# Deep Learning

Session 3

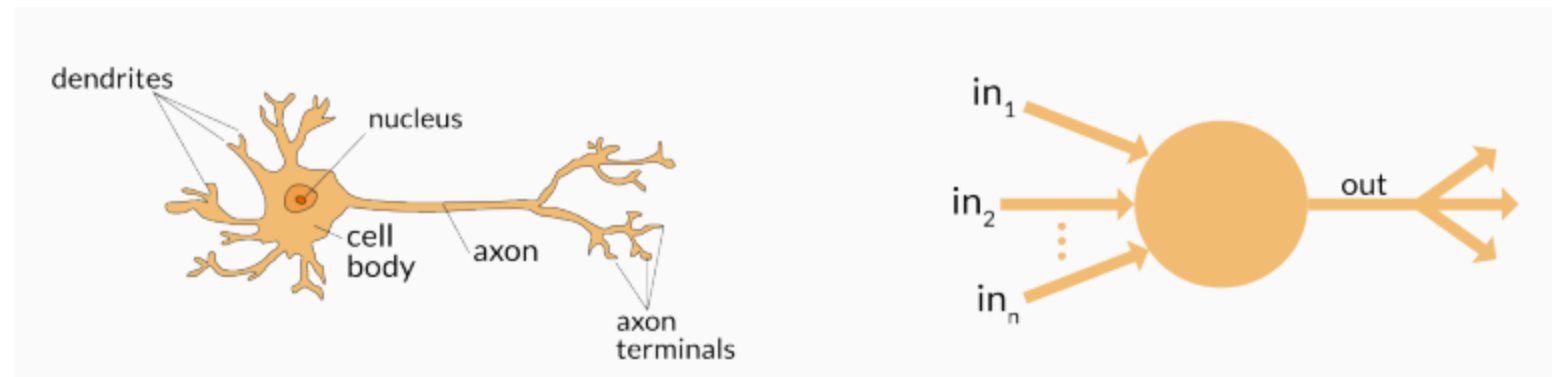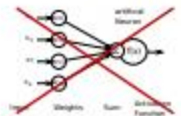## The Perceptron

**Applied Data Science**

**2024/2025**

# Perceptron



1958 Perceptron

1969 Perceptrons book

Perceptron criticized



dendrites
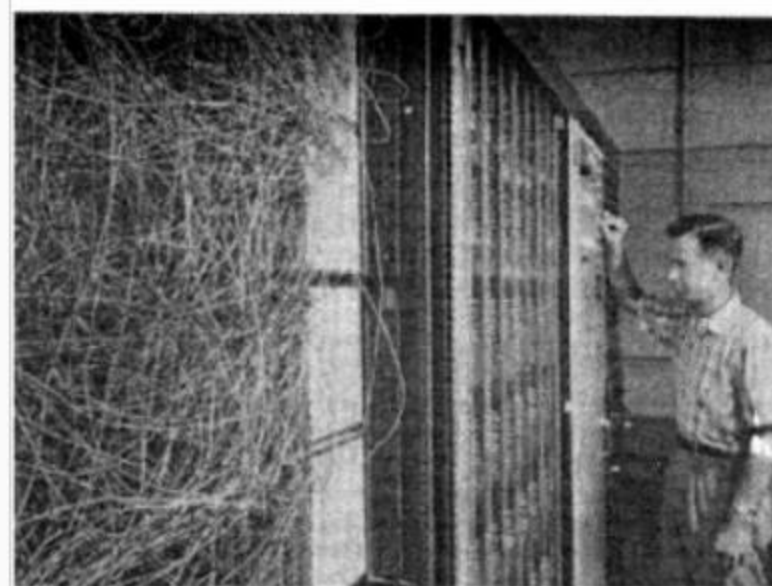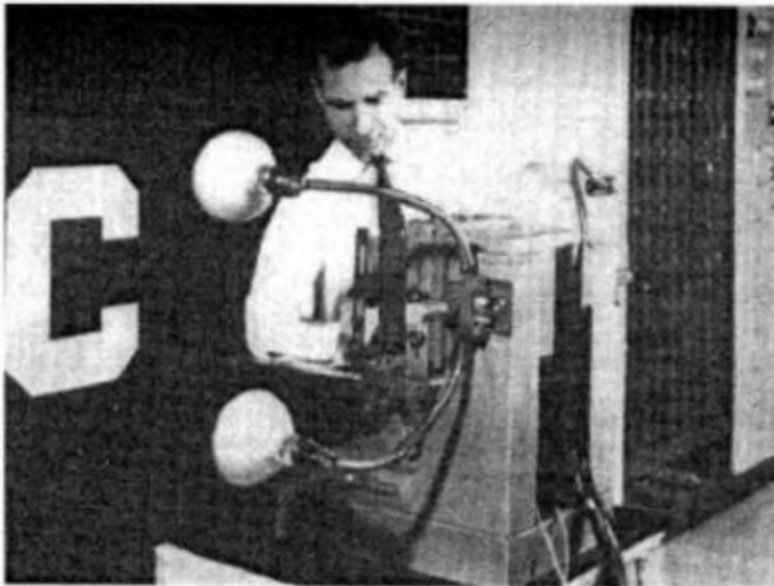
nucleus

cell body

axon

axon terminals

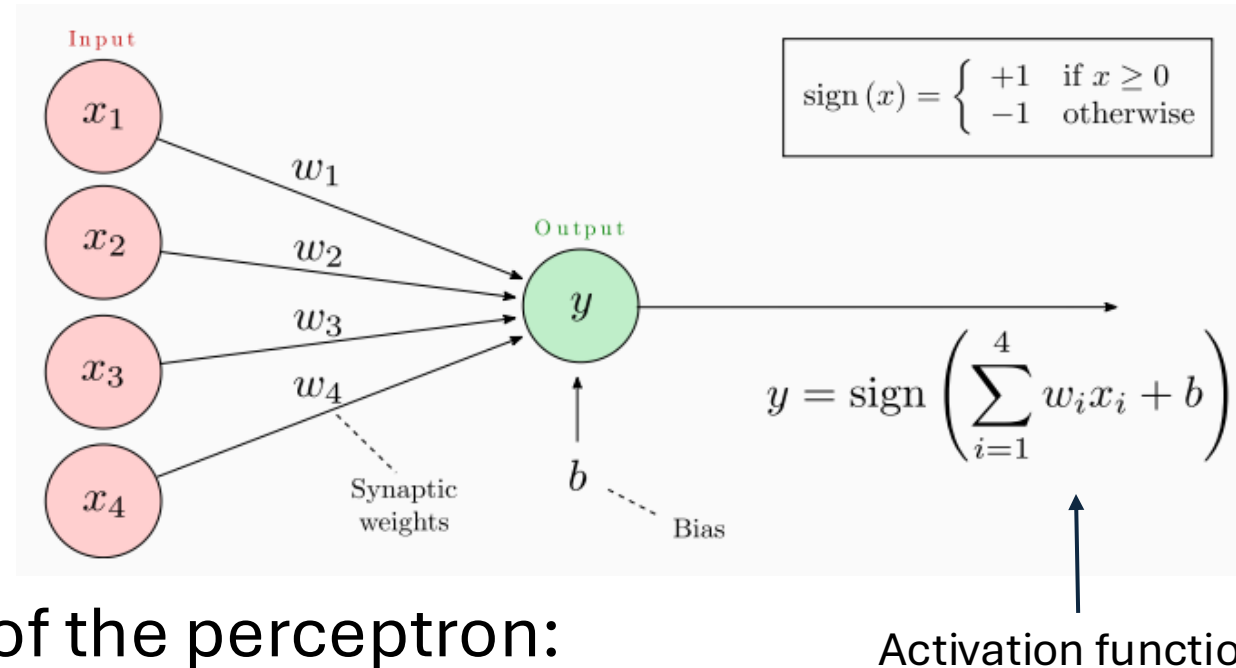$in_1$

$in_2$

$in_n$

out

# Perceptron

- First binary classifier based on supervised learning;

- Foundation of modern artificial neural networks;

