

SDSU CS 549 Spring  
2024

# Machine Learning Lecture 8: Neural Networks

---

IRFAN KHAN

# References

---

SDSU CS549 Lecture Notes by Prof Yang Xu, Spring 2023. Some slides used here

Coursera machine learning course by Dr Andrew Ng, Oct 2023

# Introduction

---

- We start with Artificial Neural Networks (ANN) used for Regression and Logistic Regression
- These ANNs fall in the category of Feed-Forward Neural Networks (FNN)
  - In FNNs, flow is uni-directional => information in the model flows in only one direction—forward—from the input nodes, through the hidden nodes and to the output nodes, without any cycles or loops
  - In contrast recurrent neural networks, which we will cover later, have a bi-directional flow.
- Many of these concepts carry over to other types of neural networks, e.g. Convolutional Neural Networks used in Object Detection

# Outline

---

Neural Network Architecture

Back Propagation Intro

Computational Graph

Computational Graph and Basic NN Code Examples

Activation Functions

Derivatives of Vectors and Matrices and Back propagation

Generalize Back propagation for Deeper ANNs Using Vectors and Matrices

Backpropagation Code Examples

Initialization

# Neural Network Architecture

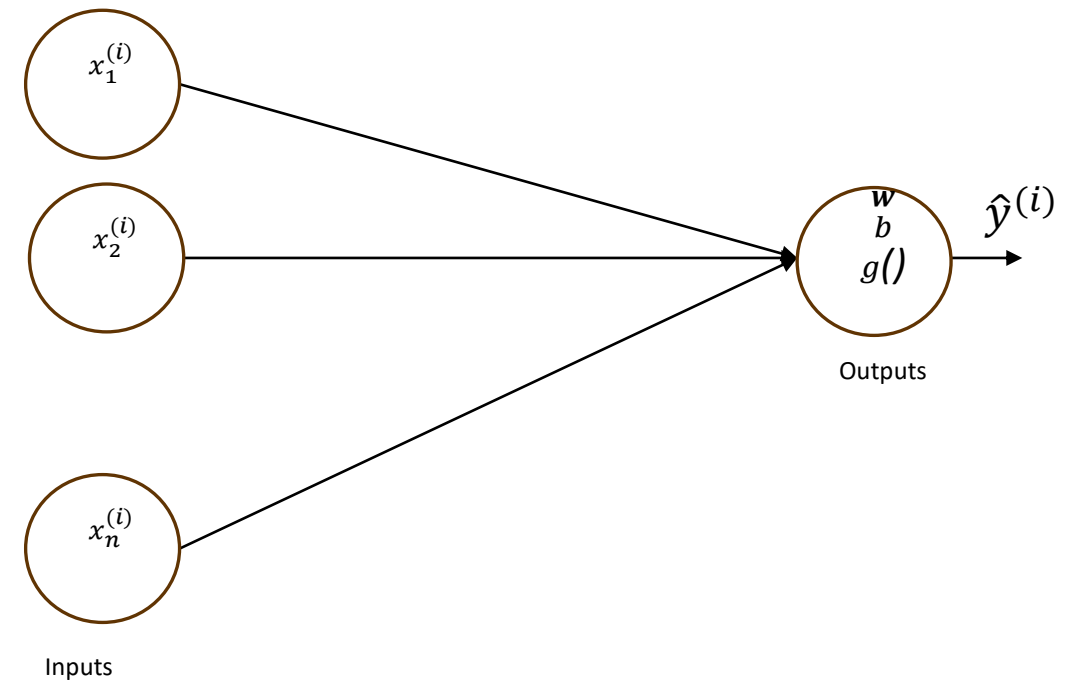
---

# Recap – Linear Regression and Logistic Regression

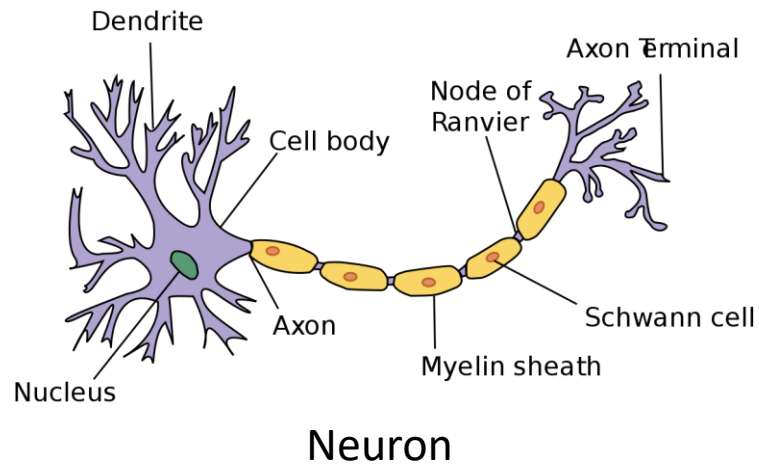
$$\begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(m)} \end{bmatrix} = g \left( \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} + \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix} \right)$$

$g()$  = pass through function for Linear Regression

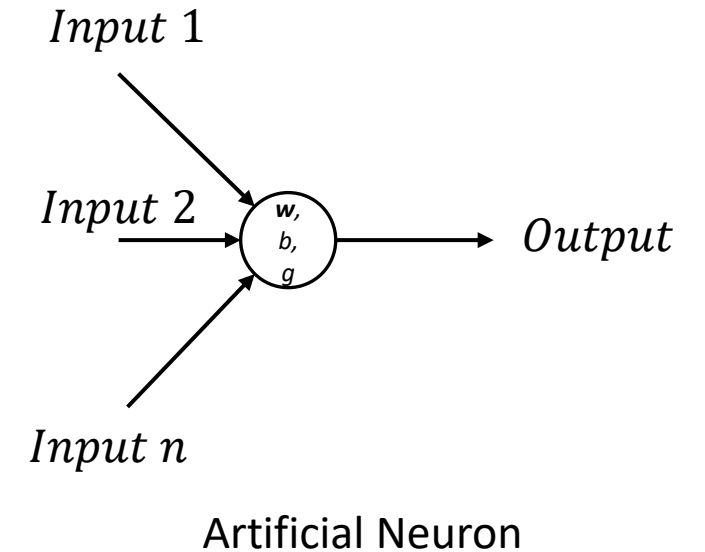
$g()$  = Sigmoid function for Logistic Regression



# Artificial Neuron

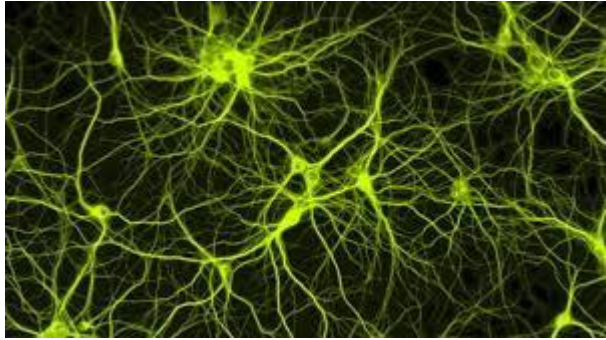


Inspired from biological neuron



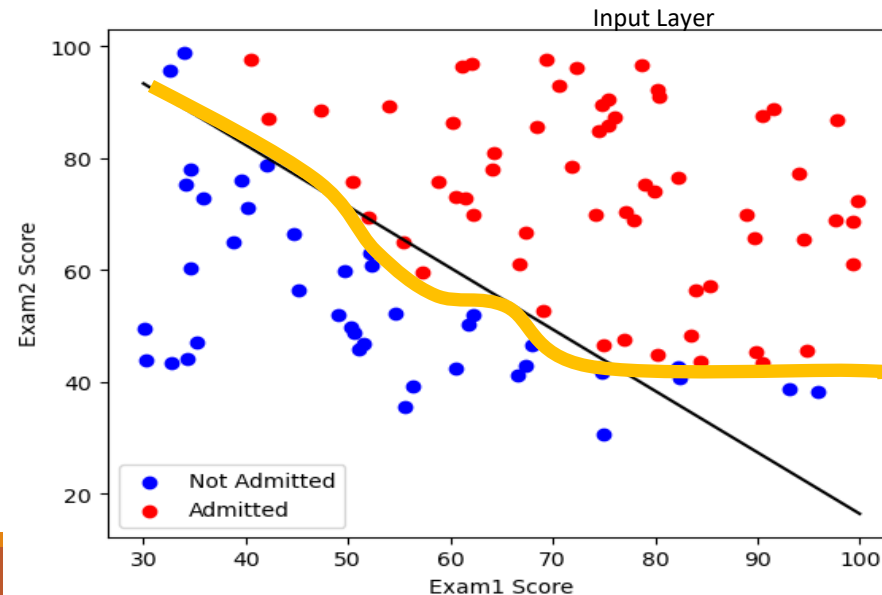
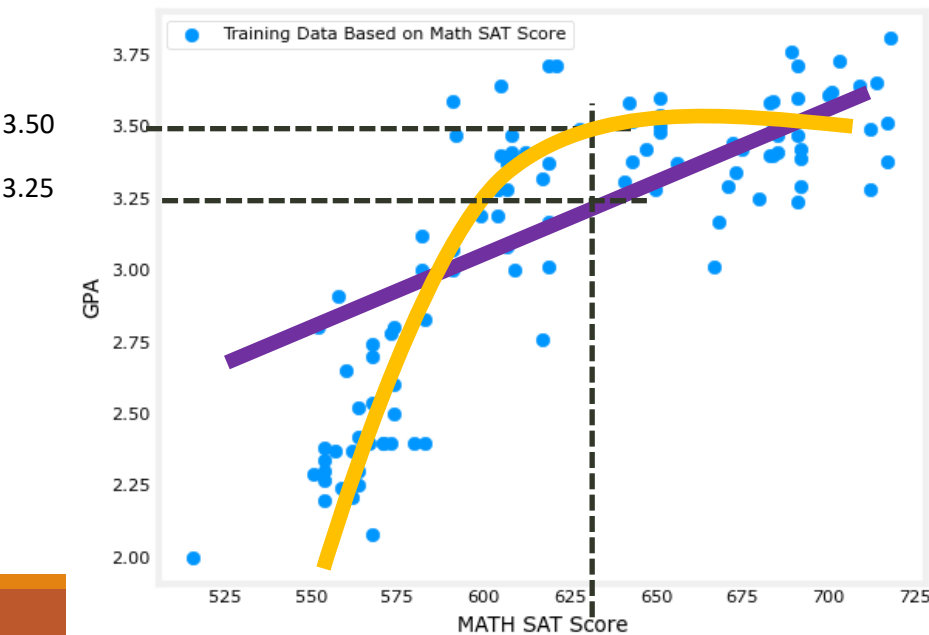
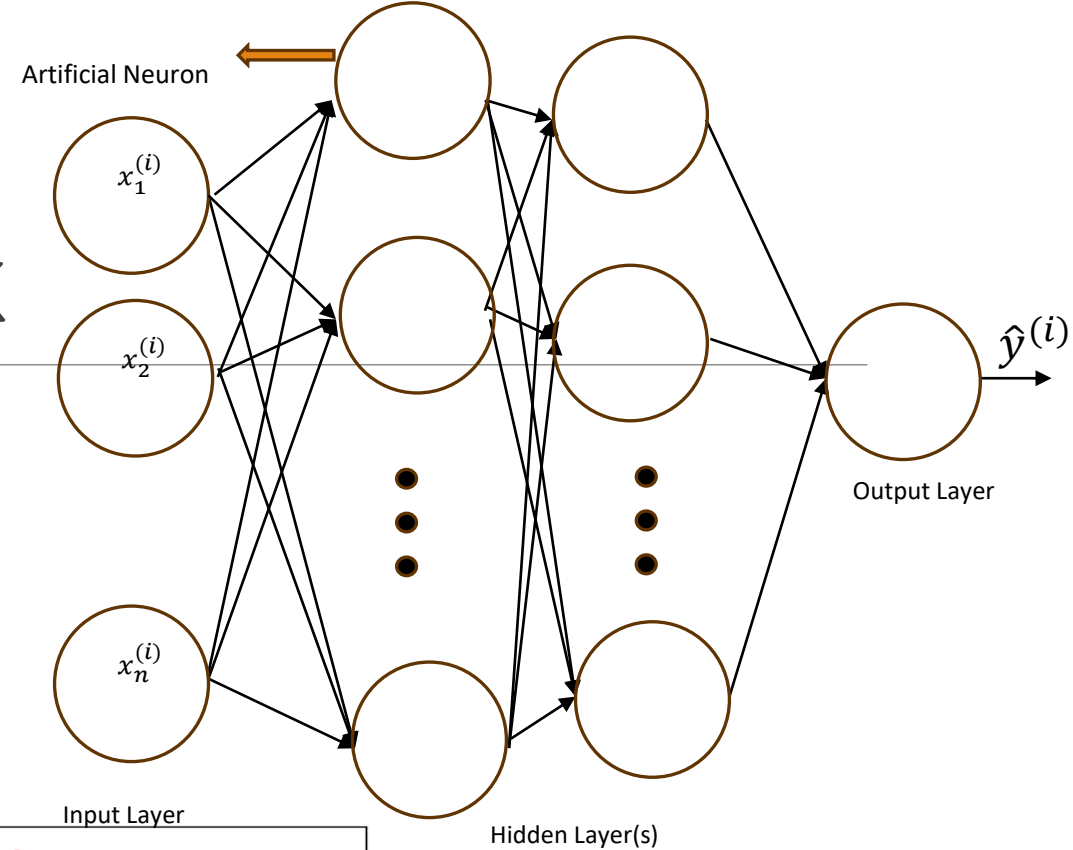
$w$ : "weight",  $b$ : "bias",  $g$ : activation function applied to  $w\mathbf{x}+b$

# Artificial Neural Network



Inspired from biological neural network => ANN

Neural Network: A dense interconnection of neurons

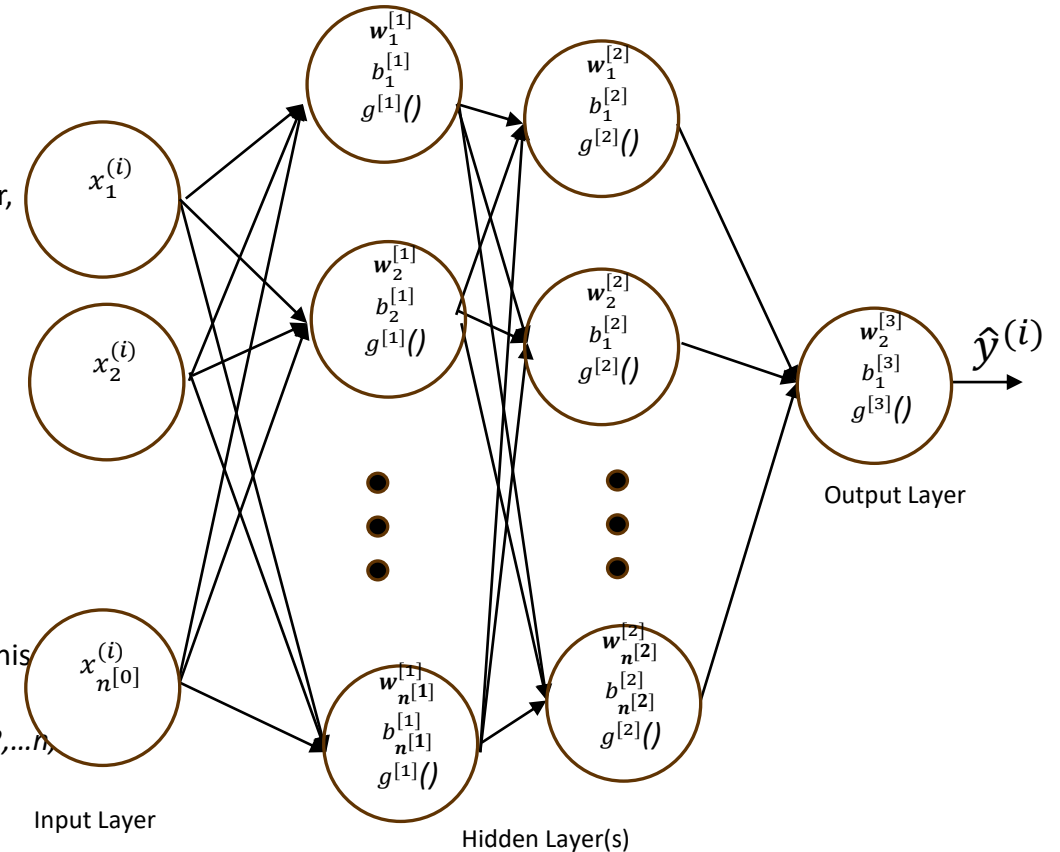


- Artificial Neural Network (ANN): A dense interconnection of neurons with *introduced non-linearities*
  - What happens if activation functions in the hidden layers are linear?
- Automatically creates more complex data fitting functions like the orange curve on the left
- The more the number of hidden layers, the “deeper” the ANN



