بسم الله الرحمن الرحيم

In the Name of Allah, the Most Gracious, the Most Merciful

قَالَ رَبِّ اشۡرَحۡ لِیۡ صَدۡرِیۡ ۙ ﴿۲۵﴾

۲۵ـ (موسیٰ علیہ السلام نے) عرض کیا: اے میرے رب! میرے لئے میرا سینہ کشادہ فرما دے۔

وَ یَسِّرۡ لِیۡۤ اَمۡرِیۡ ۙ ﴿۲۶﴾

۲۶ـ اور میرا کار (رسالت) میرے لئے آسان فرما دے۔

وَ احۡلُلۡ عُقۡدَۃً مِّنۡ لِّسَانِیۡ ﴿۲۷﴾

۲۷ـ اور میری زبان کی گرہ کھول دے۔

یَفۡقَہُوۡا قَوۡلِیۡ ۪ ﴿۲۸﴾

۲۸ـ کہ لوگ میری بات (آسانی سے) سمجھ سکیں۔

*Surah Taha with Urdu Translation*

رَبِّ اِنِّیۡ لِمَاۤ اَنۡزَلۡتَ اِلَیَّ مِنۡ خَیۡرٍ فَقِیۡرٌ
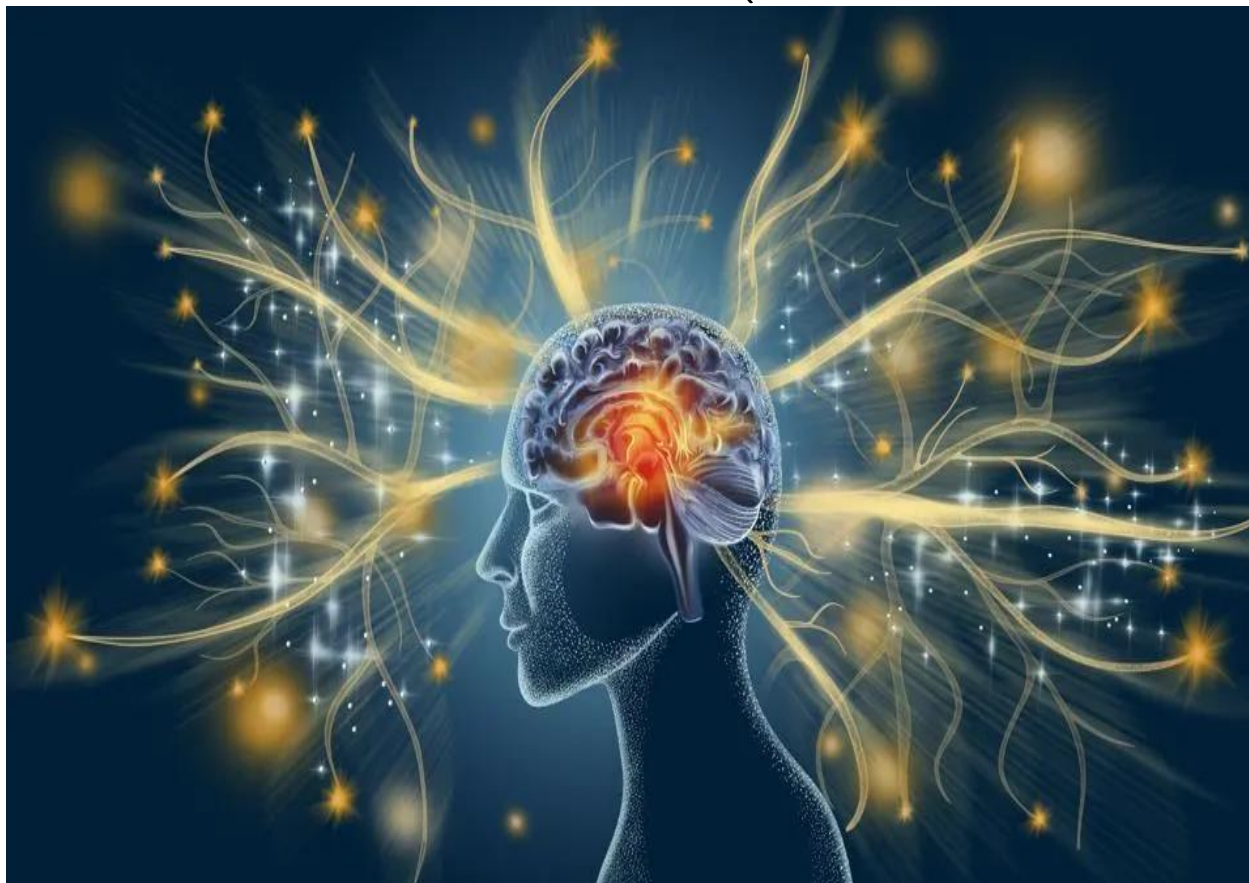
اے میرے رب! بے شک میں، جو بھلائی بھی تو میری طرف نازل فرمائے، اس کا محتاج ہوں۔

