# CIS 678 – Machine Learning

Introduction to Neural Networks

# Supervised Models

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

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

Regression

Classification

# Supervised Models

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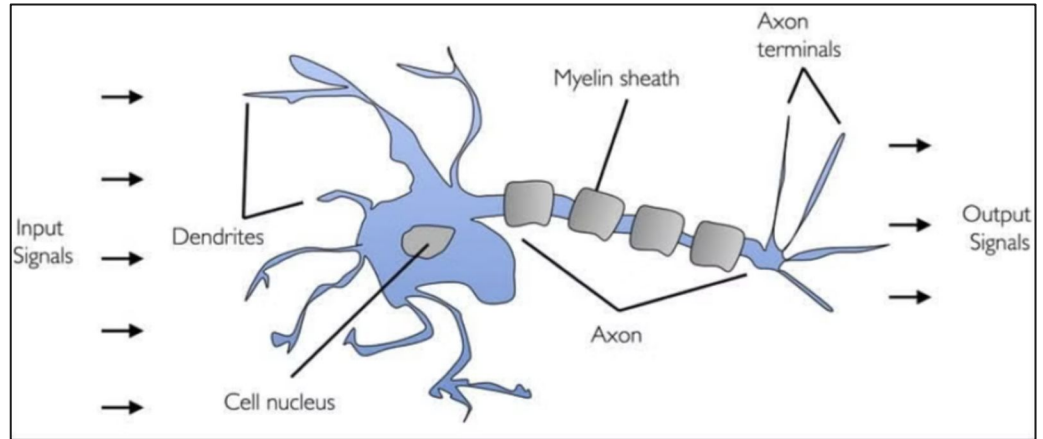- Neural Networks (NNs)