# ANN - Parameters

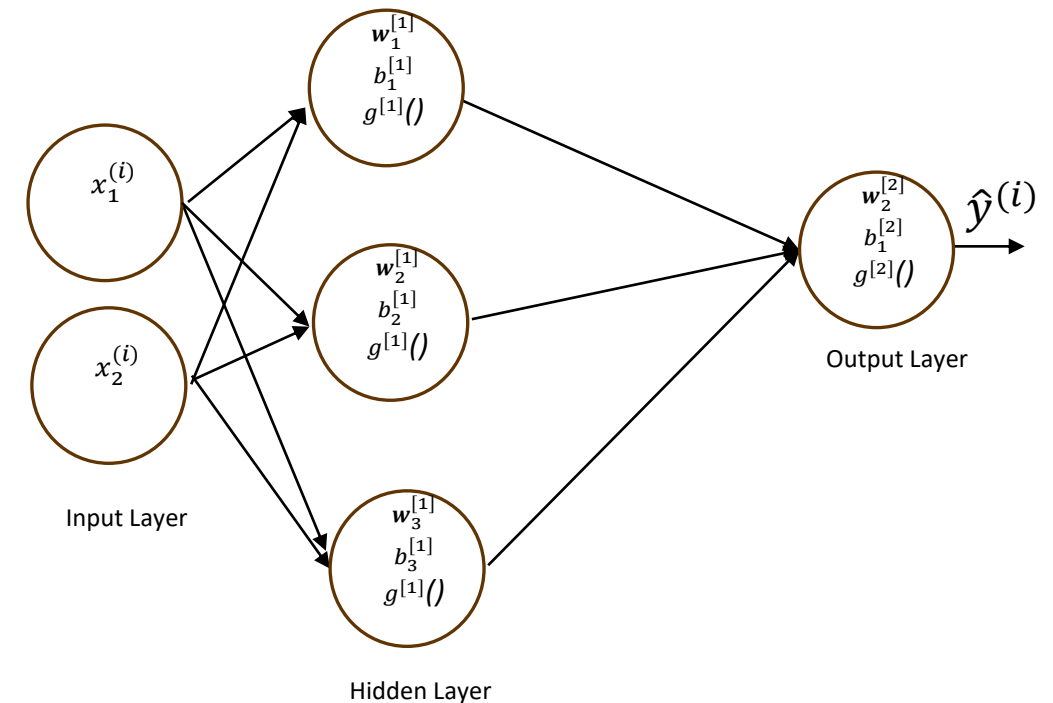
- Layers in a neural network are composed of neurons
  - $l = 0, 1, 2, \dots, L$  where  $l=0$  denotes the input layer and  $L$  denotes the output layer.
  - So total number of “hidden” layers =  $L-1$
- $n^{[l]}$  represents the number of neurons in the  $l^{th}$  layer
- $w_{ij}^{[l]}$ : weight parameter in the  $i^{th}$  neuron in layer  $l$  applied to output of the  $j^{th}$  neuron of the previous layer,  $l=1, 2, \dots, L$ 
  - Define the row vector at the  $i^{th}$  neuron in layer  $l$  as:  $\mathbf{w}_i^{[l]} = [w_{i1}^{[l]} \quad w_{i2}^{[l]} \quad \dots \quad w_{in_{l-1}}^{[l]}]$
- $b_i^{[l]}$ : bias parameter in the  $i^{th}$  neuron in layer  $l$ ,  $l=1, 2, \dots, L$
- $g^{[l]}()$ : Activation function at all neurons in layer  $l$ 
  - Activation functions help create non-linearities
  - Number of choices are available, e.g. ReLU
- For the output layer:
  - For regression, a pass-through (linear) function is used
  - For classification, sigmoid layer is used
- Layer 0 consists just of the input data and no weights or bias or activation function are associated with this layer
- As usual, input data is represented by  $x_j^{(i)}$ , where  $i = 1, 2, \dots, m$ ,  $m$  being the number of samples, and  $j = 1, 2, \dots, n$ ,  $n$  being the number features in the input data
- Layer  $L$ , the output layer, predicts the output,  $\hat{y}$ . So, only one neuron in the output layer.
  - In general,  $\hat{y}$  can be a vector with multiple dimensions we may want to predict, in which case the output layer also will have multiple neurons, one for each output dimension



$$\text{Total Parameters} = \sum_{l=1}^L n^{[l]} n^{[l-1]} + n^{[L]}$$

# Total Parameters for an ANN

- $L = 2, n = 2$  for the simple/"shallow" ANN shown on the right
- Total Parameters  $= \sum_{l=1}^L n^{[l]} n^{[l-1]} + n^{[L]} = 3*2+3 + 1*3+1 = 13$
- Number of parameters can grow rapidly with "deep" ANNs
- ResNet-50: deep CNN with over 23 million parameters, designed for image classification tasks.
  - We will cover CNNs in a later lecture
- One metric for HW hosting neural networks is TOPs (trillions of Operations per second!)
  - Entry level Nvidia A2 GPU: 36 INT8 TOPs
  - <https://www.nvidia.com/content/dam/en-zz/solutions/data-center/a2/pdf/a2-datasheet.pdf>
- How is all this possible?
  - Matrix manipulation that can be parallelized in hardware



# Matrix Operations and Parallelization Scope in Hardware

➤  $A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix}$

➤  $3 \times 2$

➤  $B = \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \end{bmatrix}$

➤  $2 \times 3$

➤  $C = A * B = \begin{bmatrix} c_{11} & c_{12} & c_{13} \\ c_{21} & c_{22} & c_{23} \\ c_{31} & c_{32} & c_{33} \end{bmatrix}$

➤ Using 'for' loops:

➤ for  $i$  in range (3):

for  $j$  in range (3):

$c[i,j]=0$

for  $k$  in range(2):

$C[i,j]=C[i,j]+A[i,k]*B[k,j]$

➤ Using matrix function in NumPy

➤  $C=np.dot(A,B)$

➤ All  $a_{ik} * b_{kj}$  simple multiplications can be executed in hardware in parallel.

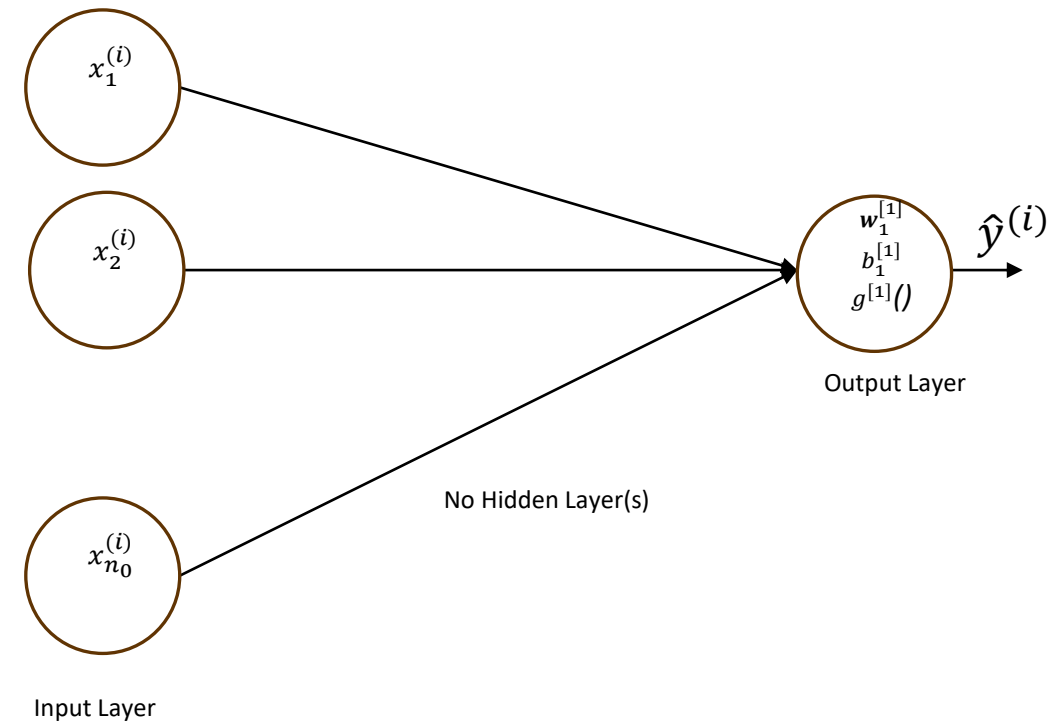
➤ All that remains next is a simple addition using terms computed above:  $c_{ij} = \sum_{k=1}^2 a_{ik} * b_{kj}$  which also can be computed in hardware in parallel for each element of  $C$

# $L=1$ corresponds to Legacy Regression Models

---

➤ Zero hidden layers

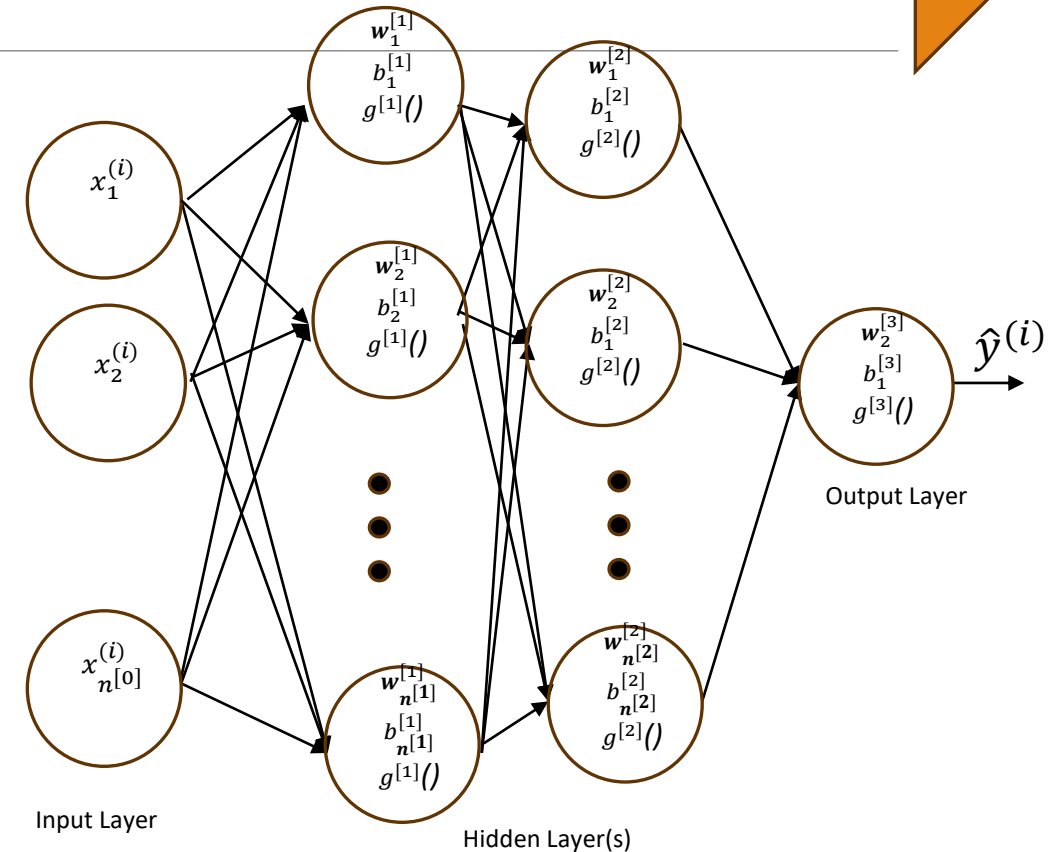
➤ Total Parameters =  $\sum_{l=1}^L n^{[l]} n^{[l-1]} + n^{[l]} = n_0 * 1 + 1 = n_0 + 1$



# ANN in Practice

- How do we train the ANN => how do we learn values for  $w_j^{[l]}$  and  $b_j^{[l]}$  which will minimize the cost function  $J$ ?
  - Gradient Descent Back Propagation
  - Needs Forward Propagation to update parameters and to compute cost function,  $J$  iteratively
- Once trained, forward propagation, allows us to compute the prediction,  $\hat{y}$
- Considerations:
  - How many hidden layers?
  - How many neurons/hidden layer?
  - Which activation functions to use?

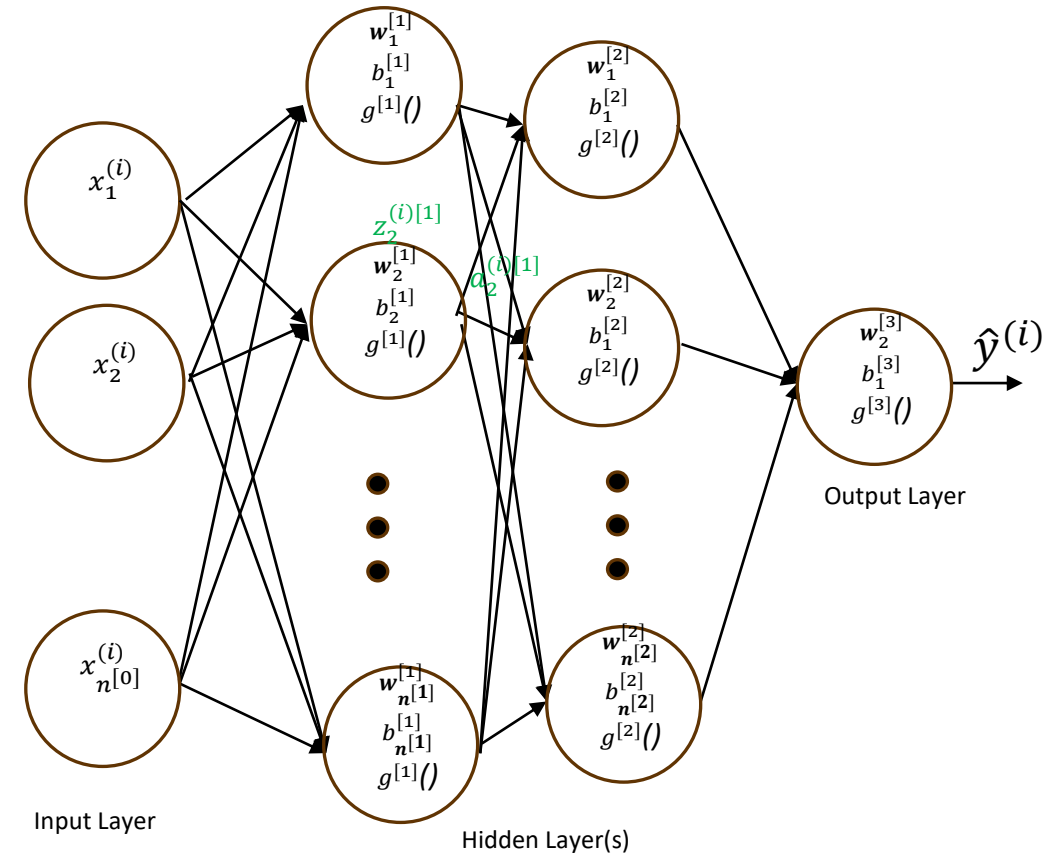
Forward Propagation for Inference/Prediction ( $\hat{y}$ ) and during training



Training: Gradient Descent Back Propagation

# Forward Propagation/Inference: Computing $\hat{y}^{(i)}$

- Introducing two new variables at neuron  $j$  in the  $l^{th}$  layer
  - $z_j^{(i)[l]}$  which computes an interim value after applying the weights and bias
  - $a_j^{(i)[l]}$  representing the output at neuron  $j$  after applying the activation function
- $z_j^{(i)[l]} = \sum_{k=1}^{n_{l-1}} a_k^{(i)[l-1]} w_{kj}^{[l]} + b_j^{[l]}, l = 1, 2, \dots, L, j = 1, 2, \dots, n^{[l]}$
- $a_j^{(i)[l]} = g^{[l]}(z_j^{(i)[l]})$
- $a_j^{(i)[0]} = x_j^{(i)} j = 1, 2, \dots, n^{[0]}$
- $\hat{y}^{(i)} = a_j^{(i)[L]} (j = 1)$ , for the single output case



# Vector/Matrix Formulation for Computing $\hat{y}^{(i)}$

➤ Input data matrix:  $X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} m \times n$

➤ Output prediction:  $\hat{y} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(m)} \end{bmatrix} m \times 1$

➤ Weight matrix at the  $l^{th}$  layer:  $W^{[l]} = \begin{bmatrix} w_{11}^{[l]} & w_{12}^{[l]} & \dots & w_{1n^{[l-1]}}^{[l]} \\ w_{21}^{[l]} & w_{22}^{[l]} & \dots & w_{2n^{[l-1]}}^{[l]} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n^{[l]}1}^{[l]} & w_{n^{[l]}2}^{[l]} & \dots & w_{n^{[l]}n^{[l-1]}}^{[l]} \end{bmatrix} n^{[l]} \times n^{[l-1]}$

➤ by convention used in PyTorch, could have been a  $n^{[l-1]} \times n^{[l]}$  matrix instead.