﴾سورۃ القصص: آیت نمبر 24 ﴿

# CS4152: Deep Learning and Neural Networks

Lecture 3  (Introduction to Neural Networks )

**Dr.  Jameel Ahmad**

Assistant Professor

Jameel.ahmad@umt.edu.p k

Department of Computer Science

School of Systems and Technology

University of Management and Technology, Lahore

# Lecture 3

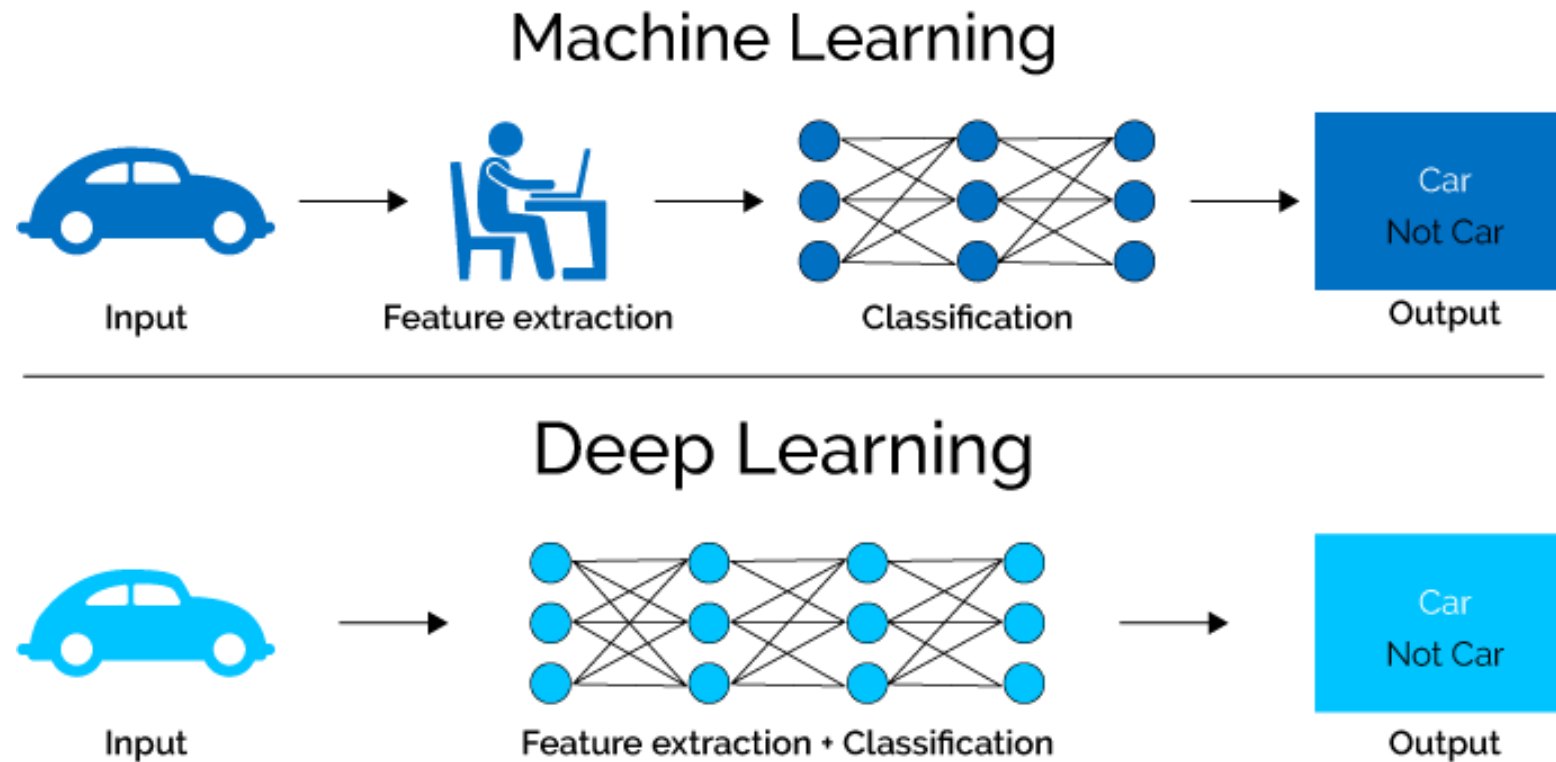**Neural Networks
An Overview
Jameel Ahmad, PhD**

