



CIS 678 - Machine Learning

Introduction to Neural Networks



Supervised Models

- kNN
- Linear Regression
- Decision Tree
- Random Forest Regressor
- Boosting Regressor
- Support Vector Regressor (SVRs)

Regression

- kNN
- Logistic Regression
- Decision Tree
- Random Forest Classifier
- Boosting Classifiers
- Support Vector Classifiers (SVCs)
- Naive Bayes

Classification



Supervised Models

- kNN
- Linear Regression
- Decision Tree
- Random Forest Regressor
- Boosting Regressor
- Support Vector Regressor (SVRs)
- Neural Networks (NNs)

Regression

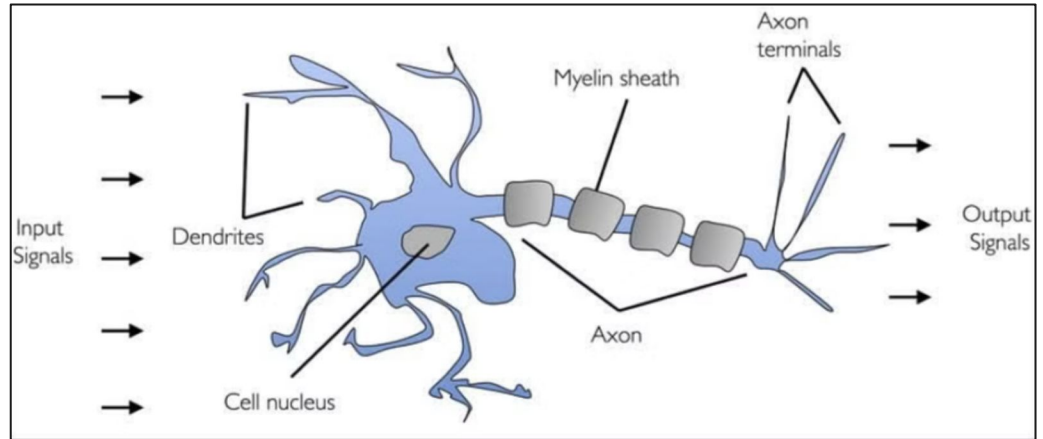
- kNN
- Logistic Regression
- Decision Tree
- Random Forest Classifier
- Boosting Classifiers
- Support Vector Classifiers (SVCs)
- Naive Bayes
- Neural Networks (NNs)

Classification

Neural Networks

Motivation src: Biological neuron

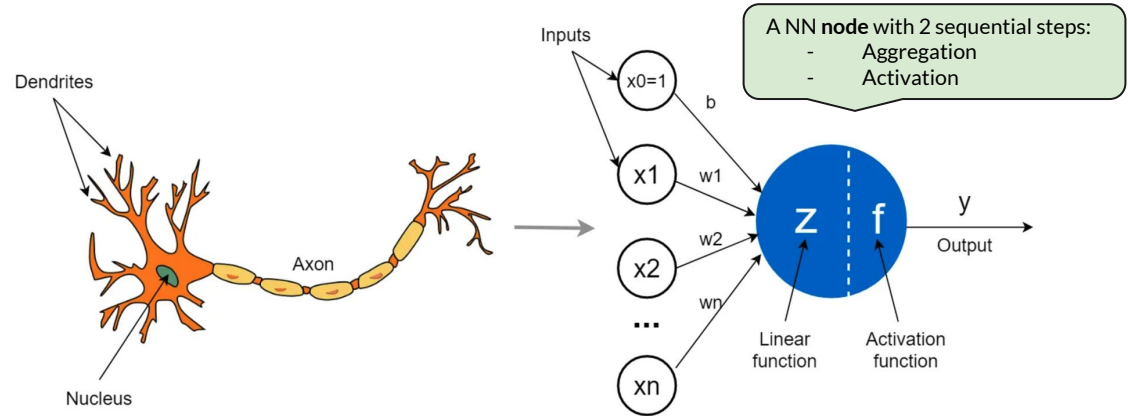
Perceptron was introduced by **Frank Rosenblatt** in 1957.



[perceptron](#)

Neural Networks

From Biological Neuron to
Artificial NN



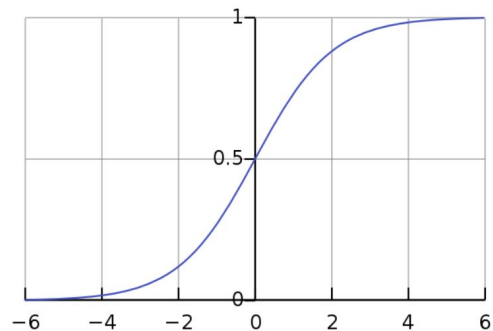
[fig-ref](#)

Logistic Regression

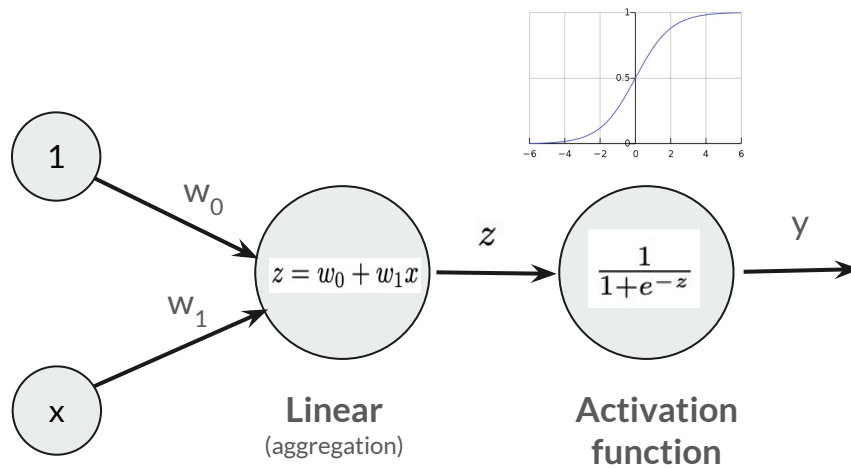
- Probabilistic classifier

$$p(x) = \frac{1}{1 + e^{-(w_0 + w_1 x)}}$$

- Sigmoid function

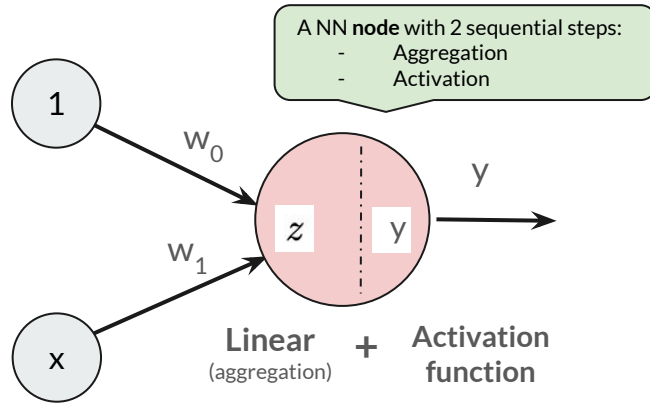


Neural Networks (Node)

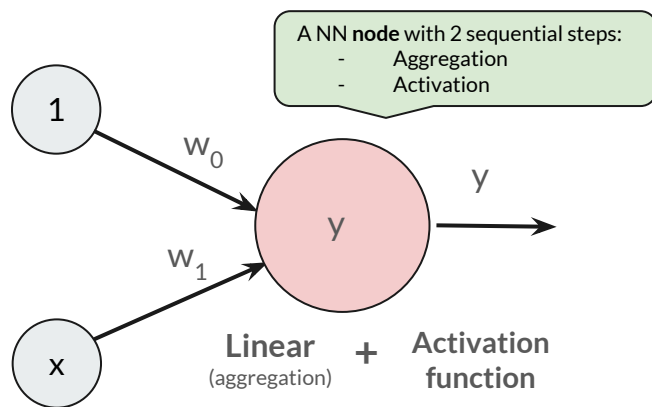


$y \in \{0, 1\}$
Logistic Regression

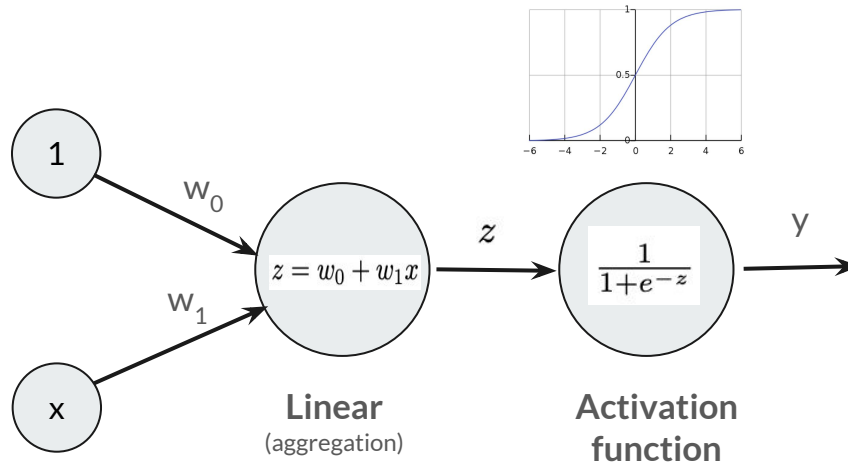
Neural Networks (Node)