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

Regression

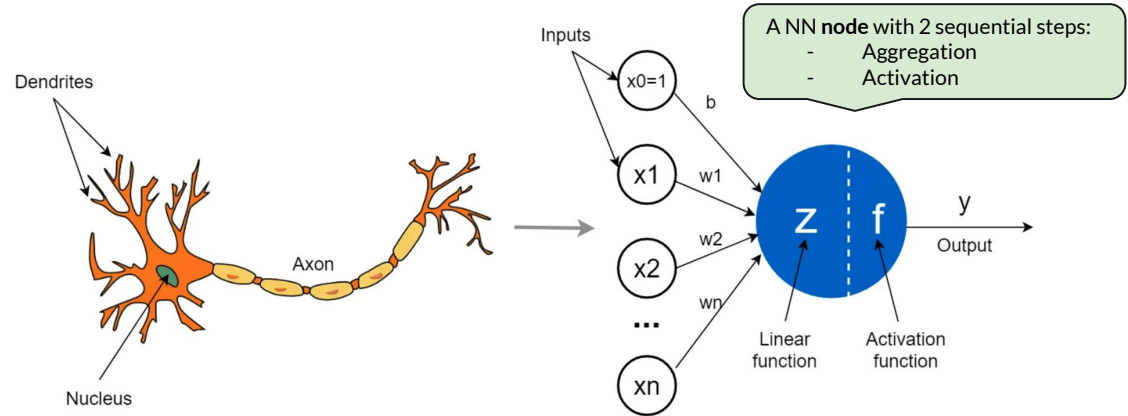Classification

# Neural Networks

Motivation src: Biological neuron

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



perceptron

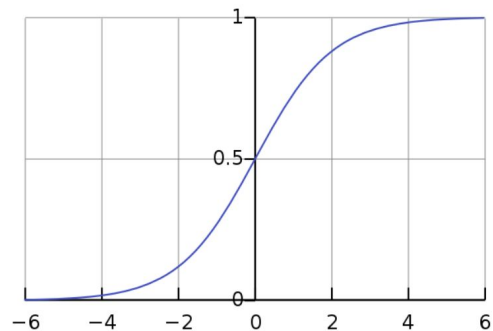# Neural Networks

From Biological Neuron to
Artificial NN



A NN **node** with 2 sequential steps:
- Aggregation
- Activation

fig-ref

# Logistic Regression

- Probabilistic classifier

- Sigmoid function

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

# Neural Networks (Node)

$$z = w_0 + w_1 x$$

1

x

$w_0$

$w_1$

$z$

**Linear**
(aggregation)

# Neural Networks (Node)



1

x

$w_0$

$w_1$

$z = w_0 + w_1 x$

$z$

$\frac{1}{1+e^{-z}}$

**Linear**
(aggregation)

**Activation function**

# Neural Networks (Node)



$$z = w_0 + w_1 x$$

$$\frac{1}{1+e^{-z}}$$

**1**

$w_0$

$z$

$y$

$w_1$

**x**

**Linear**
(aggregation)

**Activation
function**

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

**Logistic Regression**
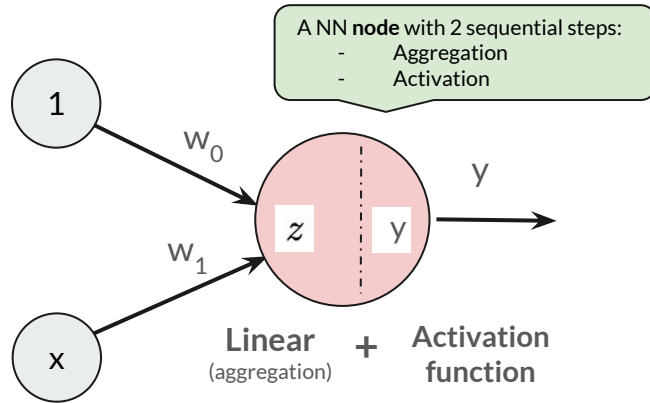
# Neural Networks (Node)

# Neural Networks (Node)

# Neural Networks (Node)



$z = w_0 + w_1 x$

$\dfrac{1}{1+e^{-z}}$

**Linear**
(aggregation)

**Activation function**

$y \in R$

**A NN with Sigmoid Activation function**

# Neural Networks (Node)



$1$    $w_0$

$x$    $w_1$

$z = w_0 + w_1 x$

**Linear**
(aggregation)

$z$

$z > 0$

**Activation
function**

$y$

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

**A NN with Step
Activation function**

**Frank Rosenblatt's
Perceptron**

# Neural Networks (Node)

A NN **node** with 2 sequential steps:
- Aggregation
- Activation

1

$w_0$

$z = w_0 + w_1 x$

???

$y$

$w_1$

x

**Linear**
(aggregation)

**Activation function**

A NN with ???
Activation function

# Neural Networks (Node)

# Neural Networks (No Hidden Layer)

A NN **node** with 2 sequential steps:
- Aggregation
- Activation

1

$w_0$

$x$

$w_1$

y

y

**Linear**
(aggregation)

**+**

**Activation
function**

A simple **one leayer NN** with 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)

# Feed-forward (FF) neural networks



$y = f(\mathbf{X}, \mathbf{Z})$

- Feedforward NN
- Perceptron

Input layer        Hidden layer(s)        Output layer

# Feed-forward (FF) neural networks

$x_1$

$x_2$

$x_3$

$x_4$

Input
(X)

# Feed-forward (FF) neural networks

$x_1 = a_1^{(1)}$

$x_2 = a_2^{(1)}$

$x_3 = a_3^{(1)}$

$x_4 = a_4^{(1)}$

A = X

# Feed-forward (FF) neural networks



$$z^{(2)} = W^{(1)}x + b^{(1)}$$

# Feed-forward (FF) neural networks



$$z^{(2)} = W^{(1)}x + b^{(1)}$$

$$a^{(2)} = f(z^{(2)})$$

mainly adapted from

# Feed-forward (FF) neural networks



$x_1 = a_1^{(1)}$

$x_2 = a_2^{(1)}$

$x_3 = a_3^{(1)}$

$x_4 = a_4^{(1)}$

$z_1^{(2)}$  $a_1^{(2)}$

$z_2^{(2)}$  $a_2^{(2)}$

$W^{(1)}$

$b^{(1)}$

Input layer

Hidden_1 layer

A = X

$$z^{(2)} = W^{(1)}x + b^{(1)}$$

$$a^{(2)} = f(z^{(2)})$$

Concept of Nodes (NN)