# Lecture Outline

- Elements of neural networks (NNs)
  - Activation functions
  - Single and Multilayer Perceptron

# ML vs. Deep Learning (DL)

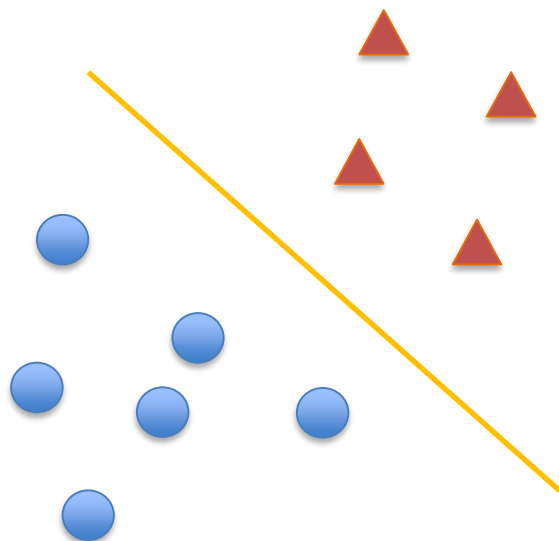*Difference between Machine Learning and Deep Learning*

- *Deep learning* (DL) is a machine learning subfield that uses multiple layers for learning data representations
  - DL is exceptionally effective at learning patterns

## Machine Learning

Input → Feature extraction → Classification → Output (Car / Not Car)

## Deep Learning

Input → Feature extraction + Classification → Output (Car / Not Car)

# Neural Network

Logistic Regression

Neural Networks

# Neural Network

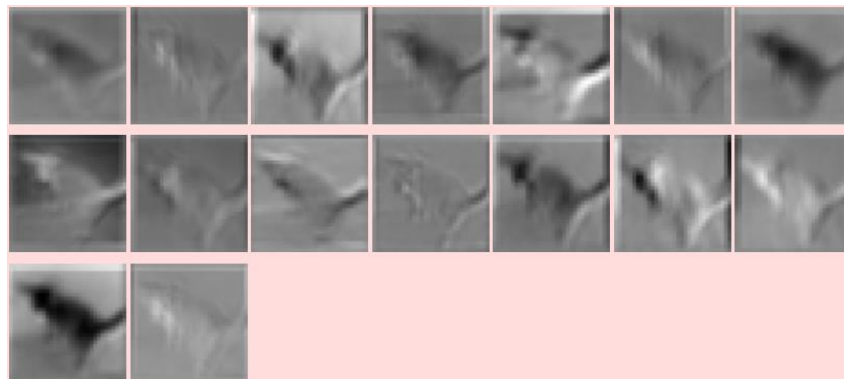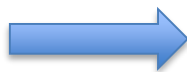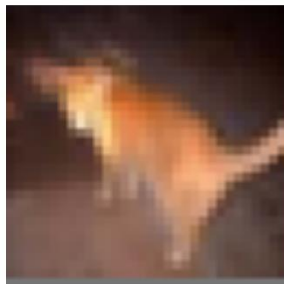- Non-linear score function $f = \ldots (\max(\mathbf{0}, \boldsymbol{W_1 x}))$
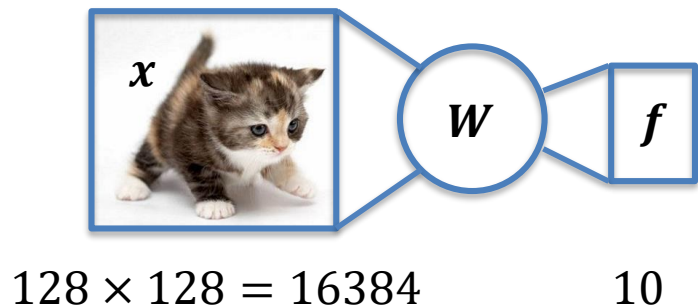