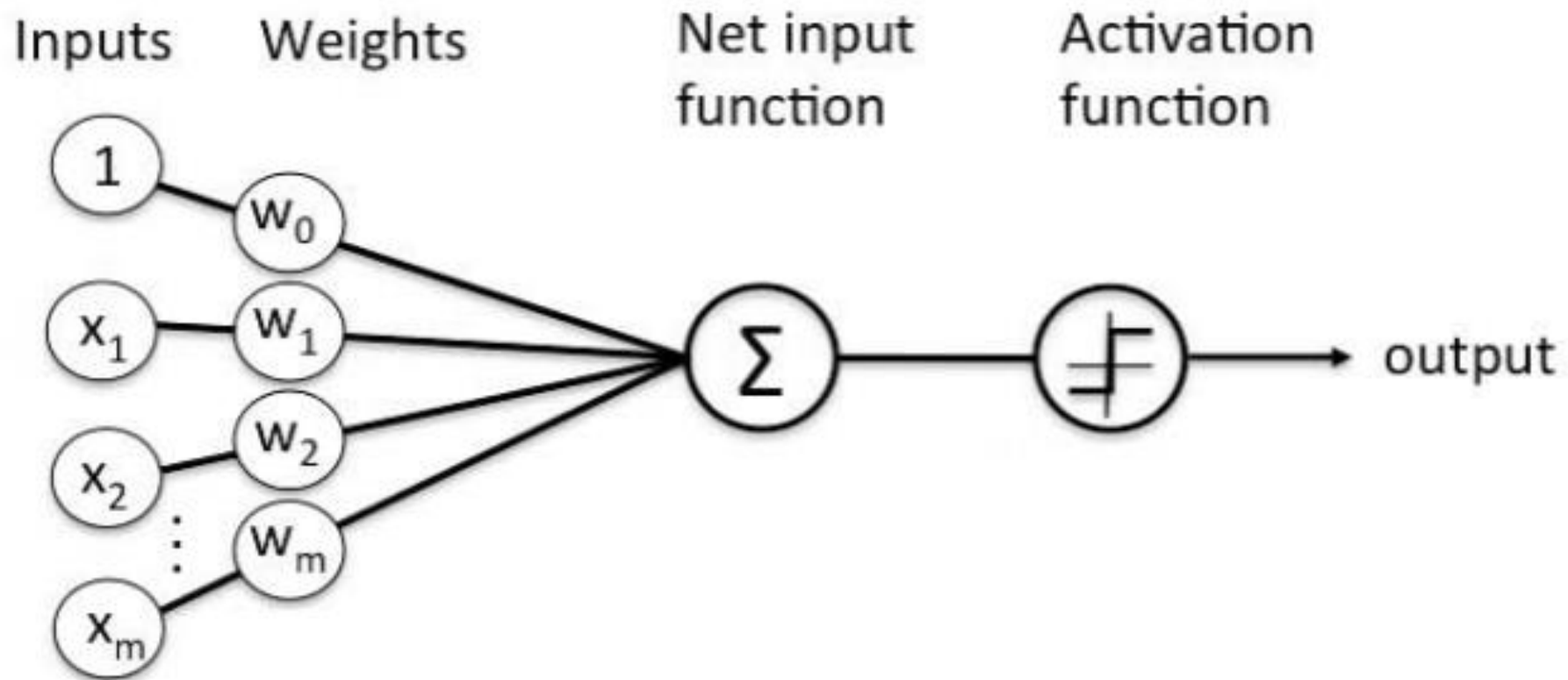

Perceptron (Frank Rosenblatt, 1958)
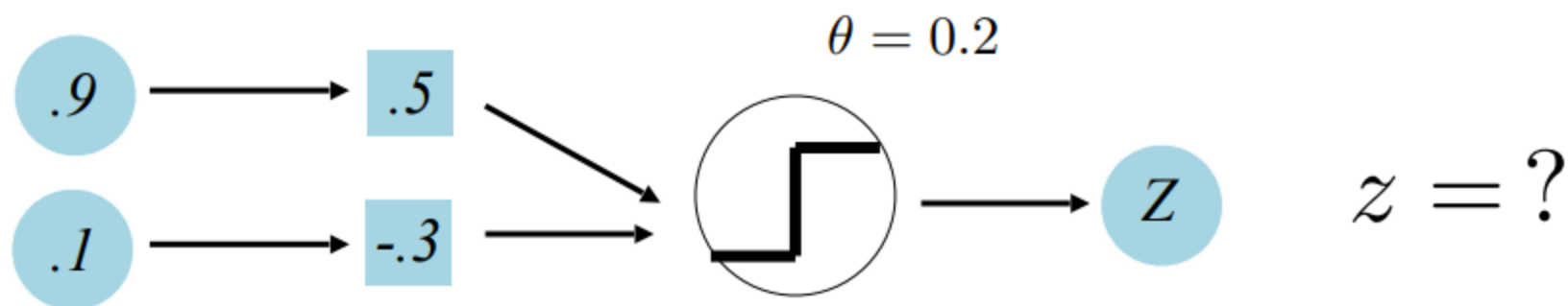
# Representation of the Perceptron

Input

$x_1$

$x_2$

$x_3$

$x_4$

$w_1$

$w_2$

$w_3$

$w_4$

Output

$y$

$$\text{sign}(x) = \begin{cases} +1 & \text{if } x \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$y = \text{sign}\left(\sum_{i=1}^{4} w_i x_i + b\right)$$

Synaptic weights

$b$

Bias

Activation function

- Parameters of the perceptron:
  - $w_k$: weights
  - b: bias

- Training → adjusting the weights and bias.

# Alternative Representation of the Perceptron
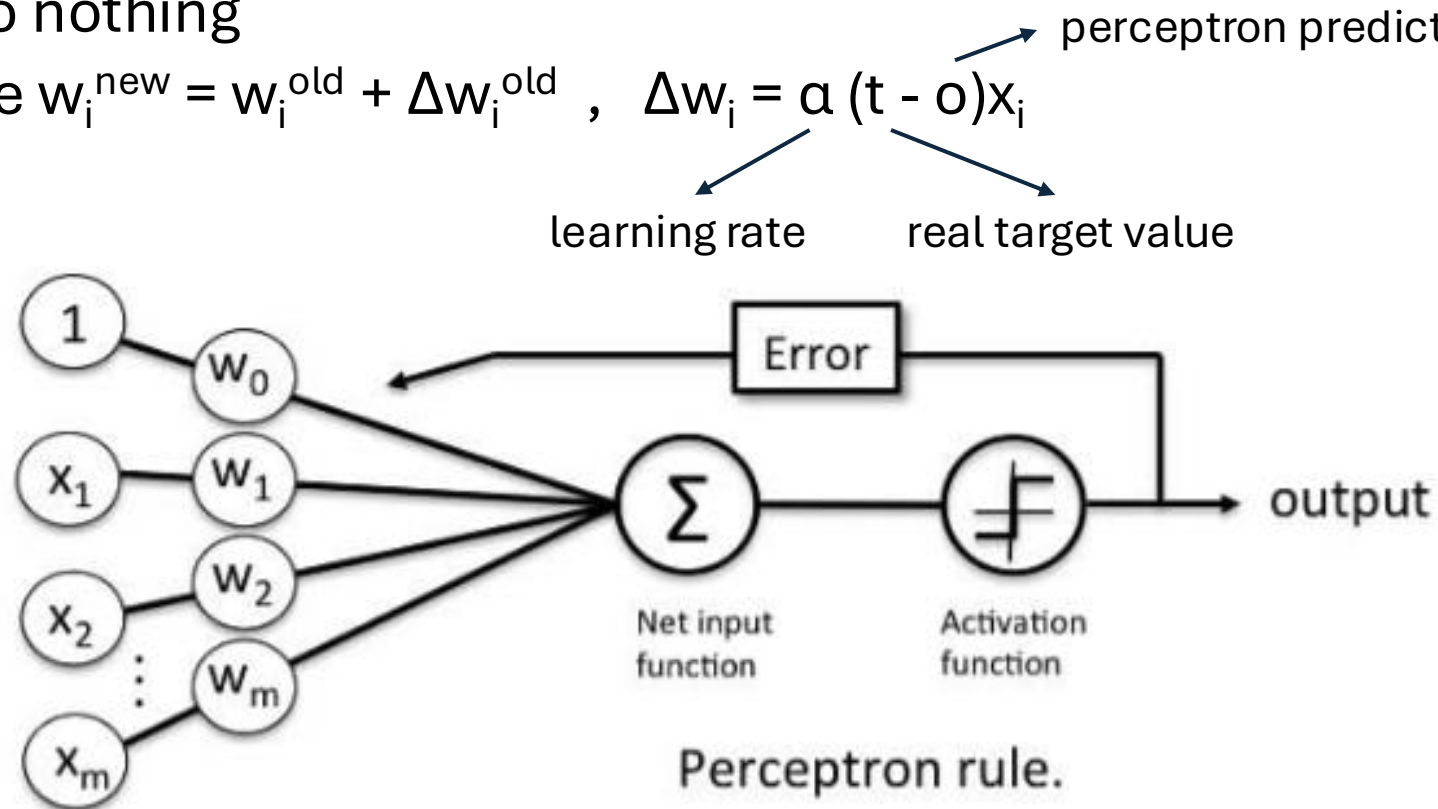
# Perceptron Examples

$$\theta = 0.2$$

$$z = ?$$

$$z = \begin{array}{l} 1, if\ w \cdot x > \theta \\ 0,\ w \cdot x <= \theta \end{array}$$

$$\theta = 0.1$$

$$z = ?$$

# Perceptron Learning Rule

- Suppose that x is a feature vector, y is the correct class label, and y' is the predicted class label computed using the current weights.
  - If y' = y, do nothing
  - Otherwise $w_i^{new} = w_i^{old} + \Delta w_i^{old}$ , $\Delta w_i = \alpha (t - o)x_i$

perceptron prediction

learning rate          real target value



Perceptron rule.