Neural Networks (Node)



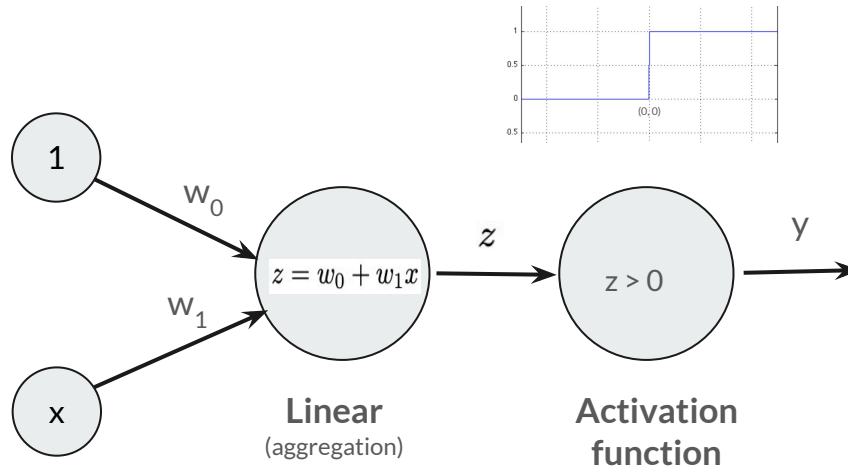
Neural Networks (Node)



$$y \in R$$

A NN with Sigmoid
Activation function

Neural Networks (Node)



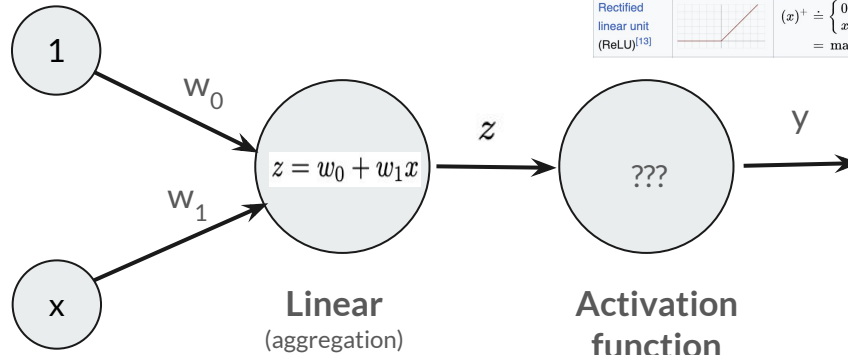
$$y \in \{0, 1\}$$

A NN with Step
Activation function

Frank Rosenblatt's
Perceptron

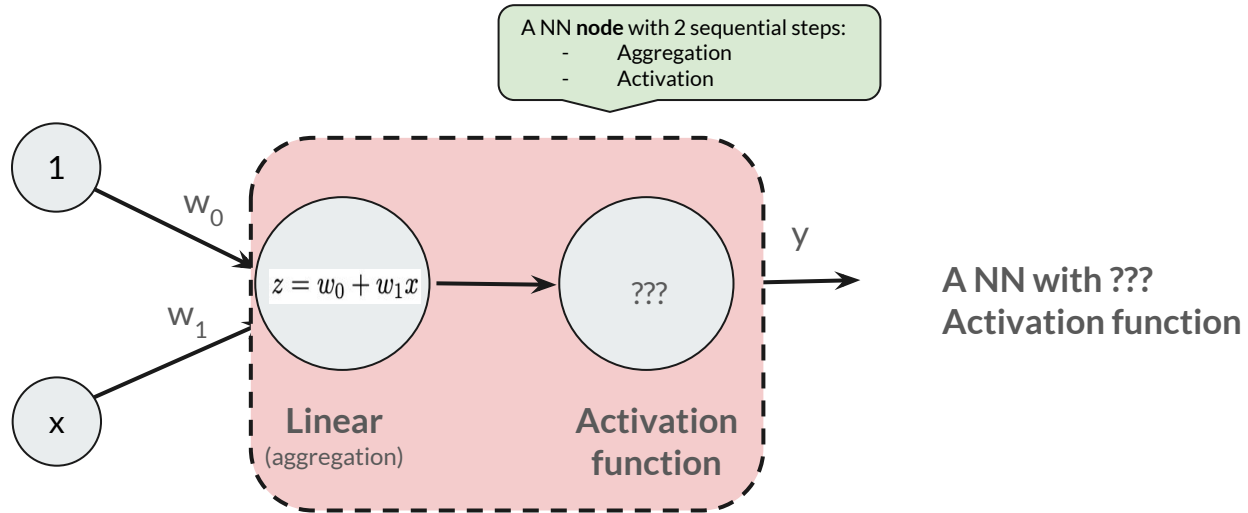
Neural Networks (Node)

Name	Plot	Function, $g(x)$
Identity		x
Binary step		$\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$
Logistic, sigmoid, or soft step		$\sigma(x) \doteq \frac{1}{1 + e^{-x}}$
Hyperbolic tangent (tanh)		$\tanh(x) \doteq \frac{e^x - e^{-x}}{e^x + e^{-x}}$
Soboleva modified hyperbolic tangent (smht)		$\text{smht}(x) \doteq \frac{e^{ax} - e^{-bx}}{e^{cx} + e^{-dx}}$
Rectified linear unit (ReLU) ^[13]		$(x)^+ \doteq \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ $= \max(0, x) = x \mathbf{1}_{x>0}$

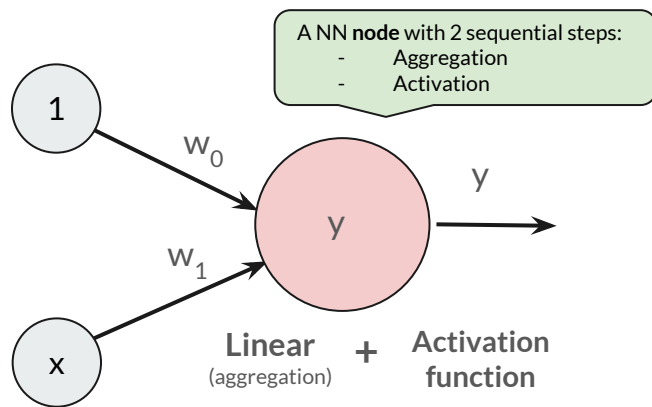


A NN with ???
Activation function

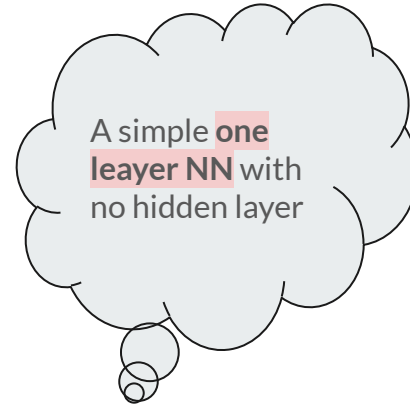
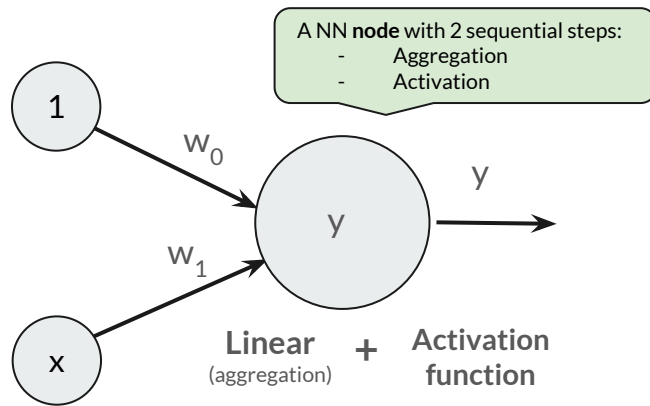
Neural Networks (Node)