On CIFAR-10
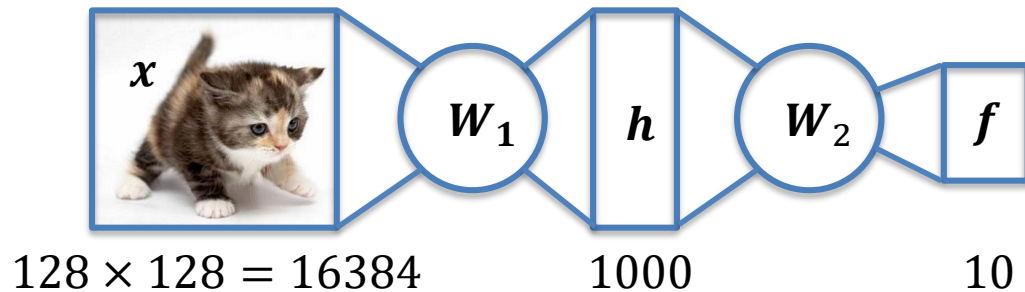


Visualizing activations of the first layer.

Source: ConvNetJS

# Neural Network

1-layer network: $f = Wx$

2-layer network: $f = W_2 \max(0, W_1 x)$



$128 \times 128 = 16384$     $10$

$128 \times 128 = 16384$     $1000$     $10$

Why is this structure useful?

# Neural Network

2-layer network: $f = W_2 \max(0, W_1 x)$



$128 \times 128 = 16384$      $1000$      $10$

Input Layer      Hidden Layer      Output Layer

input layer

hidden layer

output layer

# Net of Artificial Neurons

# Neural Network



Source: https://towardsdatascience.com/training-deep-neural-networks-9fdb1964b964

# Activation Functions

Sigmoid: $\sigma(x) = \dfrac{1}{(1+e^{-x})}$

tanh: $\tanh(x)$

ReLU: $\max(0, x)$

Leaky ReLU: $\max(0.1x, x)$

Parametric ReLU: $\max(\alpha x, x)$

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

ELU $f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$

# Neural Network

$$f = W_3 \cdot (W_2 \cdot (W_1 \cdot x)))$$

Why activation functions?

Simply concatenating linear layers would be so much cheaper...

# Neural Network

Why organize a neural network into layers?



input layer   hidden layer 1   hidden layer 2   hidden layer 3

output layer

# Biological Neurons



Credit: Stanford CS 231n

# Biological Neurons



Credit: Stanford CS 231n

# Artificial Neural Network



$x_1$

$x_2$

$x_3$

$f(W_{0,0}x + b_{0,0})$

$f(W_{0,1}x + b_{0,1})$

$f(W_{0,2}x + b_{0,2})$

$f(W_{0,3}x + b_{0,3})$

$f(W_{1,0}x + b_{1,0})$

$f(W_{1,1}x + b_{1,1})$

$f(W_{1,2}x + b_{1,2})$

$f(W_{2,0}x + b_{2,0})$

# Neural Networks

- Handwritten digit recognition (MNIST dataset)
  - The intensity of each pixel is considered an input element
  - Output is the class of the digit

**Input**

$x_1$

$x_2$

$\vdots$

$x_{256}$

16 x 16 = 256

Ink → 1

No ink → 0

**Output**

| 0.1 | is 1 |

| 0.7 | is 2 |

$\vdots$

| 0.2 | is 0 |

The image is "2"

