# AIL721: Deep Learning
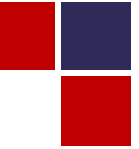
**Instructor:** James Arambam

**ScAI** | **Yardi School of Artificial Intelligence**
**Indian Institute of Technology Delhi**

# Class Announcements

ScAI

❑ **Please use your name and official email IDs in the piazza.**

  ▪ Enrollments without actual names and official email IDs will be removed.

❑ **Guidelines regarding the project topic.**

  ▪ Applications of deep learning (or neural networks) in problems related to your respective branches.

  ▪ Pick an application that interests you, and explore how best to apply learning algorithms to solve it.

  ▪ Computer Vision, Natural Language Processing, Speech Recognition, Reinforcement Learning, Healthcare etc.

❑ **Guest Lecture (online) on Training LLMs**   - Confirmed!



Dr. Maksim Tkachenko
Research Scientist,
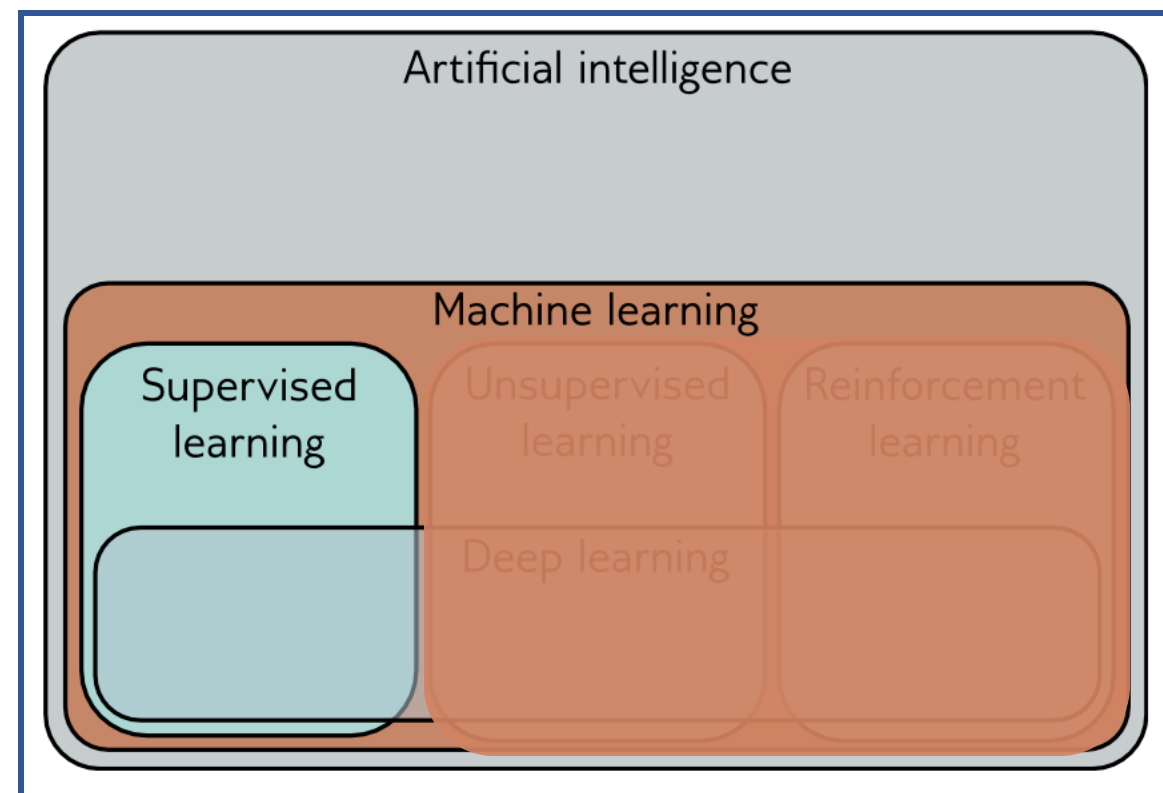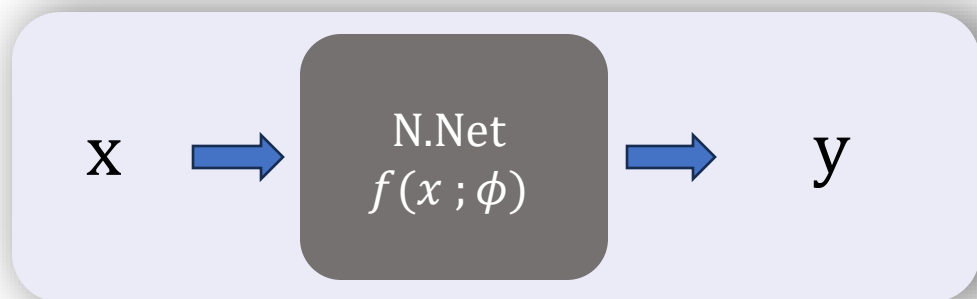Rakuten Institute of Technology, Singapore.

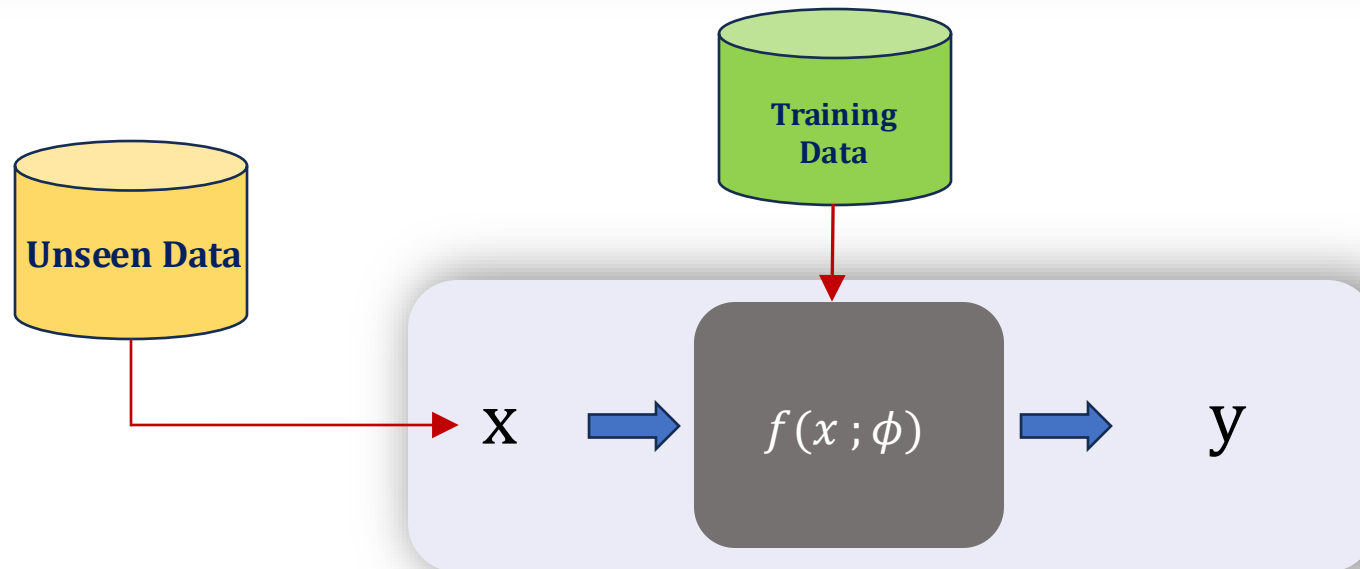❑ **IIT Delhi HPC Credit**
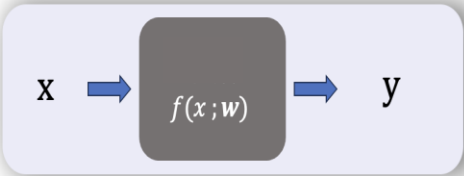
  ▪ How much for each student?

# What is Deep Learning?

❑ Deep learning is a branch of **machine learning.**

❑ A general-purpose framework for **learning from data**.

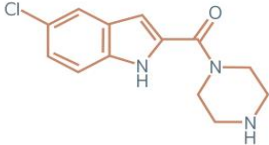❑ Based on computational models called **neural networks.**

# Supervised Learning

A type of **machine learning paradigm** that uses **labeled datasets** to train algorithms to **predict outcomes** and **recognize patterns**.



Training Data

Unseen Data

$x$ → $f(x\,;\phi)$ → $y$

# Examples



| Real world input | Model input | Model | Model output | Real world output | |
|---|---|---|---|---|---|
| 6000 square feet, 4 bedrooms, previously sold for $235K in 2005, 1 parking spot. | $\begin{bmatrix}6000\\4\\235\\2005\\1\end{bmatrix}$ | Deep learning model | $[340]$ | Predicted price is $340k | **Regression** |
| (chemical structure) | $\begin{bmatrix}1\\0\\1\\\vdots\\17\\1\\1\\\vdots\end{bmatrix}$ | Deep learning model | $\begin{bmatrix}-12.9\\56.4\end{bmatrix}$ | Freezing point is -12.9°C Boiling point is 56.4°C | **Multivariate Regression** |
| "The steak was terrible, the salad was rotten, and the soup tasted like socks" | $\begin{bmatrix}8672\\8194\\9804\\8634\\8672\\\vdots\end{bmatrix}$ | Deep learning model | $\begin{bmatrix}0.02\\0.98\end{bmatrix}$ | Positive Negative | **Binary Classification** |
| (audio waveform) | $\begin{bmatrix}125\\12054\\1253\\6178\\24\\4447\\\vdots\end{bmatrix}$ | Deep learning model | $\begin{bmatrix}0.03\\0.52\\0.18\\0.07\\0.12\\0.08\end{bmatrix}$ | Classical Electronica Hip Hop Jazz Pop Metal | **Multiclass Classification** |
| (bicycle image) | $\begin{bmatrix}124\\140\\156\\128\\142\\157\\\vdots\end{bmatrix}$ | Deep learning model | $\begin{bmatrix}0.00\\0.00\\0.01\\0.89\\0.05\\0.00\\\vdots\\0.01\end{bmatrix}$ | Aardvark Apple Bee Bicycle Bridge Clown ⋮ Zebra | |

[Image from J.D Prince 2023 DL Book]
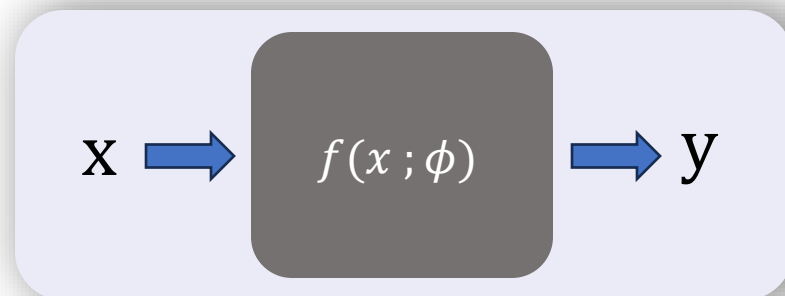
# Supervised Learning
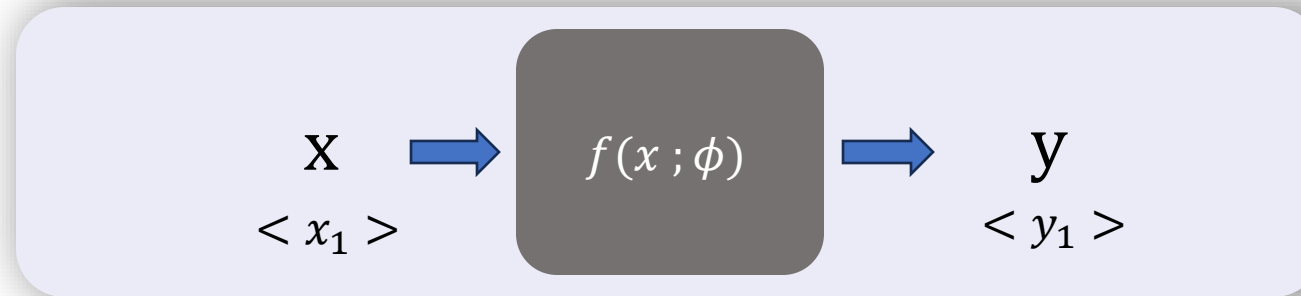
**Types of supervised learning:**

❑ Regression problem: Model predicts **real values**.

