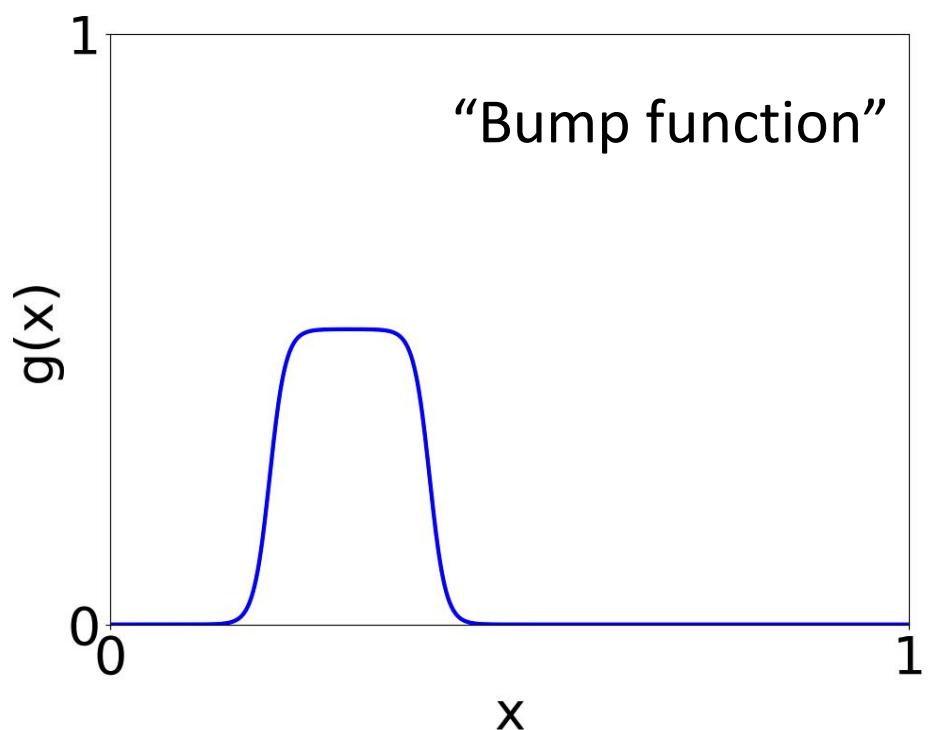
Output layer
is always assumed
linear