Each dimension represents the confidence of a digit

# Elements of Neural Networks

- NNs consist of hidden layers with neurons (i.e., computational units)
- A single neuron maps a set of inputs into an output number, or $f: R^K \rightarrow R$

$$z = a_1 w_1 + a_2 w_2 + \cdots + a_K w_K + b$$

$$a = \sigma(z)$$

$a_1$

$a_2$

$\vdots$

$a_K$

input

$w_1$

$w_2$

$\vdots$

$w_K$

weights

$+$

$z$

$b$

bias

$\sigma(z)$

Activation function

$a$

output

# Elements of Neural Networks

*Introduction to Neural Networks*

- A NN with one hidden layer and one output layer



Weights    Biases

$$hidden\ layer\ h = \sigma(W_1 x + b_1)$$

$$output\ layer\ y = \sigma(W_2 h + b_2)$$

Activation functions

4 + 2 = 6 neurons (not counting inputs)
[3 × 4] + [4 × 2] = 20 weights
4 + 2 = 6 biases

26 learnable parameters

# Elements of Neural Networks

*Introduction to Neural Networks*

- A neural network playground [link](#)

# Elements of Neural Networks

*Introduction to Neural Networks*

- Deep NNs have many hidden layers
  - Fully-connected (dense) layers (a.k.a. Multi-Layer Perceptron or MLP)
  - Each neuron is connected to all neurons in the succeeding layer

# Elements of Neural Networks

*Introduction to Neural Networks*

- A simple network, toy example

$(1 \cdot 1) + (-1) \cdot (-2) + 1 = 4$

1    4    **0.98**

**1**

-2

1

## Sigmoid Function

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

$\sigma(z)$

1

0.5

0

$z$

-1    -2    **0.12**

**-1**

1

0

$1 \cdot (-1) + (-1) \cdot 1 + 0 = -2$

# Elements of Neural Networks

*Introduction to Neural Networks*

- A simple network, toy example (cont'd)
  - For an input vector $[1 \quad -1]^T$, the output is $[0.62 \quad 0.83]^T$



$$f: R^2 \to R^2 \qquad f\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right) = \begin{bmatrix} 0.62 \\ 0.83 \end{bmatrix}$$

# Matrix Operation

- Matrix operations are helpful when working with multidimensional inputs and outputs

$$\sigma( \boxed{W} \boxed{x} + \boxed{b} ) = \boxed{a}$$

$$\sigma( \begin{bmatrix} 1 & -2 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix}) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix}$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxx}}$$

$$\begin{bmatrix} 4 \\ -2 \end{bmatrix}$$

# Matrix Operation

- Multilayer NN, matrix calculations for the first layer
  - Input vector $x$, weights matrix $W^1$, bias vector $b^1$, output vector $a^1$



$$a^1 = \sigma(W^1 x + b^1)$$

# Matrix Operation

- Multilayer NN, matrix calculations for all layers



$$\sigma(\ W^1\ x\ +\ b^1\ )$$

$$\sigma(\ W^2\ a^1\ +\ b^2\ )$$

$$\sigma(\ W^L\ X\ +\ b^L\ )$$

# Matrix Operation

- Multilayer NN, function $f$ maps inputs $x$ to outputs $y$, i.e., $y = f(x)$



$$y = f(x) = \sigma(W^L \cdots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \cdots + b^L)$$

# Softmax Layer

- In multi-class classification tasks, the output layer is typically a *softmax layer*
  - I.e., it employs a *softmax activation function*
  - If a layer with a sigmoid activation function is used as the output layer instead, the predictions by the NN may not be easy to interpret
    - Note that an output layer with sigmoid activations can still be used for binary classification

### A Layer with Sigmoid Activations

$$z_1 \xrightarrow{\quad 3 \quad} \sigma \xrightarrow{\quad 0.95 \quad} y_1 = \sigma(z_1)$$

$$z_2 \xrightarrow{\quad 1 \quad} \sigma \xrightarrow{\quad 0.73 \quad} y_2 = \sigma(z_2)$$

$$z_3 \xrightarrow{\quad -3 \quad} \sigma \xrightarrow{\quad 0.05 \quad} y_3 = \sigma(z_3)$$