Neural Networks (Node)



Neural Networks (No Hidden Layer)

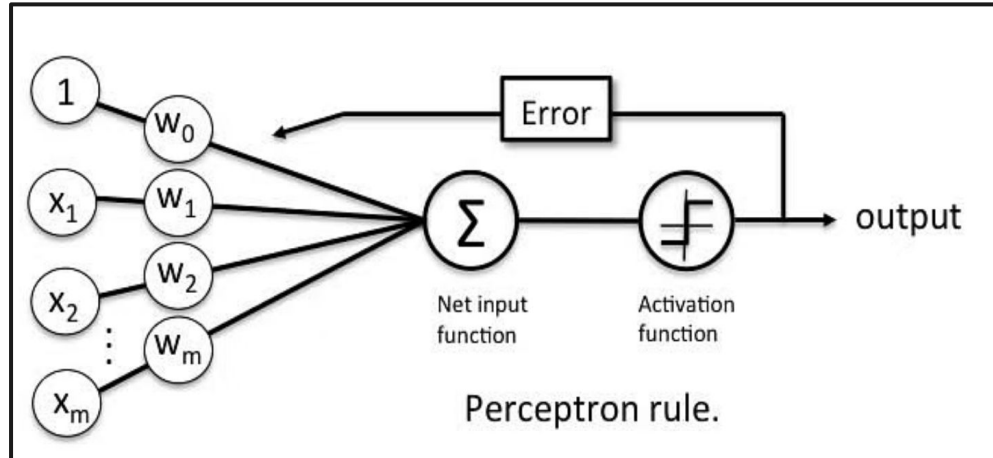


Perceptron: the first Neural Network

Motivation src: Biological neuron

Perceptron was introduced by **Frank Rosenblatt** in 1957.

A binary classifier



[perceptron](#)

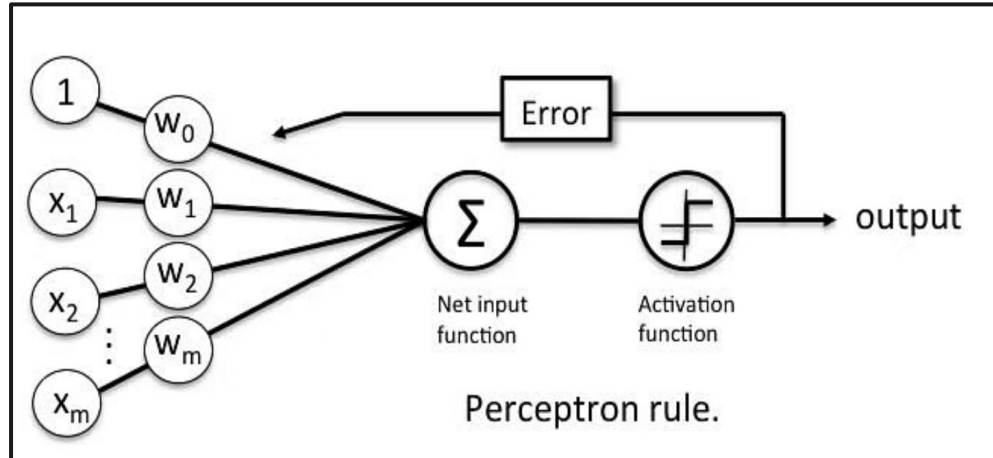
Perceptron: the first Neural Network

Motivation src: Biological neuron

Perceptron was introduced by **Frank Rosenblatt** in 1957.

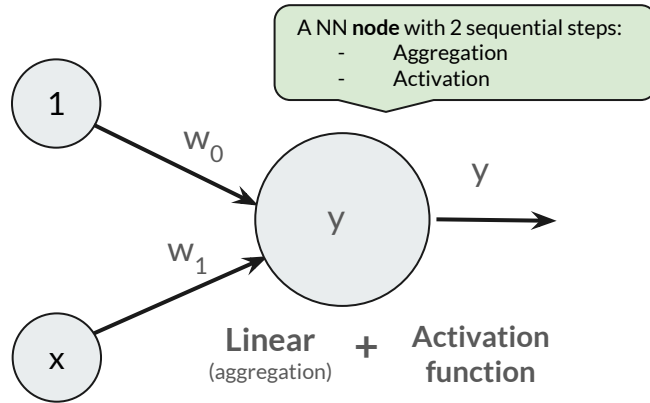
A binary classifier

[Professor's perceptron paved the way for AI – 60 years too soon](#)



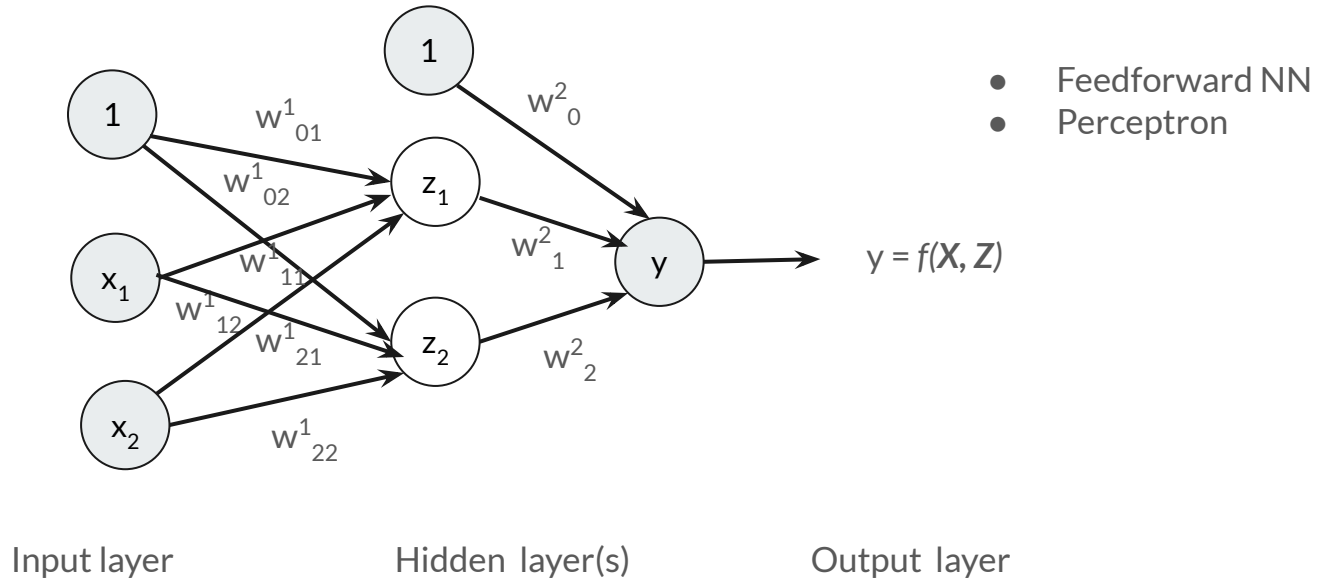
[perceptron](#)

Neural Networks (No Hidden Layer)



To take the advantage of the NN, you require at least one Hidden layer.