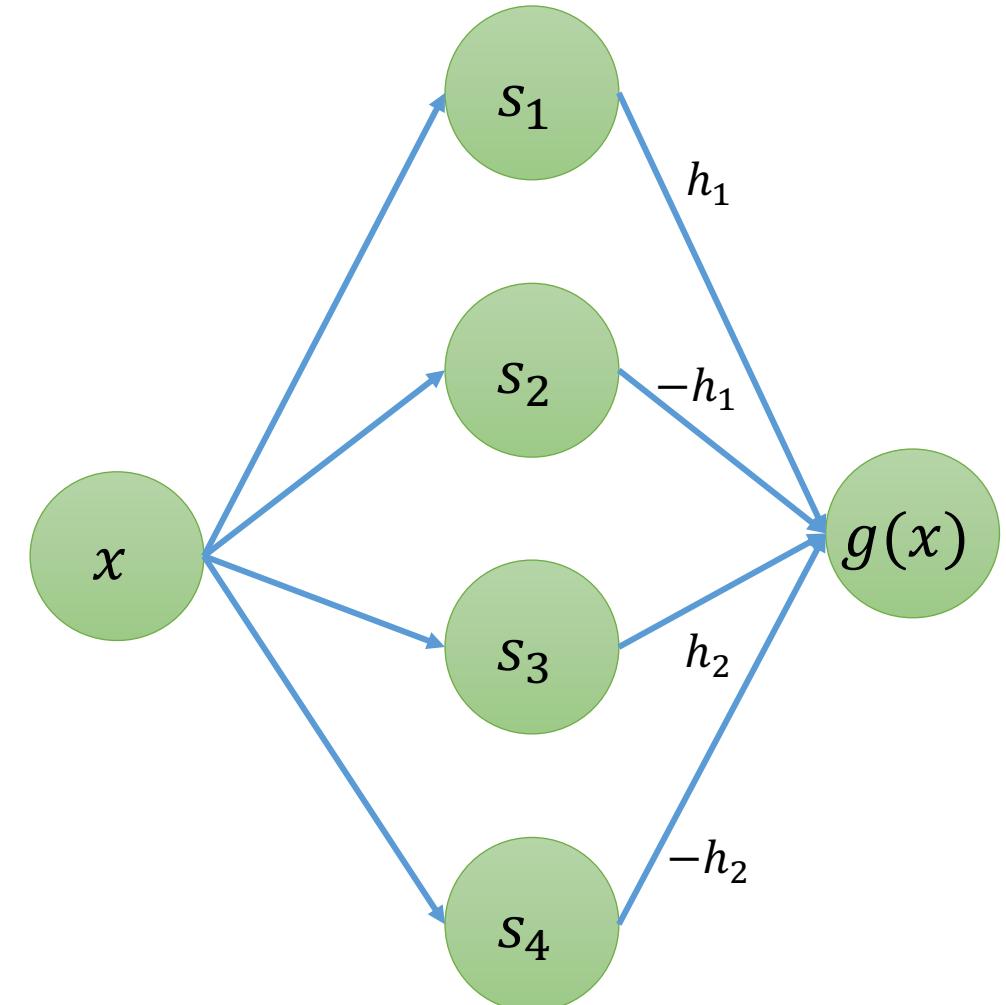
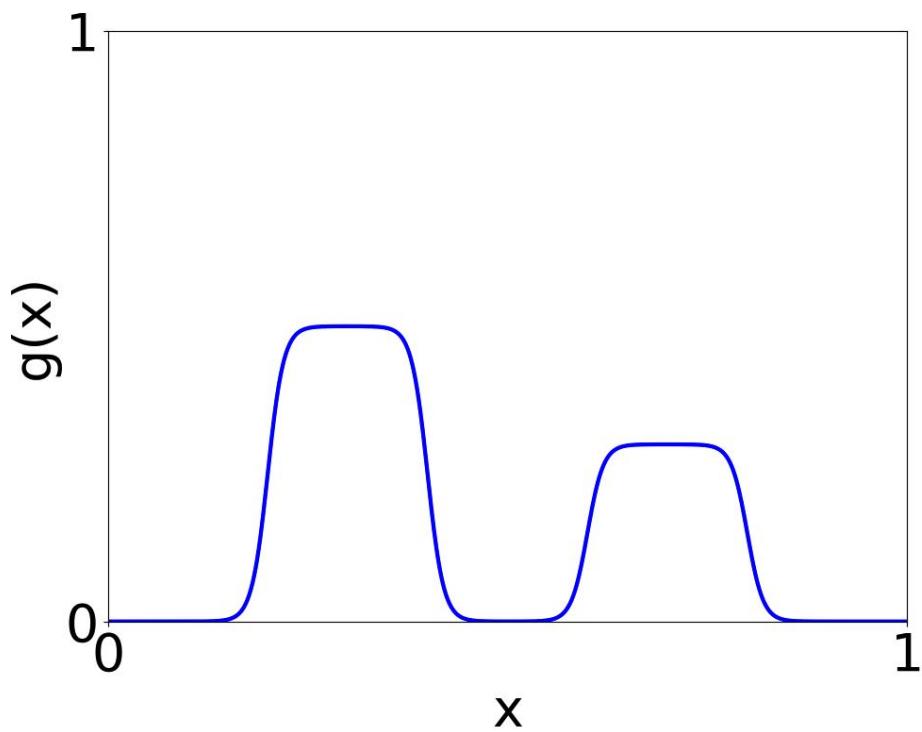
Universal approximation theorem