➤ Bias vectors at the  $l^{th}$  layer:  $B^{[l]} = \begin{bmatrix} b_1^{[l]} & b_2^{[l]} & \dots & b_{n^{[l]}}^{[l]} \end{bmatrix} 1 \times n^{[l]}$

➤ Intermediate value at the  $l^{th}$  layer:

➤  $Z^{[l]} = \begin{bmatrix} z_1^{(1)[l]} & z_2^{(1)[l]} & \dots & z_{n^{[l]}}^{(1)[l]} \\ z_1^{(2)[l]} & z_2^{(2)[l]} & \dots & z_{n^{[l]}}^{(2)[l]} \\ \vdots & \vdots & \ddots & \vdots \\ z_1^{(m)[l]} & z_2^{(m)[l]} & \dots & z_{n^{[l]}}^{(m)[l]} \end{bmatrix} m \times n^{[l]}$

➤ Output values (activations) at the  $l^{th}$  layer:  $A^{[l]} =$

$\begin{bmatrix} a_1^{(1)[l]} & a_2^{(1)[l]} & \dots & a_{n^{[l]}}^{(1)[l]} \\ a_1^{(2)[l]} & a_2^{(2)[l]} & \dots & a_{n^{[l]}}^{(2)[l]} \\ \vdots & \vdots & \ddots & \vdots \\ a_1^{(m)[l]} & a_2^{(m)[l]} & \dots & a_{n^{[l]}}^{(m)[l]} \end{bmatrix} m \times n^{[l]}$

➤  $A^{[0]} = X: m \times n = m \times n^{[0]}$

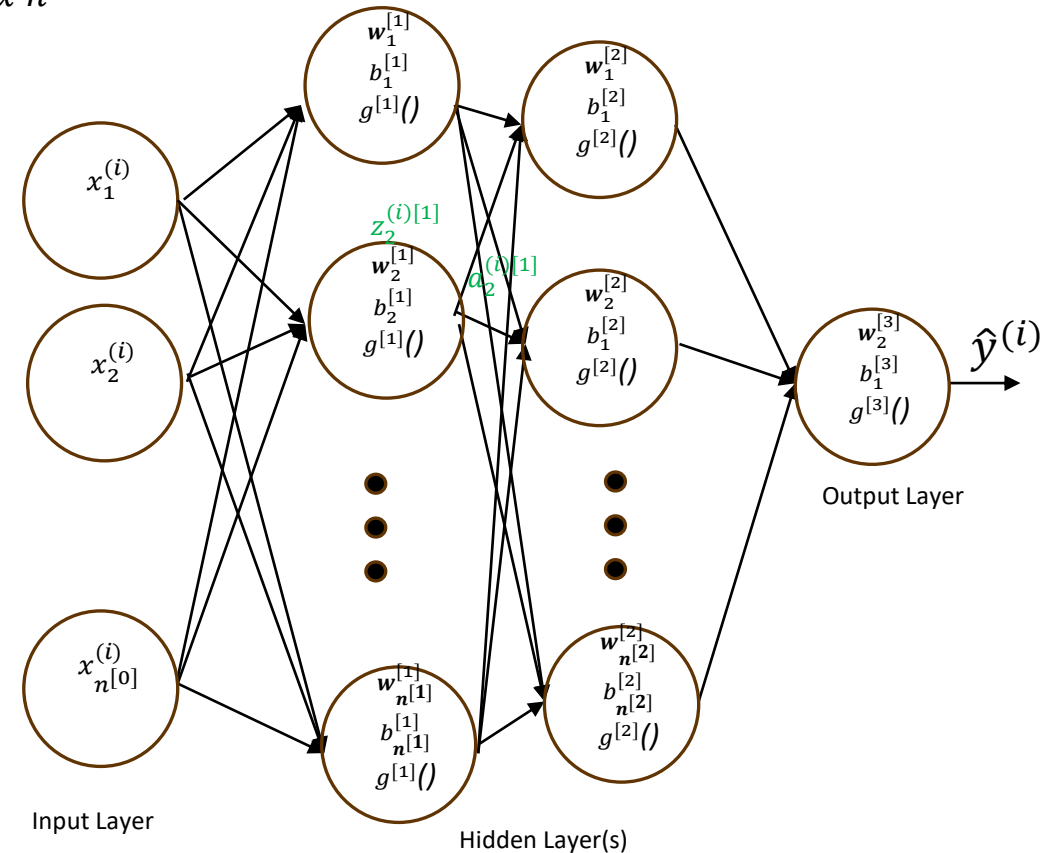
➤  $A^{[l]} = g^{[l]}(Z^{[l]}): m \times n^{[l]}$

➤  $Z^{[l]} = A^{[l-1]}W^{[l]T} + B^{[l]}: m \times n^{[l]}$

➤ Check matrix dimensions:

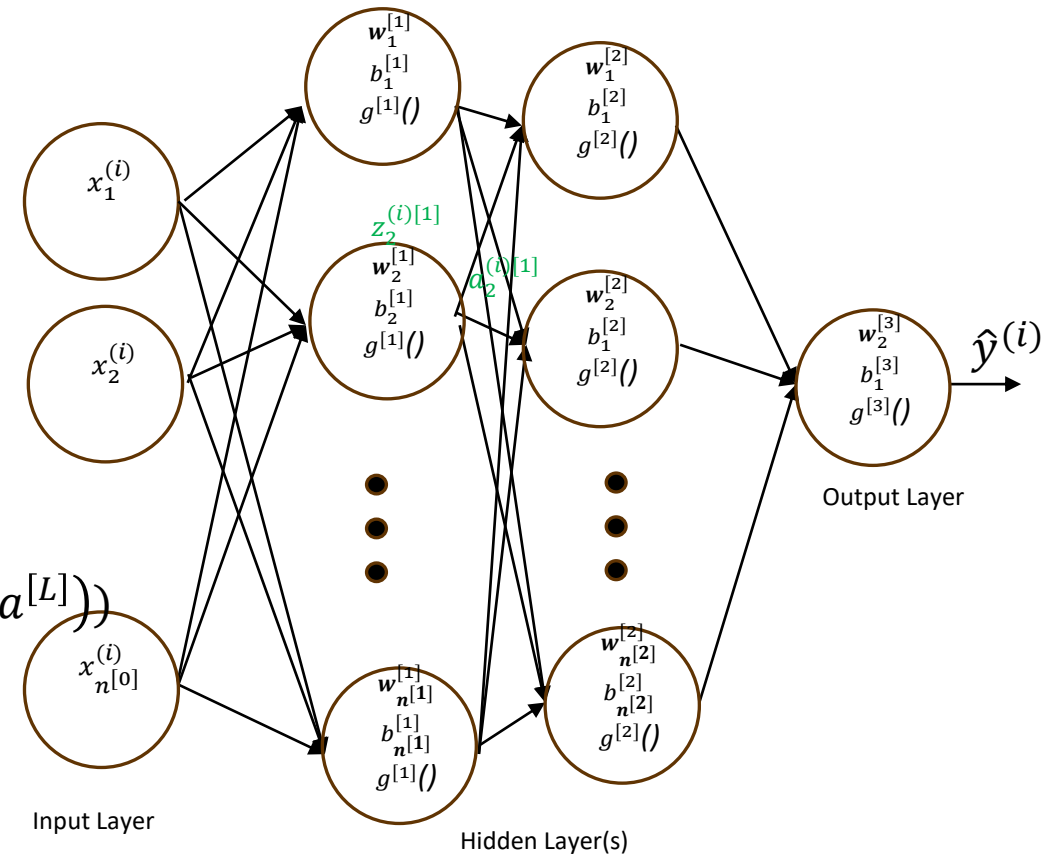
➤  $m \times n^{[l]} = m \times n^{[l-1]} * n^{[l-1]} \times n^{[l]} + m \times n^{[l]}$

➤  $\hat{y} = A^{[L]}: m \times n^{[L]}: m \times n^{[L]}: m \times 1$  when there is only one output.



# Computing the Cost Function $J$

- Let the sample data output  $\mathbf{y} = \begin{bmatrix} \mathbf{y}^{(1)} \\ \mathbf{y}^{(2)} \\ \vdots \\ \mathbf{y}^{(m)} \end{bmatrix}$
- For a regression problem,
  - $J = \frac{1}{2m} \sum_{i=1}^m (\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)})^2 = \frac{1}{2m} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}})$
- For a Logistic Regression problem,
  - $J = -\frac{1}{m} \sum_{i=1}^m \mathbf{y}^{(i)} \log(\mathbf{a}^{(i)[L]}) + (1 - \mathbf{y}^{(i)}) \log(1 - (\mathbf{a}^{(i)[L]})$
  - $J = -\frac{1}{m} (\text{Rowsum}(\mathbf{y} * \log(\mathbf{a}^{[L]} + (1 - \mathbf{y}) * \log(1 - \mathbf{a}^{[L]})))$

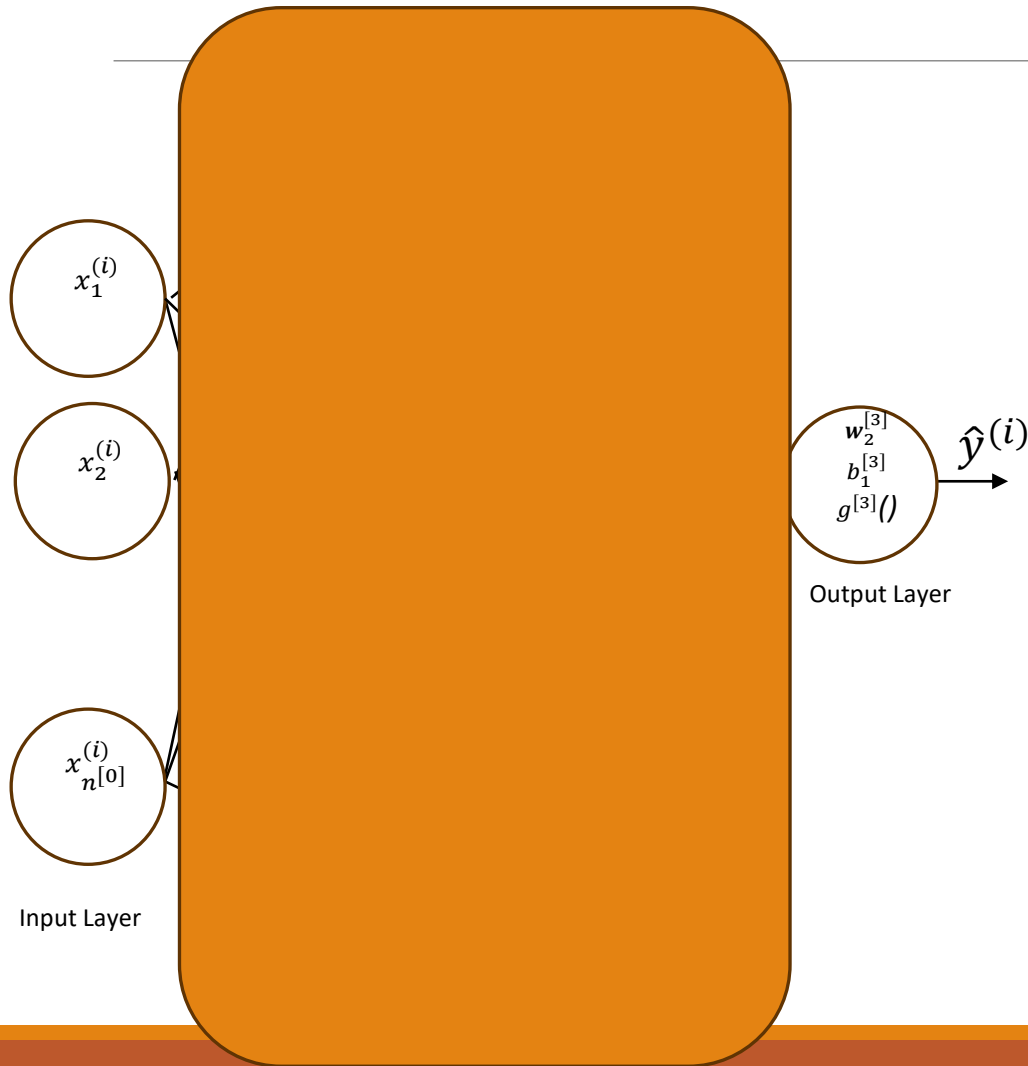




# Gradient Descent Back Propagation Intro

---

# Gradient Descent Back Propagation for an ANN



Recall Gradient Descent steps

- Take derivative of cost function wrt  $W$  and wrt  $b$
- Move  $W$  and  $b$  opposite to the direction of the gradient scaled by a learning factor of  $\alpha$
- Continue until convergence

How do we back-propagate gradients through various layers?