Feed-forward (FF) neural networks





Feed-forward (FF) neural networks

x_1

x_2

x_3

x_4

Input
(X)



Feed-forward (FF) neural networks

$$x_1 = a_1^{(1)} \quad \text{●}$$

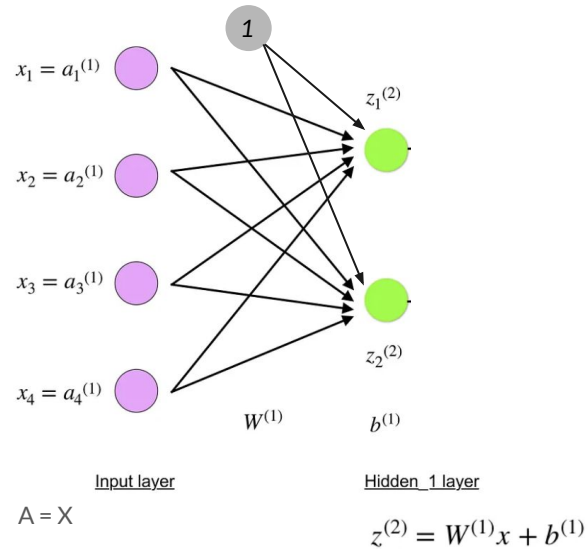
$$x_2 = a_2^{(1)} \quad \text{●}$$

$$x_3 = a_3^{(1)} \quad \text{●}$$

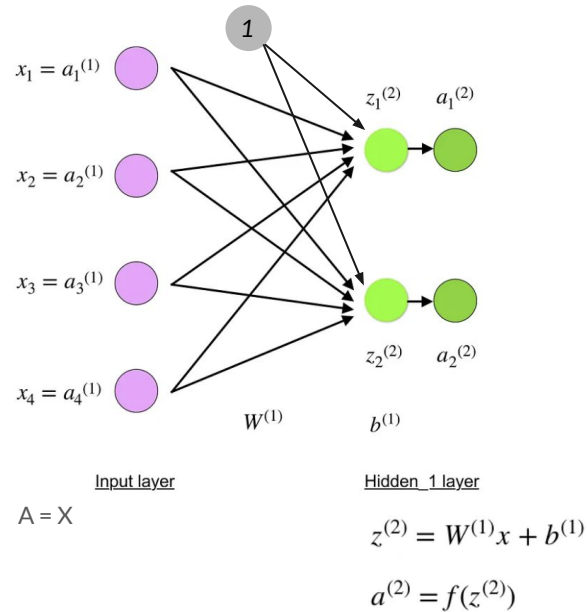
$$x_4 = a_4^{(1)} \quad \text{●}$$

$$A = X$$

Feed-forward (FF) neural networks

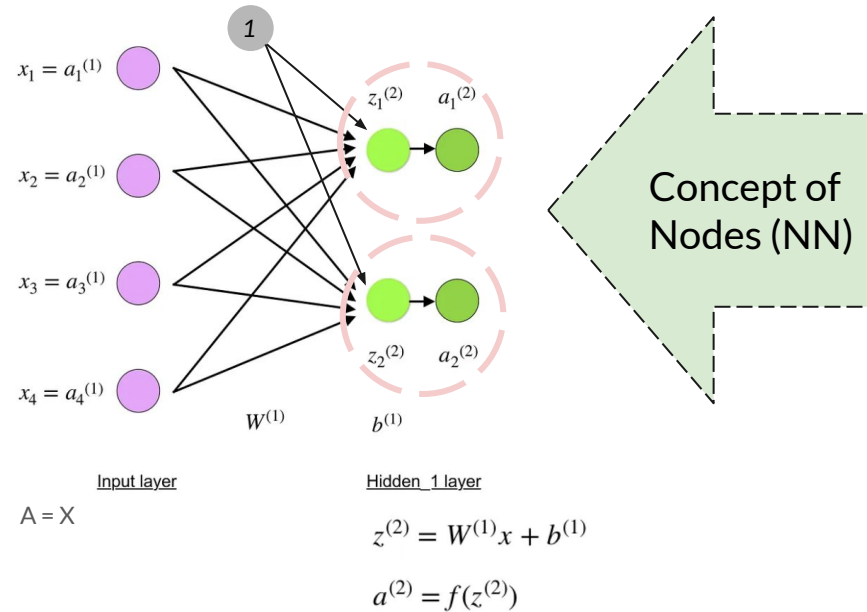


Feed-forward (FF) neural networks



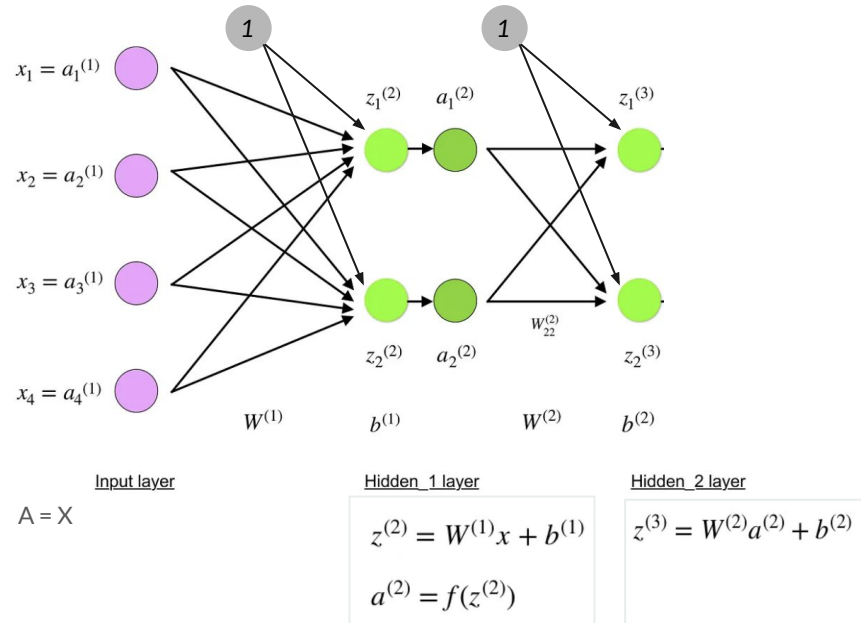
[mainly adapted from](#)

Feed-forward (FF) neural networks



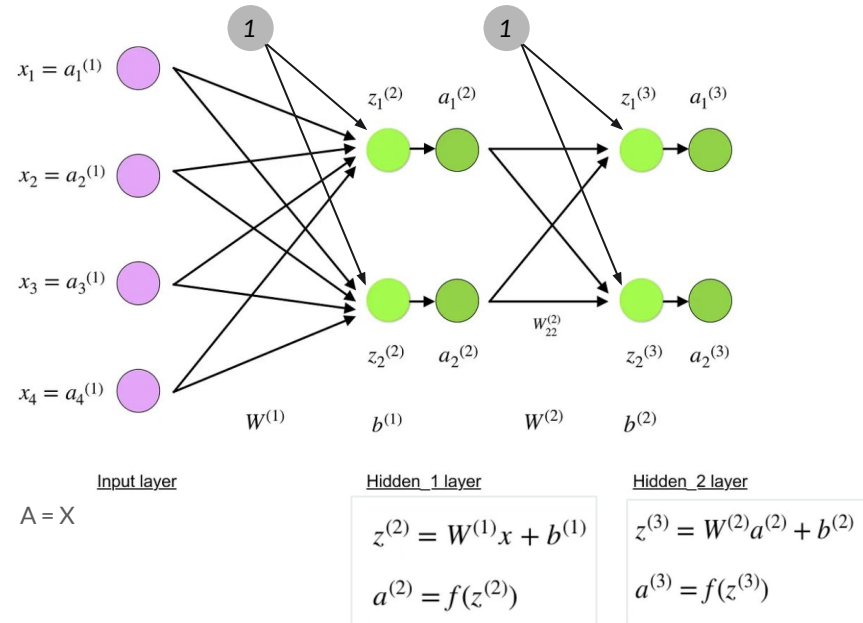
[mainly adapted from](#)