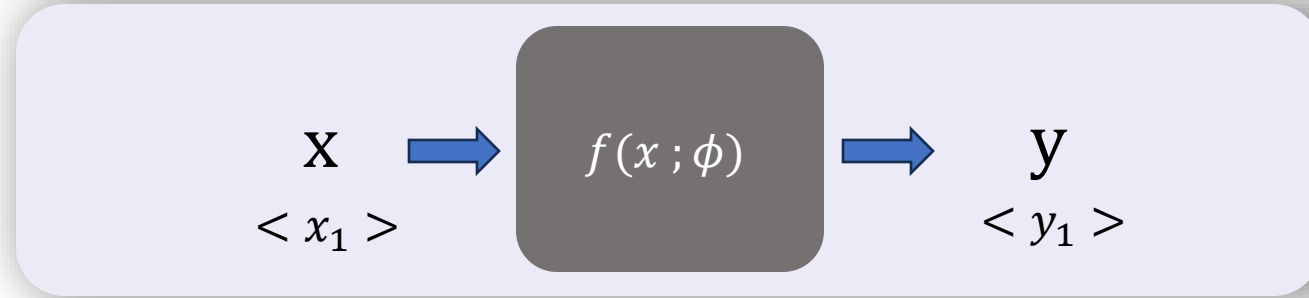❑ Classification problem: Model predicts **discrete values.**

$$x \Rightarrow \boxed{f(x\,;\phi)} \Rightarrow y$$

# Regression

$$x \Rightarrow f(x\,;\phi) \Rightarrow y$$

$$< x_1, x_2, \ldots, x_d > \qquad\qquad\qquad < y_1, y_2, \ldots, y_m >$$

**Regression**

$$x \Rightarrow f(x\,;\phi) \Rightarrow y$$

$$< x_1 > \qquad\qquad\qquad\qquad < y_1 >$$

**A Simple Regression Problem**

# Linear Regression

$$x \Rightarrow f(x\,;\phi) \Rightarrow y$$

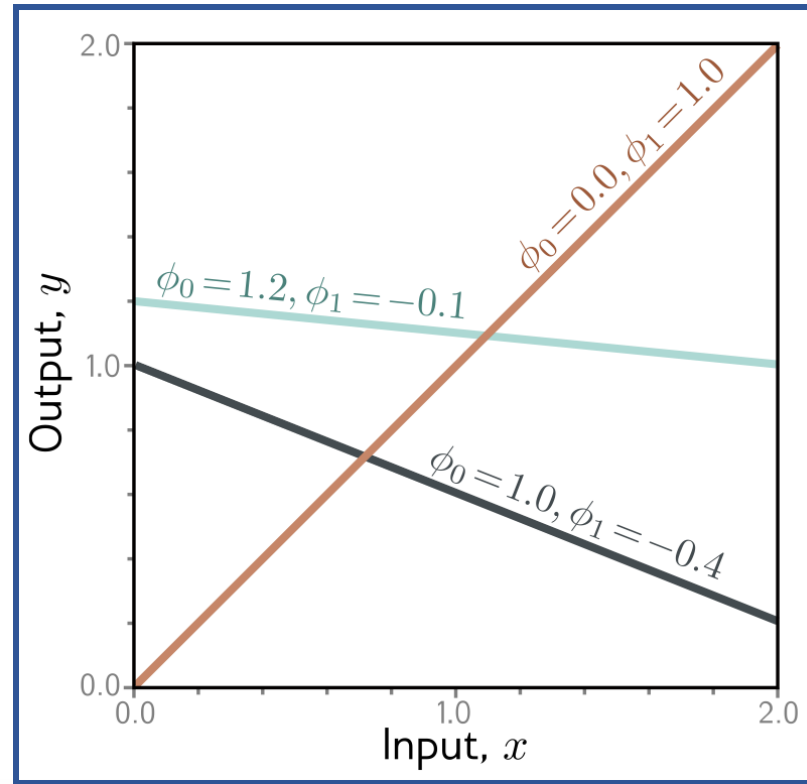$$< x_1 > \qquad\qquad\qquad < y_1 >$$

**A Simple Regression Problem**

**What is the simplest mathematical model to represent the function $f(x\,;\phi)$?**

$$f(x;\boldsymbol{\phi}) = \phi_0 + \phi_1 \cdot x$$
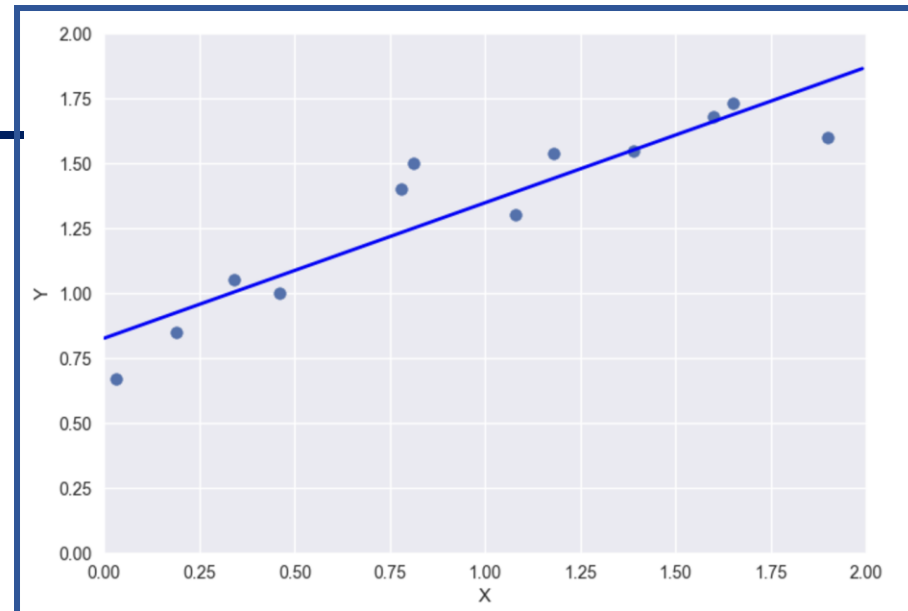
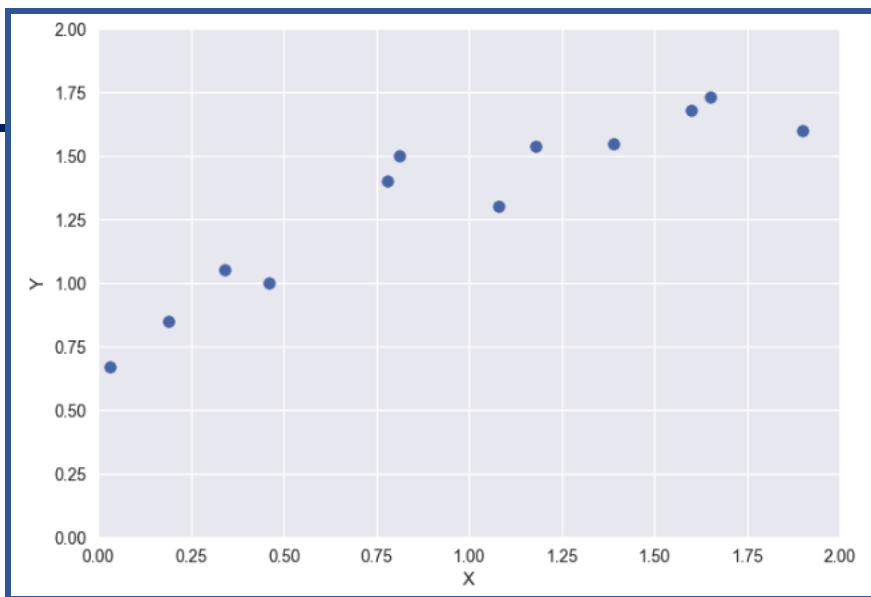$$f(x;\boldsymbol{w}) = w_0 + w_1 \cdot x$$

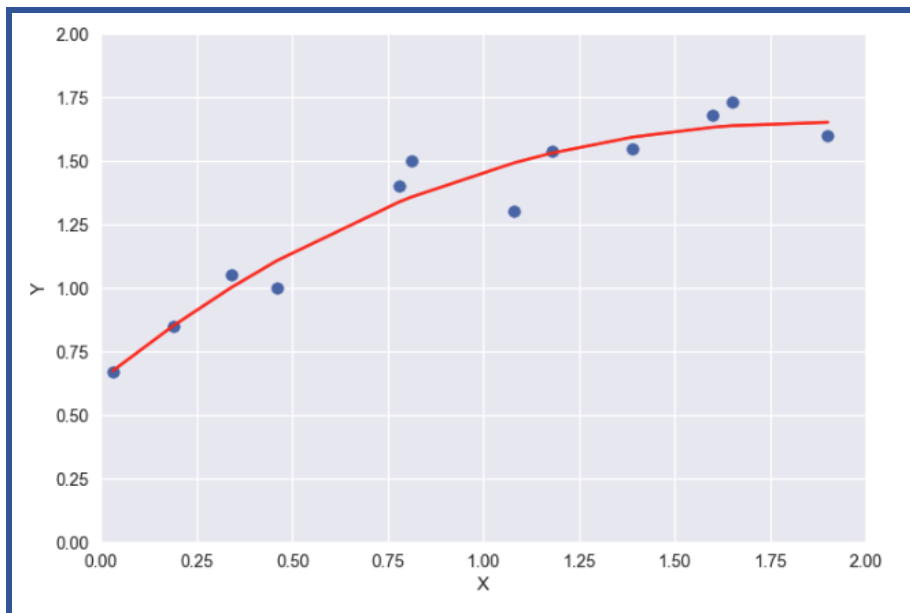$$f(x;\boldsymbol{\theta}) = \theta_0 + \theta_1 \cdot x$$

# Linear Model



What's the limitation of such a linear model?

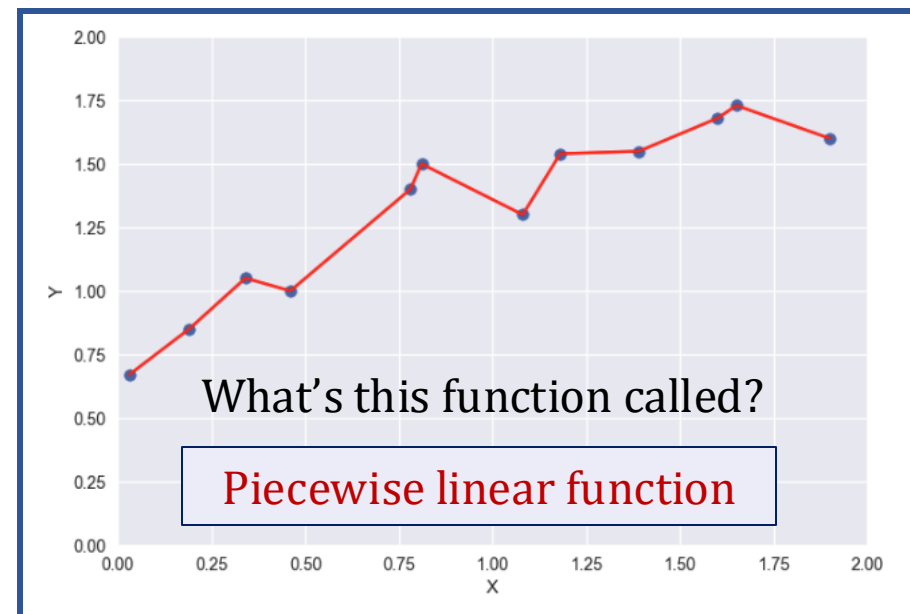Linear models can only describe the input/output relationship as a line.

[Image from J.D Prince 2023 DL Book]

Linear function

Smooth function

What's this function called?

Piecewise linear function
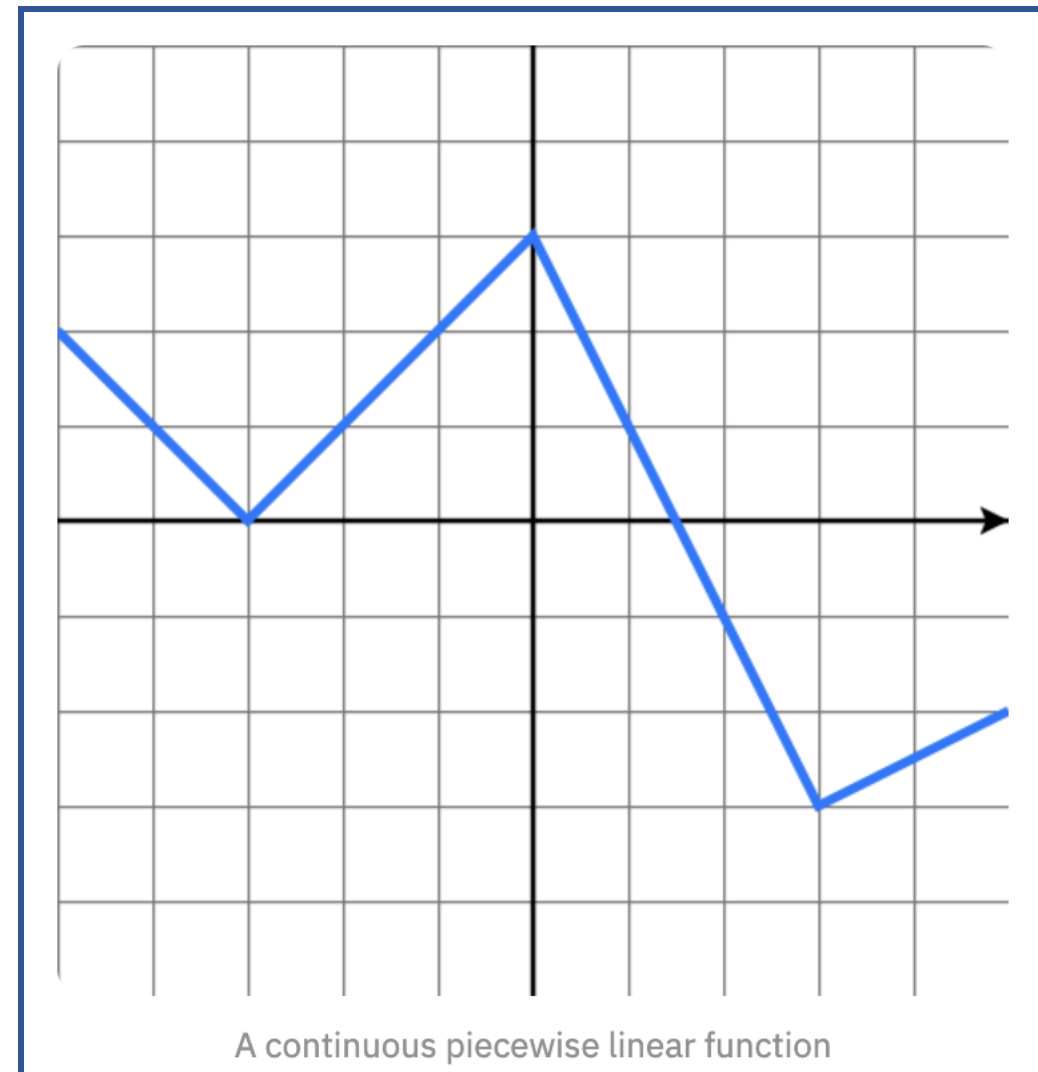
# Piecewise Linear Function

$$f(x) = \begin{cases} -x - 3 & \text{if } x \le -3 \\ x + 3 & \text{if } -3 < x < 0 \\ -2x + 3 & \text{if } 0 \le x < 3 \\ 0.5x - 4.5 & \text{if } x \ge 3 \end{cases}$$
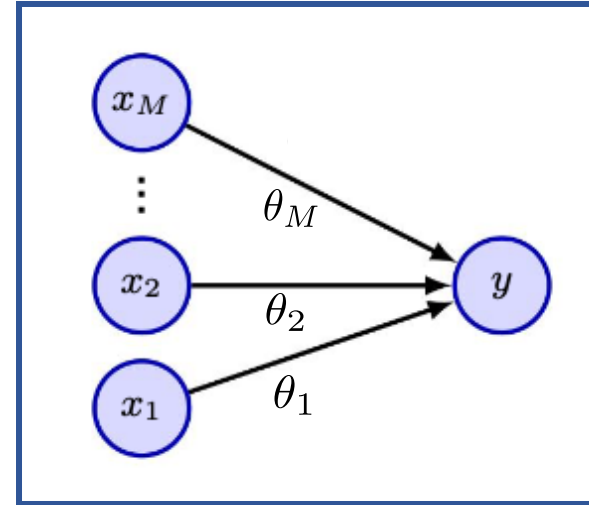