# Feed-forward (FF) neural networks



$$A = X$$

$$z^{(2)} = W^{(1)}x + b^{(1)}$$

$$a^{(2)} = f(z^{(2)})$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$

# Feed-forward (FF) neural networks



$$z^{(2)} = W^{(1)}x + b^{(1)}$$

$$a^{(2)} = f(z^{(2)})$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$

$$a^{(3)} = f(z^{(3)})$$

mainly adapted from

# Feed-forward (FF) neural networks



$$x_1 = a_1^{(1)}$$

$$x_2 = a_2^{(1)}$$

$$x_3 = a_3^{(1)}$$

$$x_4 = a_4^{(1)}$$

$z_1^{(2)}$  $a_1^{(2)}$

$z_2^{(2)}$  $a_2^{(2)}$

$W_{22}^{(2)}$

$z_1^{(3)}$  $a_1^{(3)}$

$z_2^{(3)}$  $a_2^{(3)}$

$s$

$W^{(1)}$  $b^{(1)}$  $W^{(2)}$  $b^{(2)}$  $W^{(3)}$

Input layer

A = X

Hidden_1 layer

$$z^{(2)} = W^{(1)}x + b^{(1)}$$

$$a^{(2)} = f(z^{(2)})$$

Hidden_2 layer

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$

$$a^{(3)} = f(z^{(3)})$$

Output layer

$$s = W^{(3)}a^{(3)}$$

# Feed-forward (FF) neural networks

# 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



$x_1 = a_1^{(1)}$

$x_2 = a_2^{(1)}$

$x_3 = a_3^{(1)}$

$x_4 = a_4^{(1)}$

$z_1^{(2)}$   $a_1^{(2)}$

$z_2^{(2)}$   $a_2^{(2)}$

$W_{22}^{(2)}$

$z_1^{(3)}$   $a_1^{(3)}$

$z_2^{(3)}$   $a_2^{(3)}$

$C = cost(s, y)$

$s$

$W^{(1)}$    $b^{(1)}$    $W^{(2)}$    $b^{(2)}$    $W^{(3)}$

Model parameters

Input layer

A = X

Hidden_1 layer

$$z^{(2)} = W^{(1)}x + b^{(1)}$$
$$a^{(2)} = f(z^{(2)})$$

Hidden_2 layer

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$
$$a^{(3)} = f(z^{(3)})$$

Output layer

$$s = W^{(3)}a^{(3)}$$

mainly adapted from

# Training



$x_1 = a_1^{(1)}$

$x_2 = a_2^{(1)}$

$x_3 = a_3^{(1)}$

$x_4 = a_4^{(1)}$

$z_1^{(2)}$  $a_1^{(2)}$

$z_1^{(3)}$  $a_1^{(3)}$

$C = cost(s, y)$

MSE (regression)

$s$

$W_{22}^{(2)}$

$z_2^{(2)}$  $a_2^{(2)}$

$z_2^{(3)}$  $a_2^{(3)}$

$W^{(1)}$     $b^{(1)}$     $W^{(2)}$     $b^{(2)}$     $W^{(3)}$

Model parameters

Input layer

A = X

Hidden_1 layer

$z^{(2)} = W^{(1)}x + b^{(1)}$

$a^{(2)} = f(z^{(2)})$

Hidden_2 layer

$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$

$a^{(3)} = f(z^{(3)})$

Output layer

$s = W^{(3)}a^{(3)}$

mainly adapted from

# Training



$x_1 = a_1^{(1)}$

$x_2 = a_2^{(1)}$

$x_3 = a_3^{(1)}$

$x_4 = a_4^{(1)}$

1

1

$z_1^{(2)}$   $a_1^{(2)}$

$z_2^{(2)}$   $a_2^{(2)}$

$W_{22}^{(2)}$

$z_1^{(3)}$   $a_1^{(3)}$

$z_2^{(3)}$   $a_2^{(3)}$

$C = cost(s, y)$

$s$

MSE (regression),
cross-entropy(classification)

$W^{(1)}$      $b^{(1)}$         $W^{(2)}$      $b^{(2)}$              $W^{(3)}$

Model
parameters

Input layer

A = X

Hidden_1 layer

Hidden_2 layer

Output layer

$$z^{(2)} = W^{(1)}x + b^{(1)}$$
$$a^{(2)} = f(z^{(2)})$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$
$$a^{(3)} = f(z^{(3)})$$

$$s = W^{(3)}a^{(3)}$$

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

$$\frac{\partial C}{\partial x} = [\frac{\partial C}{\partial x_1}, \frac{\partial C}{\partial x_2}, \ldots, \frac{\partial C}{\partial x_m}]$$

mainly adapted from

# 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

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

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

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

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

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

mainly adapted from

# 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*

Hidden layer

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

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

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

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

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

mainly adapted from

# 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

Hidden layer

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

constant

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

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

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

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

# 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

$$\frac{\partial C}{\partial b_j^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial b_j^l} \qquad chain\ rule$$

