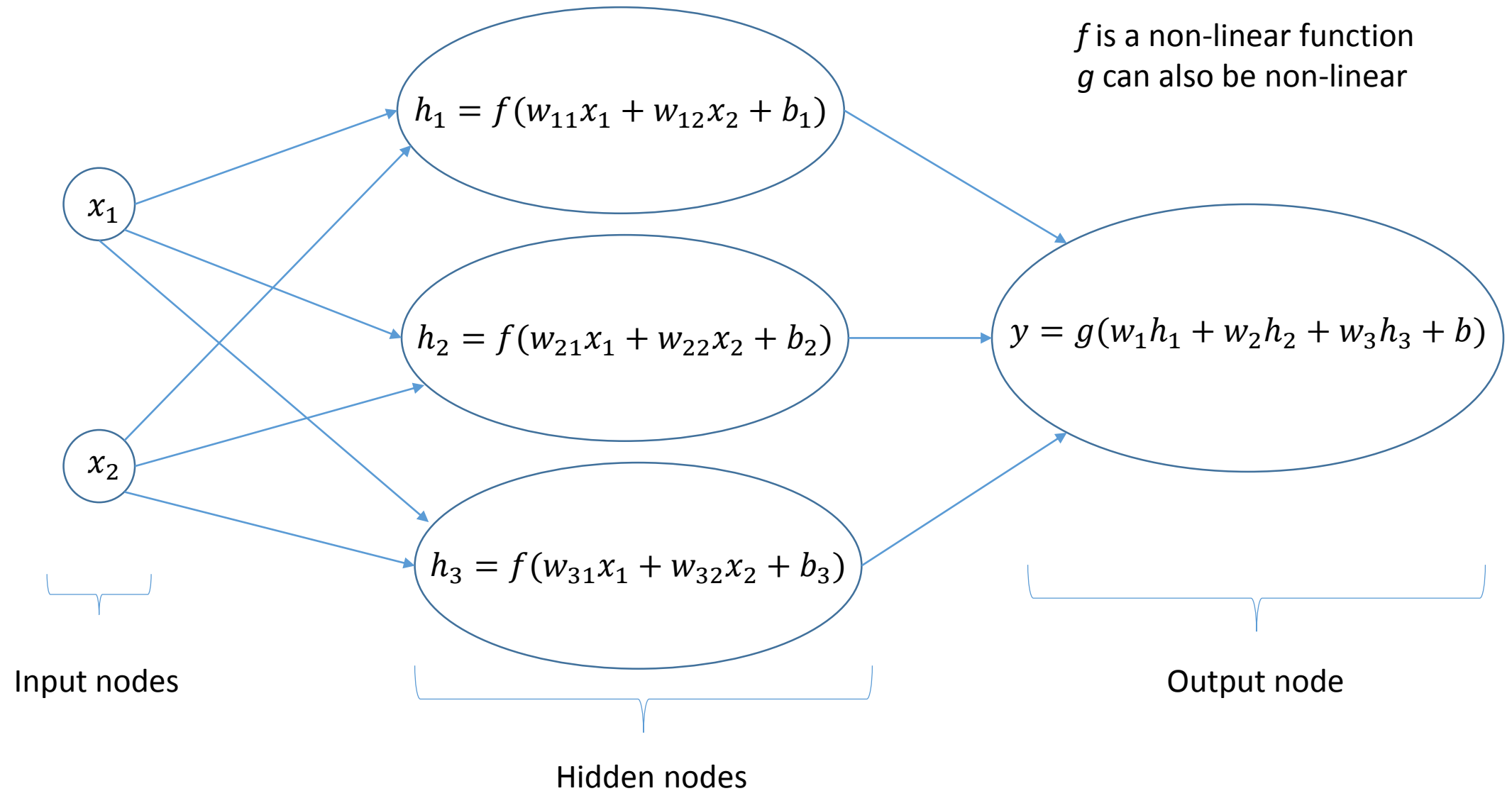# Introduction to Neural Networks

Computing Science

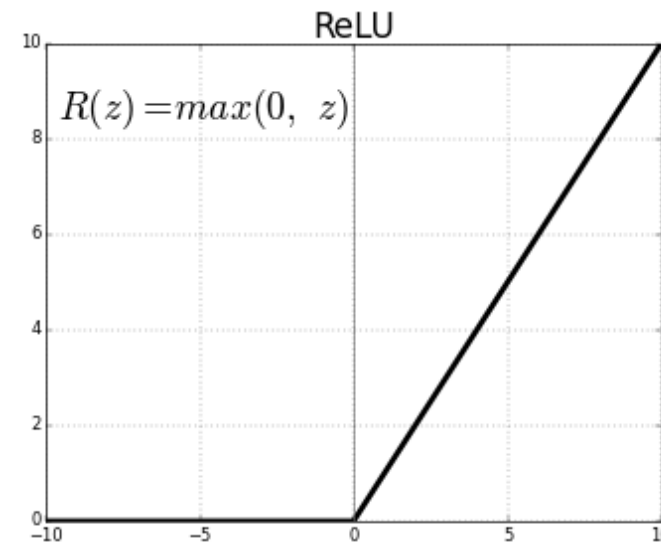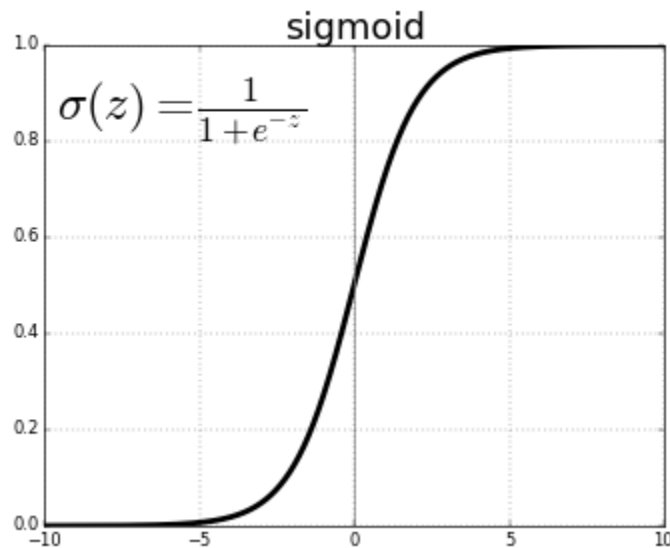University of Alberta

Nilanjan Ray

# Agenda

- What is a Neural Net?
  - Neural net as a computational graph
- Approximating "XOR" function with neural net
- Understanding backpropagation