A continuous piecewise linear function

(*Image is from Wikipedia) INDIAN INSTITUTE OF TECHNOLOGY DELHI

# How piecewise linear function relates to neural networks?

# Neural Networks

❑ Mathematically:

$$h = \sum_{i=1}^{M} \theta_i \cdot x_i \qquad (1)$$

$$y = a[h] \qquad (2)$$

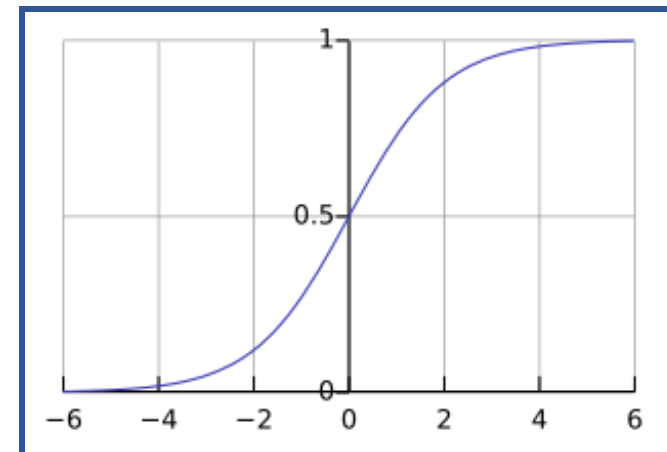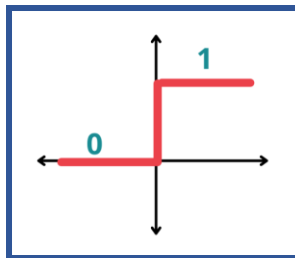$a[\cdot]$ **: Activation Function**

# Activation Function

## Earlier motivation:

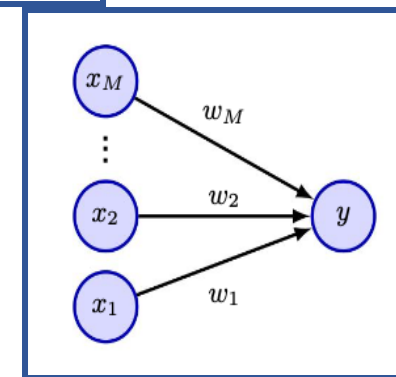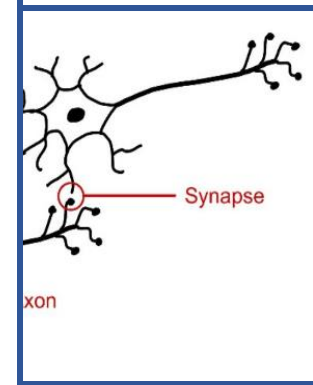❑ **Firing of neurons** depends on the **strength of the synapses**.

❑ Rosenblatt's Perceptron, 1962

Activation function a[h] is a **step function**:

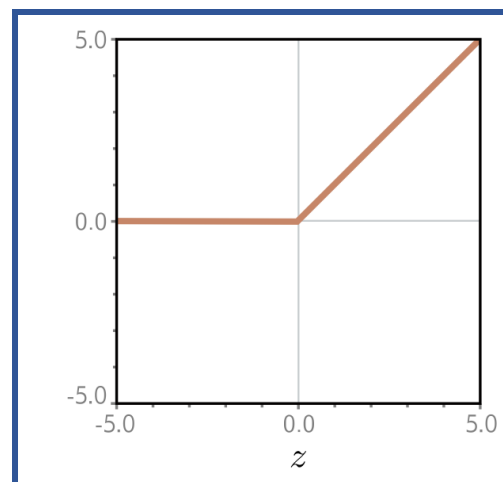$$a[h] = \begin{cases} 0, & \text{if } h \leq 0, \\ 1, & \text{if } h > 0 \end{cases}$$

❑ Depending on the type of a[h], it can introduce **non-linearity (a key property of NN).**

**Sigmoid Function**

$$a[z] = \text{ReLU}[z] = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}$$

Piecewise linear function

$$h = \sum_{i=1}^{M} \theta_i \cdot x_i$$

$$y = a[h]$$

# Neural Networks
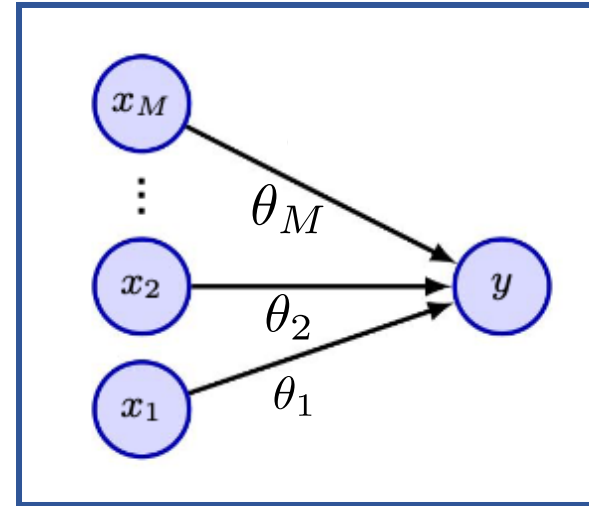
☐ Mathematically:

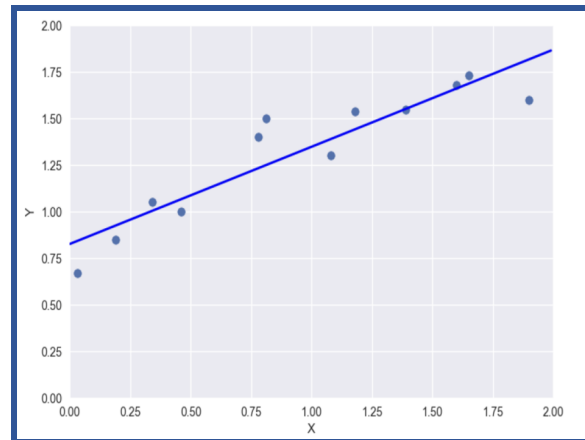$$h = \sum_{i=1}^{M} \theta_i \cdot x_i \qquad (1)$$

$$y = a[h] \qquad (2)$$



$$h = \theta_0 + \theta_1 \cdot x$$

$$y = \boxed{a[h] = h}$$

**Linear or Nonlinear activation?**



**Linear or Nonlinear Function?**

$$h = \theta_0 + \theta_1 \cdot x$$

$$y = \phi_0 + \phi_1 \cdot a[h]$$

$$y = \phi_0 + \phi_1 \cdot a[\theta_{10} + \theta_{11} \cdot x]$$

[Image is from Bishop 2024 DL Book]

# Neural Networks

ScAI

$$y = \phi_0 + \phi_1 \cdot a[\theta_{10} + \theta_{11} \cdot x]$$

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$
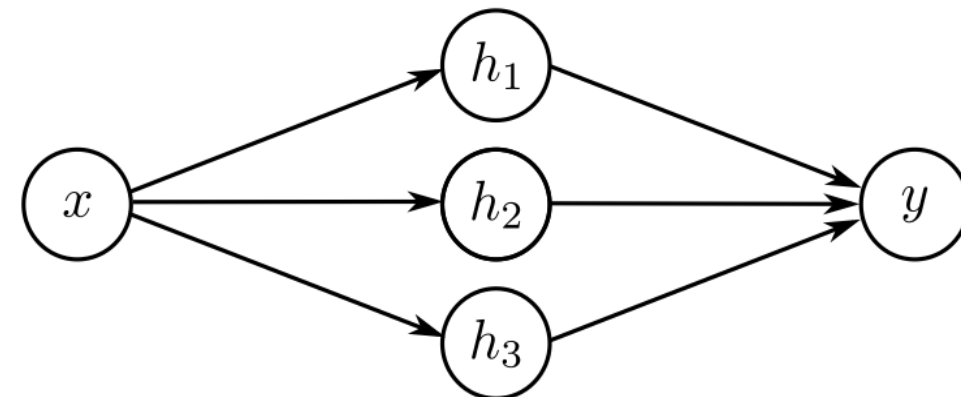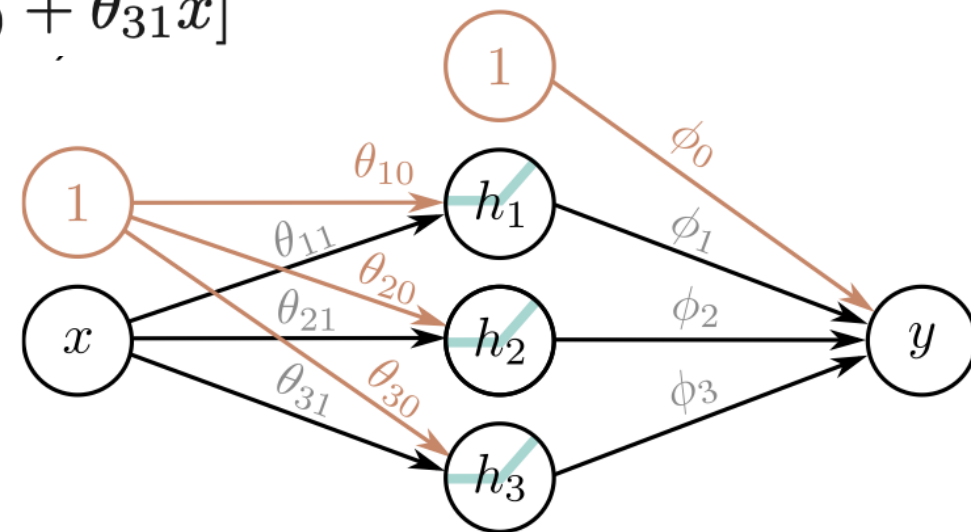
$$
\begin{aligned}
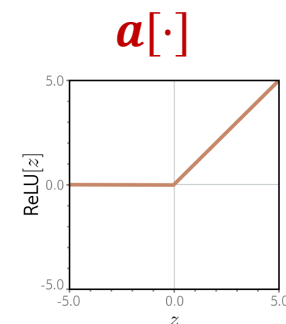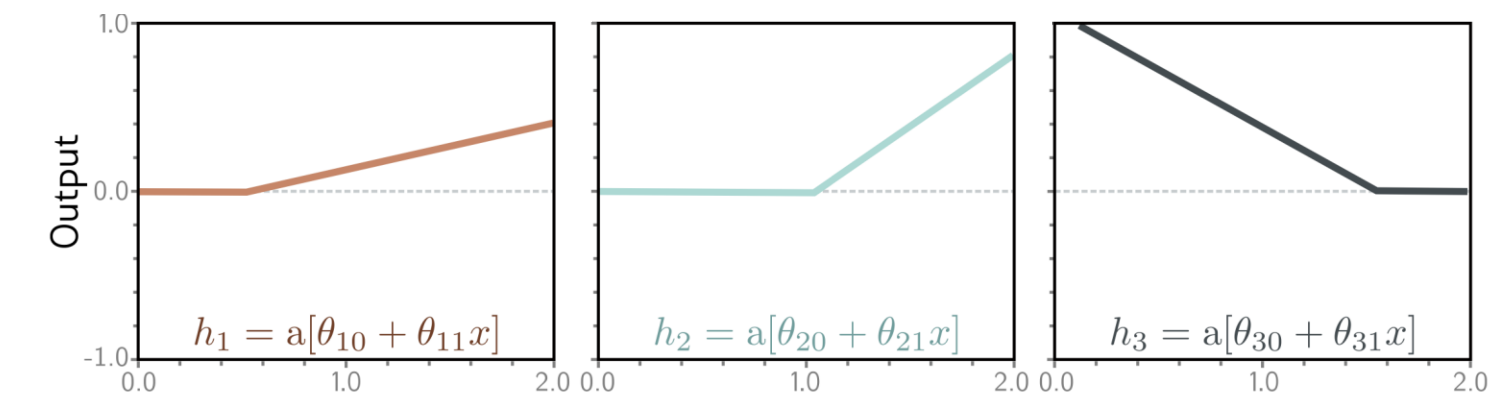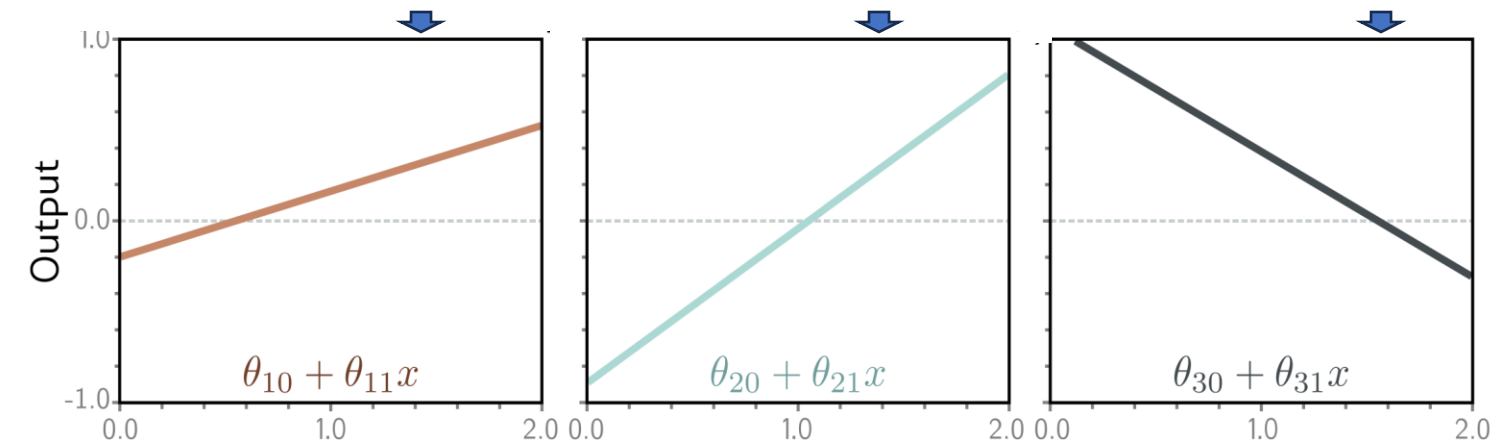h_1 &= a[\theta_{10} + \theta_{11}x] \\
h_2 &= a[\theta_{20} + \theta_{21}x] \\
h_3 &= a[\theta_{30} + \theta_{31}x]
\end{aligned}
$$