# Activation Functions

- The purpose of activation functions is to **introduce non-linearities** into the network.
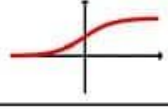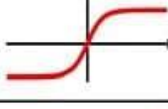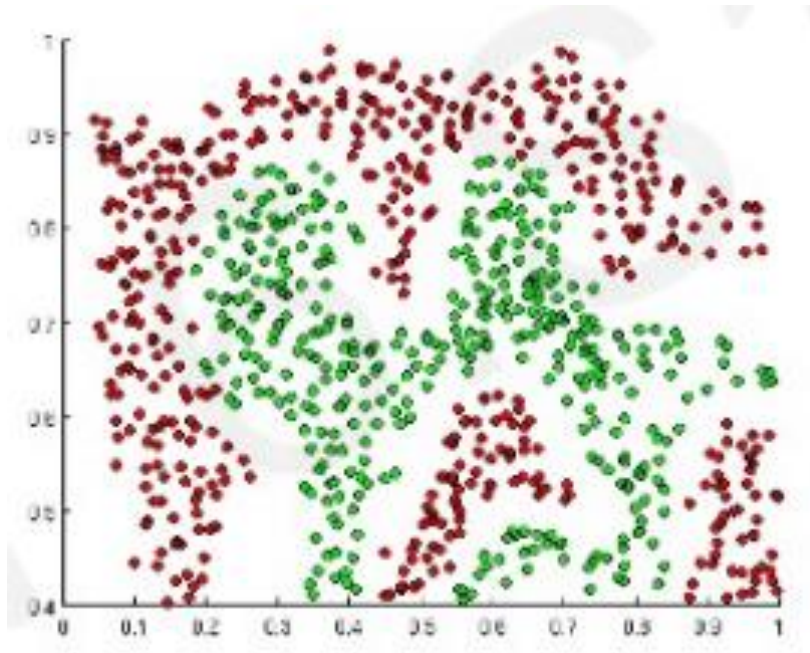
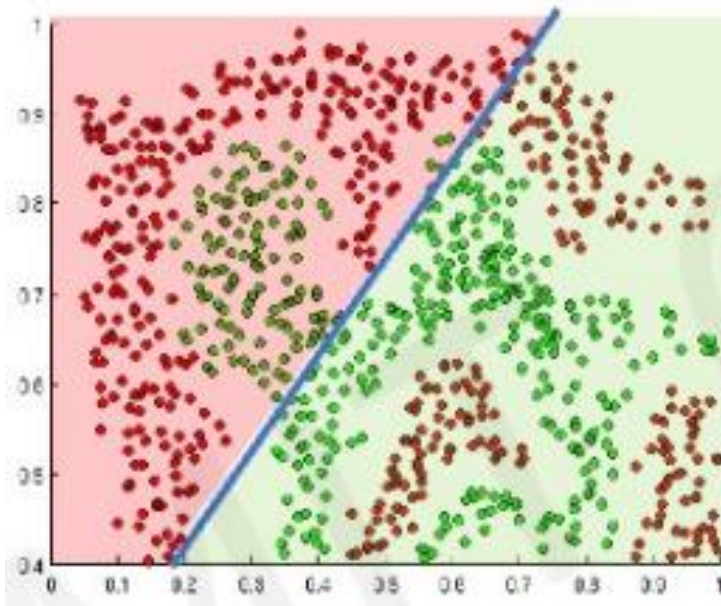| Activation Function | Equation | Example | 1D Graph |
|---|---|---|---|
| Linear | $\phi(z) = z$ | Adaline, linear regression | |
| Unit Step (Heaviside Function) | $\phi(z) = \begin{cases} 0 & z < 0 \\ 0.5 & z = 0 \\ 1 & z > 0 \end{cases}$ | Perceptron variant | |
| Sign (signum) | $\phi(z) = \begin{cases} -1 & z < 0 \\ 0 & z = 0 \\ 1 & z > 0 \end{cases}$ | Perceptron variant | |
| Piece-wise Linear | $\phi(z) = \begin{cases} 0 & z \leq -\frac{1}{2} \\ z + \frac{1}{2} & -\frac{1}{2} \leq z \leq \frac{1}{2} \\ 1 & z \geq \frac{1}{2} \end{cases}$ | Support vector machine | |
| Logistic (sigmoid) | $\phi(z) = \dfrac{1}{1 + e^{-z}}$ | Logistic regression, Multilayer NN | |
| Hyperbolic Tangent (tanh) | $\phi(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$ | Multilayer NN, RNNs | |
| ReLU | $\phi(z) = \begin{cases} 0 & z < 0 \\ z & z > 0 \end{cases}$ | Multilayer NN, CNNs | |

# Activation Functions

- The purpose of activation functions is to **introduce non-linearities** into the network.



What if we wanted to build a neural network
to distinguish green vs red points?

# Activation Functions

- The purpose of activation functions is to **introduce non-linearities** into the network.



Linear activations produce linear decisions
no matter the network size.
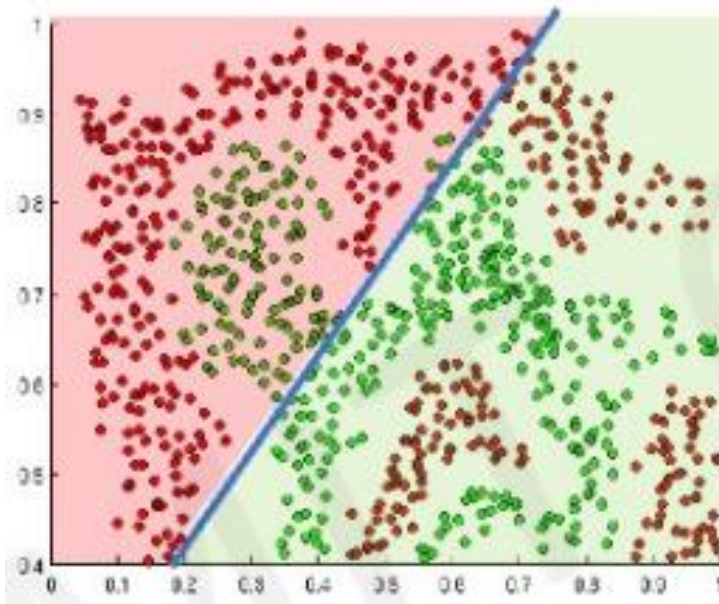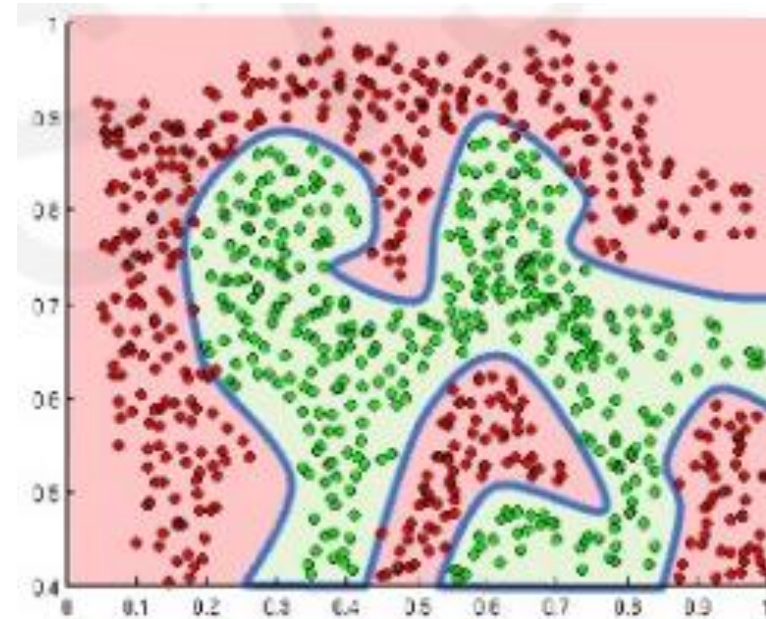
# Activation Functions

- The purpose of activation functions is to **introduce non-linearities** into the network.



Linear activations produce linear decisions no matter the network size.