Universal approximation theorem



Universal approximation theorem

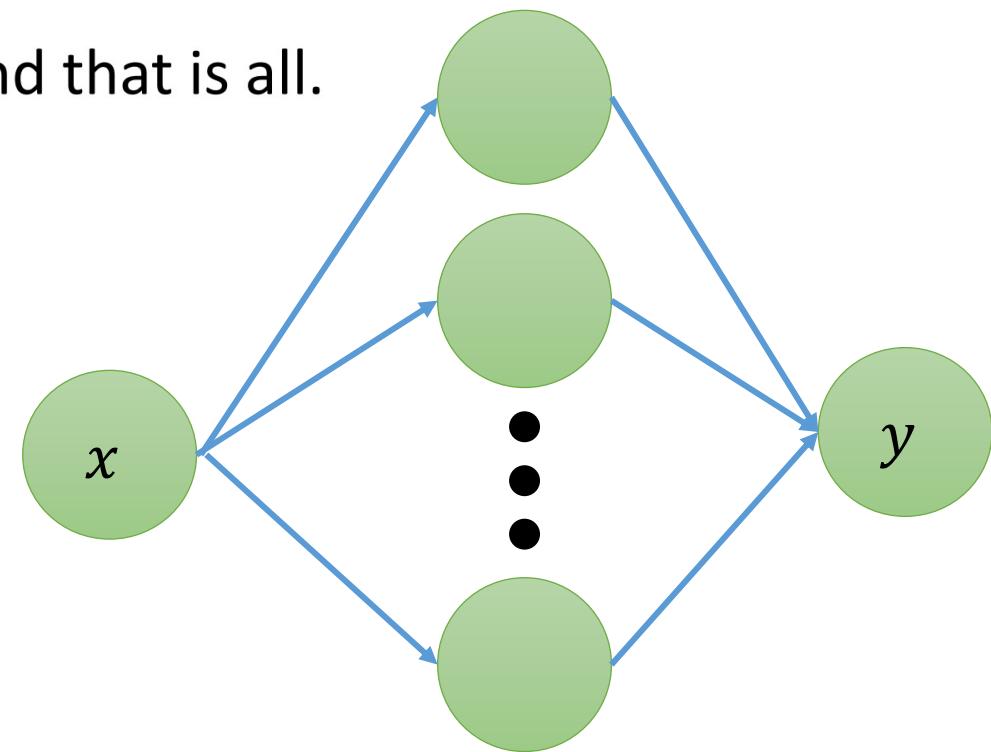
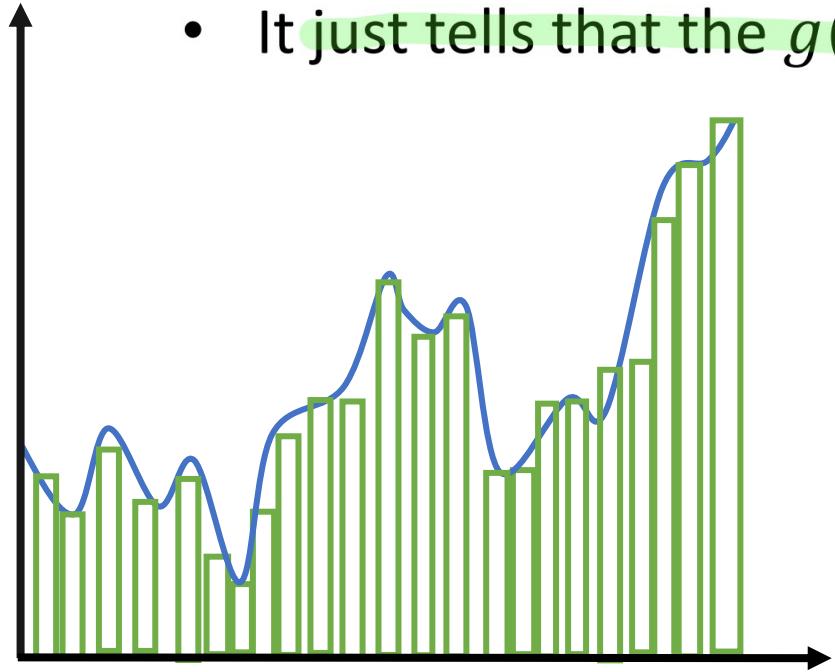


Universal approximation theorem

- Tells us that feed forward networks can approximate $f(x)$
- It *does not* tell us anything about the chances of learning the correct parameters

P₀ as long as f(x) continuous.