Feed-forward (FF) neural networks



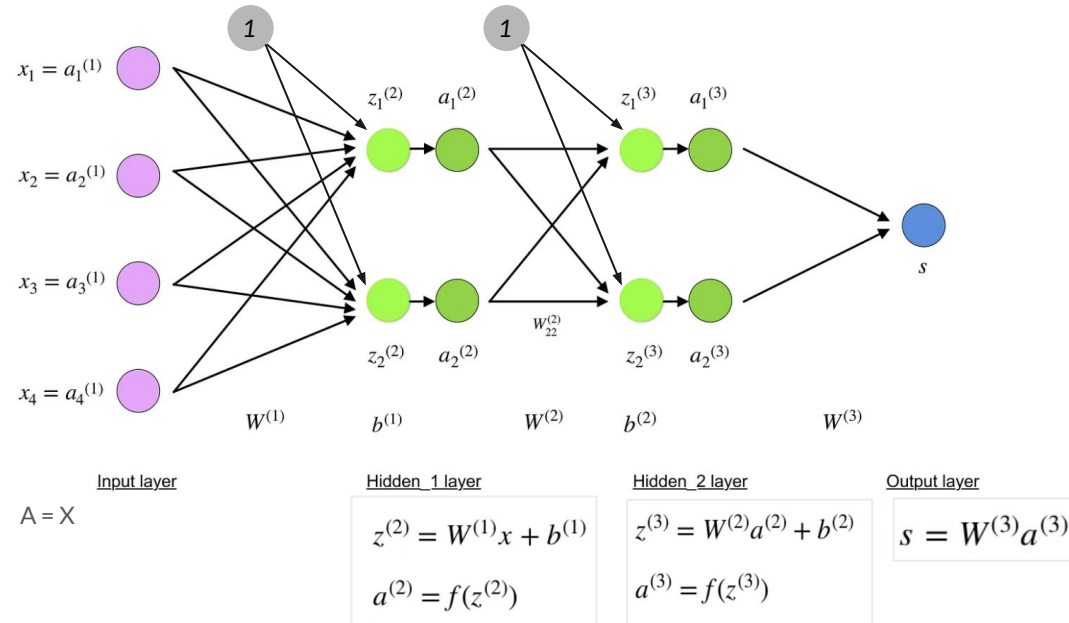
[mainly adapted from](#)

Feed-forward (FF) neural networks



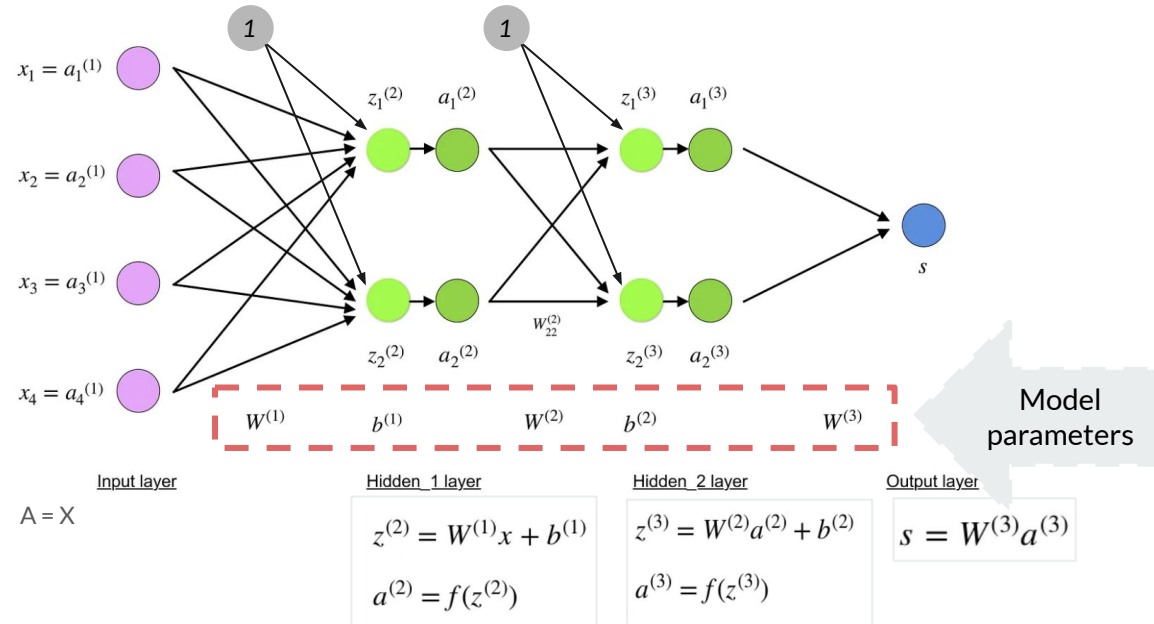
[mainly adapted from](#)

Feed-forward (FF) neural networks



[mainly adapted from](#)

Feed-forward (FF) neural networks



[mainly adapted from](#)



Question?

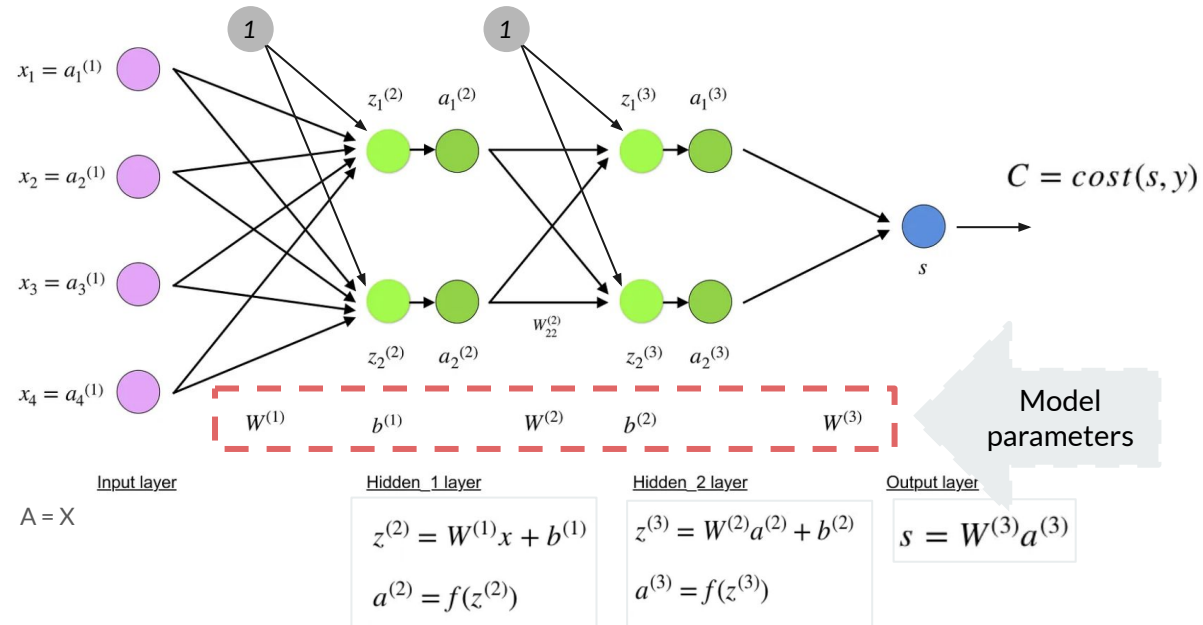
Q. Draw the diagram of a Feed Forward Neural Network with the properties given below, and estimate the minimum number of parameters your model would have:

1. **Input layer:** 3 nodes (to consume 3 input features, $\{x_1, x_2, x_3\}$)
2. **Three (3) Hidden layers** with the following configuration:
 - i) Hidden layer one: 3 nodes
 - ii) Hidden layer two: 2 nodes
 - iii) Hidden layer three: 4 nodes
1. **One bias input node** for each hidden layer in (2)
2. **Output layer:** 1 node (y)