Non-linearities allow us to approximate arbitrarly complex functions.
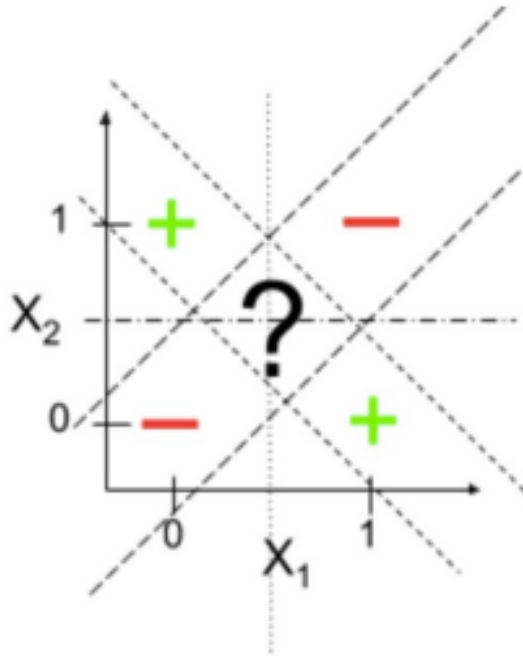
# Perceptron Limitation: XOR Problem

- XOR = "Exclusive Or"
  - Input: two binary values $x_1$ and $x_2$
  - Output:
    - 1, when exactly one input equals 1
    - 0, otherwise

| $x_1$ | $x_2$ | $x_1$ XOR $x_2$ |
|-------|-------|-----------------|
| 0 | 0 | ? |
| 0 | 1 | ? |
| 1 | 0 | ? |
| 1 | 1 | ? |

# Perceptron Limitation: XOR Problem

- Cannot solve XOR problem and so separate 1s from 0s with a perceptron (linear function):



| $x_1$ | $x_2$ | $x_1$ XOR $x_2$ |
|-------|-------|-----------------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

# Multilayer Perceptron



1958 Perceptron

1969
Perceptrons
book

~1980
Multilayer
network

1989
Universal
Approximation
Theorem

Perceptron criticized

# Mulilayer Perceptron

- AKA Artificial Neural Networks

# Mulilayer Perceptron

input layer

hidden layer 1   hidden layer 2

output layer

- How does this relate to a perceptron?



- Unit: takes as input a weighted sum and applies an activation function

# Mulilayer Perceptron

input layer

hidden layer 1    hidden layer 2

output layer

- How does this relate to a perceptron?



1

$x_1$    $w_0$

$x_2$    $w_1$

         $w_2$

$x_m$    $w_m$

Weight update    Error

$\Sigma$

Net input function

Threshold function

Output

- Unit: takes as input a weighted sum and applies an activation function

# Mulilayer Perceptron



input layer

hidden layer 1    hidden layer 2

output layer
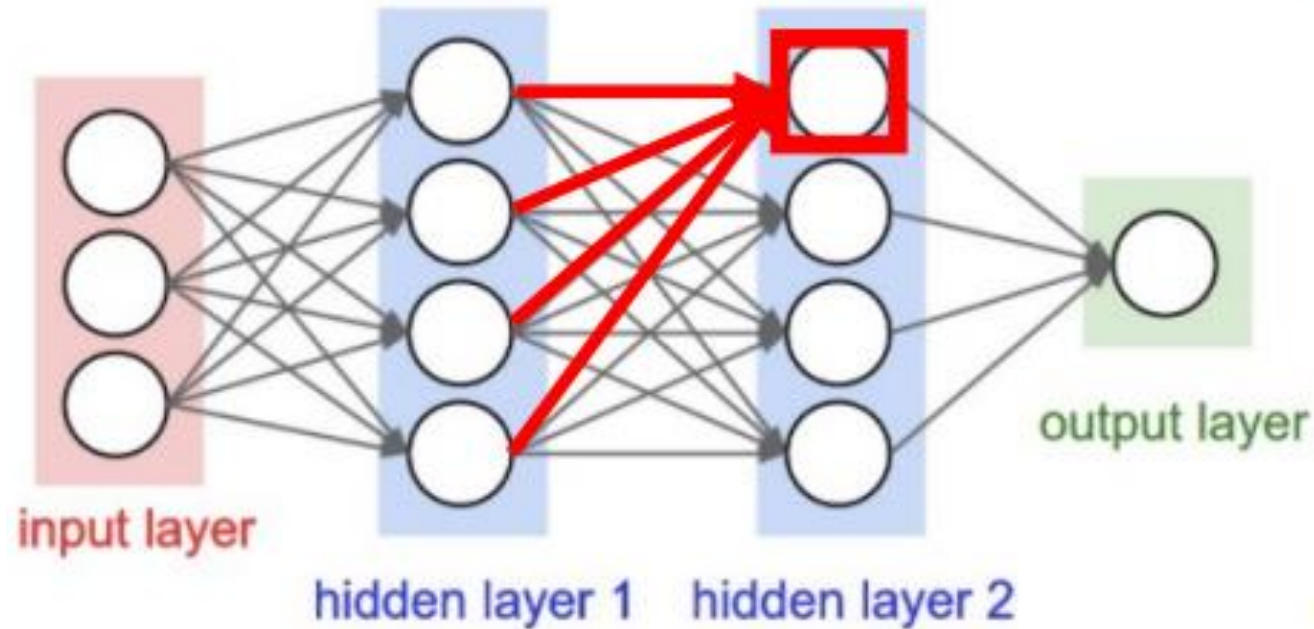
- How does this relate to a perceptron?



- Unit: takes as input a weighted sum and applies an activation function

# Artificial Neural Networks
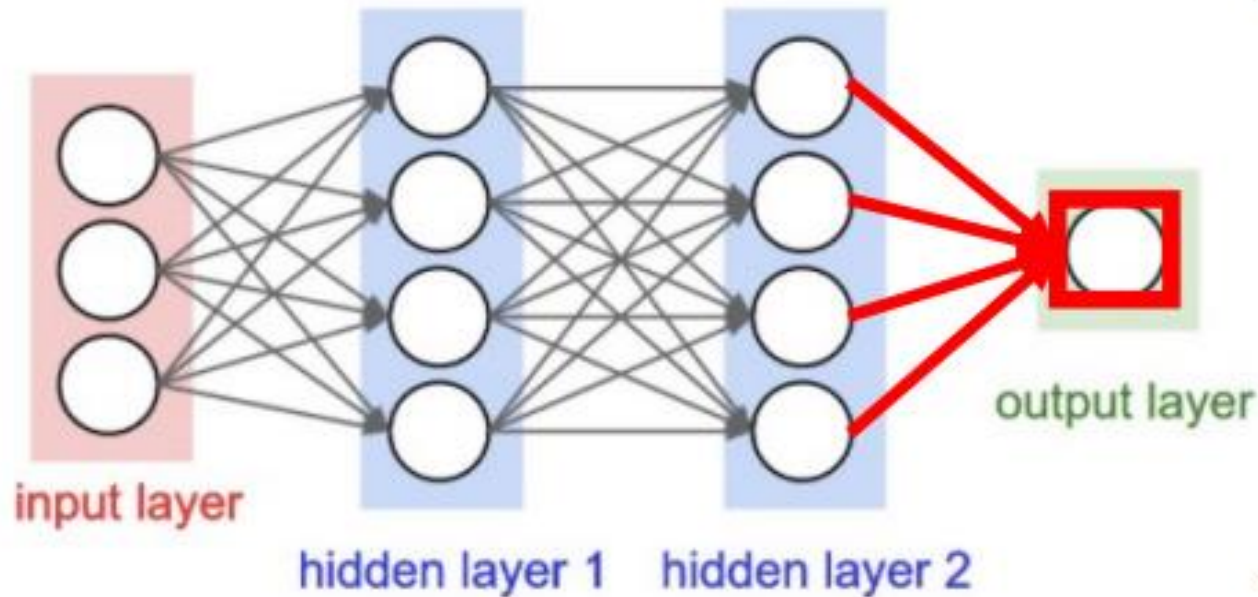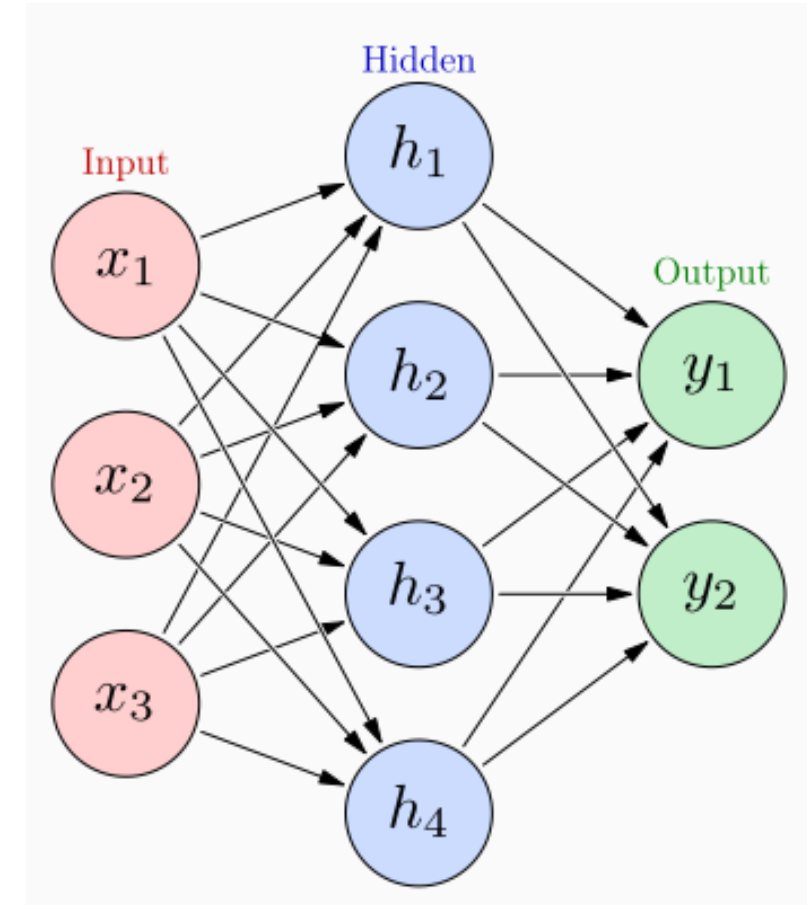
- Inter-connection of several artificial neurons (also called nodes or units);