- Universal function approximation by a neural net

# Feed forward neural network



$f$ is a non-linear function
$g$ can also be non-linear

$h_1 = f(w_{11}x_1 + w_{12}x_2 + b_1)$

$h_2 = f(w_{21}x_1 + w_{22}x_2 + b_2)$

$h_3 = f(w_{31}x_1 + w_{32}x_2 + b_3)$

$y = g(w_1h_1 + w_2h_2 + w_3h_3 + b)$
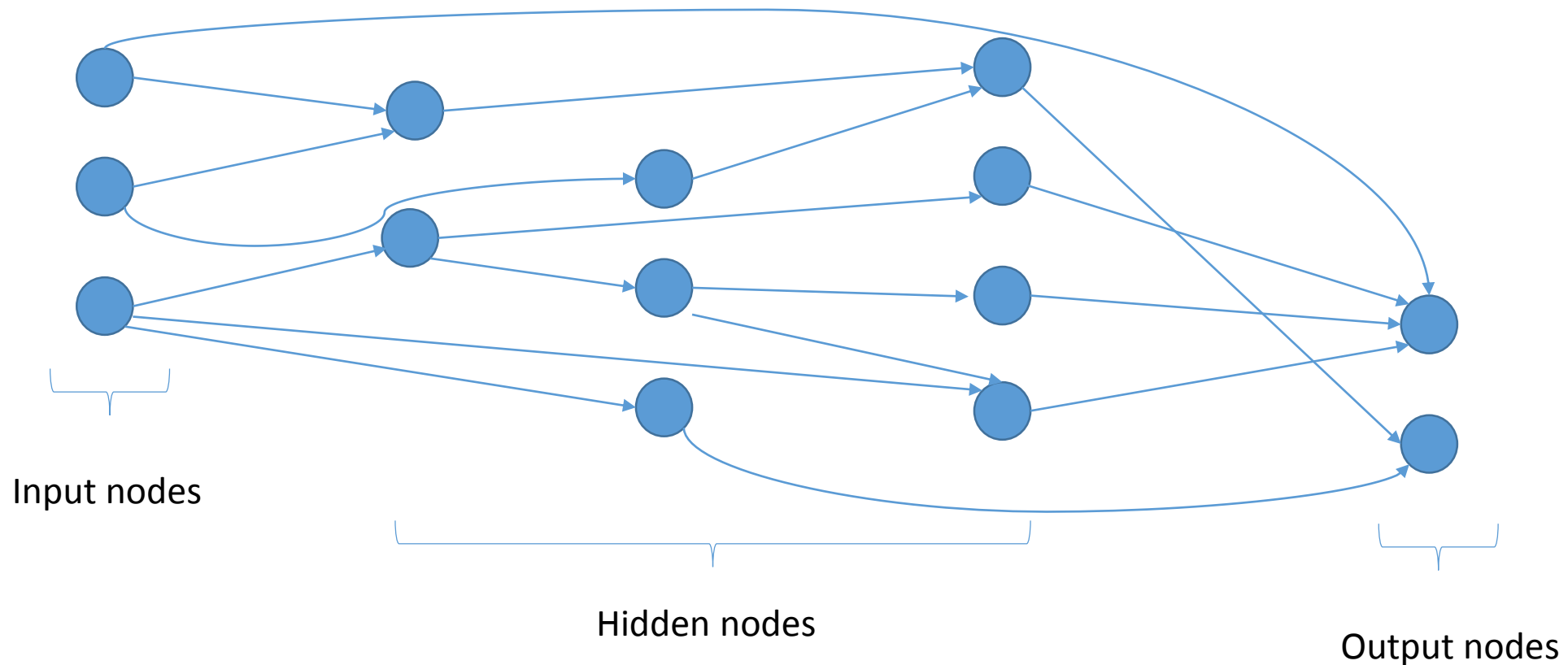
$x_1$