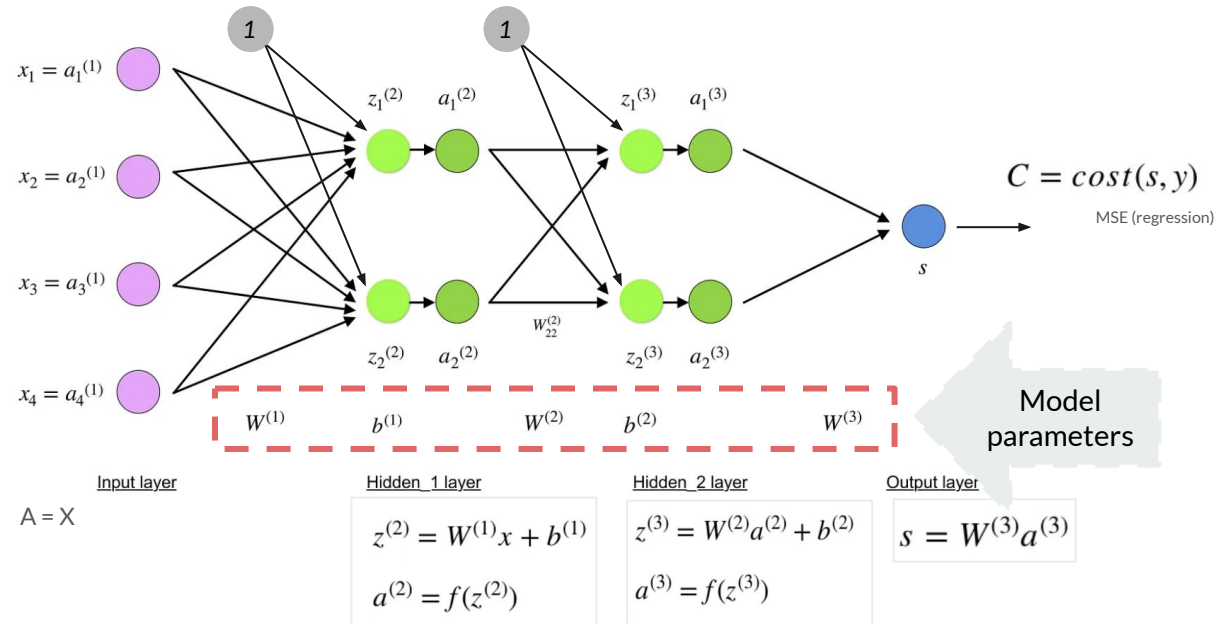
NN Training

Training



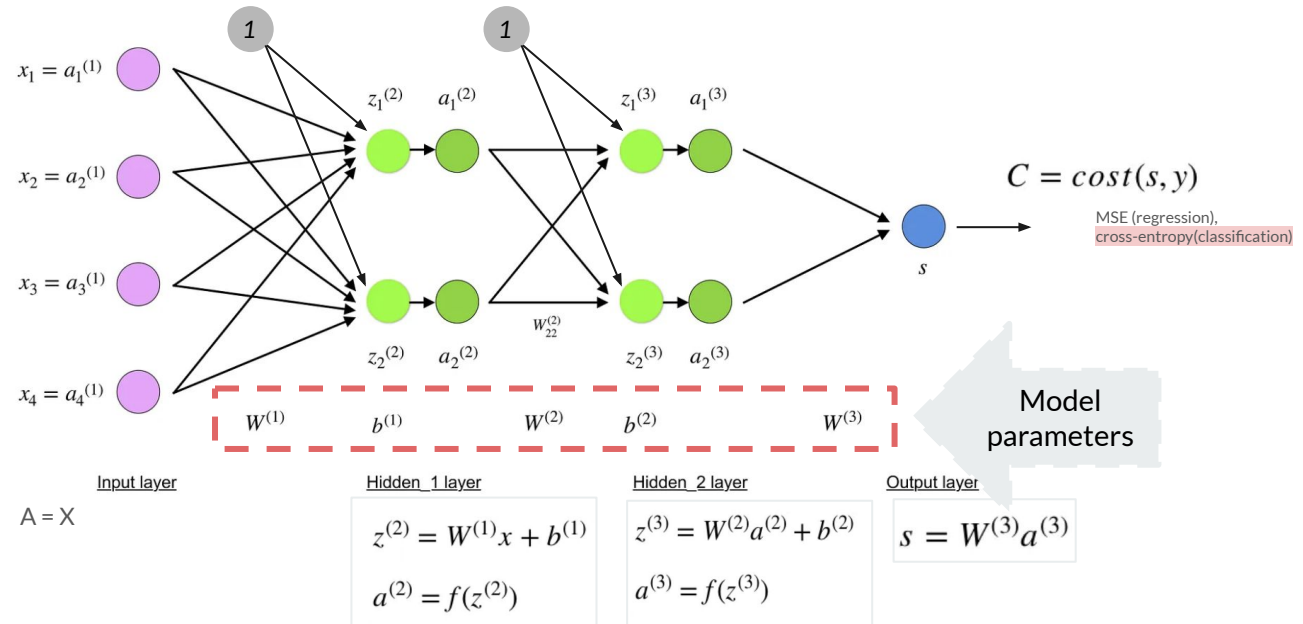
[mainly adapted from](#)

Training



[mainly adapted from](#)

Training



mainly adapted from



Gradients

- x is your parameter vector/matrix
- Partial derivatives
- Only the last (hidden) layer parameters can have direct derivatives
- Rest (including the input layer) requires to apply a chain rule

$$\frac{\partial C}{\partial x} = \left[\frac{\partial C}{\partial x_1}, \frac{\partial C}{\partial x_2}, \dots, \frac{\partial C}{\partial x_m} \right]$$



Gradients

- \mathbf{x} is your parameter vector
- Partial derivatives
- Only the last (hidden) layer parameters can have direct derivatives
- Rest (including the input layer) requires to apply a chain rule

l : layer index
 j : node index in layer l ,
 k : node index in layer $l-1$

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} \quad \text{chain rule}$$

$$z_j^l = \sum_{k=1}^m w_{jk}^l a_k^{l-1} + b_j^l \quad \text{by definition}$$

m - number of neurons in $l-1$ layer

$$\frac{\partial z_j^l}{\partial w_{jk}^l} = a_k^{l-1} \quad \text{by differentiation (calculating derivative)}$$

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} a_k^{l-1} \quad \text{final value}$$

Gradients

- \mathbf{x} is your parameter vector
- Partial derivatives
- Only the last (hidden) layer parameters can have direct derivatives
- Rest (including the input layer) requires to apply a chain rule

l : layer index
 j : node index in layer l ,
 k : node index in layer $l-1$

Hidden layer

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} \quad \text{chain rule}$$

$$z_j^l = \sum_{k=1}^m w_{jk}^l a_k^{l-1} + b_j^l \quad \text{by definition}$$

m - number of neurons in $l-1$ layer

$$\frac{\partial z_j^l}{\partial w_{jk}^l} = a_k^{l-1} \quad \text{by differentiation (calculating derivative)}$$

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} a_k^{l-1} \quad \text{final value}$$

Gradients

- \mathbf{x} is your parameter vector
- Partial derivatives
- Only the last (hidden) layer parameters can have direct derivatives
- Rest (including the input layer) requires to apply a chain rule

l : layer index
 j : node index in layer l ,
 k : node index in layer $l-1$

Hidden layer

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} \quad \text{chain rule}$$

$$z_j^l = \sum_{k=1}^m w_{jk}^l a_k^{l-1} + b_j^l \quad \text{by definition}$$

m - number of neurons in $l-1$ layer

$$\frac{\partial z_j^l}{\partial w_{jk}^l} = a_k^{l-1} \quad \text{by differentiation (calculating derivative)}$$

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} a_k^{l-1} \quad \text{final value}$$

Gradients

- x is your parameter vector
- Partial derivatives
- Only the last (hidden) layer parameters can have direct derivatives
- Rest (including the input layer) requires to apply a chain rule

l : layer index
 j : node index in layer l ,
 k : node index in layer $l-1$

$$\frac{\partial C}{\partial b_j^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial b_j^l} \quad \text{chain rule}$$
$$\frac{\partial z_j^l}{\partial b_j^l} = 1 \quad \text{by differentiation (calculating derivative)}$$

Derivative with respect to bias is 1

$$\frac{\partial C}{\partial b_j^l} = \frac{\partial C}{\partial z_j^l} 1 \quad \text{final value}$$



Gradient descent

- Can you recall our gradient descent Linear Regression model training?