- Each "level" in the graph is called a lauer:
  - Input layer;
  - Hidden layer(s);
  - Output layer.

- Each neuron in the hidden layers acts as a classifier / feature detector;

- Fedforward neural network (no cycles):
  - First na simplest type of neural network;
  - Information moves in one direction.

# Artificial Neural Networks

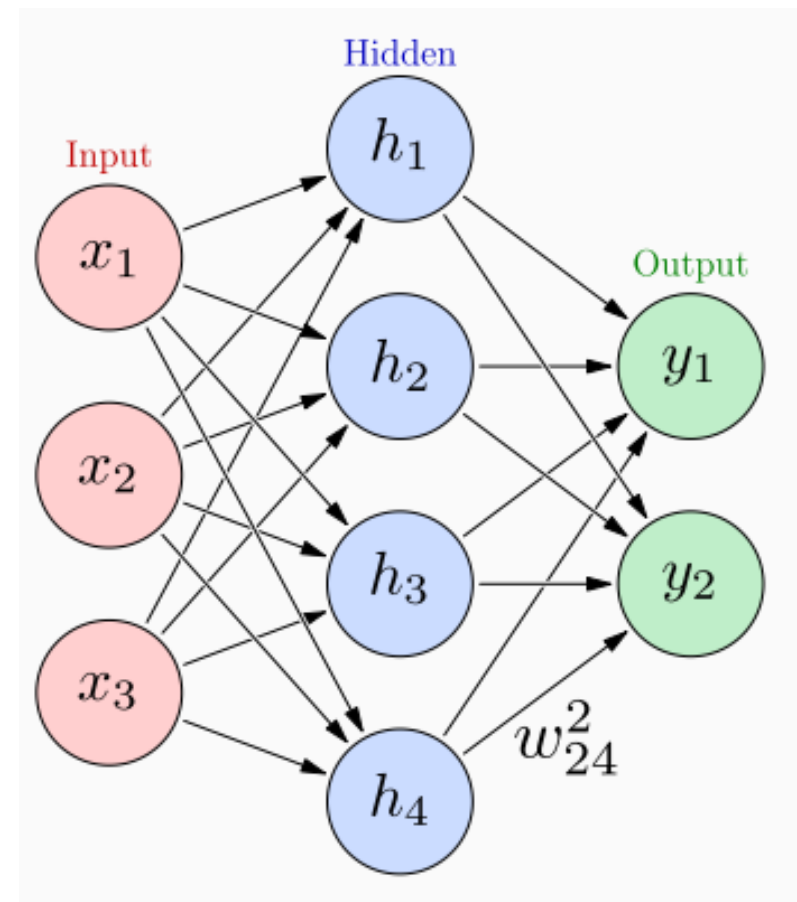$$h_1 = g_1 \left( w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1 \right)$$

$$h_2 = g_1 \left( w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1 \right)$$

$$h_3 = g_1 \left( w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1 \right)$$

$$h_4 = g_1 \left( w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1 \right)$$

$$y_1 = g_2 \left( w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2 \right)$$

$$y_2 = g_2 \left( w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2 \right)$$



- $w_{ij}^k$ weight between previous node j and next node i at layer k;
- $g_k$ is any activation function applied to each its input vector

# Artificial Neural Networks

$$h_1 = g_1\left(w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1\right)$$

$$h_2 = g_1\left(w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1\right)$$

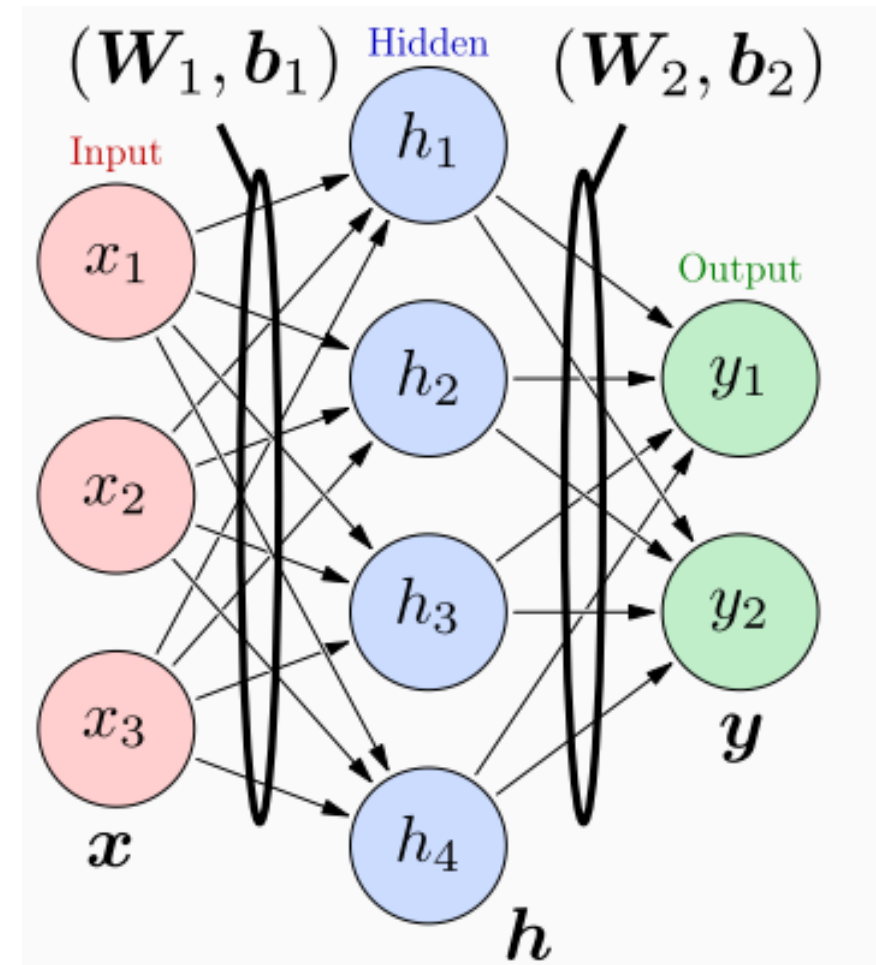$$h_3 = g_1\left(w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1\right)$$

$$h_4 = g_1\left(w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1\right)$$

$$\boldsymbol{h} = g_1\left(\boldsymbol{W}_1 \boldsymbol{x} + \boldsymbol{b}_1\right)$$

$$y_1 = g_2\left(w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2\right)$$

$$y_2 = g_2\left(w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2\right)$$

$$\boldsymbol{y} = g_2\left(\boldsymbol{W}_2 \boldsymbol{h} + \boldsymbol{b}_2\right)$$