Derivative with respect to bias is 1

$$\frac{\partial z_j^l}{\partial b_j^l} = 1 \qquad by\ differentiation\ (calculating\ derivative)$$

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

# Gradient descent

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

$$while \ (termination \ condition \ not \ met)$$

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

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

$$end$$

mainly adapted from

# Training



$x_1 = a_1^{(1)}$

$x_2 = a_2^{(1)}$

$x_3 = a_3^{(1)}$

$x_4 = a_4^{(1)}$

$z_1^{(2)}$  $a_1^{(2)}$  $z_1^{(3)}$  $a_1^{(3)}$

$z_2^{(2)}$  $a_2^{(2)}$  $z_2^{(3)}$  $a_2^{(3)}$

$W_{22}^{(2)}$

$C = cost(s, y)$

MSE,
cross-entropy

$s$

$W^{(1)}$  $b^{(1)}$  $W^{(2)}$  $b^{(2)}$  $W^{(3)}$

Model parameters

Input layer

Hidden_1 layer

Hidden_2 layer

Output layer

A = X

$$z^{(2)} = W^{(1)}x + b^{(1)}$$
$$a^{(2)} = f(z^{(2)})$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$
$$a^{(3)} = f(z^{(3)})$$

$$s = W^{(3)}a^{(3)}$$

mainly adapted from

# Training

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



$x_1 = a_1^{(1)}$

$x_2 = a_2^{(1)}$

$x_3 = a_3^{(1)}$

$x_4 = a_4^{(1)}$

$z_1^{(2)}$   $a_1^{(2)}$   $z_1^{(3)}$   $a_1^{(3)}$

$z_2^{(2)}$   $a_2^{(2)}$

$W_{22}^{(2)}$

$z_2^{(3)}$   $a_2^{(3)}$

$C = cost(s, y)$

MSE,
cross-entropy

$s$

$W^{(1)}$   $b^{(1)}$   $W^{(2)}$   $b^{(2)}$   $W^{(3)}$

Model
parameters

Input layer

A = X

Hidden_1 layer

$z^{(2)} = W^{(1)}x + b^{(1)}$

$a^{(2)} = f(z^{(2)})$

Hidden_2 layer

$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$

$a^{(3)} = f(z^{(3)})$

Output layer

$s = W^{(3)}a^{(3)}$

mainly adapted from

# Training

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



$$\frac{\partial C}{\partial w^{(2)}_{22}} = \frac{\partial C}{\partial z^{(3)}_2} \cdot \frac{\partial z^{(3)}_2}{\partial w^{(2)}_{22}}$$