**Hidden Units**

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

[Image from J.D Prince 2023 DL Book]

INDIAN INSTITUTE OF TECHNOLOGY DELHI

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$
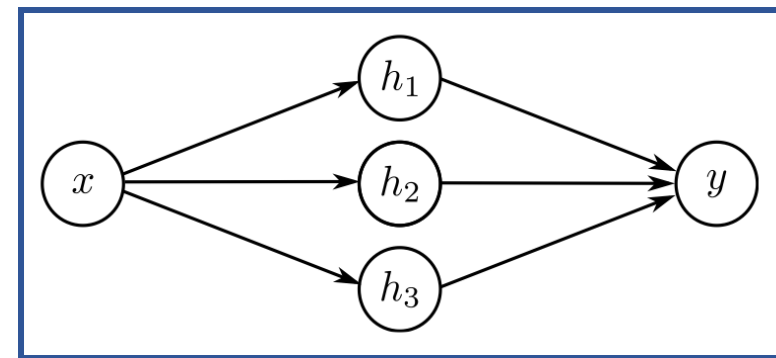


$a[\cdot]$

ReLU Activation

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

[Image from J.D Prince 2023 DL Book] INDIAN INSTITUTE OF TECHNOLOGY DELHI

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$



a)

Output

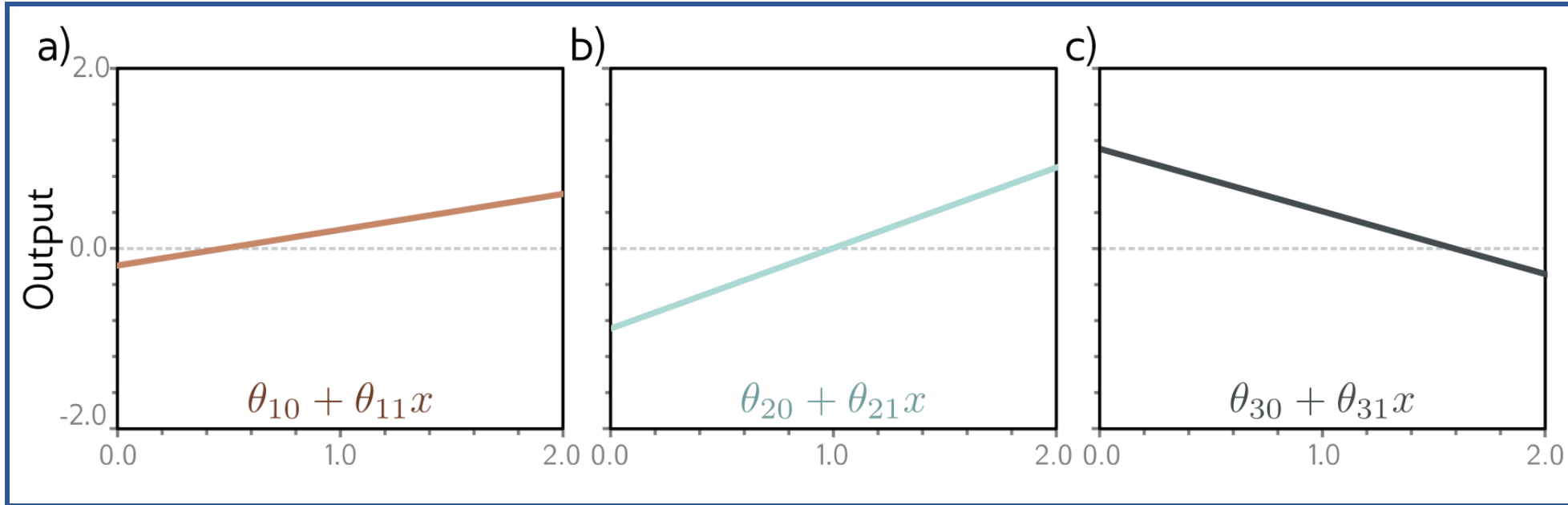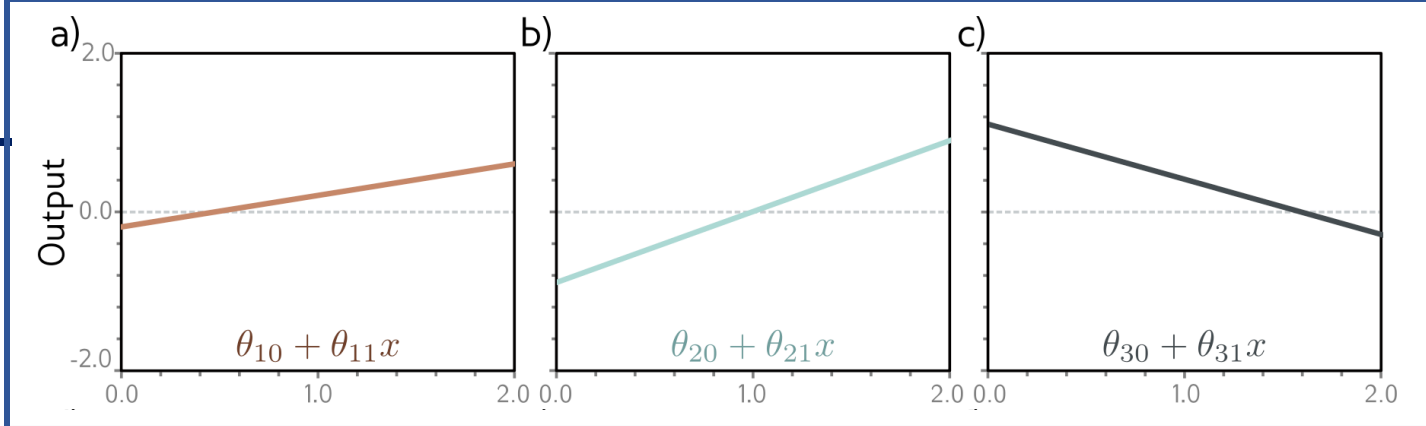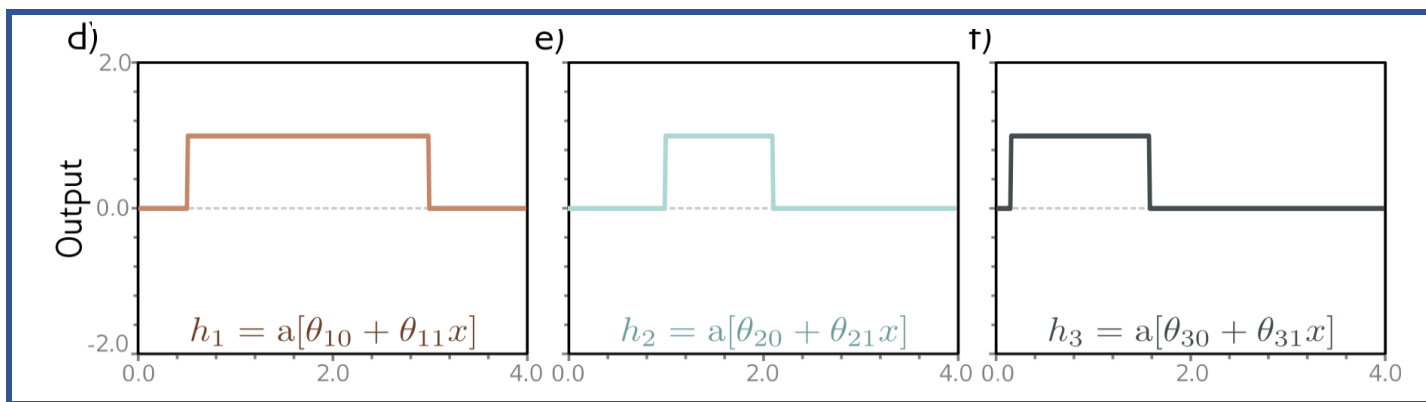$\theta_{10} + \theta_{11}x$

b)

$\theta_{20} + \theta_{21}x$

c)

$\theta_{30} + \theta_{31}x$

**Activation Fn:**

$$\text{rect}[z] = \begin{cases} 0 & z < 0 \\ 1 & 0 \leq z \leq 1 \\ 0 & z > 1 \end{cases}$$

a) $\theta_{10} + \theta_{11}x$

b) $\theta_{20} + \theta_{21}x$

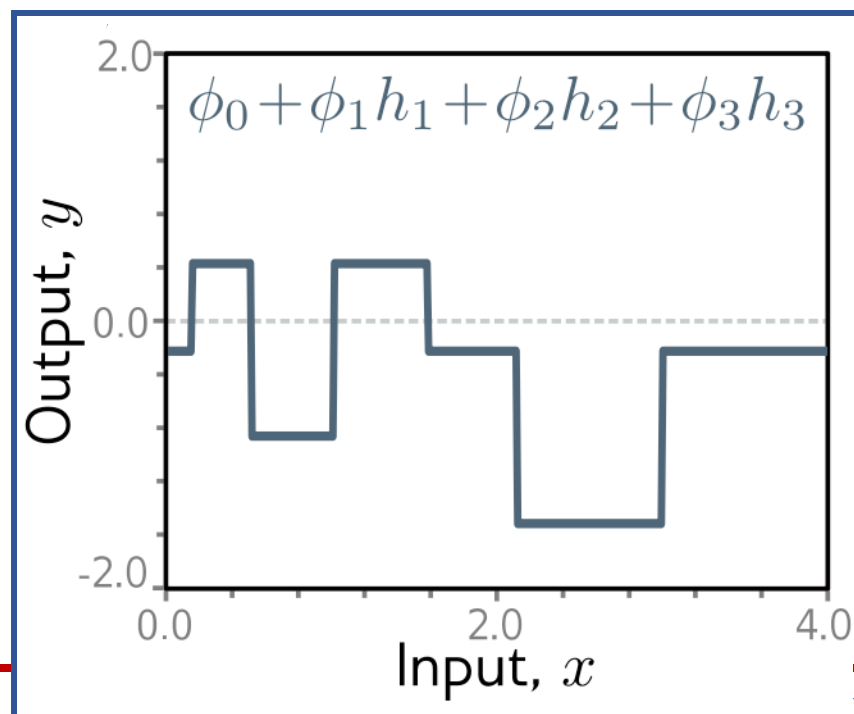c) $\theta_{30} + \theta_{31}x$

**Activation Fn:**

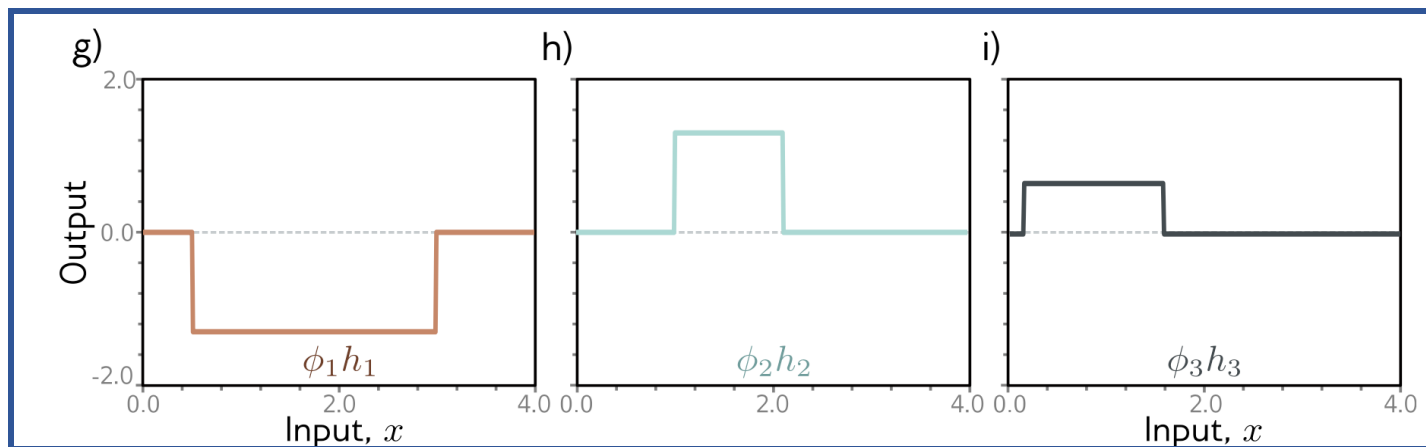$$\text{rect}[z] = \begin{cases} 0 & z < 0 \\ 1 & 0 \leq z \leq 1 \\ 0 & z > 1 \end{cases}$$

d) $h_1 = \text{a}[\theta_{10} + \theta_{11}x]$

e) $h_2 = \text{a}[\theta_{20} + \theta_{21}x]$

f) $h_3 = \text{a}[\theta_{30} + \theta_{31}x]$

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

g) $\phi_1 h_1$

h) $\phi_2 h_2$

i) $\phi_3 h_3$

$\phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$

# Activation Pattern

Unit 1



Unit 2



Unit 3



$\phi_1 h_1$

$\phi_2 h_2$

$\phi_3 h_3$

What kind of function is it?