- It just tells that the $g(x)$ exists, and that is all.



More precisely...

[Cybenko, G. "[Approximations by superpositions of sigmoidal functions](#)", Mathematics of Control, Signals, and Systems, 1989.]

Let $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ be a non-constant, bounded and continuous (activation) function

I_m denotes the m -dimensional unit hypercube $[0,1]^m$ and the space of real-valued functions on I_m is denoted with $C(I_m)$.

Then any function $f \in C(I_m)$ can be approximated given any $\epsilon > 0$, integer N , real constants $v_i, b_i \in \mathbb{R}$ and real vectors $w_i \in \mathbb{R}^m$ for $i = 1, \dots, N$:

$$f(x) \approx g(x) = \sum_{i=1}^N v_i \sigma(w_i^T x + b_i)$$

And: $|g(x) - f(x)| < \epsilon$, for all $x \in I_m$

However...

Networks with a single-hidden layer, need to have exponential width

In practice, deeper networks work better.

Lot's of ongoing work to provide theory on what network properties lead to which approximation capabilities.

E.g., [Lu et al. "[The Expressive Power of Neural Networks: A View from the Width](#)". NeurIPS 2017]

So far: NNs are universal approximators

Next: How do we find the network parameters?

2019 notes

- Previous: What types of functions can be approximated by neural networks?
any kind of $f(x)$ with continuity
- Next: How do we train neural networks?

General procedure

Iterative gradient descent to find parameters Θ :

Initialize weights with small random values

Initialize biases with 0 or small positive values

Compute gradients

Update parameters with SGD

from experience

SGD in a nutshell:

Compute the negative gradient at θ^0

$$\rightarrow -\nabla C(\theta^0)$$

Times the learning rate η

$$\rightarrow -\eta \nabla C(\theta^0)$$

The nonlinearity of a neural network causes most interesting loss functions to become non-convex

Chain rule

$$y = g(x) \text{ and } z = f(g(x)) = f(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

For vector types

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

Two ways to compute gradients

Consider the scalar function:

$$f = \exp(\exp(x) + \exp(x)^2) + \sin(\exp(x) + \exp(x)^2)$$

Symbolic differentiation gives us:

$$\frac{df}{dx} =$$

$[\exp(x^2)]' = 2\exp(x)$
 $\cdot \exp(x)$

$$\begin{aligned} & \exp(\exp(x) + \exp(x)^2)(\exp(x) + 2\exp(x)^2) \\ & + \cos(\exp(x) + \exp(x)^2)(\exp(x) + 2\exp(x)^2) \end{aligned}$$

$= 2\exp(x)^2$

If we were to write a program...

Consider the scalar function:

$$f = \exp(\exp(x) + \exp(x)^2) + \sin(\exp(x) + \exp(x)^2)$$

We would define and compute intermediate variables:

$$a = \exp(x)$$

$$b = a^2$$

$$c = a + b$$

$$d = \exp(c)$$

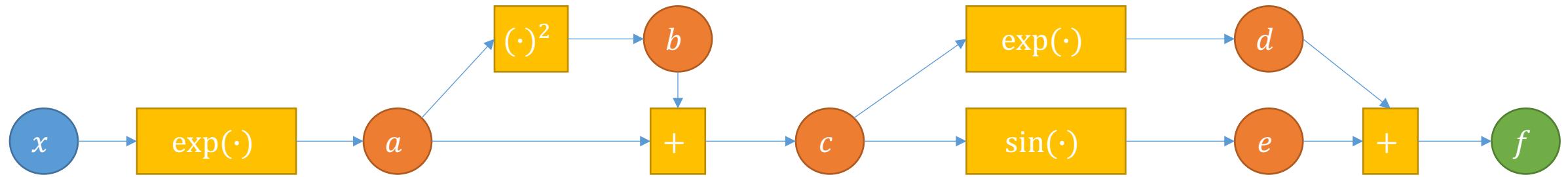
$$e = \sin(c)$$

$$f = d + e.$$

We can draw this as graph

$$f = \exp(\exp(x) + \exp(x)^2) + \sin(\exp(x) + \exp(x)^2)$$

$$\begin{aligned}a &= \exp(x) \\b &= a^2 \\c &= a + b \\d &= \exp(c) \\e &= \sin(c) \\f &= d + e.\end{aligned}$$



Mechanically writing down derivatives

$$\frac{df}{dd} = 1$$

$$\frac{df}{de} = 1$$

$$\frac{df}{dc} = \underbrace{\frac{df}{dd} \frac{dd}{dc}}_{=1} + \underbrace{\frac{df}{de} \frac{de}{dc}}_{=1}$$

$$\frac{df}{db} = \underbrace{\frac{df}{dc} \frac{dc}{db}}_{=1}$$

$$\frac{df}{da} = \underbrace{\frac{df}{dc} \frac{dc}{da}}_{c \text{ is a constant}} + \frac{df}{db} \frac{db}{da}$$

$$\frac{df}{dx} = \frac{df}{da} \frac{da}{dx}$$

$$\frac{df}{dd} = 1$$

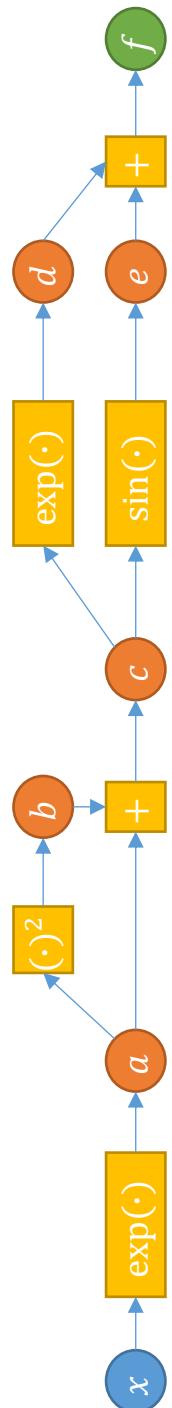
$$\frac{df}{de} = 1$$

$$\frac{df}{dc} = \frac{df}{dd} \exp(c) + \frac{df}{de} \cos(c)$$

$$\frac{df}{db} = \frac{df}{dc}$$

$$\frac{df}{da} = \frac{df}{dc} + \frac{df}{db} 2a$$

$$\frac{df}{dx} = \frac{df}{da} \exp(x)$$



Backpropagation in Neural Networks

Example & Derivation on the blackboard

① Backpropagation - single unit

$$* f^{(3)} = \|z^{(2)} - y\|^2 = z^{(2)T} z^{(2)} - 2z^{(2)}y + y^T y$$

$$z^{(3)} = C(\theta)$$

$\frac{\partial C}{\partial C} = \frac{\partial z^{(3)}}{\partial z^{(3)}} = 1,$

$$\boxed{f^{(3)} = \|z^{(2)} - y\|^2}$$

$$\frac{\partial z^{(3)}}{\partial z^{(2)}} \cdot 1 \quad < \quad \dots$$

$$z^{(2)} \uparrow$$

$$\boxed{f^{(2)} = \sigma(z^{(1)})}$$

$$\frac{\partial z^{(2)}}{\partial z^{(1)}} \cdot \frac{\partial z^{(3)}}{\partial z^{(2)}} \cdot 1 \quad < \quad \dots$$