# Computation Graph

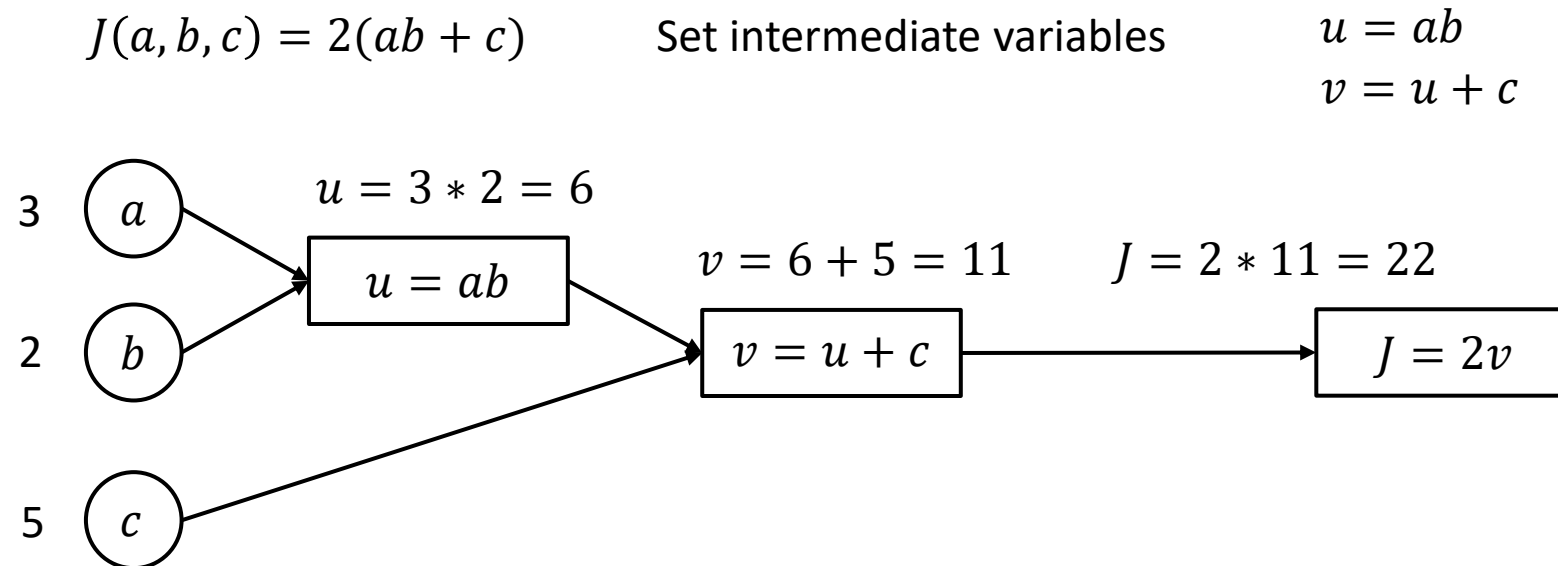
---

# Computation Graph

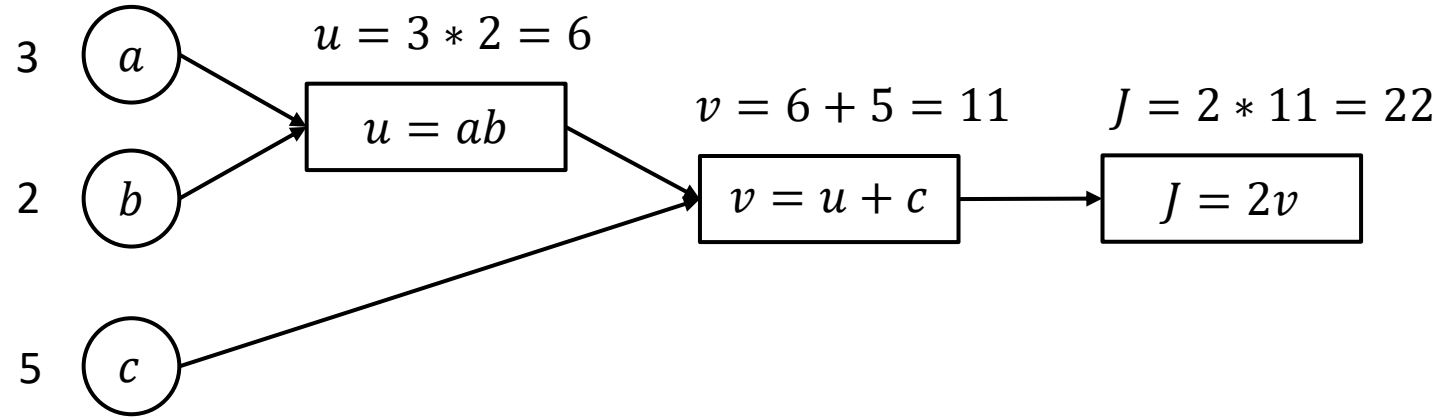
What is it?

A useful tool to understand/visualize ML models

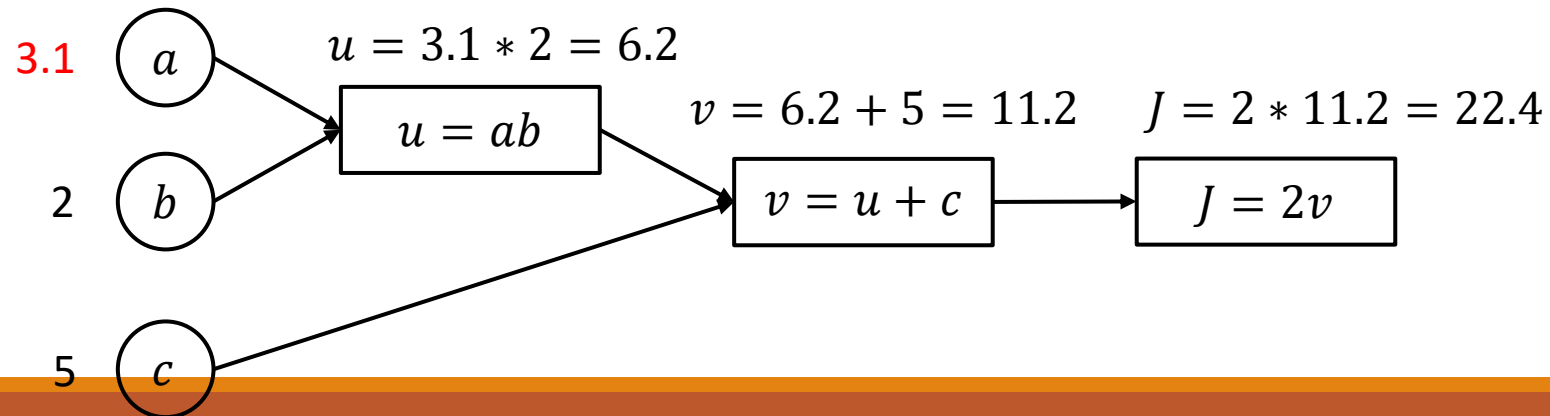
Simple arithmetic example



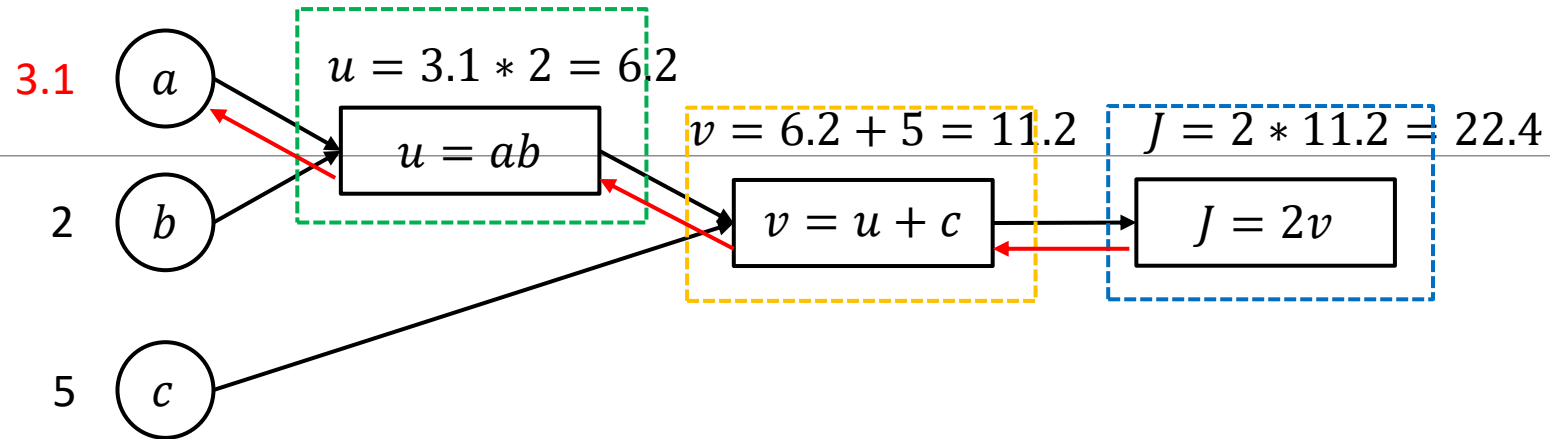
# Input change leads to output change



How much will  $J$  change if  $a$  changes a little bit?



# Derivative of $J$ w.r.t $a$



$$\frac{\partial J}{\partial a} = \frac{22.4 - 22}{3.1 - 3} = 4 \quad J \text{ will change by 4 units if } a \text{ changes by 1}$$

What about going Backward **step by step**?

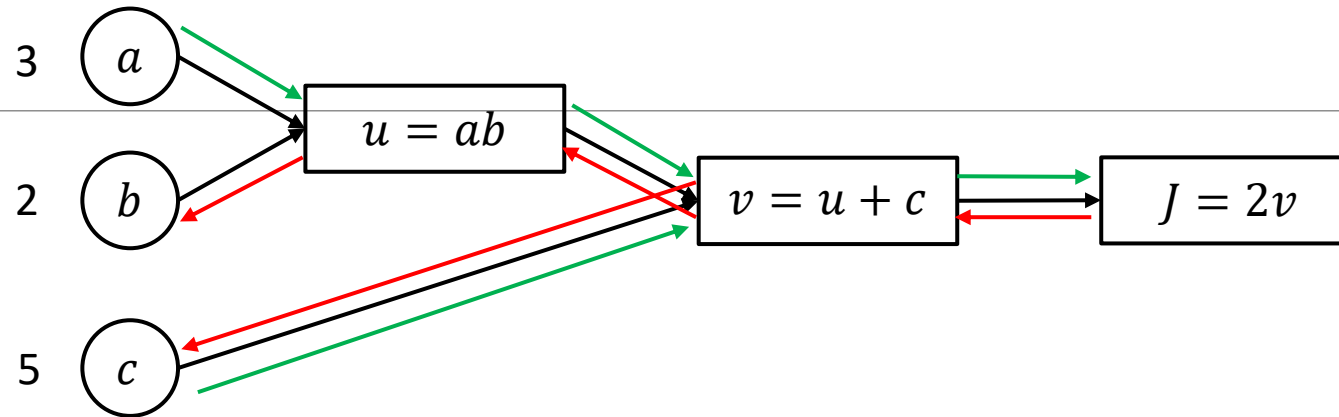
$$\frac{\partial u}{\partial a} = \frac{6.2 - 6}{3.1 - 3} = 2$$

$$\frac{\partial v}{\partial u} = \frac{11.2 - 11}{6.2 - 6} = 1$$

$$\frac{\partial J}{\partial v} = \frac{22.4 - 22}{11.2 - 11} = 2$$

Apply chain rule:  $\frac{\partial J}{\partial v} \frac{\partial v}{\partial u} \frac{\partial u}{\partial a} = 2 * 1 * 2 = 4 = \frac{\partial J}{\partial a}$  Same result

# Derivative of $J$ w.r.t all inputs



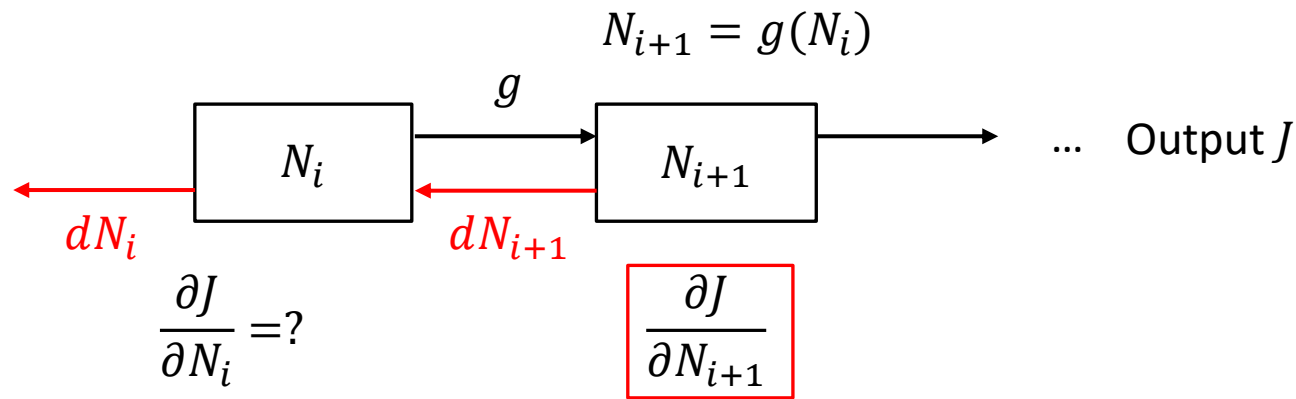
$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial v} \frac{\partial v}{\partial u} \frac{\partial u}{\partial b} = 2 * 1 * 3 = 6$$

$$\frac{\partial J}{\partial c} = \frac{\partial J}{\partial v} \frac{\partial v}{\partial c} = 2 * 1 = 2$$

Forward pass (green), compute the cost  $J$

Backward pass (red), compute the derivative  $\frac{\partial J}{\partial ?}$

# Chain rule between nodes



$$\frac{\partial J}{\partial N_i} = ?$$

$$= \frac{\partial J}{\partial N_{i+1}} \cdot \frac{\partial N_{i+1}}{\partial N_i} \quad \text{dot-product for vectors/matrices}$$

$$\boxed{dN_i = \frac{\partial J}{\partial N_{i+1}} g'(N_i)} = dN_{i+1} g'(N_i)$$

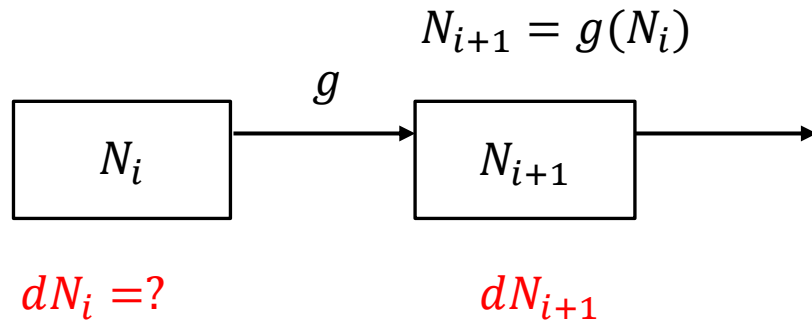
The derivative of one node can be inferred from that of the adjacent node, and the **function** that connects them

Notation for **derivatives**:

$$dN_{i+1} \text{ is for } \frac{\partial J}{\partial N_{i+1}}$$



# Practice

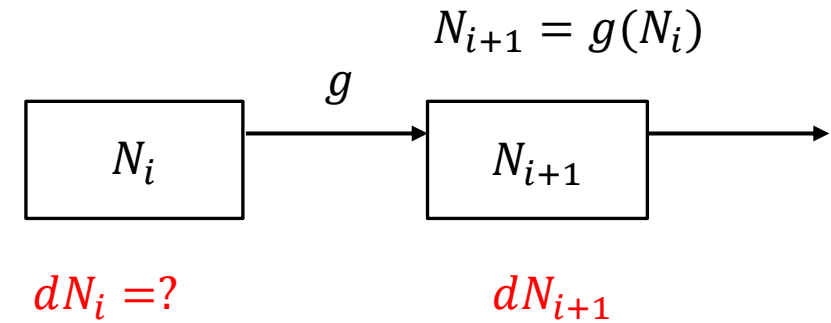


Sigmoid  $g(z) = \frac{1}{1+e^{-z}}$

$$dN_i = dN_{i+1}g'(N_i) = ?$$

Answer:  $g'(z) = g(z)(1 - g(z))$

$$dN_i = dN_{i+1}g(N_i)(1 - g(N_i)) = dN_{i+1}N_{i+1}(1 - N_{i+1})$$



Tanh  $g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

$$dN_i = dN_{i+1}g'(N_i) = ?$$

Answer:  $g'(z) = 1 - g^2(z)$

$$dN_i = dN_{i+1}(1 - g^2(N_i)) = dN_{i+1}(1 - N_{i+1}^2)$$

# Computation graph for logistic regression

## ■ Review:

$$z = x\mathbf{w}^T + b$$

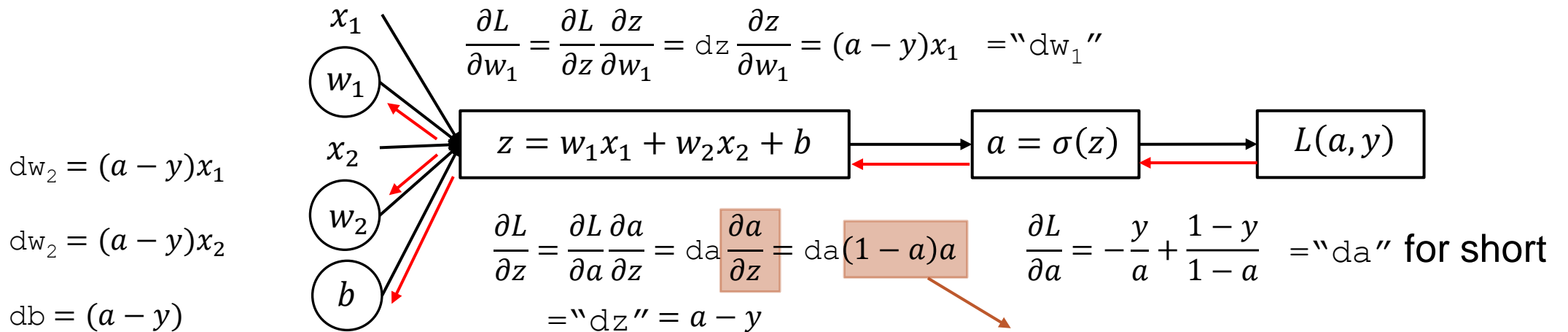
$$\hat{y} = \underline{a} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\underline{L(a, y)} = -(y \log(a) + (1 - y) \log(1 - a))$$