$\phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$
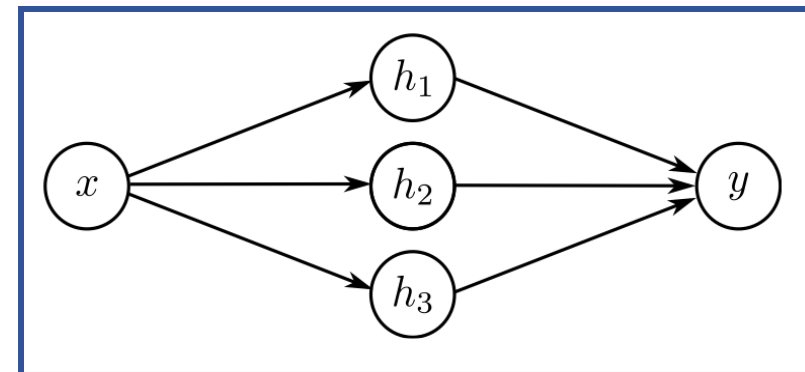
Shaded region:
- Unit 1 active
- Unit 2 inactive
- Unit 3 active



Piecewise linear function

How piecewise linear function relates to neural networks?

[Image from J.D Prince 2023 DL Book]

INDIAN INSTITUTE OF TECHNOLOGY DELHI

# Neural Networks

ScAI

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$



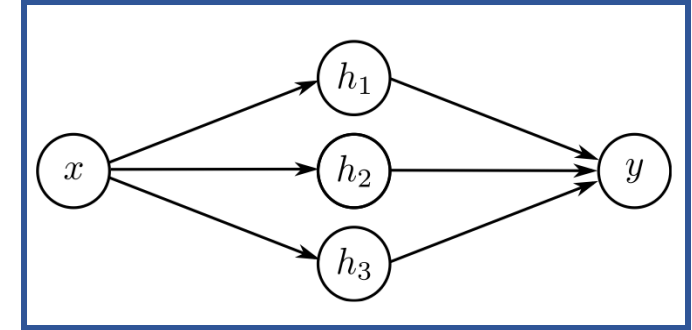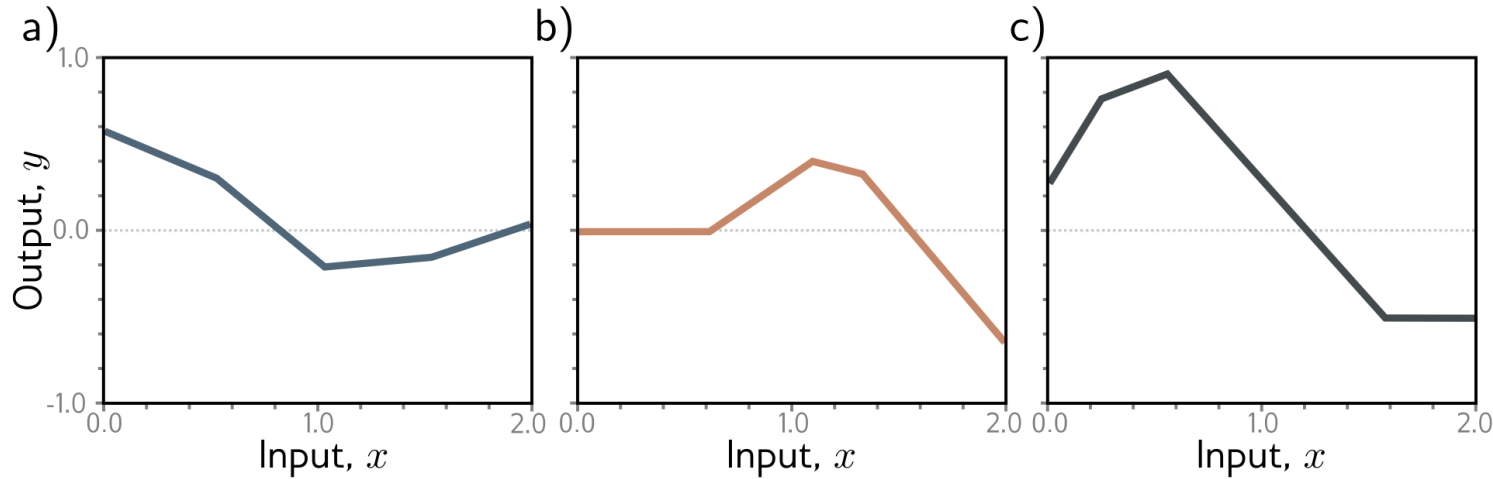$$\phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

**What factors influence the curve's shape?**

❑ **Number of model parameters/weights:** $\phi_0, \phi_1, \phi_2, \phi_3, \theta_{10}, \theta_{11}, \theta_{20}, \theta_{21}, \theta_{30}, \theta_{31}$

❑ **Type of activation function:** ReLU, Sigmoid, tanh, etc.

**Data**

**One more key factor remaining. Anyone?**

[Image from J.D Prince 2023 DL Book] INDIAN INSTITUTE OF TECHNOLOGY DELHI

# Neural Networks

**How many linear regions (or segments) in the above piecewise linear function?**

❑ **Neural Network with three hidden units:**

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$

❑ **Neural Network with "D" hidden units:**

$$h_d = a[\theta_{d0} + \theta_{d1}x] \qquad y = \phi_0 + \sum_{d=1}^{D} \phi_d h_d$$

How many linear regions (or segments) with "D" hidden units under ReLU activation?

**D+1**

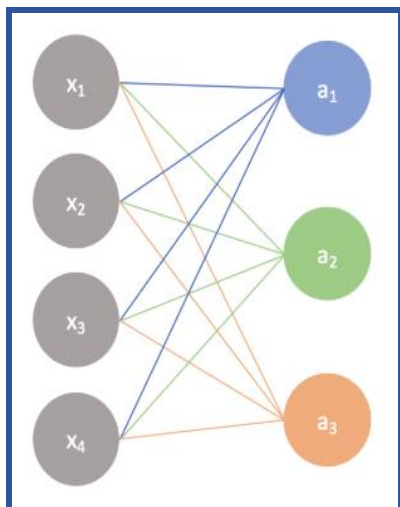# General Form of Neural Networks

❑ **Neural Network with three hidden units:**

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$

❑ **Neural Network with "D" hidden units:**

$$h_d = a[\theta_{d0} + \theta_{d1}x] \qquad y = \phi_0 + \sum_{d=1}^{D} \phi_d h_d$$
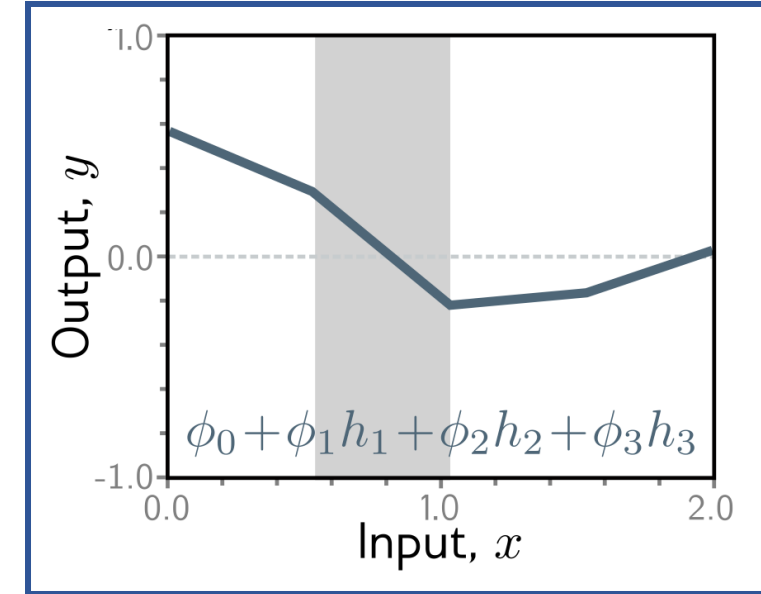
❑ **Matrix Form:**
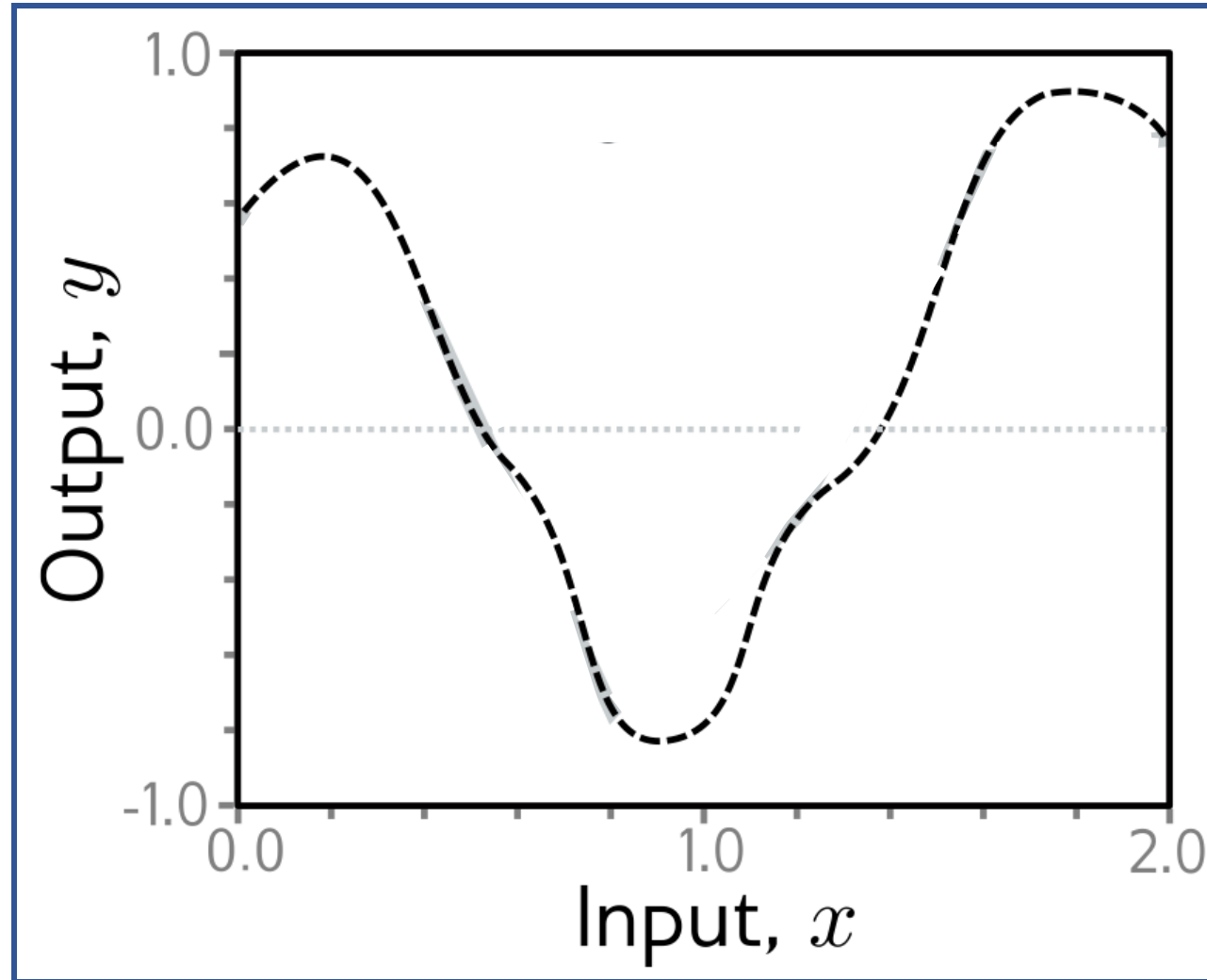
# Neural Networks