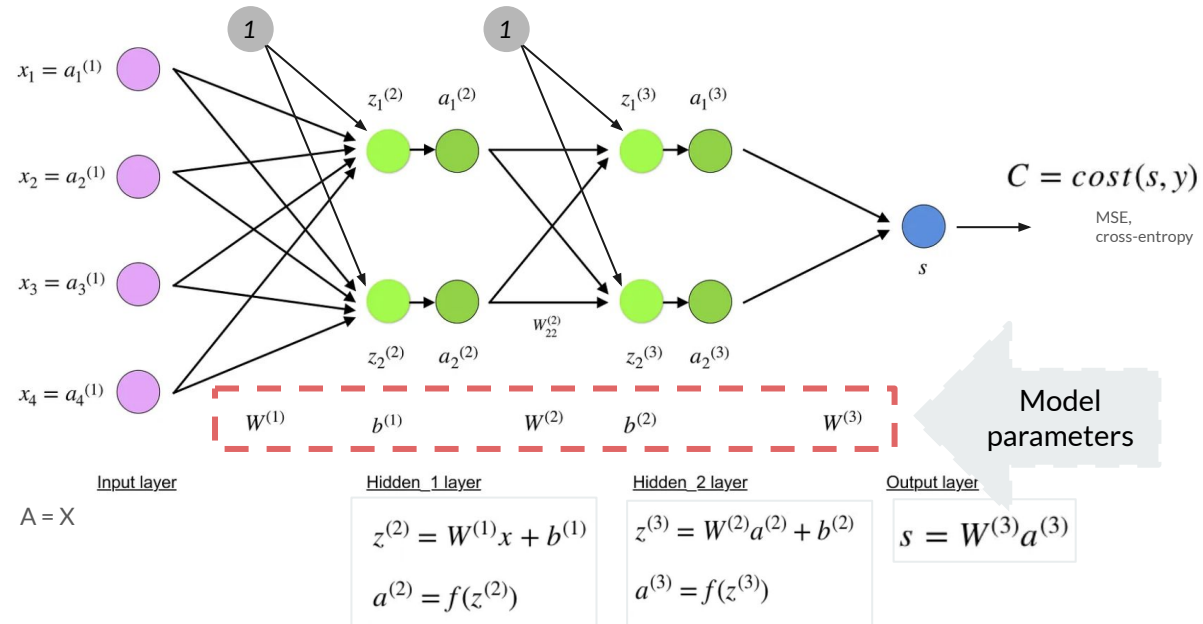
while (termination condition not met)

$$w := w - \epsilon \frac{\partial C}{\partial w}$$

$$b := b - \epsilon \frac{\partial C}{\partial b}$$

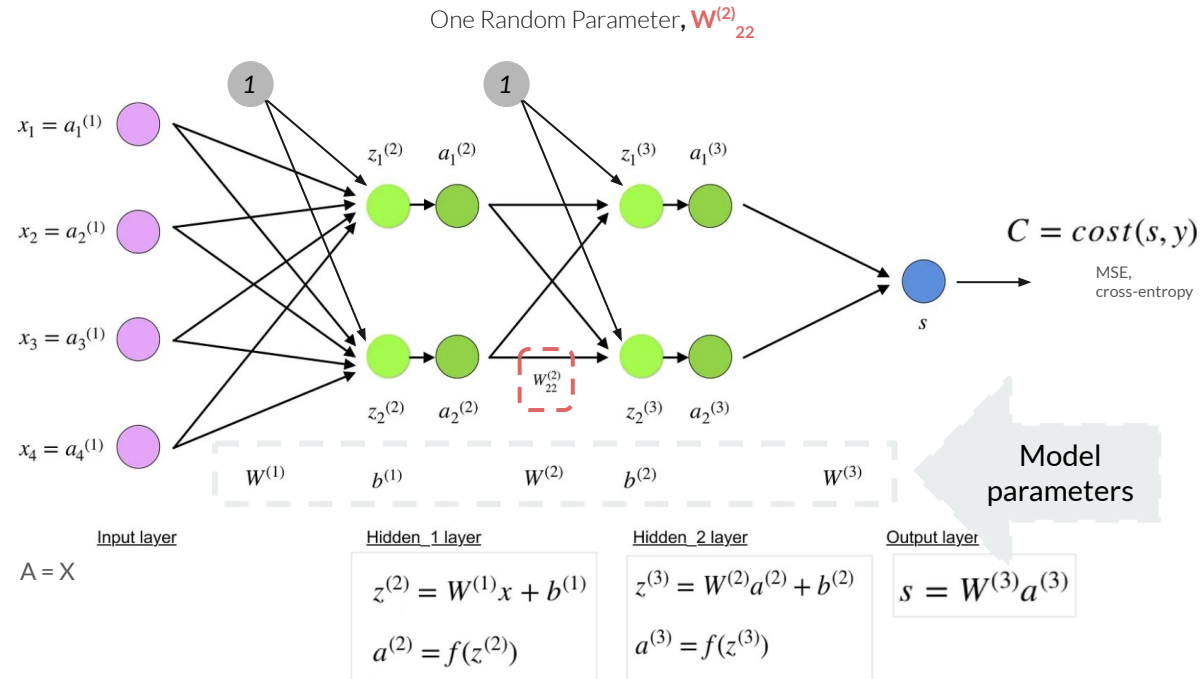
end

Training



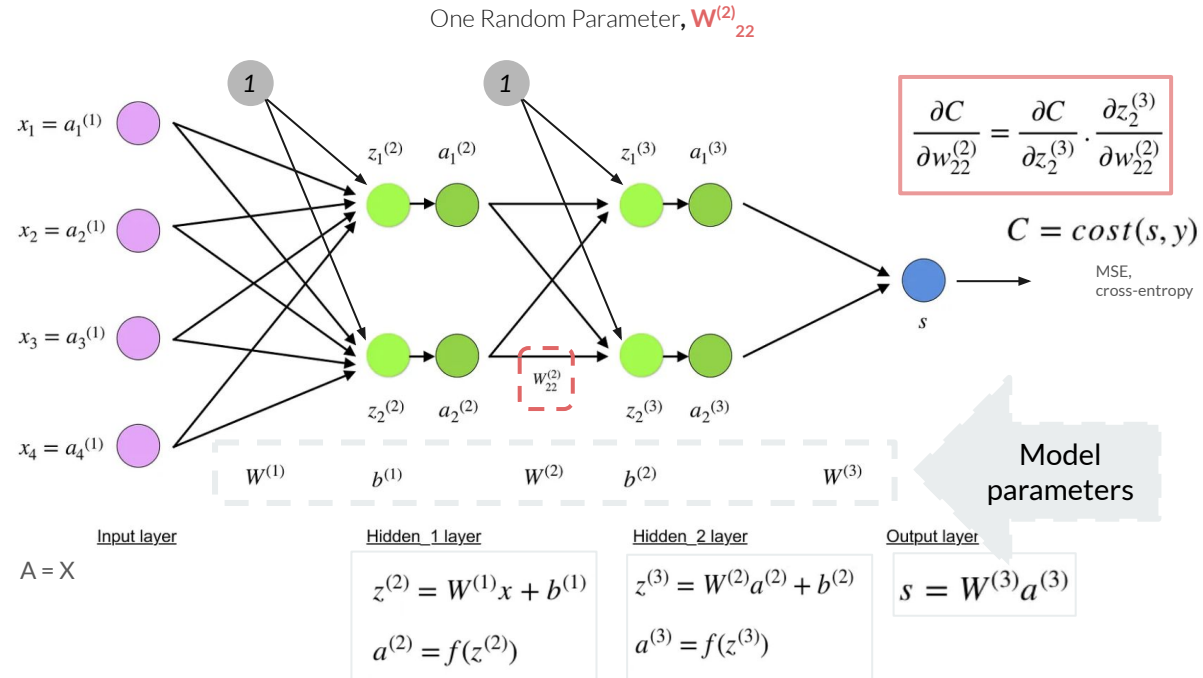
[mainly adapted from](#)

Training



[mainly adapted from](#)