$a$  for “activation”

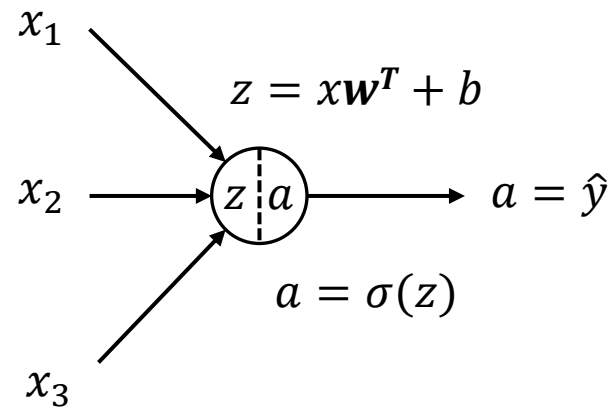
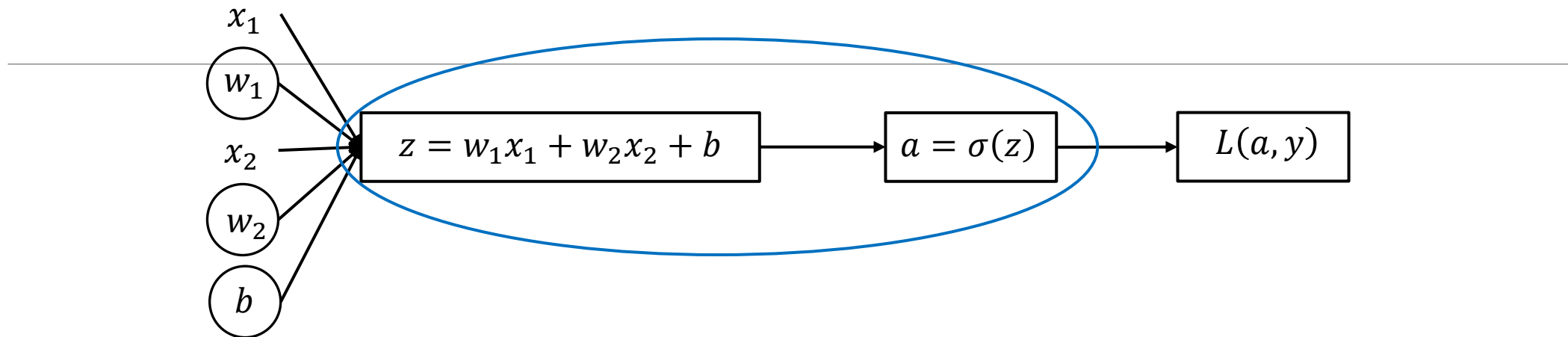
**loss** function  $L(a, y)$

= The **cost** of a single training example

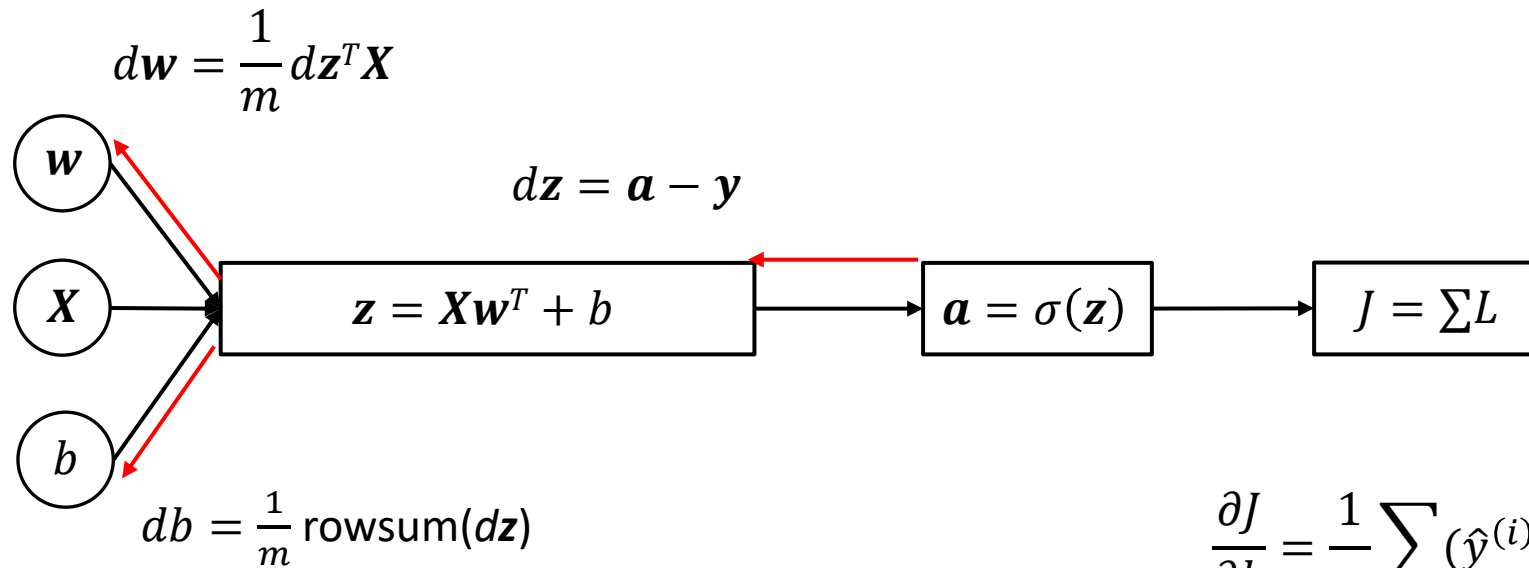


The derivative of sigmoid function (previous slide)

# Wrapped into a single unit



# Vectorizing logistic regression (recap)



$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)})$$

$$\frac{\partial J}{\partial w_k} = \frac{1}{m} \sum_i (\hat{y}^{(i)} - y^{(i)}) x_k^{(i)}, k=1,2,\dots,n$$

# Code Examples

---

# Computational graph in PyTorch

Computational graph is implemented by `torch.autograd` in PyTorch, and automatic differentiation engine

```
import torch
```

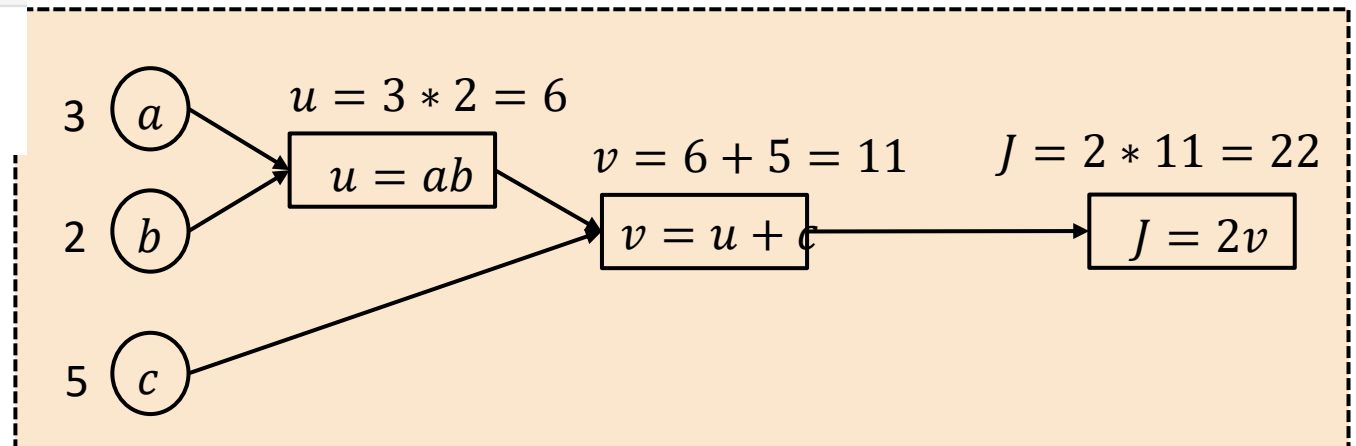
```
a = torch.tensor([3.], requires_grad=True)
b = torch.tensor([2.], requires_grad=True)
c = torch.tensor([5.], requires_grad=True)
print(a)
print(b)
print(c)
```

```
tensor([3.], requires_grad=True)
tensor([2.], requires_grad=True)
tensor([5.], requires_grad=True)
```

“requires\_grad=True” to indicate they are leaf tensor of the computational graph

```
# define the output function
J = 2 * (a * b + c)
print(J)
```

```
tensor([22.], grad_fn=<MulBackward0>)
```



# What is a Tensor?

---

- A tensor is a mathematical object that generalizes the concept of scalars, vectors, and matrices to higher-dimensional spaces. In the context of machine learning and deep learning, tensors are fundamental data structures used to represent and store multi-dimensional data.
- Key concepts: Dimension/Rank, Shape, Dtype
- Dimension/Rank
  - `[[[1,2,3],[4,5,6]],[[7,8,9],[10,11,12]]]` has three dimensions
- Here are some types of tensors:
  - **Scalar:** A scalar is a single numerical value. In tensor notation, a scalar is considered a tensor of rank 0.
  - **Vector:** A vector is an ordered collection of scalars. It is a one-dimensional tensor. For example, a vector with three elements could represent a point in 3D space.
  - **Matrix:** A matrix is a two-dimensional array of scalars. It is a two-dimensional tensor. Matrices are often used to represent linear transformations, often used to carry image data or tabular data
  - **Tensor:** A tensor is a multi-dimensional array of numerical values. It is a generalization of scalars, vectors, and matrices. Tensors can have any number of dimensions, and each dimension is called a mode or an axis. In deep learning, tensors with three or more dimensions are commonly used to represent complex data structures, such as images, sequences, and volumes. A tensor with rank 3 can be visualized as a cube or a stack of matrices
- Shape of a tensor is related to its dimensions. The number of entries in each dimension
  - `[[[1,2,3],[4,5,6]],[[7,8,9],[10,11,12]]]` has shape (2,2,3)
- Dtype refers to the data type of the elements, e.g. int32, float32, and boolean
- In libraries like PyTorch, tensors provide a convenient way to perform operations on multi-dimensional data, making them a foundational concept in the field of deep learning.

# Computational graph in PyTorch

Let the output J change by 1.0, and let the gradient “flow back”

Remember to clear the gradients before you run the next backward step

```
external_grad = torch.tensor([1.0])  
J.backward(gradient = external_grad)
```

```
# clear all gradients  
a.grad.data.zero_()  
b.grad.data.zero_()  
c.grad.data.zero_()
```

The gradients for all variables in the graph are automatically computed

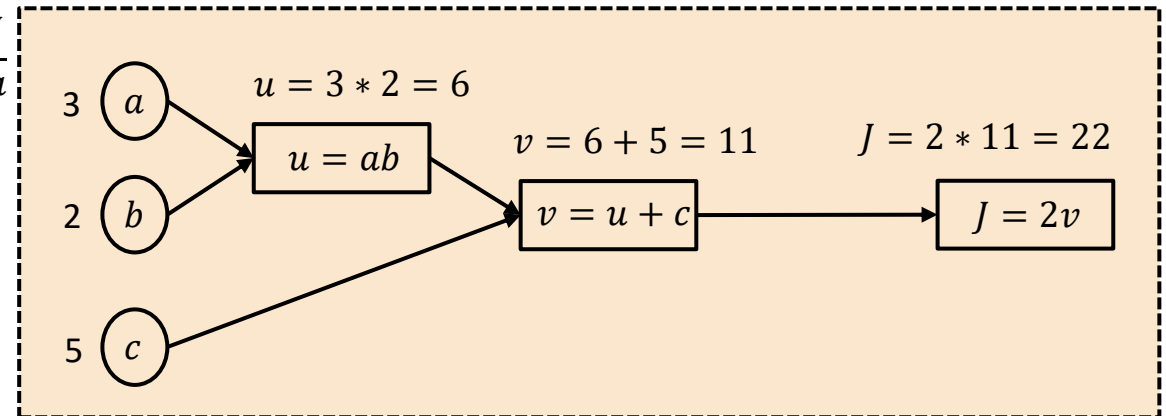
```
print(a.grad)  
print(b.grad)  
print(c.grad)
```

```
tensor([4.])  
tensor([6.])  
tensor([2.])
```

$$\frac{\partial J}{\partial v} \frac{\partial v}{\partial u} \frac{\partial u}{\partial a} = 2 * 1 * 2 = 4 = \frac{\partial J}{\partial a}$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial v} \frac{\partial v}{\partial u} \frac{\partial u}{\partial b} = 2 * 1 * 3 = 6$$

$$\frac{\partial J}{\partial c} = \frac{\partial J}{\partial v} \frac{\partial v}{\partial c} = 2 * 1 = 2$$





# Basic neural network in PyTorch

torch.nn.Linear (pass-through activation fn)

1. Use the torch.nn module

```
import torch
import torch.nn as nn
import numpy as np
torch.manual_seed(0)
```

2. Define a Linear layer

```
m = nn.Linear(4, 3)
# Equivalent to nn.Linear(in_features = 4, out_features = 3, bias = True)
```

3. Let's look at the initial weight and bias:

```
print(m.weight)
print(m.bias)
```

```
Parameter containing:
tensor([[ 0.3742, -0.0806,  0.0529,  0.4527],
        [-0.4638, -0.3148, -0.1266, -0.1949],
        [ 0.4320, -0.3241, -0.2302, -0.3493]], requires_grad=True)
Parameter containing:
tensor([-0.4683, -0.2919,  0.4298], requires_grad=True)
```

*Recall wt matrix:  $n_l \times n_{l-1}$*

*Bias vector:  $1 \times n_l$*

# Basic neural network in PyTorch

---

## 4. Set some dummy values

```
# Manually set weight and bias
m.weight.data = torch.tensor(np.arange(1,13).reshape((3,4))).float()
print(m.weight)
m.bias.data = torch.ones_like(m.bias.data).float()
print(m.bias)
```

```
Parameter containing:
tensor([[ 1.,  2.,  3.,  4.],
        [ 5.,  6.,  7.,  8.],
        [ 9., 10., 11., 12.]], requires_grad=True)
```

```
Parameter containing:
tensor([1., 1., 1.], requires_grad=True)
```

# Basic neural network in PyTorch

---

5. Prepare a single input data example

```
x = torch.randn(1, 4)
print(x.shape)
```

torch.Size([1, 4])

Get the output by *calling* the model

```
out = m(x)
print(out.shape)
```

torch.Size([1, 3])

x is  $1 \times 4$ :

By default, the first dimension of a tensor stands for batch size

- out is  $1 \times 3$ :
- Pytorch calculates  $XW^T+B$ . Check the dimensions:
- X is  $1 \times 4$ ,  $W^T$  is  $4 \times 3$ , B is  $1 \times 3$ , hence the output is  $1 \times 3$

# Basic neural network in PyTorch

---

## 6. Prepare a batch of data examples

```
x1 = torch.randn(100, 4)
out1 = m(x1)
print(out1.shape)

torch.Size([100, 3])
```

x is  $100 \times 4$ :

By default, the first dimension of a tensor stands for batch size

- out is  $100 \times 3$ :
- Pytorch calculates  $XW^T+B$ . Check the dimensions:
- X is  $100 \times 4$ ,  $W^T$  is  $4 \times 3$ ,  $B$  is  $1 \times 3$ , hence the output is  $100 \times 3$

# Basic neural network in PyTorch

## 7. Let's manually check what computation is done

Prepare an input vector: [1,1,1,1]

```
x2 = torch.ones(1,4)
print(x2)
out2 = m(x2)
print(out2)
```

```
tensor([[1., 1., 1., 1.]])
tensor([[11., 27., 43.]], grad_fn=<AddmmBackward>)
```

```
weight = m.weight.data.numpy()
bias = m.bias.data.numpy()
x = x2.numpy().squeeze()

print(np.dot(weight[0,:], x) + bias[0])
print(np.dot(weight[1,:], x) + bias[1])
print(np.dot(weight[2,:], x) + bias[2])

out2_manual = np.sum(weight, axis=1) + bias

print(np.equal(out2_manual, out2.data.numpy()))

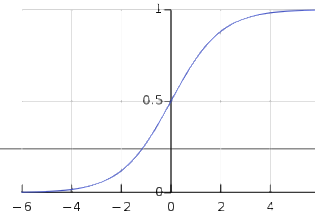
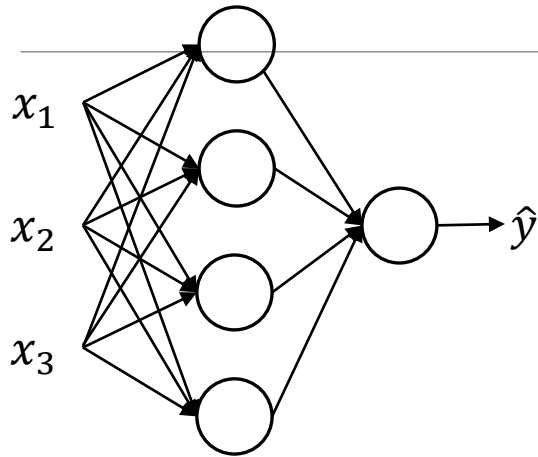
11.0
27.0
43.0
[[ True  True  True]]
```

$$\text{out2} = \mathbf{XW}^T + \mathbf{B} = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 5 & 9 \\ 2 & 6 & 10 \\ 3 & 7 & 11 \\ 4 & 8 & 12 \end{bmatrix} + \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 11 & 27 & 43 \end{bmatrix}$$

# Activation Functions

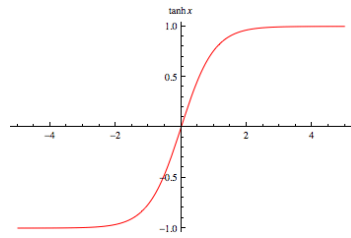
---

# Activation functions



Sigmoid

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



$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

“hyperbolic tangent”. Almost always works better for hidden layers than sigmoid.