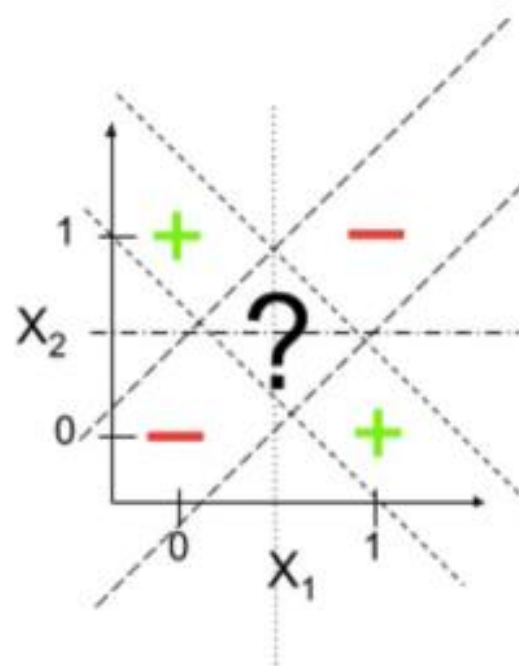- The matrices $W_k$ and biases $b_k$ are learned from labeled training data.
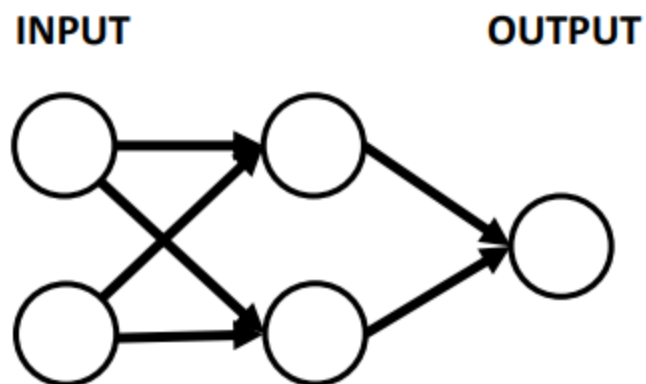
# Artificial Neural Networks



- It can have 1 hidden layer only (shallow network);
- It can have more than 1 hidden layer (deep network);
- Each layer can have a different size, and hidden and output layers often have different activation functions.
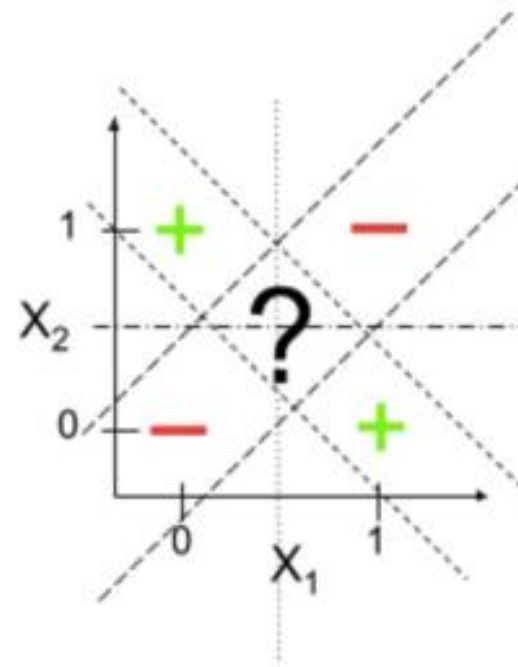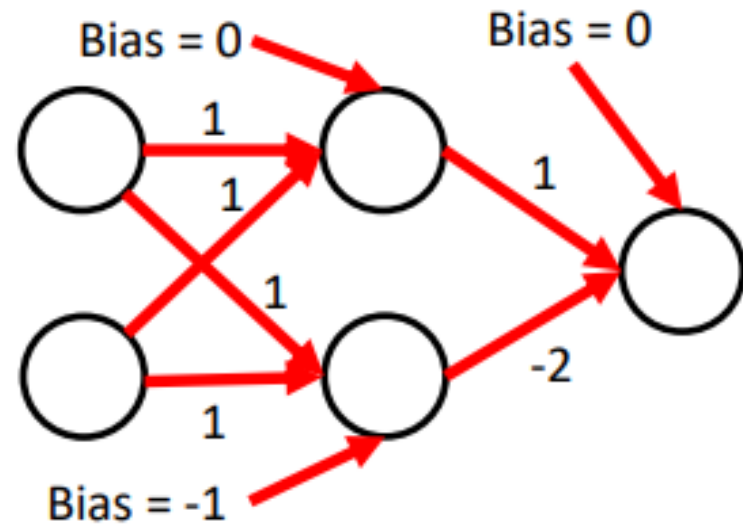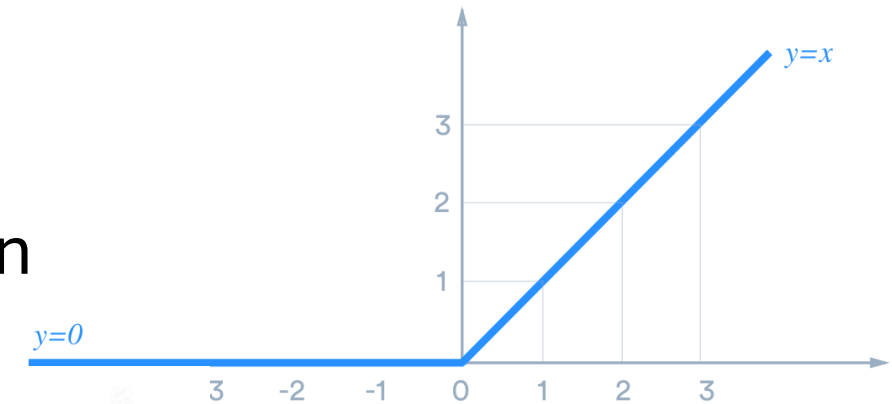
# Revisiting the XOR Problem

- Non-linear function: separate 1s from 0s



| $x_1$ | $x_2$ | $x_1$ XOR $x_2$ |
|-------|-------|-----------------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

# Revisiting the XOR Problem

- Non-linear function: separate 1s from 0s

- Approach: use the ReLU activation function

# Revisiting the XOR Problem

- Non-linear function: separate 1s from 0s

- Approach: use the ReLU activation function

# Revisiting the XOR Problem

- Non-linear function: separate 1s from 0s

- Approach: use the ReLU activation function

# Revisiting the XOR Problem

- Non-linear function: separate 1s from 0s

- Neural networks can solve the XOR problem!

And so model non-linear functions!



| $x_1$ | $x_2$ | $x_1$ XOR $x_2$ |
|-------|-------|------------------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

# Output Activations

- Desired output driven by task



Regression (predict **continuous** value)

Classification (predict **discrete** value)

# Output Activations

- Linear (No Activation Function)

# Output Activations

- Desired output driven by task



Regression
(predict **continuous** value)

Classification
(predict **discrete** value)

# Output Activations

- Sigmoid for Binary Classification



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

If >= 0.5, output 1;

Else, outputs 0

# Output Activations

- Softmax for Multiclass Classification


Converts vector of scores into a probability distribution that sums to 1; e.g.,

# Output Activations

- Softmax for Multiclass Classification

Converts vector of scores into a probability distribution that sums to 1; e.g.,



softmax

$\rightarrow max(p_i)$

Ideal prediction; 1 hot vector of one 1 and the rest 0s

# Output Activations

- Softmax for Multiclass Classification

Converts vector of scores into a probability distribution that sums to 1

Get rid of negative values while preserving original order of scores; e causes negative scores to become slightly larger than 0 while positive values grow exponentially (choosing $e$ rather than another exponent base simplifies math during training)

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$

$i = 1, ..., K$

Number of classes

Want to divide each node's score by sum of all entries to make them sum to 1 (normalization)

# Output Activations

- Softmax for Multiclass Classification

Converts vector of scores into a probability distribution that sums to 1; e.g.,

# Output Activations

- Softmax for Multiclass Classification

Converts vector of scores into a probability distribution that sums to 1; e.g.,

|  | Scoring Function |
|---|---|
| Dog | -3.44 |
| Cat | 1.16 |
| Boat | -0.81 |
| Airplane | 3.91 |