$$z^{(1)} \uparrow$$

$$\boxed{f^{(1)} = \sum \theta^{(1)} z^{(0)}}$$

$$\frac{\partial z^{(1)}}{\partial \theta^{(1)}} \cdot \frac{\partial z^{(2)}}{\partial z^{(1)}} \cdot \frac{\partial z^{(3)}}{\partial z^{(2)}} \cdot 1 \quad < \quad = \frac{\partial C}{\partial \theta} = \frac{\partial z^{(3)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial \theta^{(1)}}$$

$$z^{(0)} = \bar{x}$$

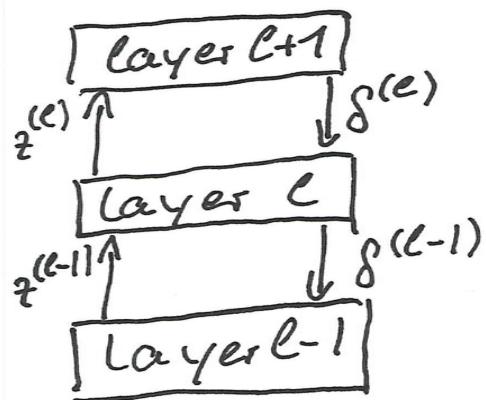
$$\frac{\partial z^{(3)}}{\partial z^{(2)}}^* = 2z^{(2)} - 2y = 2(z^{(2)} - y)$$

$$\frac{\partial z^{(2)}}{\partial z^{(1)}}^* = \sigma(z^{(1)})(1 - \sigma(z^{(1)}))$$

$$\frac{\partial z^{(1)}}{\partial \theta^{(1)}}^* = z^{(0)} = \bar{x}$$

$$= 2(z^{(2)} - y)\sigma(z^{(1)})(1 - \sigma(z^{(1)}))\bar{x}$$

⑪ BP layer wise

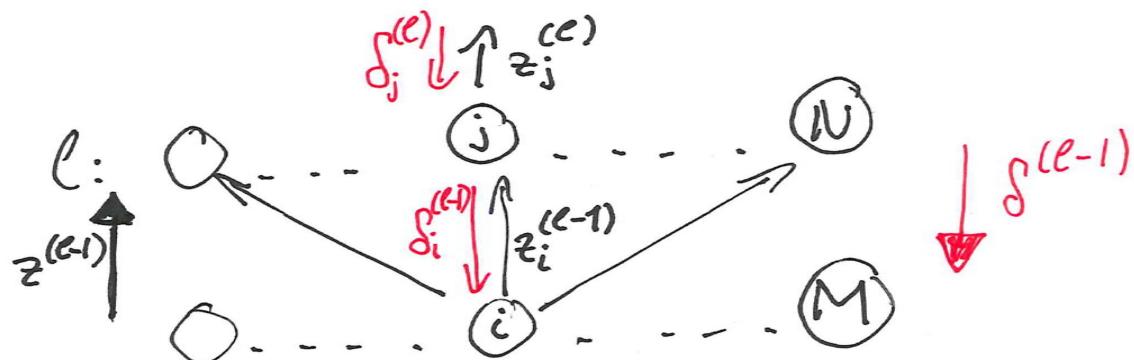


$\delta^{(l-1)}$: Gradient that flows from layer (l) to $(l-1)$

$\delta_i^{(l-1)}$: Gradient wrt to the i^{th} unit in the $(l-1)^{\text{th}}$ -layer

Example: "Layer 2 receives inputs $z^{(1)}$ and passes down gradients $\delta^{(1)}$ ".