# Softmax Layer

- The softmax layer applies softmax activations to output a probability value in the range [0, 1]
  - The values $z$ inputted to the softmax layer are referred to as *logits*

**Probability**:
- $0 < y_i < 1$
- $\sum_i y_i = 1$

### A Softmax Layer



**3**

$z_1$   $e$   $e^{z_1}$   **20**   $\div$   **0.88**   $y_1 = e^{z_1} \Big/ \sum_{j=1}^{3} e^{z_j}$

**1**

$z_2$   $e$   $e^{z_2}$   **2.7**   $\div$   **0.12**   $y_2 = e^{z_2} \Big/ \sum_{j=1}^{3} e^{z_j}$

**-3**

$z_3$   $e$   $e^{z_3}$   **0.05**   $\div$   **≈0**   $y_3 = e^{z_3} \Big/ \sum_{j=1}^{3} e^{z_j}$

$+ \;\; \sum_{j=1}^{3} e^{z_j}$

# Activation Functions

- <span style="color:red">Non-linear activations</span> are needed to learn complex (non-linear) data representations
  - Otherwise, NNs would be just a linear function (such as $W_1 W_2 x = W x$)
  - NNs with large number of layers (and neurons) can approximate more complex functions
    - Figure: more neurons improve representation (but, may overfit)



3 hidden neurons      6 hidden neurons      20 hidden neurons

# Activation: Sigmoid

- *Sigmoid function* σ: takes a real-valued number and "squashes" it into the range between 0 and 1
    - The output can be interpreted as the firing rate of a biological neuron
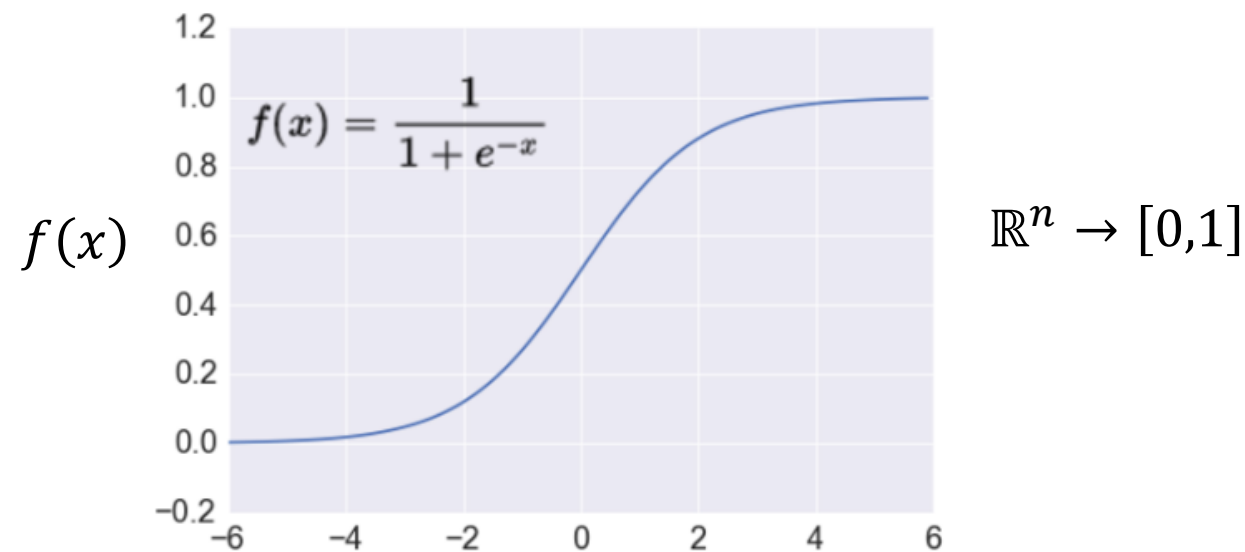        - Not firing = 0; Fully firing = 1
    - When the neuron's activation are 0 or 1, sigmoid neurons saturate
        - Gradients at these regions are almost zero (almost no signal will flow)
    - Sigmoid activations are less common in modern NNs