$$C = cost(s, y)$$

MSE,
cross-entropy

Model parameters

$x_1 = a_1^{(1)}$

$x_2 = a_2^{(1)}$

$x_3 = a_3^{(1)}$

$x_4 = a_4^{(1)}$

$z_1^{(2)} \quad a_1^{(2)}$

$z_2^{(2)} \quad a_2^{(2)}$

$W^{(2)}_{22}$

$z_1^{(3)} \quad a_1^{(3)}$

$z_2^{(3)} \quad a_2^{(3)}$

$s$

$W^{(1)} \qquad b^{(1)} \qquad W^{(2)} \qquad b^{(2)} \qquad W^{(3)}$

Input layer

A = X

Hidden_1 layer

$$z^{(2)} = W^{(1)}x + b^{(1)}$$
$$a^{(2)} = f(z^{(2)})$$

Hidden_2 layer

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$
$$a^{(3)} = f(z^{(3)})$$

Output layer

$$s = W^{(3)}a^{(3)}$$

mainly adapted from

# Error Backpropagation

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

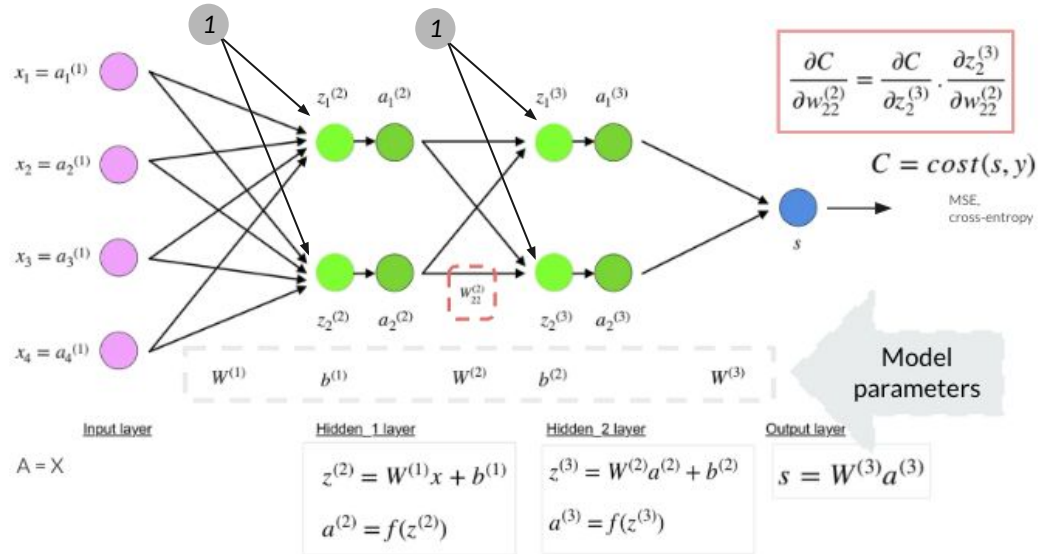$$\frac{\partial C}{\partial w^{(2)}_{22}} = \frac{\partial C}{\partial z^{(3)}_{2}} \cdot \frac{\partial z^{(3)}_{2}}{\partial w^{(2)}_{22}}$$

# Error Backpropagation

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

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

$$= \frac{\partial C}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot a_2^{(2)}$$



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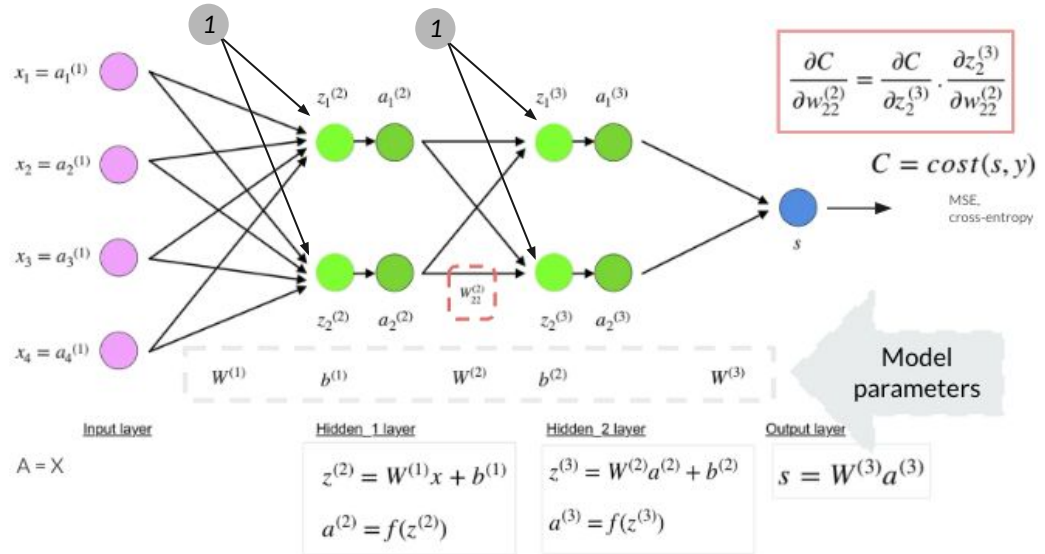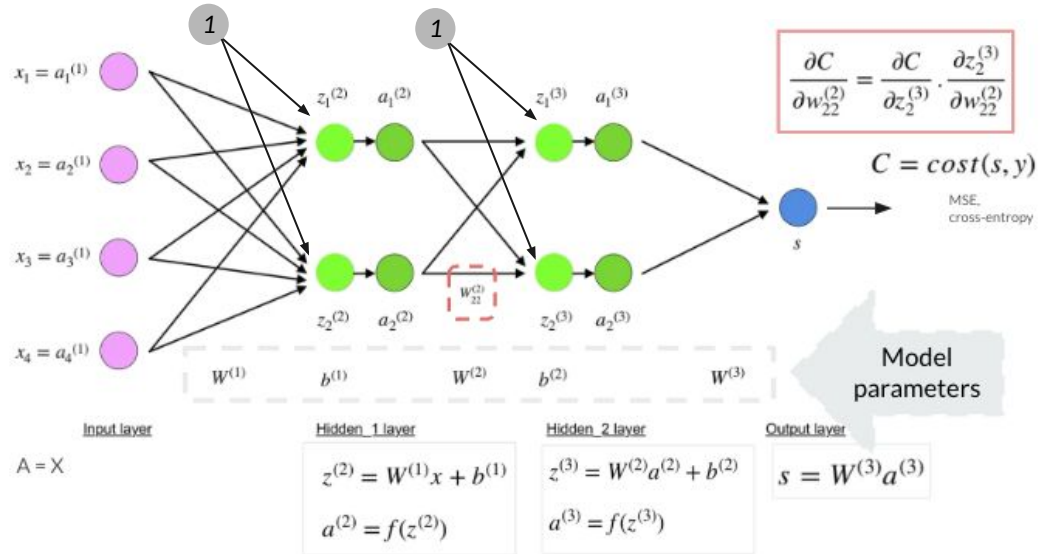
$C = cost(s, y)$