❑ **Neural Network with three hidden units:**

$$y = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}x] + \phi_2 a[\theta_{20} + \theta_{21}x] + \phi_3 a[\theta_{30} + \theta_{31}x]$$
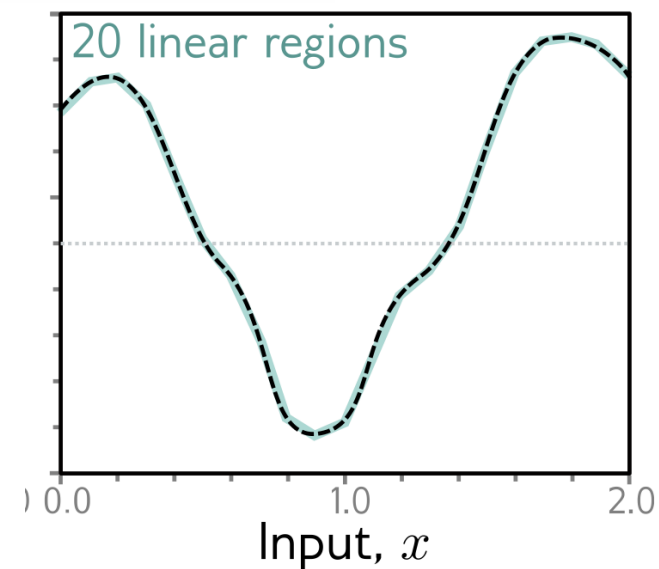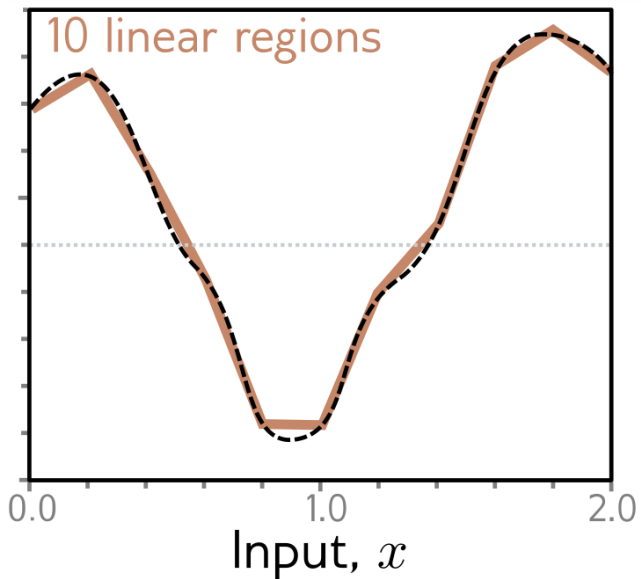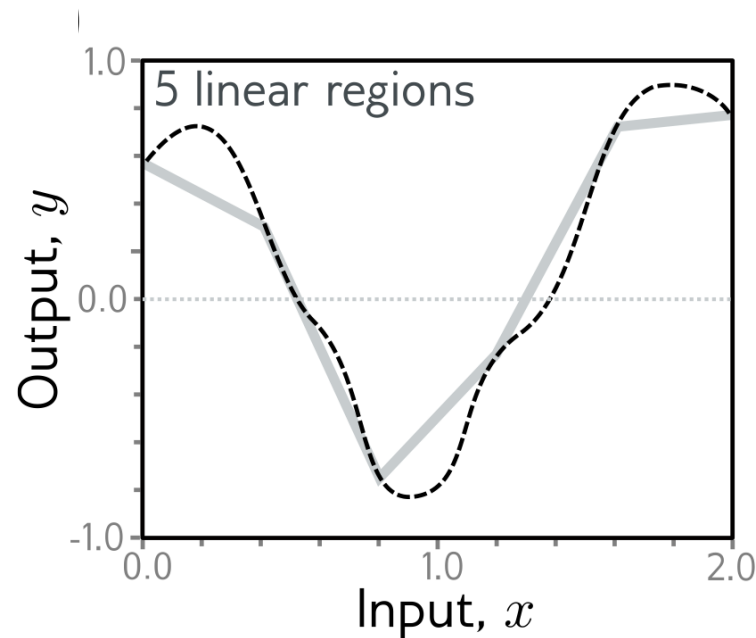
What the function be like if a[z] = $\psi_0 + \psi_1 \cdot z$?

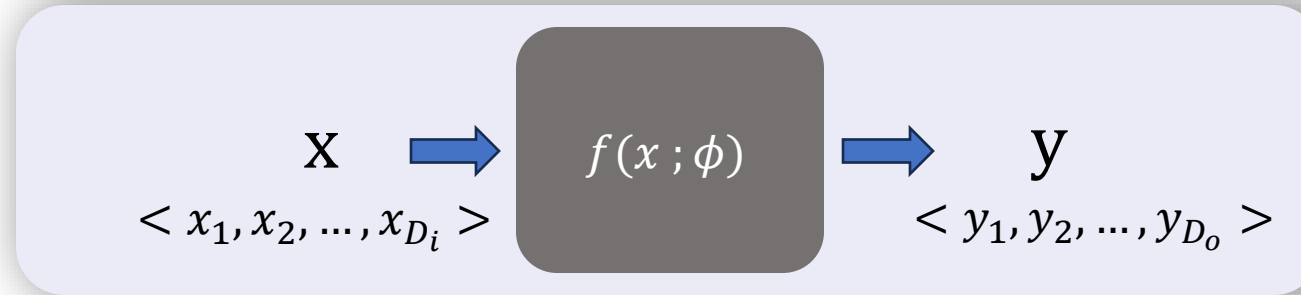# Neural Networks

**Any arbitrary continuous function.**

# Neural Networks

## Power of Neural Network



With enough hidden units (linear regions), we can describe any 1D function with arbitrary accuracy.

[Image from J.D Prince 2023 DL Book]

# Neural Networks

## Universal Approximation Theorem

"a formal proof that, with enough hidden units, a shallow neural network can describe any continuous function"

# Neural Networks

$$x \implies f(x\,;\phi) \implies y$$

$$< x_1, x_2, \ldots, x_{D_i} > \qquad\qquad\qquad < y_1, y_2, \ldots, y_{D_o} >$$

**Di-input/Do-output**

$$x \implies f(x\,;\phi) \implies y$$

$$< x_1 > \qquad\qquad\qquad\qquad < y_1 >$$

**1-input/ 1-output**

# Neural Networks

☐ Two Outputs

- 1 input, 4 hidden units, 2 outputs



$x \Rightarrow f(x\,;\phi) \Rightarrow y$

$< x_1 >$      $< y_1, y_2 >$

**1-input/ 2-output**

$$h_1 = \mathrm{a}[\theta_{10} + \theta_{11}x]$$
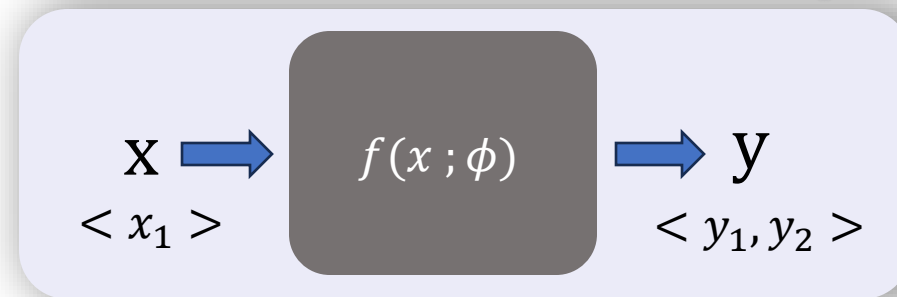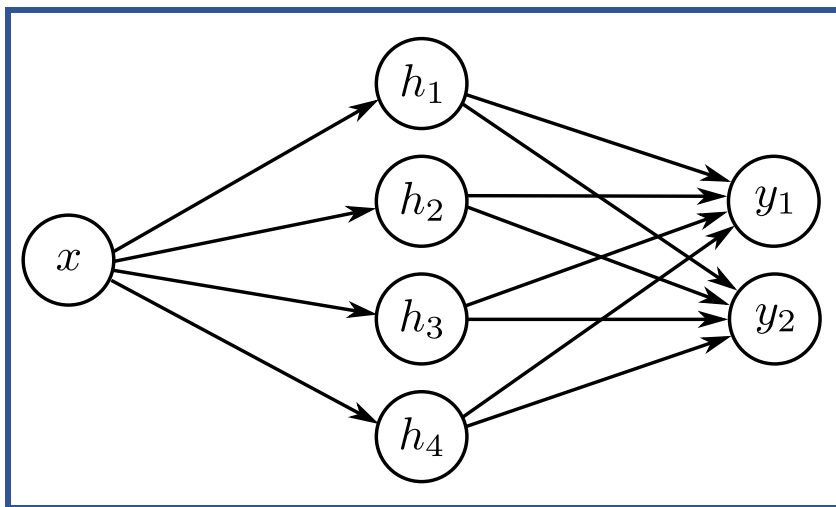$$h_2 = \mathrm{a}[\theta_{20} + \theta_{21}x]$$
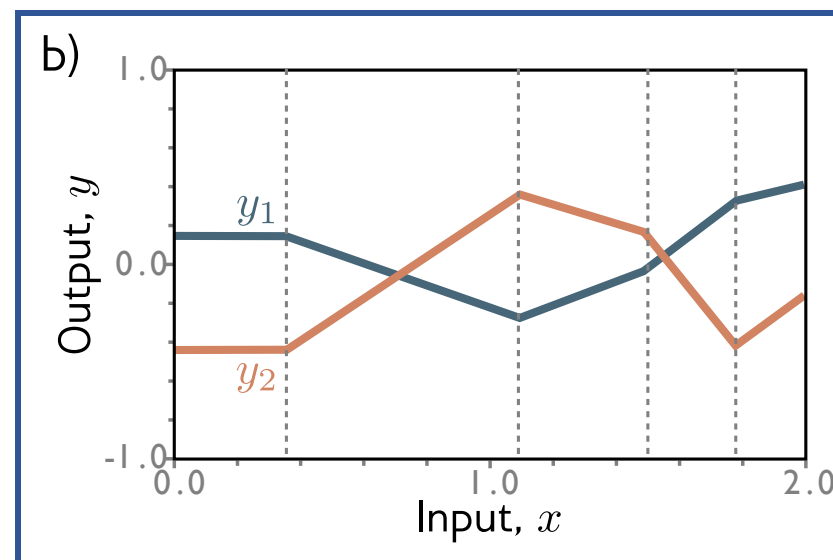$$h_3 = \mathrm{a}[\theta_{30} + \theta_{31}x]$$
$$h_4 = \mathrm{a}[\theta_{40} + \theta_{41}x]$$

$$y_1 = \phi_{10} + \phi_{11}h_1 + \phi_{12}h_2 + \phi_{13}h_3 + \phi_{14}h_4$$
$$y_2 = \phi_{20} + \phi_{21}h_1 + \phi_{22}h_2 + \phi_{23}h_3 + \phi_{24}h_4$$



b)

[Image from J.D Prince 2023 DL Book]
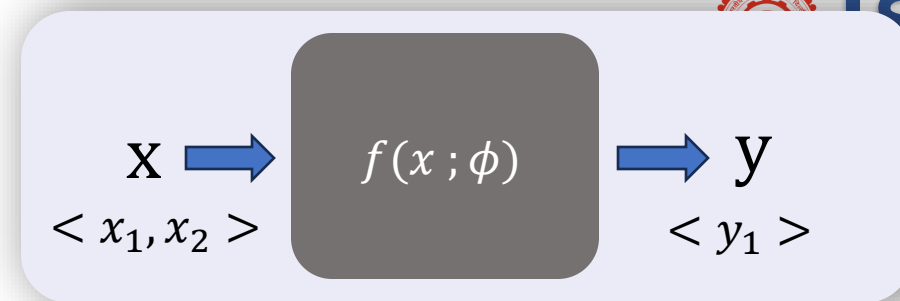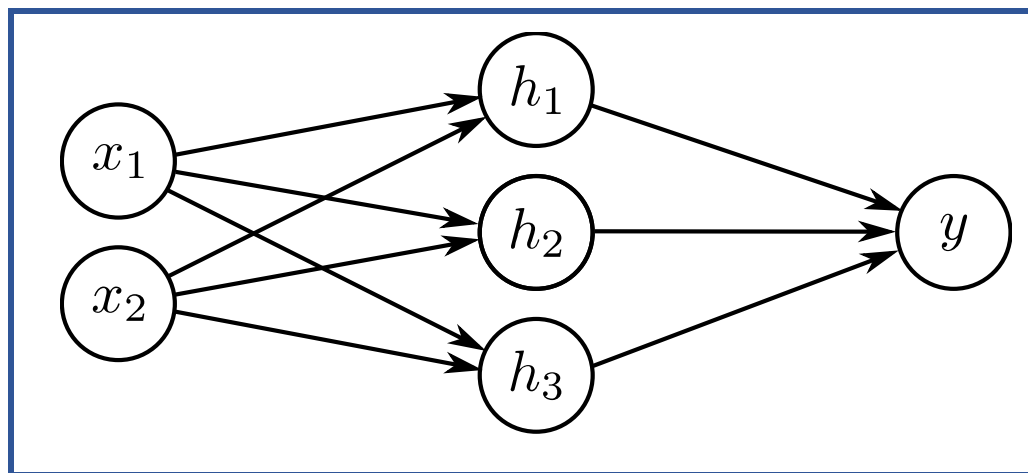
# Neural Networks