$f(x)$

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

$\mathbb{R}^n \rightarrow [0,1]$

# Activation: Tanh

- *Tanh function*: takes a real-valued number and "squashes" it into range between -1 and 1
  - Like sigmoid, tanh neurons saturate
  - Unlike sigmoid, the output is zero-centered
    - It is therefore preferred than sigmoid
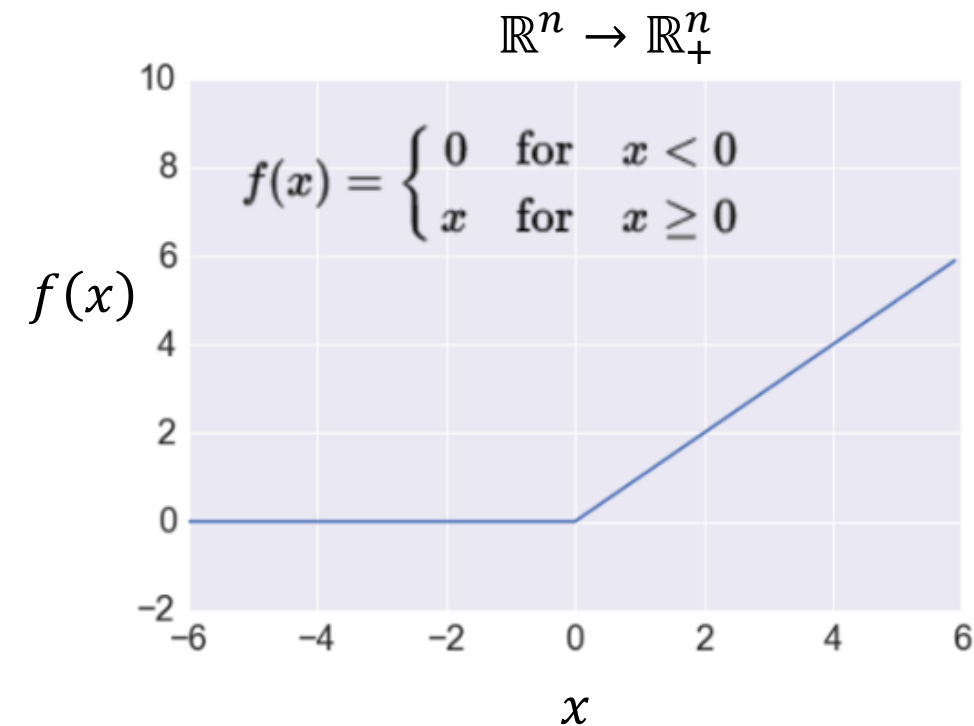  - Tanh is a scaled sigmoid: $\tanh(x) = 2 \cdot \sigma(2x) - 1$

$f(x)$

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

$$\mathbb{R}^n \rightarrow [-1,1]$$

# Activation: ReLU

- **ReLU** (Rectified Linear Unit): takes a real-valued number and thresholds it at zero

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

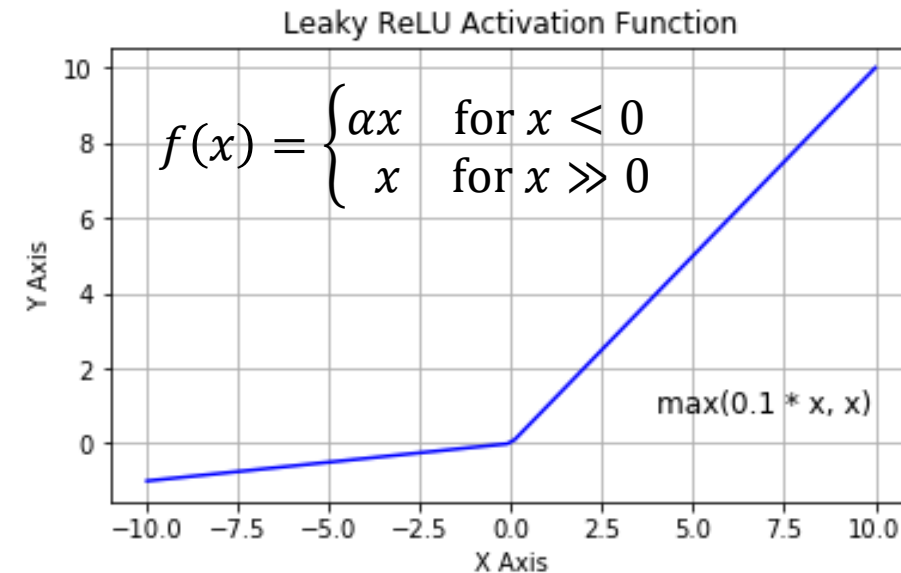$$\mathbb{R}^n \rightarrow \mathbb{R}^n_+$$