Downside for both: small gradient for extreme z values

Single input example  $x$

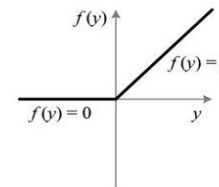
$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = \cancel{\sigma(z^{[1]})} \quad g^{[1]}(z^{[1]})$$

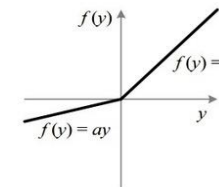
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = \cancel{\sigma(z^{[2]})} \quad g^{[2]}(z^{[2]})$$

ReLU



LeakyReLU



Rectified linear unit (ReLU)

$$g(z) = \max(0, z)$$

Leaky ReLU

$$g(z) = \max(\alpha z, z)$$

# Rule-of-thumb for choosing activation functions

---

- Binary classification,  $y \in \{0,1\}$ , use sigmoid for output layer.
- Regression, use a linear pass-through function for output layer
- For all the other units, use ReLU (increasingly the default choice) or tanh
  - ReLU learns much faster than tanh.



# Why do we need non-linear activation functions?

Why not linear activation?

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = z^{[1]}$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

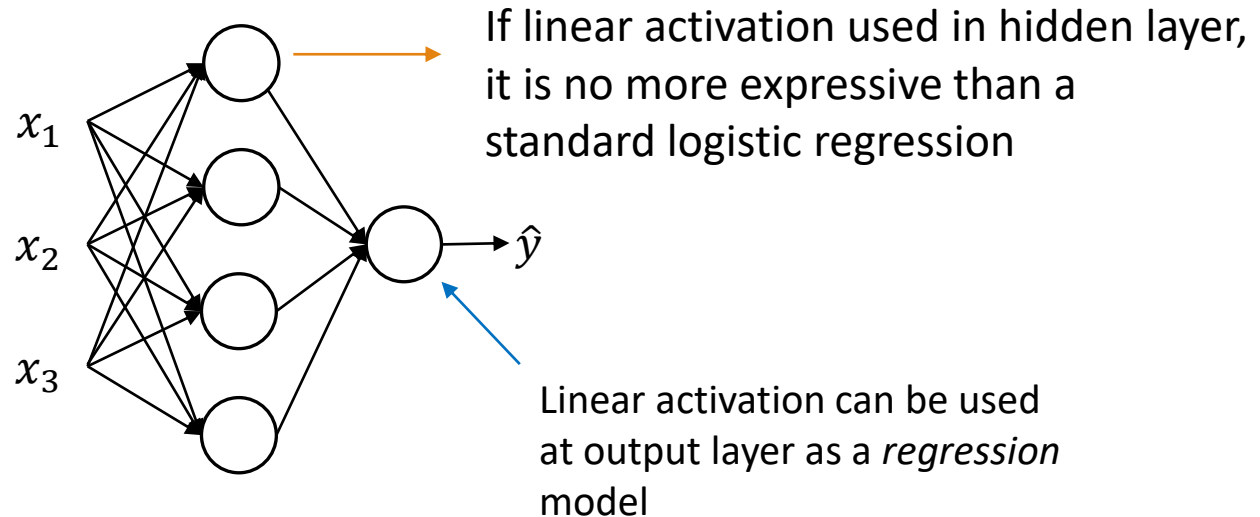
$$a^{[2]} = z^{[2]}$$



$$a^{[2]} = W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]}$$

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

$$= W'x + b'$$



# Derivatives of activation functions

---

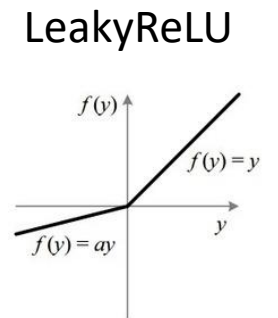
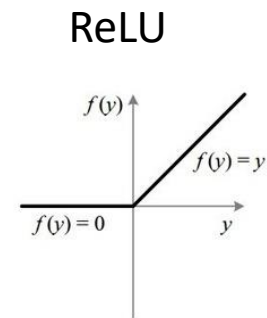
Sigmoid	$g(z) = \frac{1}{1 + e^{-z}} = a$	$g'(z) = g(z)(1 - g(z)) = a(1 - a)$
---------	-----------------------------------	-------------------------------------

Tanh	$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = a$	$g'(z) = 1 - (1 - g(z))^2 = 1 - a^2$
------	--	--------------------------------------

ReLU	$g(z) = \max(0, z)$	$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & z < 0 \end{cases}$
------	---------------------	--

LeakyReLU	$g(z) = \max(\alpha z, z)$	$g'(z) = \begin{cases} 1, & z > 0 \\ \alpha, & z < 0 \end{cases}$
-----------	----------------------------	---

$\alpha$  is a tiny positive value



# Activation functions in PyTorch

---

`torch.nn.Sigmoid`, `torch.nn.LogSigmoid`

`torch.nn.Tanh`, `torch.nn.HardTanh`

`torch.nn.ReLU`, `torch.nn.LeakyReLU`, `torch.nn.ReLU6`, `torch.nn.ReLU`, ...

# Derivatives of Vectors and Matrices/Backpropagation

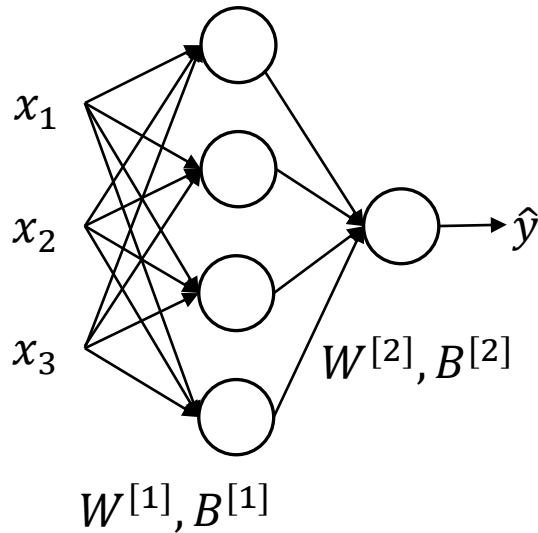
---

# Gradient descent for neural networks

Input  $X$  dimension:  $n^{[0]}$

# of units in layer 1:  $n^{[1]}$

# of units in layer 2:  $n^{[2]}$



$W^{[1]}$  dimension:  $n^{[1]} \times n^{[0]}$       $B^{[1]}$  dimension:  $1 \times n^{[1]}$

$W^{[2]}$  dimension:  $n^{[2]} \times n^{[1]}$       $B^{[2]}$  dimension:  $1 \times n^{[2]}$

Binary classification case

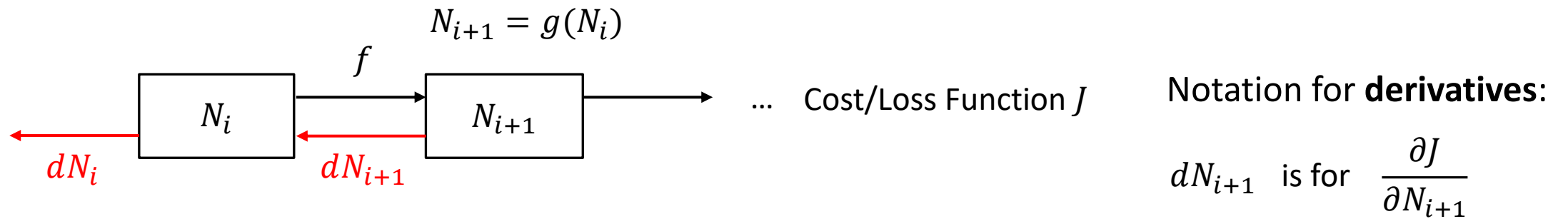
$$J(W^{[1]}, B^{[1]}, W^{[2]}, B^{[2]}) = \frac{1}{m} \sum_i L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{m} \sum_i L(\mathbf{a}^{(i)[2]}, y^{(i)})$$

$$dW^{[1]} = \frac{\partial J}{\partial W^{[1]}} \quad dB^{[1]} = \frac{\partial J}{\partial B^{[1]}} \quad dW^{[2]} = \frac{\partial J}{\partial W^{[2]}} \quad dB^{[2]} = \frac{\partial J}{\partial B^{[2]}}$$

How to compute the **gradients** (derivatives)?

# Recall – Chain Rule Between Nodes

---



$$dN_i = dN_{i+1} \cdot f'(N_i)$$

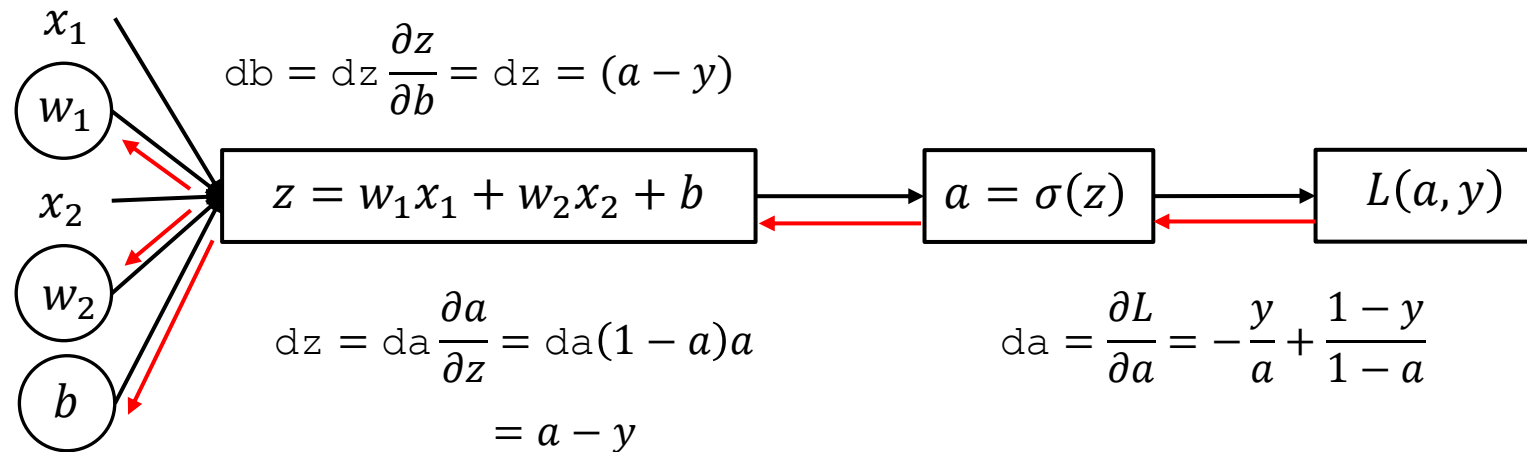
dot-product for matrices/vectors

The derivative of one node can be inferred from that of the adjacent node, and the **function** that connects them

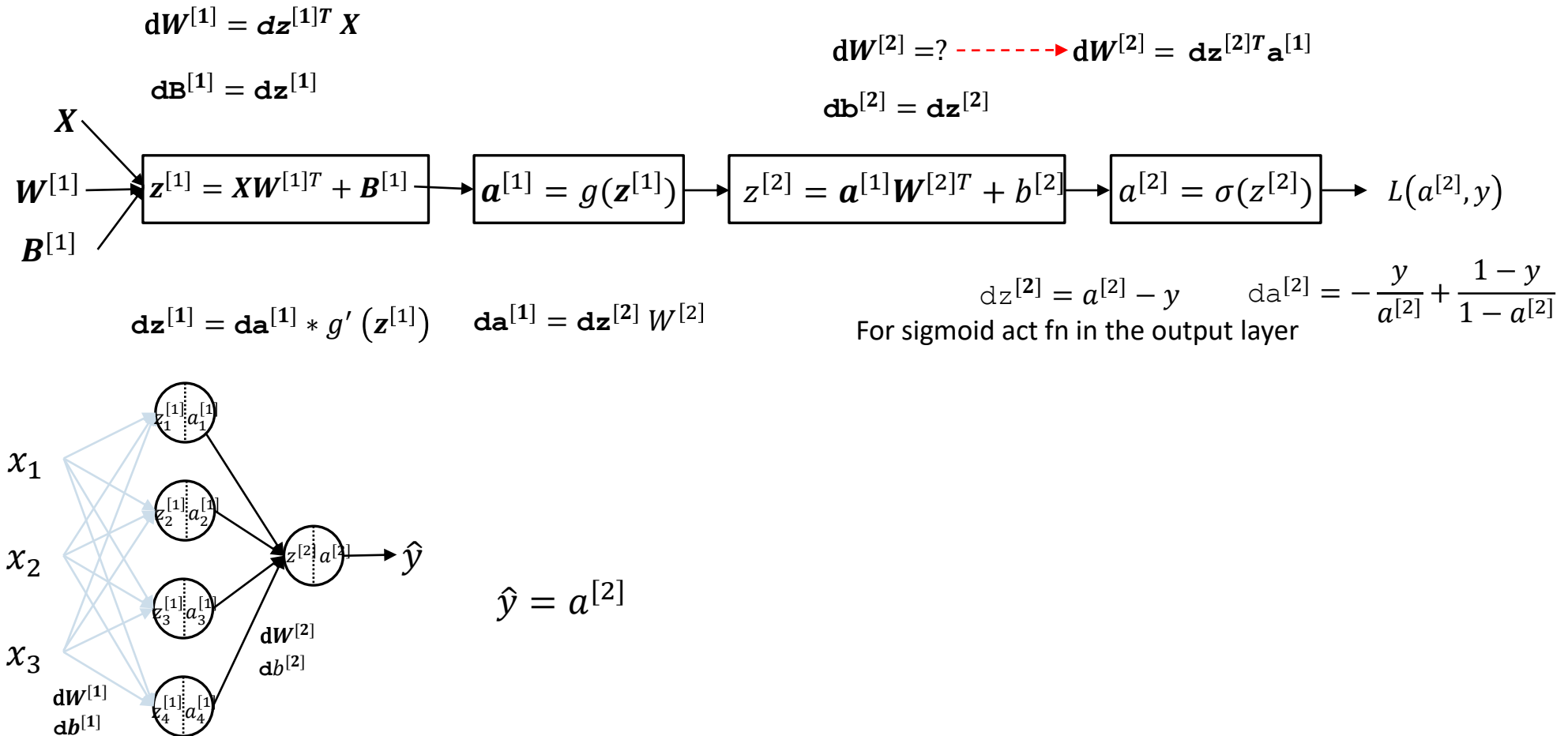
# Backpropagation

Recap of logistic regression: Loss Function:  $L(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$ ;  $a = \hat{y}$

$$\left. \begin{aligned} dw_1 &= dz \frac{\partial z}{\partial w_1} = dz x_1 = (a - y)x_1 \\ dw_2 &= dz \frac{\partial z}{\partial w_2} = dz x_2 = (a - y)x_2 \end{aligned} \right\} \mathbf{dw} = (dw_1, dw_2) = dz \cdot (x_1, x_2) = dz \cdot \mathbf{x}$$