$x_2$

Input nodes

Hidden nodes

Output node

# Feed forward net: non-linear functions

- Non-linear functions at hidden nodes are known as "activation function"
  - Sigmoid, ReLU, ELU, ….

## sigmoid

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

## ReLU

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

Why activation functions are non-linear?

# Feedforward net in general: Directed acyclic graph



Input nodes

Hidden nodes

Output nodes

# What's the big deal about neural net?

- Mathematically very rich: it can approximate any function

- It is biologically inspired: (loosely) resembles brain connections

- Computationally:
  - Simple: matrix-vector multiplication and point-wise non-linear function
  - Highly paralleizable: cuBLAS, GEMM, Batched GEMM!

- Excellent <span style="color:red">empirical</span> results on "generalization capability" over variety of applications!

# Neural network as a computational graph

Input ($4 \times 2$ matrix)

X

($4 \times 8$ matrix) $Z_1 = XW_1$  
($4 \times 8$ matrix) $Z_2 = Z_1 + b_1$  
($4 \times 8$ matrix) $Z_3 = \sigma(Z_2)$  
($4 \times 1$ vector) $Z_4 = Z_3 W_2$  
($4 \times 1$ vector) $Z_5 = Z_4 + b_2$

*   +   $\sigma$   *   +   $\sigma$

Output ($4 \times 1$ vector)

$Y^p = \sigma(Z_5)$

Loss

$W_1$  
Parameter ($2 \times 8$ matrix)

$b_1$  
Parameter ($1 \times 8$ vector)

$W_2$  
Parameter ($8 \times 1$ vector)

$b_2$  
Parameter (scalar)

Y

Ideal output ($4 \times 1$ vector)

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

Sigmoid function; applied pointwise to a vector or matrix input

This network is trying to learn XOR function

| X | | Y |
|---|---|---|
| 1 | 0 | 1 |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 1 | 0 |

# How does PyTorch optimize parameters?

- By gradient descent PyTorch adjusts network parameters to reduce the value of the loss function.

- But how?
  - Answer: Backpropagation

- Let us learn how to do backpropagation on a computational graph!