# Output Activations

- Softmax for Multiclass Classification
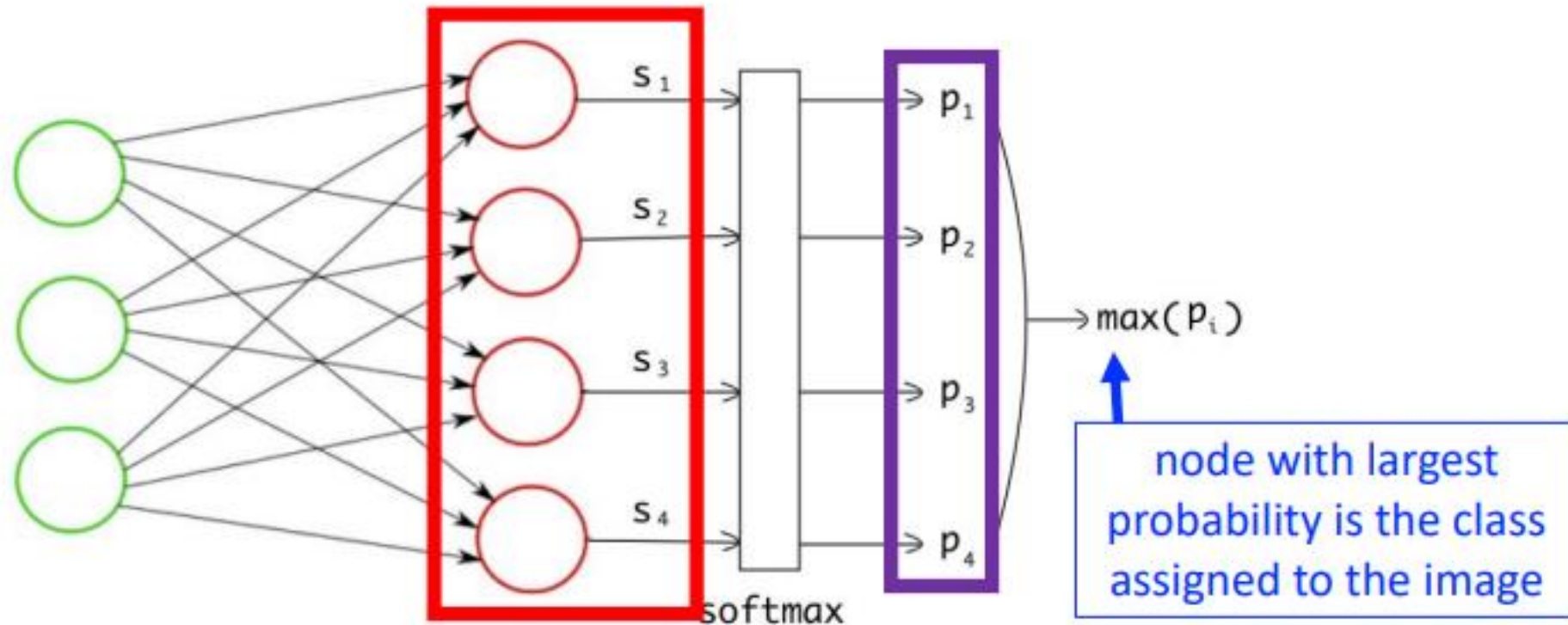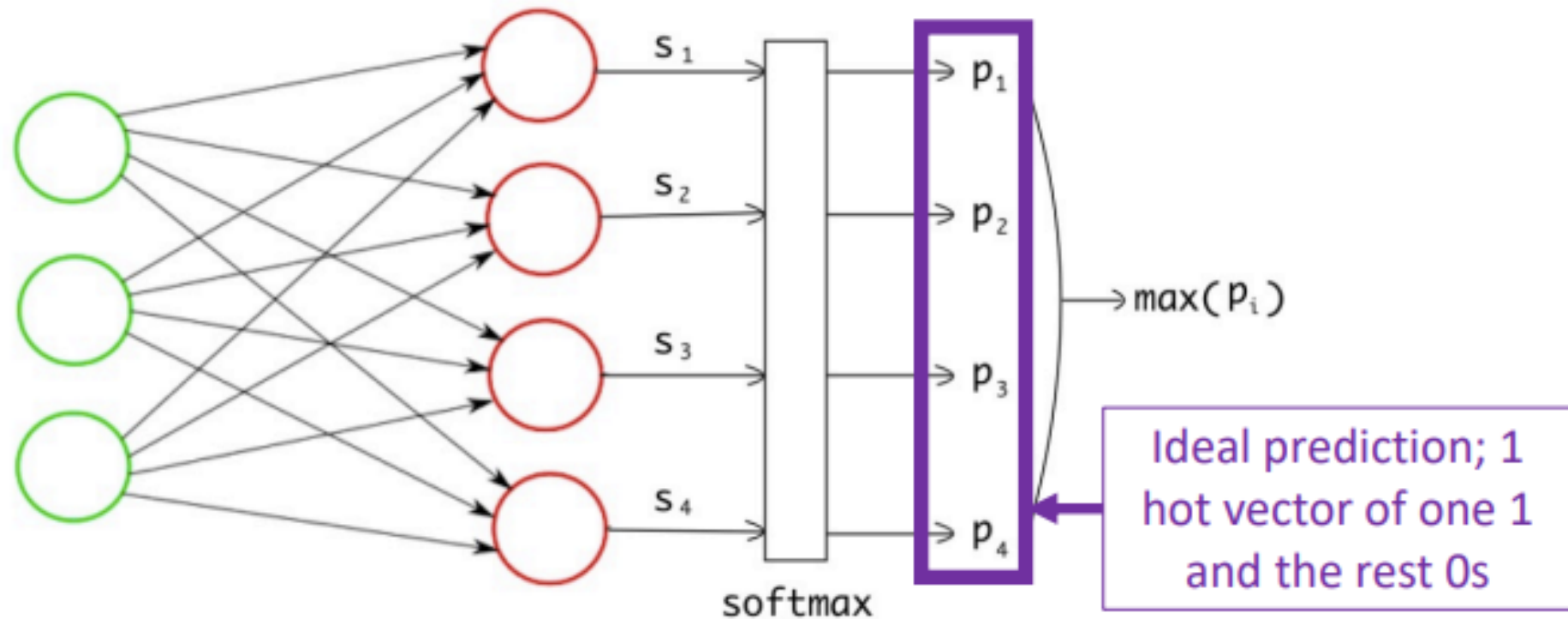
Converts vector of scores into a probability distribution that sums to 1; e.g.,

# Output Activations

- Softmax for Multiclass Classification

Converts vector of scores into a probability distribution that sums to 1; e.g.,

$$e^{z_i} \qquad \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$

| | Scoring Function | Unnormalized Probabilities | Normalized Probabilities |
|---|---|---|---|
| Dog | -3.44 | 0.0321 | 0.0006 |
| Cat | 1.16 | 3.1899 | 0.0596 |
| Boat | -0.81 | 0.4449 | 0.0083 |
| Airplane | 3.91 | 49.8990 | 0.9315 |