## Backpropagation – 1 sample data point (two-layer neural network)





---

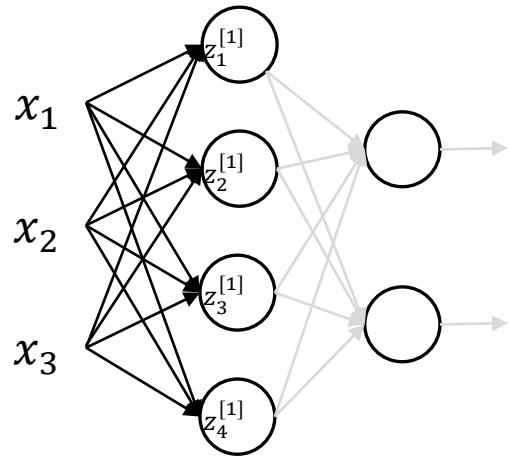
Generalize Back Propagation for Deeper  
ANNs Using Vectors and Matrices ==➡

# Derivatives of the Loss function w.r.t. $W$ and $B$

$$\boxed{\mathbf{z}^{(i)[1]} = \mathbf{X}^{(i)} W^{[1]T} + B^{[1]}} \quad \text{In order to compute } d\mathbf{W}^{[1]}, \text{ need to know } \frac{\partial \mathbf{z}^{(i)[1]}}{\partial W^{[1]}} \quad \text{How?}$$

➤ The partial derivatives of **each** component of  $\mathbf{z}^{(i)[1]}$  w.r.t. **each** element of  $W^{[1]}$

$$\mathbf{z}^{(i)[1]} = [z_1^{(i)[1]} \quad z_2^{(i)[1]} \quad z_3^{(i)[1]} \quad z_4^{(i)[1]}] = [x_1^{(i)} \quad x_2^{(i)} \quad x_3^{(i)}] \begin{bmatrix} w_{11}^{[1]} & w_{21}^{[1]} & w_{31}^{[1]} & w_{41}^{[1]} \\ w_{12}^{[1]} & w_{22}^{[1]} & w_{32}^{[1]} & w_{42}^{[1]} \\ w_{13}^{[1]} & w_{23}^{[1]} & w_{33}^{[1]} & w_{43}^{[1]} \end{bmatrix} + [b_1^{[1]} \quad b_2^{[1]} \quad b_3^{[1]} \quad b_4^{[1]}]$$



$$\frac{\partial \mathbf{z}^{(i)[1]}}{\partial W^{[1]}} = \begin{bmatrix} \frac{\partial z_1^{(i)[1]}}{\partial w_{11}^{[1]}} & \frac{\partial z_1^{(i)[1]}}{\partial w_{12}^{[1]}} & \frac{\partial z_1^{(i)[1]}}{\partial w_{13}^{[1]}} \\ \frac{\partial z_2^{(i)[1]}}{\partial w_{21}^{[1]}} & \frac{\partial z_2^{(i)[1]}}{\partial w_{22}^{[1]}} & \frac{\partial z_2^{(i)[1]}}{\partial w_{23}^{[1]}} \\ \frac{\partial z_3^{(i)[1]}}{\partial w_{31}^{[1]}} & \frac{\partial z_3^{(i)[1]}}{\partial w_{32}^{[1]}} & \frac{\partial z_3^{(i)[1]}}{\partial w_{33}^{[1]}} \\ \frac{\partial z_4^{(i)[1]}}{\partial w_{41}^{[1]}} & \frac{\partial z_4^{(i)[1]}}{\partial w_{42}^{[1]}} & \frac{\partial z_4^{(i)[1]}}{\partial w_{43}^{[1]}} \end{bmatrix} = \begin{bmatrix} x_1^{(i)} & 0 & 0 \\ 0 & x_2^{(i)} & 0 \\ 0 & 0 & x_3^{(i)} \\ 0 & 0 & 0 \end{bmatrix}$$

Collapse into 1x3

➤ Derivative chain rule:

$$\text{➤ } d\mathbf{W}^{(i)[1]} = d\mathbf{z}^{(i)[1]T} \frac{\partial \mathbf{z}^{(i)[1]}}{\partial W^{[1]}} = d\mathbf{z}^{(i)[1]T} \mathbf{x}^{(i)}$$

$$\text{➤ } 4 \times 3 \qquad \qquad \qquad 4 \times 1 \quad 1 \times 3$$

➤ In general, at each layer  $l$ ,

$$\text{➤ } d\mathbf{W}^{(i)[l]} = d\mathbf{z}^{(i)[l]T} \frac{\partial \mathbf{z}^{(i)[l]}}{\partial W^{[l]}} = d\mathbf{z}^{(i)[l]T} \mathbf{a}^{(i)[l-1]}$$

$$\text{➤ } n^{[l]} \times n^{[l-1]} \qquad \qquad \qquad n^{[l]} \times 1 \quad 1 \times n^{[l-1]}$$

➤ How about  $d\mathbf{B}^{(i)[l]}$ ?

$$\text{➤ } d\mathbf{B}^{(i)[l]} = d\mathbf{z}^{(i)[l]} \frac{\partial \mathbf{z}^{(i)[l]}}{\partial B^{[l]}}$$

$$\text{➤ } d\mathbf{B}^{(i)[l]} = d\mathbf{z}^{(i)[l]}$$

$$\text{➤ } 1 \times n^{[l]}$$

# Derivatives of the Loss Function w.r.t. $\mathbf{a}$ and $\mathbf{z}$

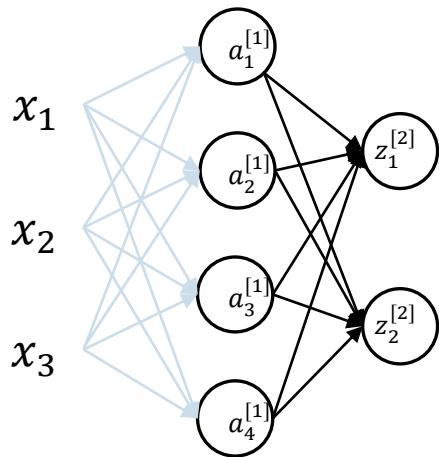
Hidden to output

$$\mathbf{z}^{(i)[2]} = \mathbf{a}^{(i)[1]} \mathbf{W}^{[2]T} + \mathbf{B}^{[2]}$$

Want to know  $d\mathbf{a}^{[1]}$

Need to compute  $\frac{\partial \mathbf{z}^{(i)[2]}}{\partial \mathbf{a}^{[1]}}$

$$\mathbf{z}^{(i)[2]} = [z_1^{(i)[2]} \quad z_2^{(i)[2]}] = \begin{bmatrix} a_1^{(i)[1]} & a_2^{(i)[1]} & a_3^{(i)[1]} & a_4^{(i)[1]} \end{bmatrix} \begin{bmatrix} w_{11}^{[2]} & w_{21}^{[2]} \\ w_{12}^{[2]} & w_{22}^{[2]} \\ w_{13}^{[2]} & w_{23}^{[2]} \\ w_{14}^{[2]} & w_{24}^{[2]} \end{bmatrix} + \begin{bmatrix} b_1^{[2]} & b_2^{[2]} \end{bmatrix}$$

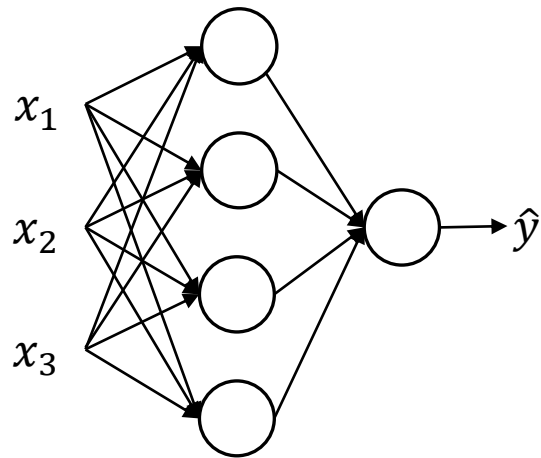


$$\frac{\partial \mathbf{z}^{(i)[2]}}{\partial \mathbf{a}^{(i)[1]}} = \begin{bmatrix} \frac{\partial z_1^{(i)[2]}}{\partial a_1^{(i)[1]}} & \frac{\partial z_1^{(i)[2]}}{\partial a_2^{(i)[1]}} & \frac{\partial z_1^{(i)[2]}}{\partial a_3^{(i)[1]}} & \frac{\partial z_1^{(i)[2]}}{\partial a_4^{(i)[1]}} \\ \frac{\partial z_2^{(i)[2]}}{\partial a_1^{(i)[1]}} & \frac{\partial z_2^{(i)[2]}}{\partial a_2^{(i)[1]}} & \frac{\partial z_2^{(i)[2]}}{\partial a_3^{(i)[1]}} & \frac{\partial z_2^{(i)[2]}}{\partial a_4^{(i)[1]}} \end{bmatrix} = \begin{bmatrix} w_{11}^{[2]} & w_{12}^{[2]} & w_{13}^{[2]} & w_{14}^{[2]} \\ w_{21}^{[2]} & w_{22}^{[2]} & w_{23}^{[2]} & w_{24}^{[2]} \end{bmatrix} = \mathbf{W}^{[2]}$$

- $d\mathbf{a}^{(i)[1]} = d\mathbf{z}^{(i)[2]} \mathbf{W}^{[2]}$
- In general:  $d\mathbf{a}^{(i)[l-1]} = d\mathbf{z}^{(i)[l]} \mathbf{W}^{[l]}$
- How about  $d\mathbf{z}^{(i)[1]}$ , or in general,  $d\mathbf{z}^{(i)[l-1]}$ ?
- $\mathbf{a}^{(i)[l-1]} = g(\mathbf{z}^{(i)[l-1]})$
- $d\mathbf{z}^{(i)[l-1]} = d\mathbf{a}^{(i)[l-1]} * g'(\mathbf{z}^{(i)[l-1]})$
- $d\mathbf{z}^{(i)[l-1]} = d\mathbf{z}^{(i)[l]} \mathbf{W}^{[l]} * g'(\mathbf{z}^{(i)[l-1]})$
- For a sigmoid activation function in the output layer:
  - $d\mathbf{z}^{(i)[L]} = \mathbf{a}^{(i)[L]} - \mathbf{y}^{(i)}$

# Vectorizing backprop across $m$ training examples

$$d\mathbf{W}^{[l]} = \frac{\partial}{\partial \mathbf{W}^{[l]}} J = \frac{\partial}{\partial \mathbf{W}^{[l]}} \left( \frac{1}{m} \sum_i L^{(i)} \right) = \frac{1}{m} \sum_{i=1}^m d\mathbf{W}^{(i)[l]} \quad d\mathbf{B}^{[l]} = \frac{1}{m} \sum_{i=1}^m d\mathbf{B}^{(i)[l]} \quad d\mathbf{z}^{[l]} = \frac{1}{m} \sum_{i=1}^m d\mathbf{z}^{(i)[l]} \quad d\mathbf{a}^{[l]} = \frac{1}{m} \sum_{i=1}^m d\mathbf{a}^{(i)[l]}$$



Single  
example

$$d\mathbf{z}^{(i)[l-1]} = d\mathbf{z}^{(i)[l]} \mathbf{W}^{[l]} * g'(\mathbf{z}^{(i)[l-1]})$$

$$d\mathbf{a}^{(i)[l-1]} = d\mathbf{z}^{(i)[l]} \mathbf{W}^{[l]}$$

$$d\mathbf{W}^{(i)[l]} = d\mathbf{z}^{(i)[l]T} \mathbf{a}^{(i)[l-1]}$$

$$d\mathbf{B}^{(i)[l]} = d\mathbf{z}^{(i)[l]}$$

$$d\mathbf{z}^{(i)[L]} = \mathbf{a}^{(i)[L]} - \mathbf{y}^{(i)} \text{ for sigmoid activation function at the output layer}$$

$m$  examples

$$d\mathbf{z}^{[l-1]} = d\mathbf{z}^{[l]} \mathbf{W}^{[l]} * g'(\mathbf{z}^{[l-1]})$$

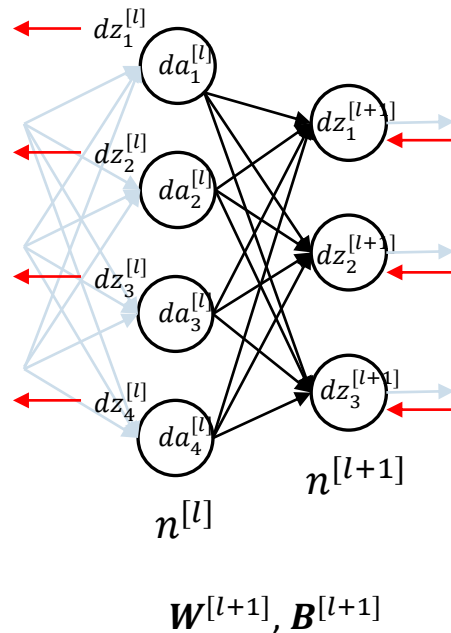
$$d\mathbf{a}^{[l-1]} = d\mathbf{z}^{[l]} \mathbf{W}^{[l]}$$

$$d\mathbf{W}^{[l]} = \frac{1}{m} d\mathbf{z}^{[l]T} \mathbf{a}^{[l-1]}$$

$$d\mathbf{B}^{[l]} = \frac{1}{m} \text{rowsum}(d\mathbf{z}^{[l]})$$

$$d\mathbf{z}^{[L]} = \frac{1}{m} \text{rowsum}(\mathbf{a}^{[L]} - \mathbf{y}) \text{ for sigmoid activation function at the output layer}$$

# Back Propagation Process



$dz^{[L]} = \frac{1}{m} \text{rowsum}(\mathbf{a}^{[L]} - \mathbf{y})$  for sigmoid activation function at the output layer

$$\mathbf{z}^{[l+1]} = \mathbf{a}^{[l]} \mathbf{W}^{[l+1]T} + \mathbf{B}^{[l+1]}$$

$$\mathbf{a}^{[l]} = g'(\mathbf{z}^{[l]})$$

$$\mathbf{z}^{[l+1]}: 1 \times n^{[l+1]}, \mathbf{a}^{[l]}: 1 \times n^{[l]}, \mathbf{W}^{[l+1]}: n^{[l+1]} \times n^{[l]}, \mathbf{B}^{[l+1]}: 1 \times n^{[l+1]}$$

$$d\mathbf{W}^{[l+1]} = d\mathbf{z}^{[l+1]T} \mathbf{a}^{[l]}$$

$$n^{[l+1]} \times n^{[l]} \quad n^{[l+1]} \times 1 \quad 1 \times n^{[l]}$$

$$d\mathbf{B}^{[l+1]} = d\mathbf{z}^{[l+1]}$$

$$1 \times n^{[l+1]} \quad 1 \times n^{[l+1]}$$

$$d\mathbf{a}^{[l]} = d\mathbf{z}^{[l+1]} \mathbf{W}^{[l+1]}$$

$$\triangleright d\mathbf{z}^{[l]} = d\mathbf{a}^{[l]} * g'(\mathbf{z}^{[l]}) = d\mathbf{z}^{[l+1]} \mathbf{W}^{[l+1]} * g'(\mathbf{z}^{[l]})$$

Usually combined

$$1 \times n^{[l]}$$

$$1 \times n^{[l]}$$

$$1 \times n^{[l+1]} \quad n^{[l+1]} \times n^{[l]}$$

$d\mathbf{z}^{[l]}$  is then passed to layer  $l - 1$  to compute  $d\mathbf{W}^{[l]}$  and  $d\mathbf{B}^{[l]}$