MSE, cross-entropy

Model parameters

$x_1 = a_1^{(1)}$  $x_2 = a_2^{(1)}$  $x_3 = a_3^{(1)}$  $x_4 = a_4^{(1)}$

$z_1^{(2)}$  $a_1^{(2)}$  $z_1^{(3)}$  $a_1^{(3)}$

$z_2^{(2)}$  $a_2^{(2)}$  $z_2^{(3)}$  $a_2^{(3)}$

$W^{(2)}_{22}$

$W^{(1)}$  $b^{(1)}$  $W^{(2)}$  $b^{(2)}$  $W^{(3)}$

Input layer  Hidden 1 layer  Hidden 2 layer  Output layer

$A = X$

$z^{(2)} = W^{(1)}x + b^{(1)}$

$a^{(2)} = f(z^{(2)})$

$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$

$a^{(3)} = f(z^{(3)})$

$s = W^{(3)}a^{(3)}$

mainly adapted from

# Error Backpropagation

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

$$\frac{\partial C}{\partial w^{(2)}_{22}} = \frac{\partial C}{\partial z^{(3)}_2} \cdot \frac{\partial z^{(3)}_2}{\partial w^{(2)}_{22}}$$

$$= \frac{\partial C}{\partial a^{(3)}_2} \cdot \frac{\partial a^{(3)}_2}{\partial z^{(3)}_2} \cdot a^{(2)}_2$$

$$= \frac{\partial C}{\partial a^{(3)}_2} \cdot f'(z^{(3)}_2) \cdot a^{(2)}_2$$

# Error Backpropagation

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

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

$$= \frac{\partial C}{\partial a_2^{(3)}} \cdot \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \cdot a_2^{(2)}$$

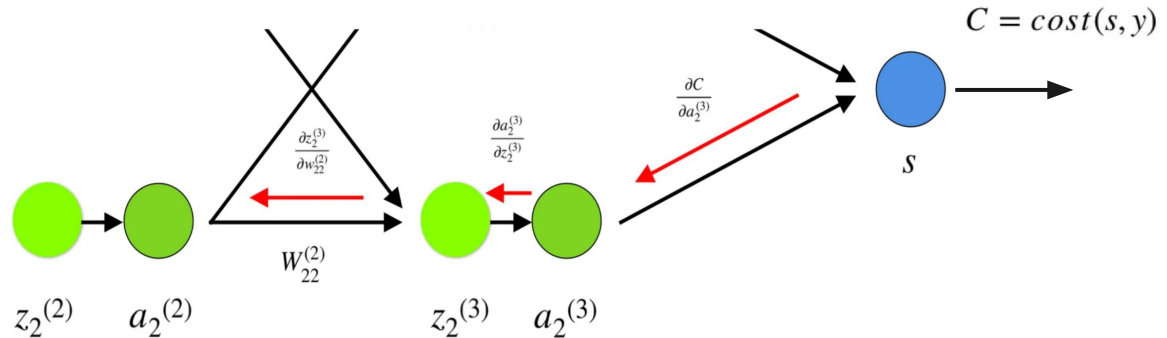$$= \frac{\partial C}{\partial a_2^{(3)}} \cdot f'(z_2^{(3)}) \cdot a_2^{(2)}$$



$C = cost(s, y)$

# QA