# Output Activations

- Softmax for Multiclass Classification

Converts vector of scores into a probability distribution that sums to 1; e.g.,

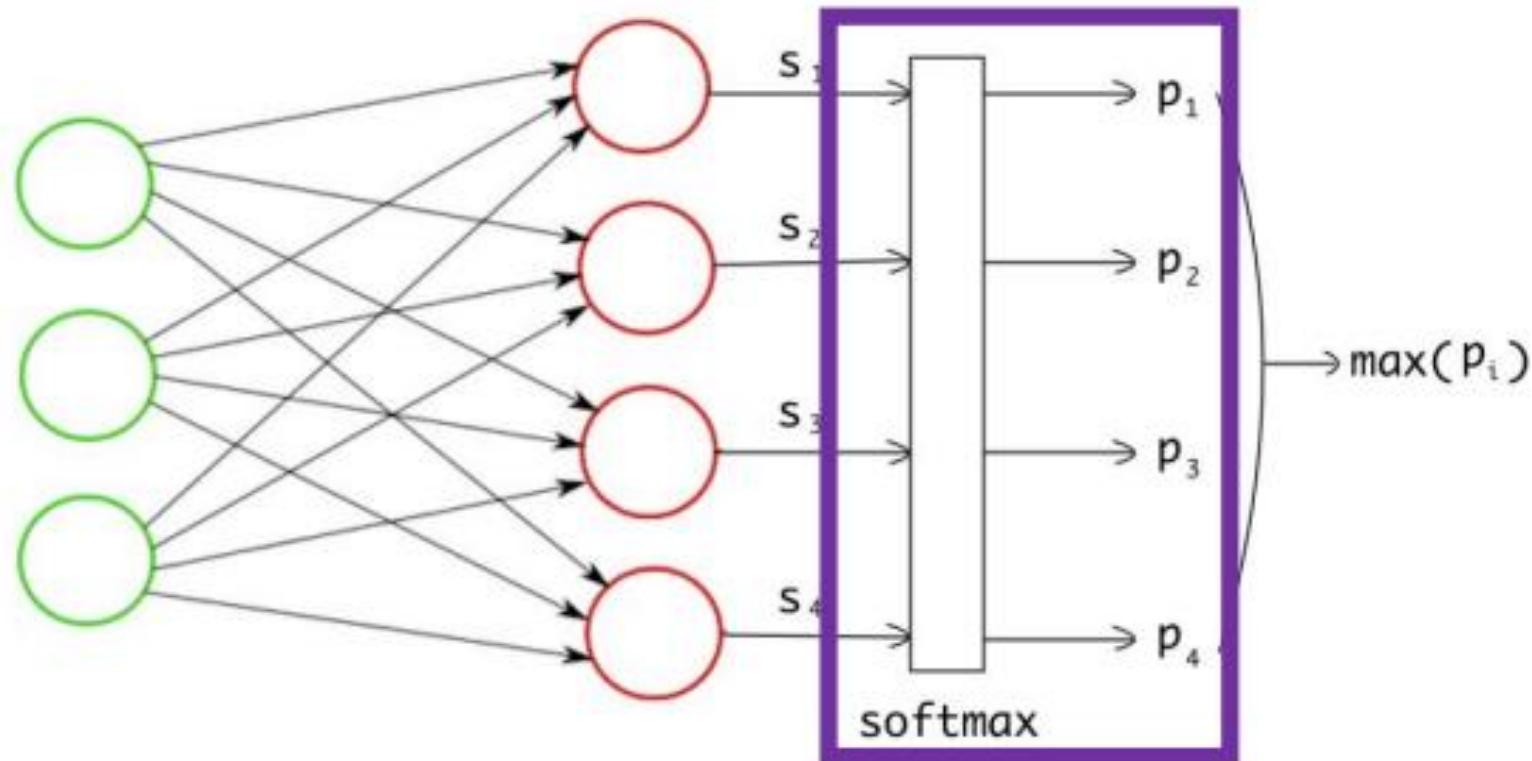# Output Activations

- Softmax for Multiclass Classification

Converts vector of scores into a probability distribution that sums to 1; e.g.,

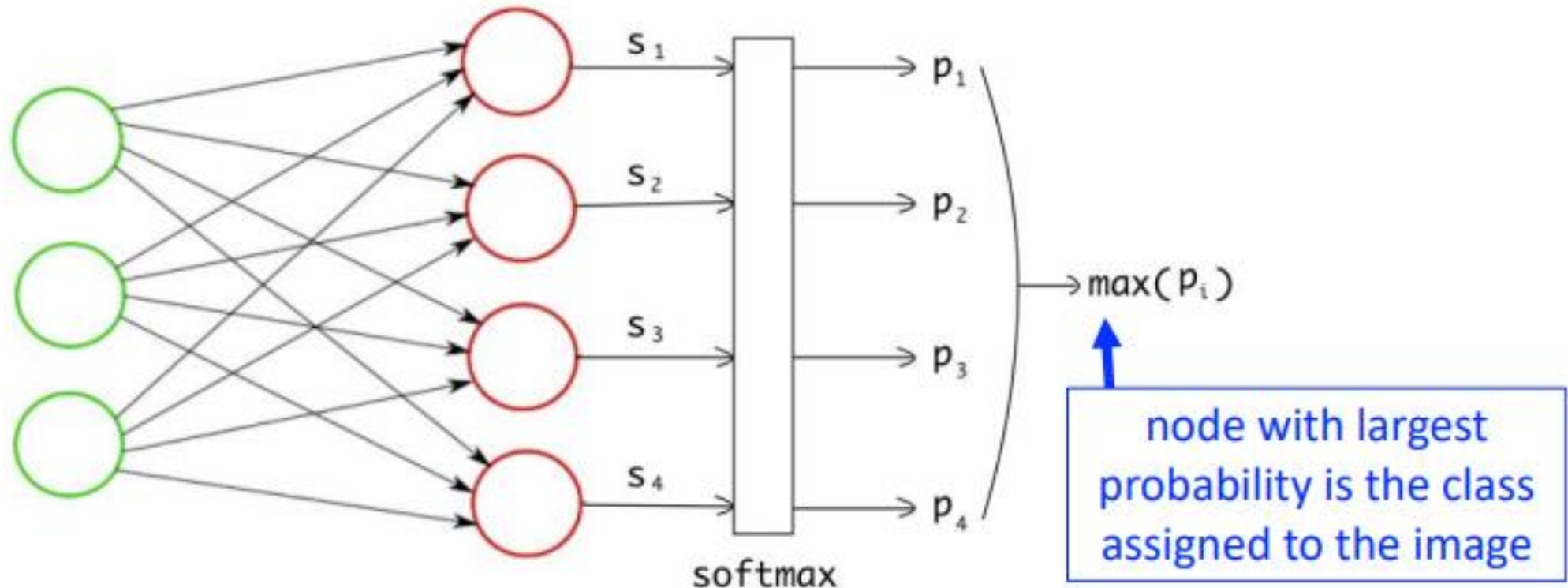| | Scoring Function | Unnormalized Probabilities | Normalized Probabilities |
|---|---|---|---|
| Dog | -3.44 | 0.0321 | 0.0006 |
| Cat | 1.16 | 3.1899 | 0.0596 |
| Boat | -0.81 | 0.4449 | 0.0083 |
| Airplane | 3.91 | 49.8990 | 0.9315 |

# Output Activations

# Example Problem

**Will I pass this class?**

Let's start with a simple two feature model

$x_1$ = Numbers of lectures you attend

$x_2$ = Hours spent on the final project

# Example Problem: Will I pass this class?

# Example Problem: Will I pass this class?

$x_2$ = Hours spent on the final project

?

$\begin{bmatrix} 4 \\ 5 \end{bmatrix}$

Legend

● Pass

● Fail

$x_1$ = Number of lectures you attend

# Example Problem: Will I pass this class?

$$x^{(1)} = [4, 5]$$

$x_1$

$x_2$

$z_1$

$z_2$

$z_3$

$\hat{y}_1$

Predicted: 0.1

# Example Problem: Will I pass this class?



$$x^{(1)} = [4, 5]$$

$x_1$

$x_2$

$z_1$

$z_2$

$z_3$

$\hat{y}_1$

Predicted: 0.1
Actual: 1

# Quantifying Loss



The **loss** of our network measures the cost incurred from incorrect predictions

$x^{(1)} = [4, 5]$

$x_1$

$x_2$

$z_1$

$z_2$

$z_3$

$\hat{y}_1$

Predicted: 0.1
Actual: 1

$$\mathcal{L}\left(\underline{f\left(x^{(i)}; W\right)}, \underline{y^{(i)}}\right)$$

Predicted        Actual

# Empirical Loss

The **empirical loss** measures the total loss over our entire dataset

$$X = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$

$$f(x) = \begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \\ \vdots \end{bmatrix} \qquad y = \begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$$

Also known as:
- Objective function
- Cost function
- Empirical Risk

$$J(W) = \frac{1}{n}\sum_{i=1}^{n}\mathcal{L}\left(f\left(x^{(i)};W\right), y^{(i)}\right)$$

Predicted        Actual

# Binary Cross Entropy Loss
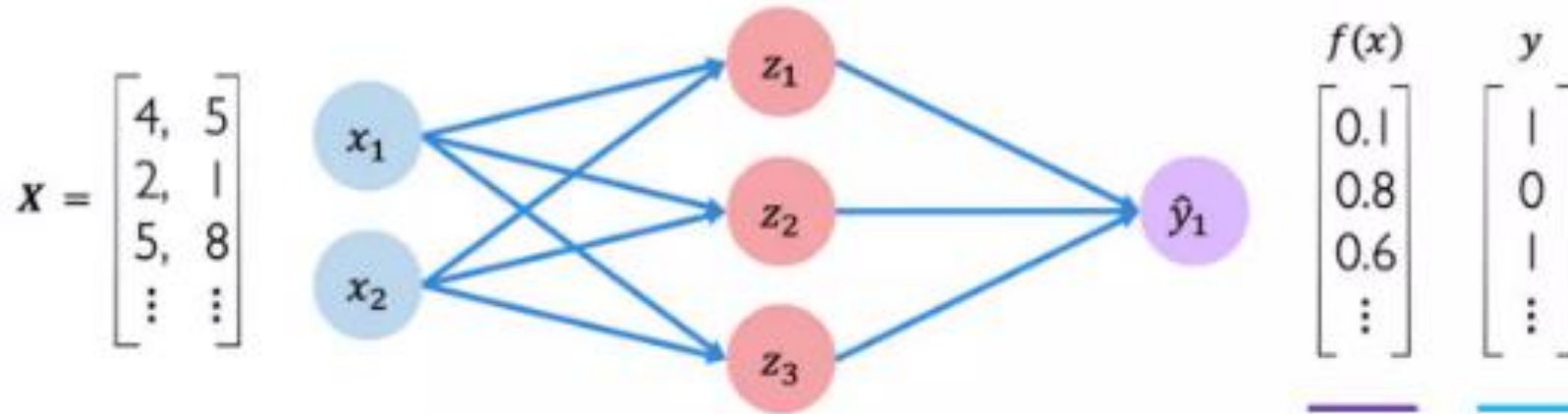


Cross entropy loss can be used with models that output a probability between 0 and 1

$$X = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$

$f(x)$

$$\begin{bmatrix} 0.1 \\ 0.8 \\ 0.6 \\ \vdots \end{bmatrix}$$

$y$

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$$

$$J(W) = \frac{1}{n} \sum_{i=1}^{n} y^{(i)} \log\left(f\left(x^{(i)}; W\right)\right) + \left(1 - y^{(i)}\right) \log\left(1 - f\left(x^{(i)}; W\right)\right)$$

Actual · Predicted · Actual · Predicted

# Mean Squared Error Loss

Mean squared error loss can be used with regression models that output continuous real numbers

$$J(W) = \frac{1}{n}\sum_{i=1}^{n}\left(y^{(i)} - f(x^{(i)}; W)\right)^2$$

Actual    Predicted

Final Grades (percentage)

# The Perceptron from Scratch

- Let's implement the perceptron architecture from scratch using only numpy.

- Use the *perceptron.py* script as a starting point.