- Most modern deep NNs use ReLU activations
- ReLU is fast to compute
  - Compared to sigmoid, tanh
  - Simply threshold a matrix at zero
- Accelerates the convergence of gradient descent
  - Due to linear, non-saturating form
- Prevents the gradient vanishing problem

$$f(x) = \begin{cases} 0 & \text{for} \quad x < 0 \\ x & \text{for} \quad x \geq 0 \end{cases}$$

# Activation: Leaky ReLU

*Introduction to Neural Networks*

- The problem of ReLU activations: they can "die"
  - ReLU could cause weights to update in a way that the gradients can become zero and the neuron will not activate again on any data
  - E.g., when a large learning rate is used

- *Leaky ReLU* activation function is a variant of ReLU
  - Instead of the function being 0 when $x < 0$, a leaky ReLU has a small negative slope (e.g., $\alpha = 0.01$, or similar)

  - This resolves the dying ReLU problem
  - Most current works still use ReLU
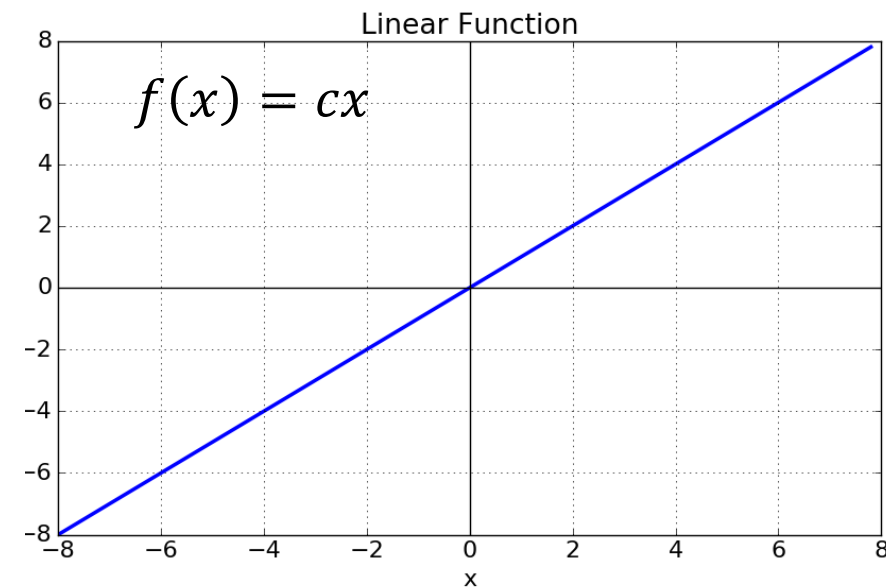    - With a proper setting of the learning rate, the problem of dying ReLU can be avoided

Leaky ReLU Activation Function

$$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \gg 0 \end{cases}$$

max(0.1 * x, x)

Y Axis

X Axis

# Activation: Linear Function

- *Linear function* means that the output signal is proportional to the input signal to the neuron

$$\mathbb{R}^n \to \mathbb{R}^n$$

  - If the value of the constant $c$ is 1, it is also called <span style="color:red">identity activation function</span>
  - This activation type is used in regression problems
    - E.g., the last layer can have linear activation function, in order to output a real number (and not a class membership)



Linear Function

$$f(x) = cx$$

# References

1. Hung-yi Lee – Deep Learning Tutorial
2. Ismini Lourentzou – Introduction to Deep Learning
3. CS231n Convolutional Neural Networks for Visual Recognition (Stanford CS course) ([link](#))
4. James Hays, Brown – Machine Learning Overview
5. Param Vir Singh, Shunyuan Zhang, Nikhil Malik – Deep Learning
6. Sebastian Ruder – An Overview of Gradient Descent Optimization Algorithms ([link](#))