# Backpropagation Code Examples

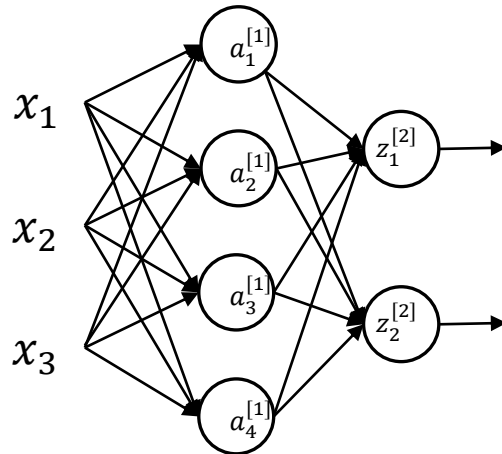
---

# Backprop in PyTorch

## Necessary imports

```
import torch
import torch.nn as nn
import numpy as np
torch.manual_seed(0)
```

<torch.\_C.Generator at 0x11b7ddb0>



## Define linear and activation layers

```
linear1 = nn.Linear(3, 4)
act1 = nn.ReLU()
linear2 = nn.Linear(4, 2)
act2 = nn.Sigmoid()
```

## Initialize to values easier to read:

```
linear1.weight.data = torch.tensor(np.arange(1,13).reshape((4,3))).float()
print(linear1.weight)
linear1.bias.data = torch.zeros_like(linear1.bias.data).float()
print(linear1.bias)
```

```
Parameter containing:
tensor([[ 1.,  2.,  3.],
        [ 4.,  5.,  6.],
        [ 7.,  8.,  9.],
        [10., 11., 12.]], requires_grad=True)
Parameter containing:
tensor([0., 0., 0., 0.], requires_grad=True)
```

# Backprop in PyTorch (layer 1)

Prepare some toy data, and watch the output from layer 1

```
x = torch.tensor([[1.,2.,3.]], requires_grad=True)
print('x: ', x)
z1 = linear1(x)
a1 = act1(z1)
```

```
print(z1)
print(a1)
print(a1.shape)
```

```
x:  tensor([[1., 2., 3.]], requires_grad=True)
tensor([[14., 32., 50., 68.]], grad_fn=<AddmmBackward>)
tensor([[14., 32., 50., 68.]], grad_fn=<ReluBackward0>)
torch.Size([1, 4])
```

Forward:

$$z1 = \begin{matrix} & \text{x} \\ & [1 & 2 & 3] \end{matrix} \begin{matrix} & \text{Weight.T} \\ \begin{bmatrix} 1 & 4 & 7 & 10 \\ 2 & 5 & 8 & 11 \\ 3 & 6 & 9 & 12 \end{bmatrix} \end{matrix} + \begin{matrix} & \text{Bias} \\ [0 & 0 & 0 & 0] \end{matrix}$$



# Backprop in PyTorch (layer 1)

Create some pseudo gradients that are in the same shape as a1

```
external_grad = torch.ones_like(a1) * 0.5
a1.backward(gradient=external_grad)
```

$$dz = da = \begin{bmatrix} .5 & .5 & .5 & .5 \end{bmatrix} \quad W = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix}$$

Gradients computed:

```
print(x.grad)
print(linear1.weight.grad)
print(linear1.bias.grad)

tensor([[11., 13., 15.]])
tensor([[0.5000, 1.0000, 1.5000],
        [0.5000, 1.0000, 1.5000],
        [0.5000, 1.0000, 1.5000],
        [0.5000, 1.0000, 1.5000]])
tensor([0.5000, 0.5000, 0.5000, 0.5000])
```

$$dx = dzW = \begin{bmatrix} .5 & .5 & .5 & .5 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix} = \begin{bmatrix} 11 & 13 & 15 \end{bmatrix}$$

$$dW = dz^T \cdot x = \begin{bmatrix} .5 \\ .5 \\ .5 \\ .5 \end{bmatrix} \begin{bmatrix} 1. & 2. & 3. \end{bmatrix} =$$

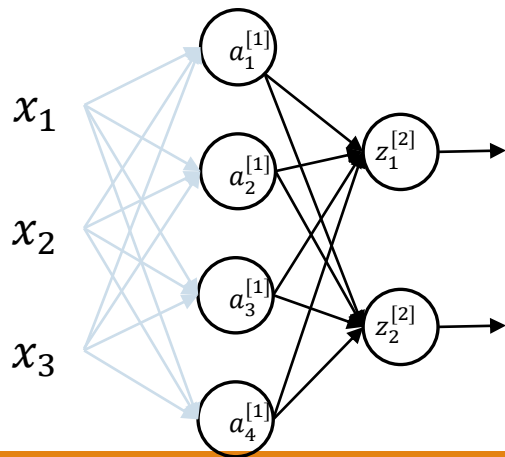
$$dB = dz = \begin{bmatrix} .5 & .5 & .5 & .5 \end{bmatrix}$$

# Backprop in PyTorch (layer 2)

Initialize the weight of linear2 to zeros (0),  
because sigmoid function saturates very fast

Forward pass:

```
x = torch.tensor([[1., 2., 3.]])
z1 = linear1(x)
a1 = act1(z1)
z2 = linear2(a1)
a2 = act2(z2)
```



```
linear2.weight.data = torch.zeros_like(linear2.weight.data).float()
linear2.bias.data = torch.zeros_like(linear2.bias.data).float()
print(linear2.weight)
print(linear2.bias)
```

Parameter containing:  
tensor([[0., 0., 0., 0.],  
[0., 0., 0., 0.]], requires\_grad=True)

Parameter containing:  
tensor([0., 0.], requires\_grad=True)

Output:

```
z2: tensor([[0., 0.]], grad_fn=<AddmmBackward>)
a2: tensor([[0.5000, 0.5000]], grad_fn=<SigmoidBackward>)
```

$$\sigma(0) = \frac{1}{1 + e^0} = 0.5$$

# Backprop in PyTorch (layer 2)

Backward pass:

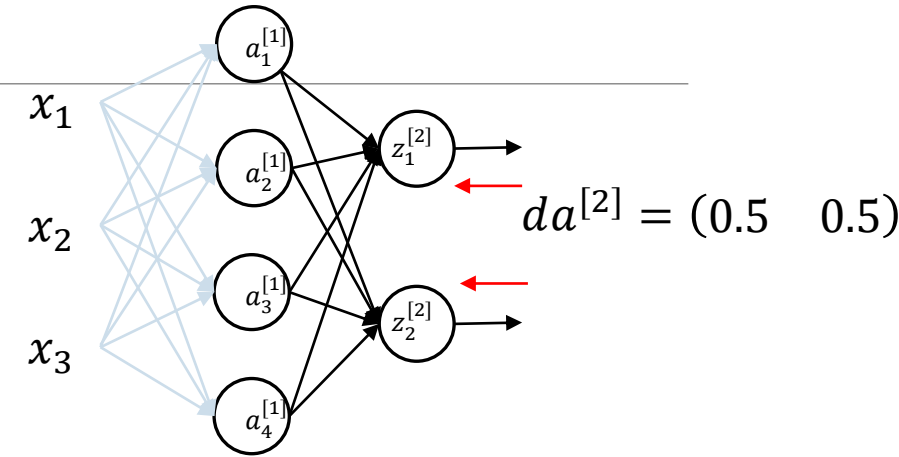
Tell PyTorch to save the gradients for intermediate variables

```
z2.retain_grad()
a1.retain_grad()
external_grad = torch.ones_like(a2) * 0.5
a2.backward(grad=external_grad)
```

Manually check  $dz_2$

```
# dz = a*(1 - a) * da
a2_np = a2.data.numpy()
dz2_manual = a2_np * (1 - a2_np) * external_grad.numpy()
dz2 = z2.grad.numpy()
print(np.equal(dz2_manual, dz2))
print('dz2: ', dz2)
```

Output:   
 `[[ True True]]`   
 `dz2: [[0.125 0.125]]`



$$dz^{[2]} = da^{[2]} * \frac{\partial a^{[2]}}{\partial z^{[2]}}$$

$$= da^2 * g'(z^{[2]})$$

$$= da^2 * \sigma'(z^{[2]})$$

Derivatives of sigmoid function:  
 $\sigma'(z^{[2]})$   
 $= \sigma(z^{[2]}) \cdot (1 - \sigma(z^{[2]}))$   
 $= a^{[2]} \cdot (1 - a^{[2]})$

$$(0.5 \ 0.5) * (0.5 \ 0.5)(1 - (.5 \ .5)))$$

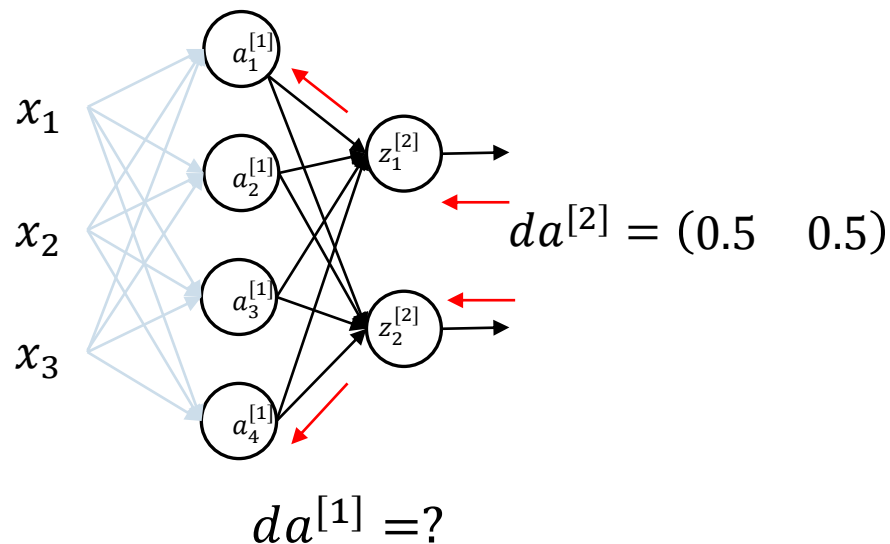
$$= (.25 \ .25) = (.125 \ .125)$$

# Backprop in PyTorch (layer 2) contd.

$$W^{[2]} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

We know:  $dz^{[2]} = (.125 \quad .125)$  so  $db^{[2]} = dz^{[2]} = (.125 \quad .125)$

$$da^{[1]} = dz^{[2]}W^{[2]} = (.125 \quad .125) \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} = [0. \quad 0. \quad 0. \quad 0.]$$



```
print('dW2: ', linear2.weight.grad)
print('db2: ', linear2.bias.grad)
print('dz2: ', z2.grad)
print('da1: ', a1.grad)
```

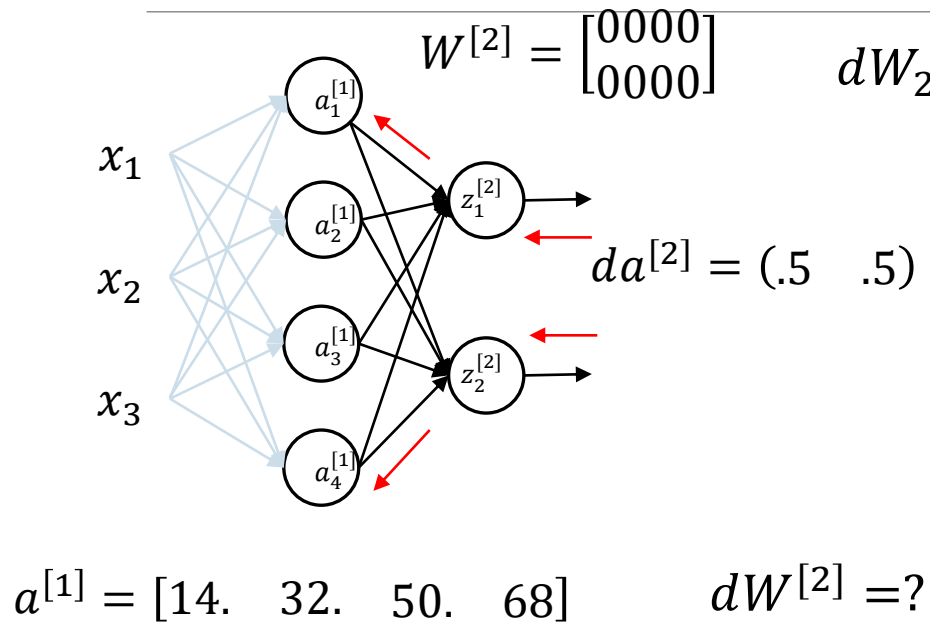
```
dW2: tensor([[1.7500, 4.0000, 6.2500, 8.5000],
             [1.7500, 4.0000, 6.2500, 8.5000]])
```

```
db2: tensor([0.1250, 0.1250])
```

```
dz2: tensor([[0.1250, 0.1250]])
```

```
da1: tensor([[0., 0., 0., 0.]])
```

# Backprop in PyTorch (layer 2) contd.



$$dW_2 = dz^{[2]T} a^{[1]}$$

$$= \begin{pmatrix} 0.125 \\ 0.125 \end{pmatrix} (14 \ 32 \ 50 \ 68)$$

$$= \begin{pmatrix} 1.75 & 4. & 6.25 & 8.5 \\ 1.75 & 4. & 6.25 & 8.5 \end{pmatrix}$$

Manually check  $dW_2$

```
a1_np = a1.data.numpy()
print('a1: ', a1_np)
dW2 = np.dot(dz2.T, a1_np)
print('dW2: ', dW2)
```

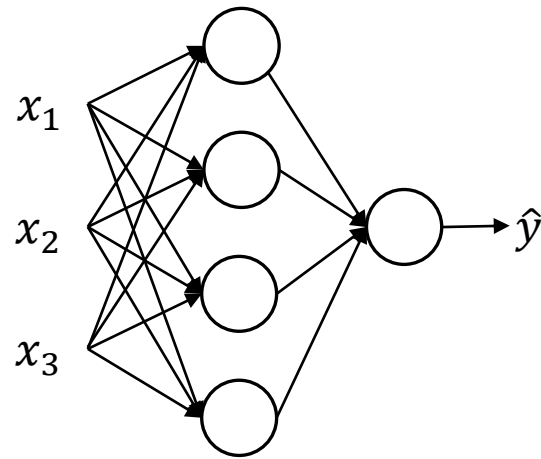
Output:

```
dW2:  [[1.75  4.   6.25  8.5 ]
       [1.75  4.   6.25  8.5 ]]
```

# Initializing $W$ and $B$ ?

---

# How to initialize $W$ parameters? To zeros?



$$W^{[1]} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\Rightarrow a_1^{[1]} = a_2^{[1]} = a_3^{[1]} = a_4^{[1]}$$

$$dW^{[2]} = dz^{[2]T} a^{[1]}$$

All same

$$W^{[2]} = [0, 0, 0, 0]$$

All same



$$da^{[1]} = W^{[2]T} dz^{[2]}$$

$$dz^{[1]} \text{ All same}$$



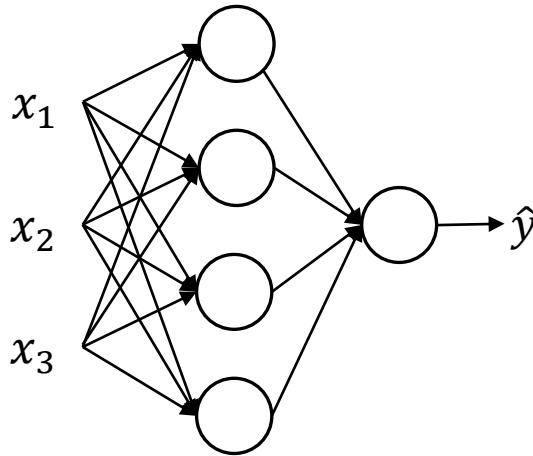
$$dW^{[1]} = dz^{[1]T} x$$

Update to all rows are same

All hidden units computing the same function!  
Pointless!

# Use Random initialization for $W$ Instead

---



$$W^{[1]} = \text{np.random.randn}(4, 3) * \underline{0.01}$$

$$B^{[1]} = \text{np.zeros}((4, 1))$$

$$W^{[2]} = \text{np.random.randn}(1, 4) * \underline{0.01}$$

$$B^{[2]} = \text{np.zeros}((1, 1))$$