Training

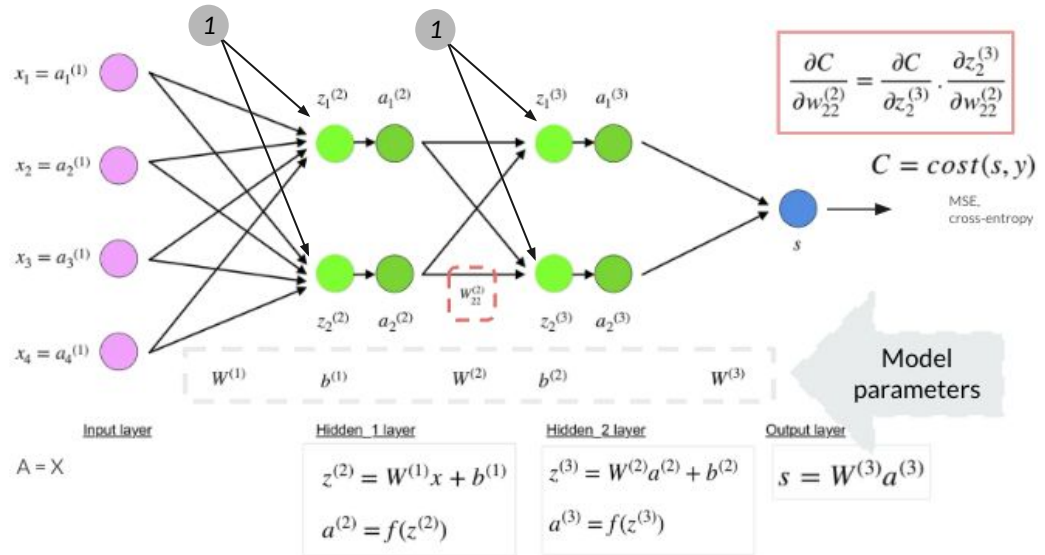


mainly adapted from

Error Backpropagation

One Random Parameter, $W_{22}^{(2)}$

$$\frac{\partial C}{\partial w_{22}^{(2)}} = \frac{\partial C}{\partial z_2^{(3)}} \cdot \frac{\partial z_2^{(3)}}{\partial w_{22}^{(2)}}$$

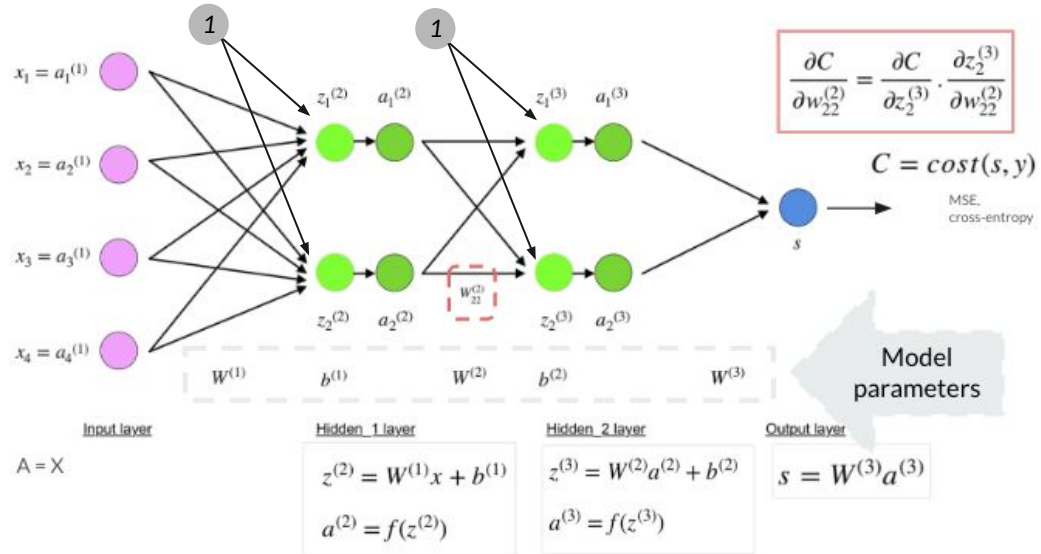


[mainly adapted from](#)

Error Backpropagation

One Random Parameter, $W_{22}^{(2)}$

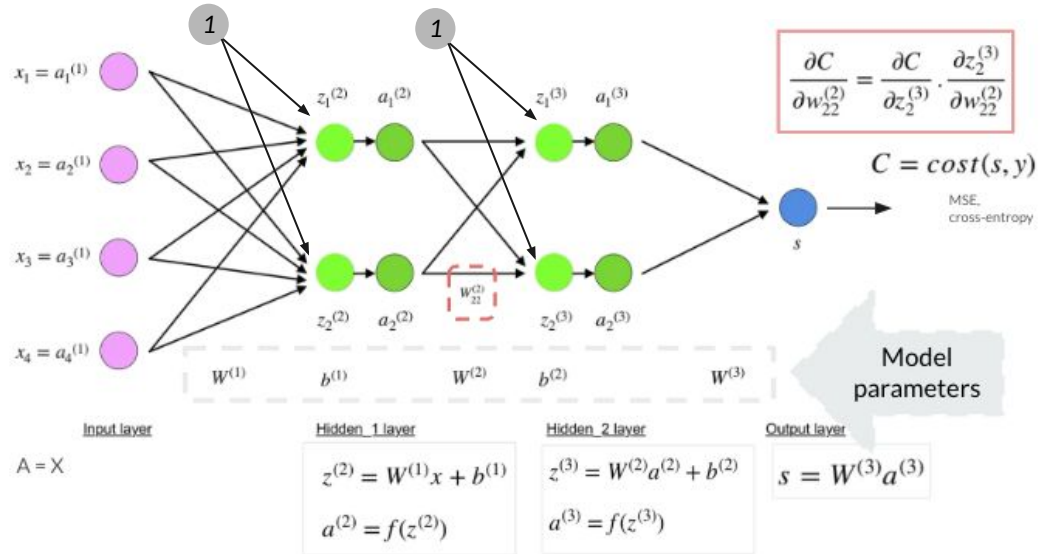
$$\begin{aligned}\frac{\partial C}{\partial w_{22}^{(2)}} &= \frac{\partial C}{\partial z_2^{(3)}} \cdot \frac{\partial z_2^{(3)}}{\partial w_{22}^{(2)}} \\ &= \frac{\partial C}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot a_2^{(2)}\end{aligned}$$



Error Backpropagation

One Random Parameter, $W_{22}^{(2)}$

$$\begin{aligned}\frac{\partial C}{\partial w_{22}^{(2)}} &= \frac{\partial C}{\partial z_2^{(3)}} \cdot \frac{\partial z_2^{(3)}}{\partial w_{22}^{(2)}} \\ &= \frac{\partial C}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot a_2^{(2)} \\ &= \frac{\partial C}{\partial a_2^{(3)}} \cdot f'(z_2^{(3)}) \cdot a_2^{(2)}\end{aligned}$$

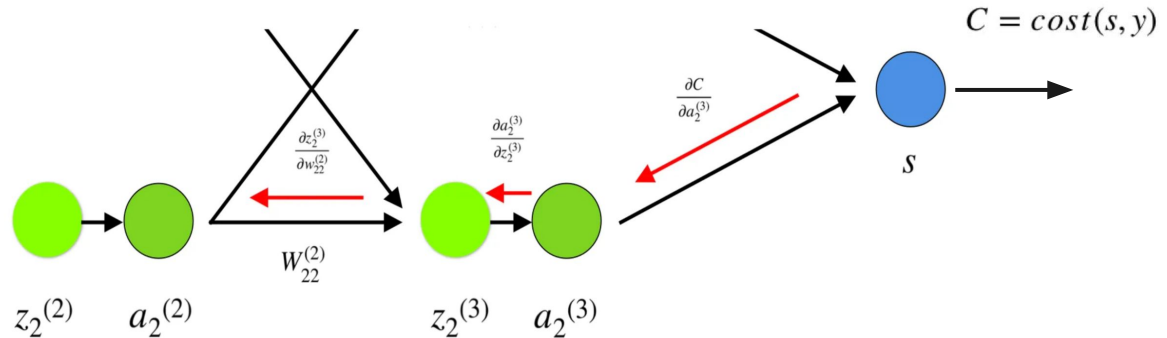


mainly adapted from

Error Backpropagation

One Random Parameter, $\mathbf{W}_{22}^{(2)}$

$$\begin{aligned}\frac{\partial C}{\partial w_{22}^{(2)}} &= \frac{\partial C}{\partial z_2^{(3)}} \cdot \frac{\partial z_2^{(3)}}{\partial w_{22}^{(2)}} \\ &= \frac{\partial C}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot a_2^{(2)} \\ &= \frac{\partial C}{\partial a_2^{(3)}} \cdot f'(z_2^{(3)}) \cdot a_2^{(2)}\end{aligned}$$





QA