## ☐ Two Inputs

- 2 inputs, 3 hidden units, 1 outputs

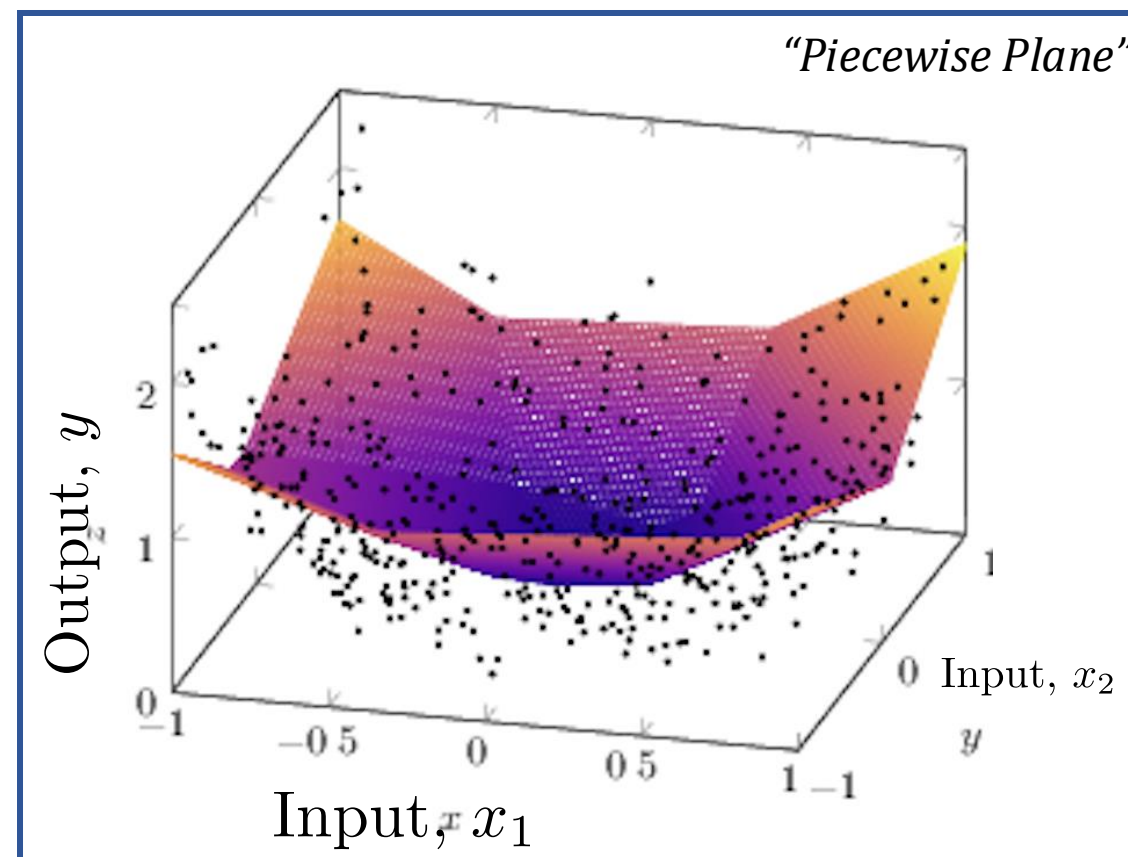$$h_1 = a[\theta_{10} + \theta_{11}x_1 + \theta_{12}x_2]$$
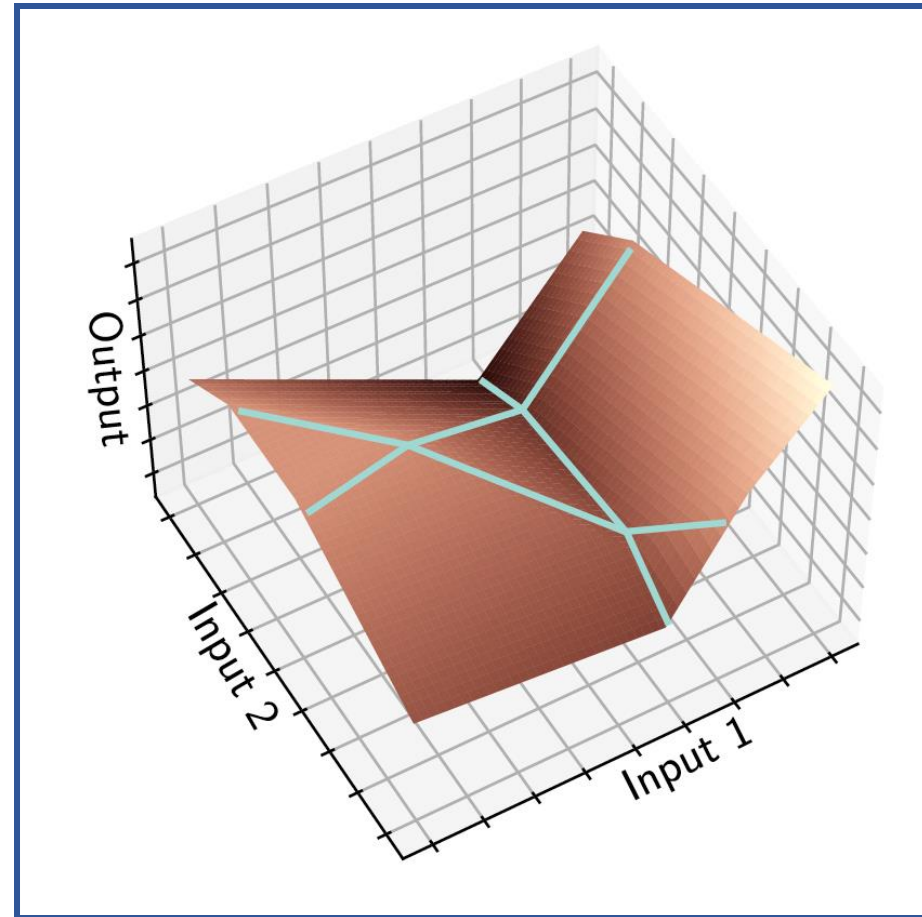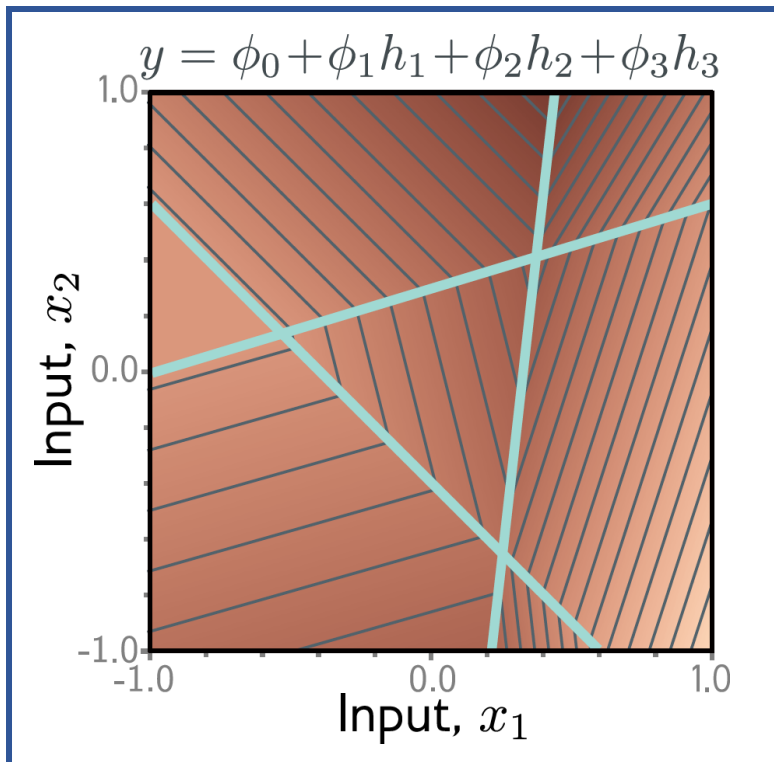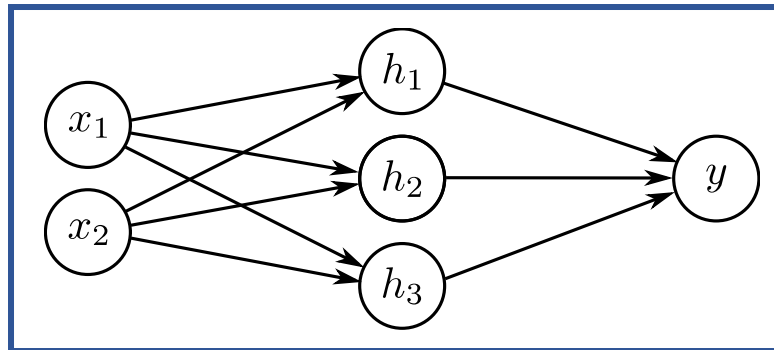$$h_2 = a[\theta_{20} + \theta_{21}x_1 + \theta_{22}x_2]$$
$$h_3 = a[\theta_{30} + \theta_{31}x_1 + \theta_{32}x_2]$$

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

$x$ ➡ $f(x\,;\phi)$ ➡ $y$

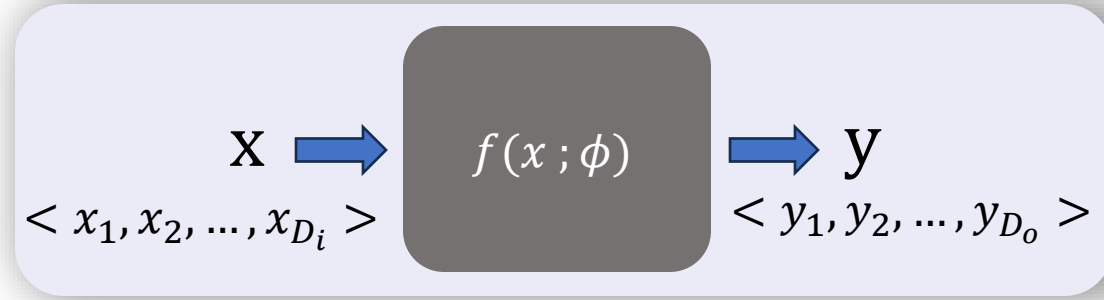$< x_1, x_2 >$          $< y_1 >$

**2-input/ 1-output**



"Piecewise Plane"

[Image from J.D Prince 2023 DL Book]

# Neural Network

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

[Image from J.D Prince 2023 DL Book]

# Neural Networks

$x$ ⟹ $f(x\,;\phi)$ ⟹ $y$

$< x_1, x_2, \dots, x_{D_i} >$    $< y_1, y_2, \dots, y_{D_o} >$
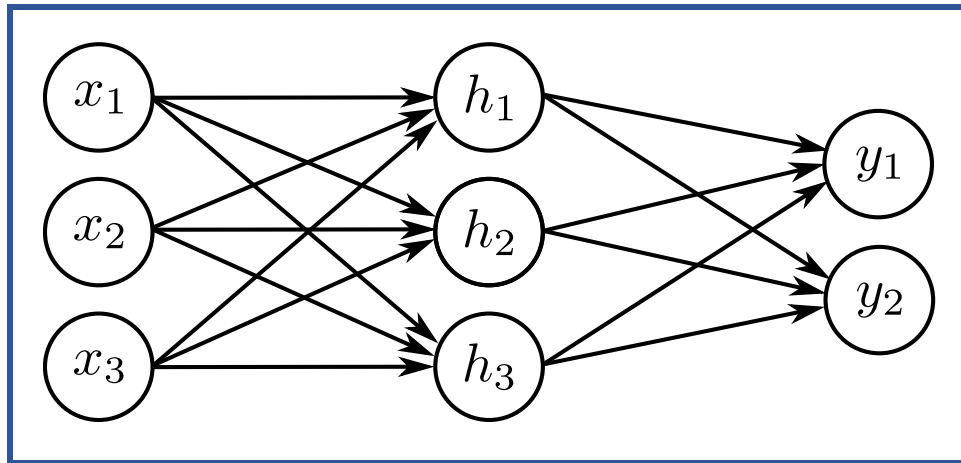
**Di-input/Do-output**

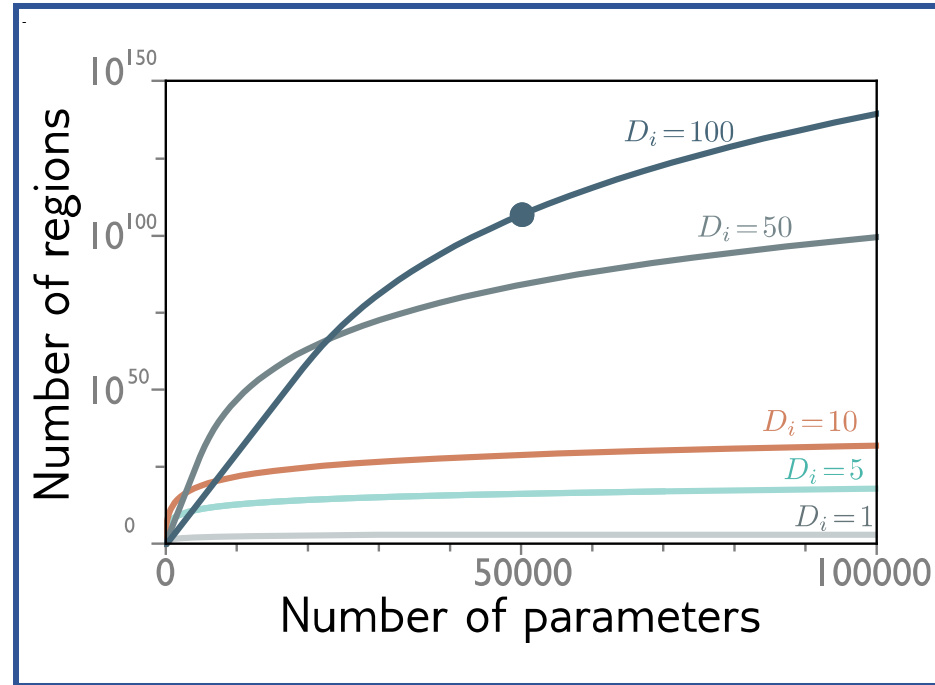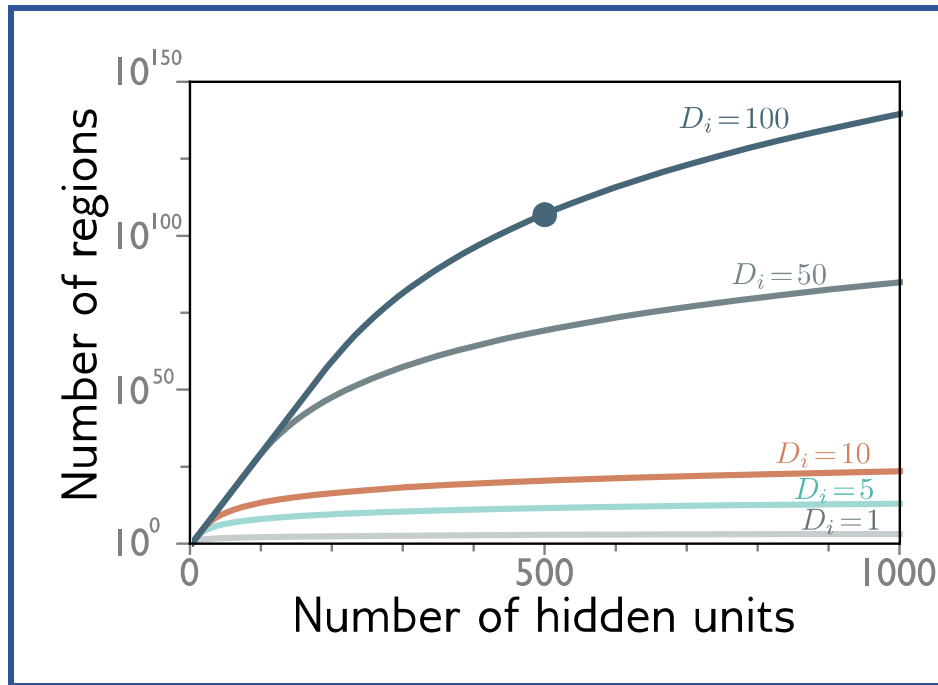❑ D$_i$-Inputs, D hidden units, D$_o$-Outputs:

$$h_d = a \left[ \theta_{d0} + \sum_{i=1}^{D_i} \theta_{di} x_i \right]$$

$$y_j = \phi_{j0} + \sum_{d=1}^{D} \phi_{jd} h_d$$

- e.g., Three inputs, three hidden units, two outputs



[Image from J.D Prince 2023 DL Book]

# Neural Networks

❏ #output regions vs #hidden units vs #parameters:

# Neural Networks

❑ Nomenclature:

<div style="background:pink">What is this neural network called?</div>



Input layer

Hidden layer

Output layer

$x_1$ $x_2$ $x_3$

$h_1$ $h_2$ $h_3$ $h_4$

$y_1$ $y_2$

Weight or parameter

Neuron or hidden unit
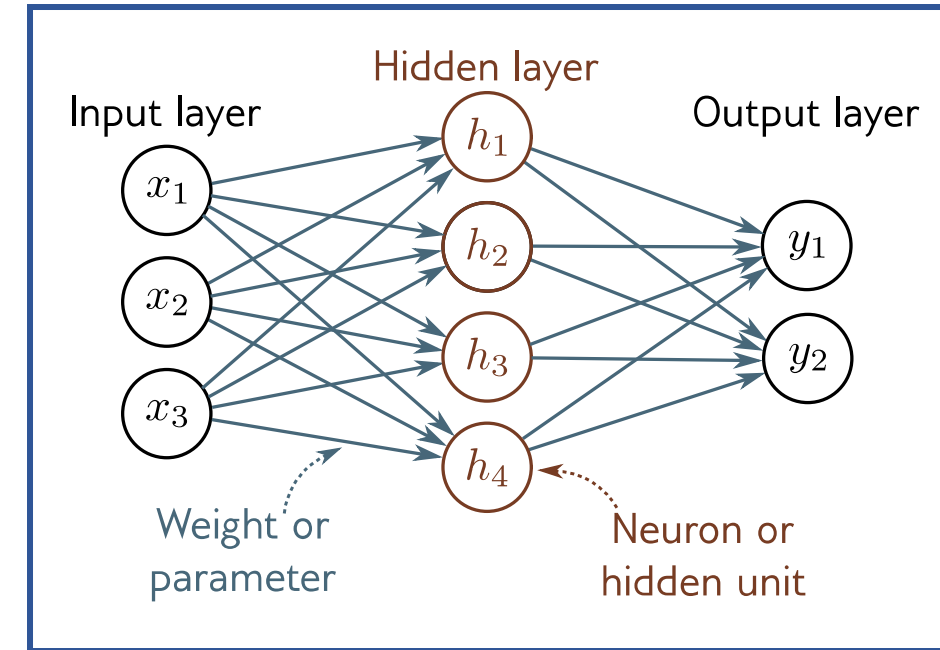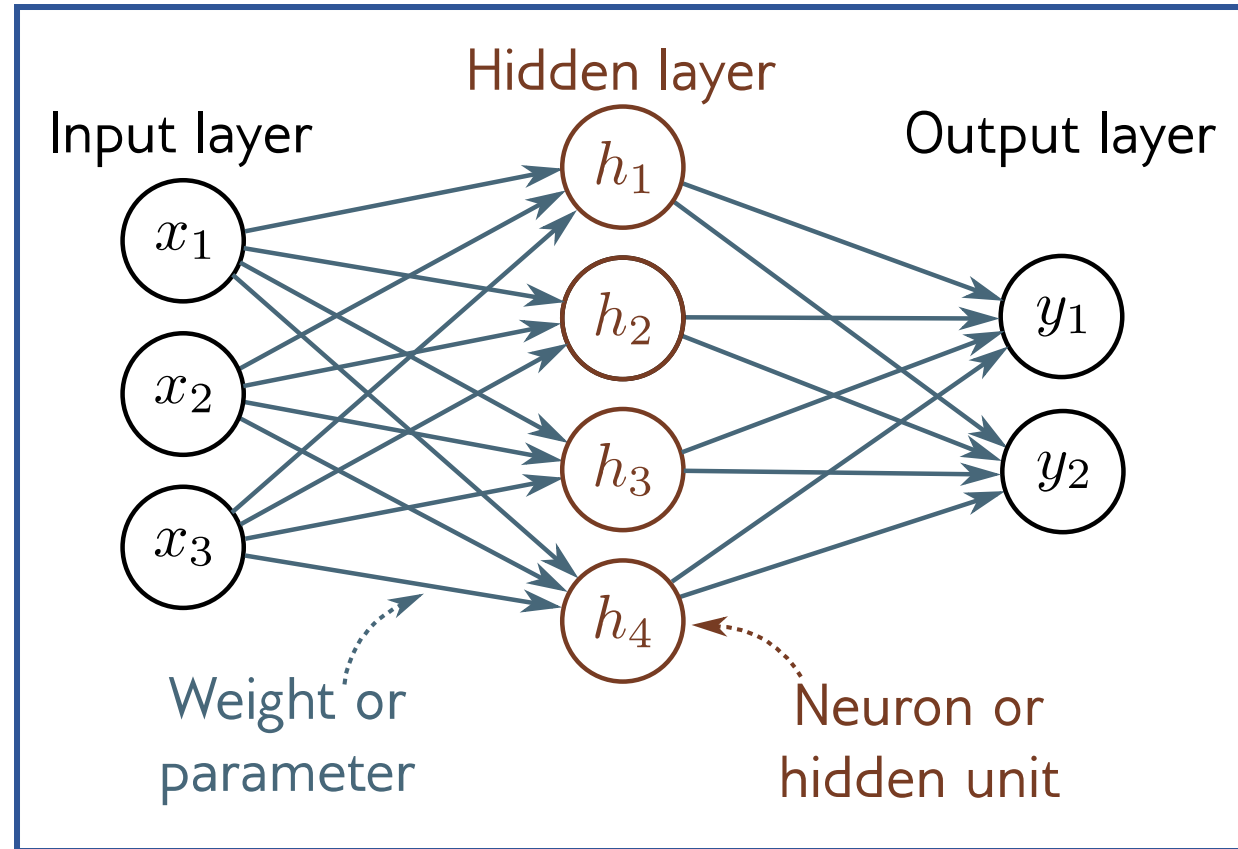
1. Single-layered neural network.

2. Two-layered neural network.

Two layers of learnable parameters (as per the Bishop 2024 DL Book)
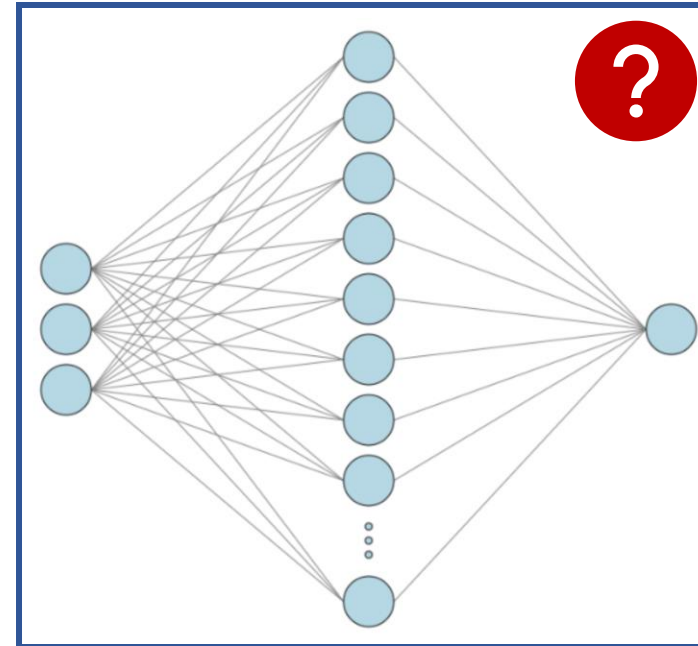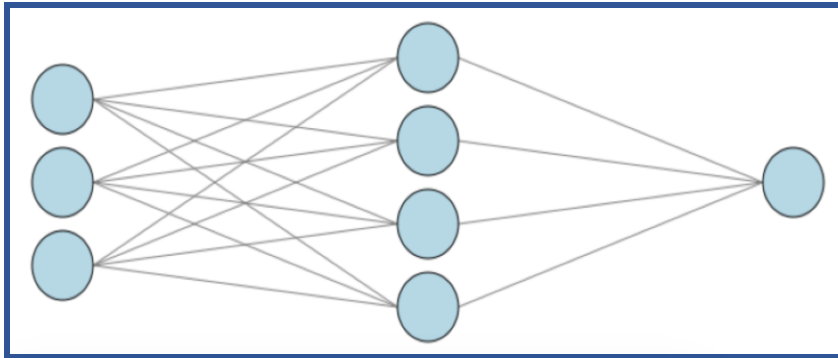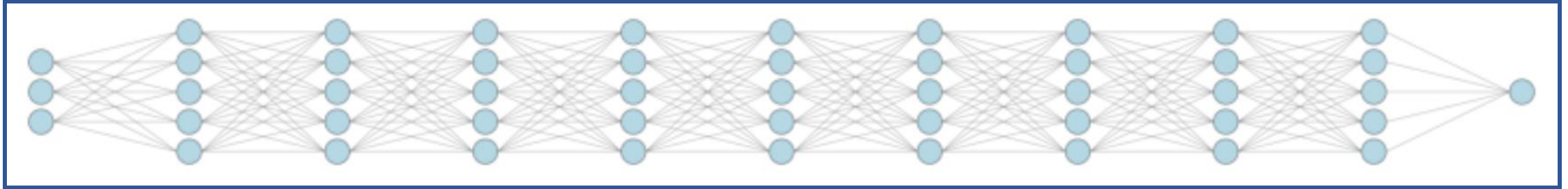
3. Three-layered neural network.

[Image from J.D Prince 2023 DL Book]

# Shallow Neural Networks



Input layer

Hidden layer

Output layer

$x_1$   $x_2$   $x_3$

$h_1$   $h_2$   $h_3$   $h_4$

$y_1$   $y_2$

Weight or parameter

Neuron or hidden unit

**What happens if we add more layers?**

# Deep Neural Network
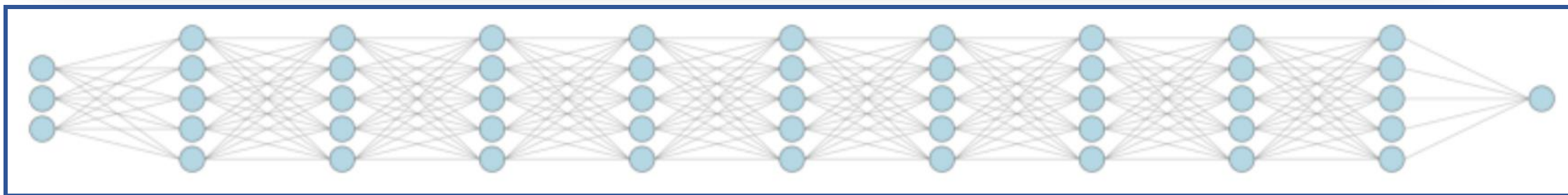


**What's different?**

**Number of parameters**

INDIAN INSTITUTE OF TECHNOLOGY DELHI
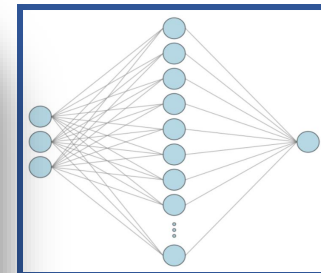
# Deep Neural Network

- ❑ A network with **two layers of learnable parameters - universal approximation** capabilities.

- ❑ A network with **more than two layers** – can represent a given function with **far fewer parameters**.



**Paper:** On the Number of Linear Regions of Deep Neural Networks. Mont´ufar et al. NeurIPS-2014

INDIAN INSTITUTE OF TECHNOLOGY DELHI