⑪a



$$\delta_i^{(l-1)} = \frac{\partial C}{\partial z_i^{(l-1)}} = \sum_{j=1}^N \underbrace{\frac{\partial C}{\partial z_j^{(l)}}}_{\delta_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial z_i^{(l-1)}}$$

⑪b

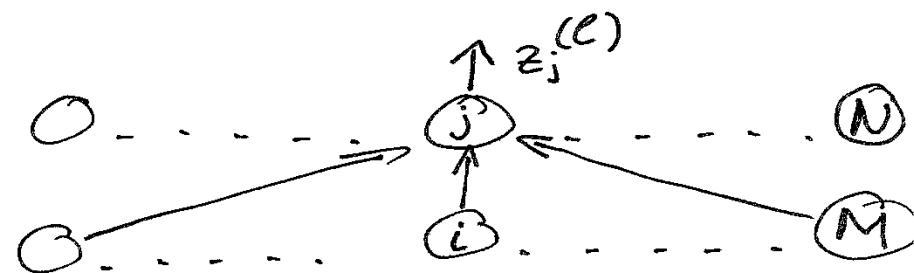
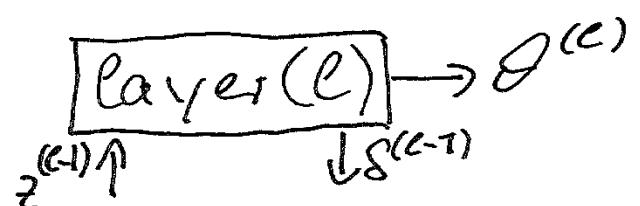
$$\begin{aligned} \delta^{(l-1)} &= \sum_{j=1}^N \delta_j^{(l)} \frac{\partial z_j^{(l)}}{\partial z^{(l-1)}} \\ &= \frac{\partial C}{\partial z^{(l)}} \frac{\partial z^{(l)}}{\partial z^{(l-1)}} \end{aligned}$$

$$\begin{aligned} <1> \quad &<1, M> <M, 1> \\ z_j^{(l)} &= \sigma(\delta_j^{(l)} z^{(l-1)}) \end{aligned}$$

$$\begin{aligned} <N, 1> \quad &<N, M> <M, 1> \\ z^{(l)} &= \sigma(\delta^{(l)} z^{(l-1)}) \end{aligned}$$

$$\sigma \left| \left(\delta_1, \delta_2, \dots \right) \middle/ \begin{pmatrix} z_1 \\ z_2 \\ \vdots \end{pmatrix} \right|$$

(III) BP - parameter update



$$\frac{\partial C}{\partial z_j^{(\ell)}} \cdot \frac{\partial z_j^{(\ell)}}{\partial \theta_j^{(\ell)}} \quad \left. \right\} \text{ weight update for } j^{\text{th}} \text{ unit}$$

$$\frac{\partial C}{\partial \theta^{(\ell)}} = \sum_{j=1}^N \frac{\partial C}{\partial z_j^{(\ell)}} \cdot \frac{\partial z_j^{(\ell)}}{\partial \theta^{(\ell)}} \quad \begin{matrix} j=1 \\ \vdots \\ j=2 \end{matrix} \rightarrow \begin{matrix} N \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_M \end{matrix} \quad \left. \right\}$$

$$= \sum_{j=1}^N \delta_j^{(\ell)} \frac{\partial z_j^{(\ell)}}{\partial \theta^{(\ell)}}$$

$$= \frac{\partial C}{\partial z^{(\ell)}} \frac{\partial z^{(\ell)}}{\partial \theta^{(\ell)}}$$

$$N \begin{bmatrix} \alpha_1 & \cdots & \alpha_M \\ \alpha_1 & \cdots & \alpha_M \\ \vdots & \ddots & \vdots \\ \alpha_1 & \cdots & \alpha_M \\ \alpha_1 & \cdots & \alpha_M \end{bmatrix} \quad \left. \right\}$$

$$N \begin{bmatrix} \alpha_1^1 & \cdots & \alpha_1^M \\ \alpha_2^1 & \cdots & \alpha_2^M \\ \vdots & \ddots & \vdots \\ \alpha_N^1 & \cdots & \alpha_N^M \end{bmatrix} \quad \left. \right\}$$

Next week

Convolutional Neural Networks