# Chain rule of derivative for a computational node

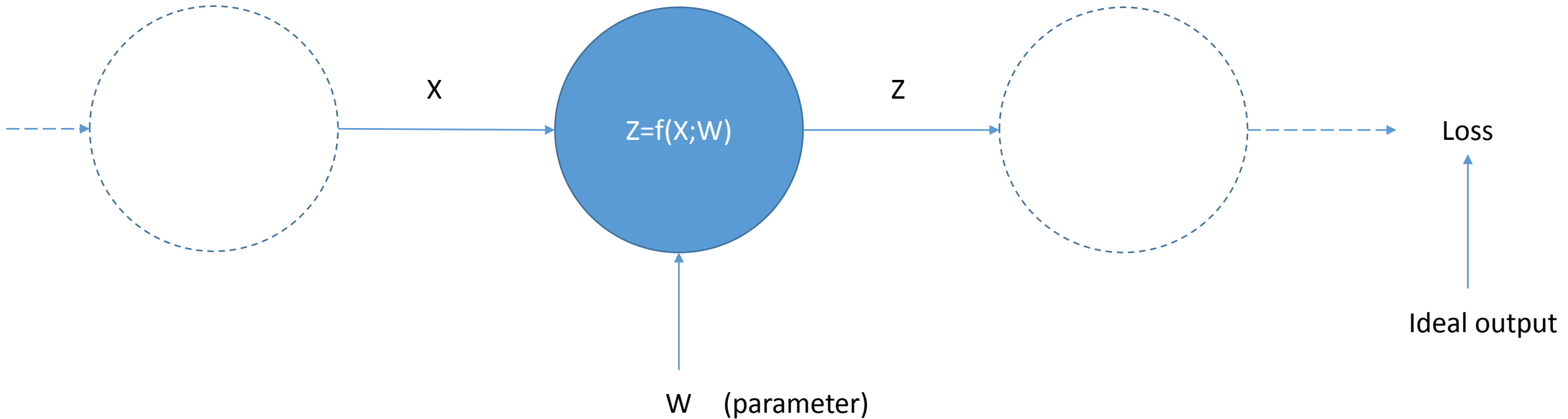X      Z=f(X;W)      Z      Loss

Ideal output

W    (parameter)

If X, Z, W are all scalars, then usual chain rule of derivative applies:

$$\frac{\partial(\text{Loss})}{\partial X} = \frac{\partial Z}{\partial X}\frac{\partial(\text{Loss})}{\partial Z}$$

$$\frac{\partial(\text{Loss})}{\partial W} = \frac{\partial Z}{\partial W}\frac{\partial(\text{Loss})}{\partial Z}$$
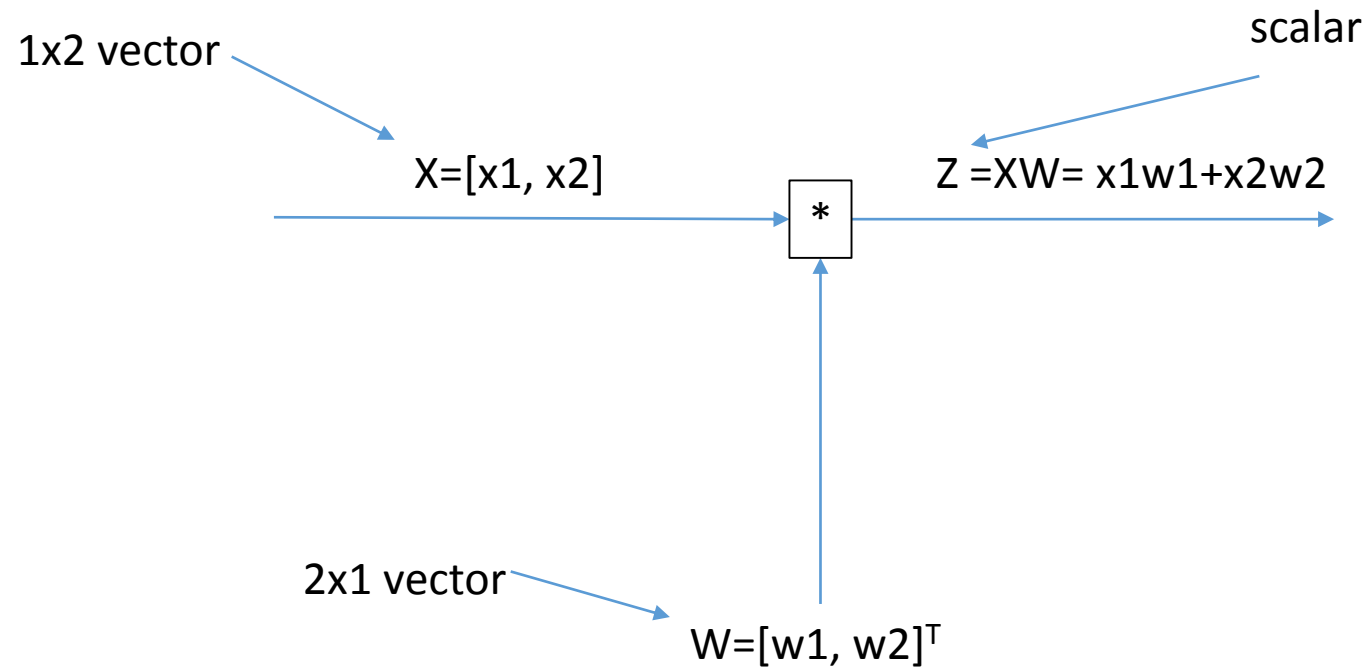
# Chain rule of derivative...



X

Z=f(X;W)

Z

Loss

Ideal output

W    (parameter)

If X, Z, W are matrices or vectors, then :

$$\nabla_X(\text{Loss}) = \left(\frac{\partial Z}{\partial X}\right) * \nabla_Z(\text{Loss})$$
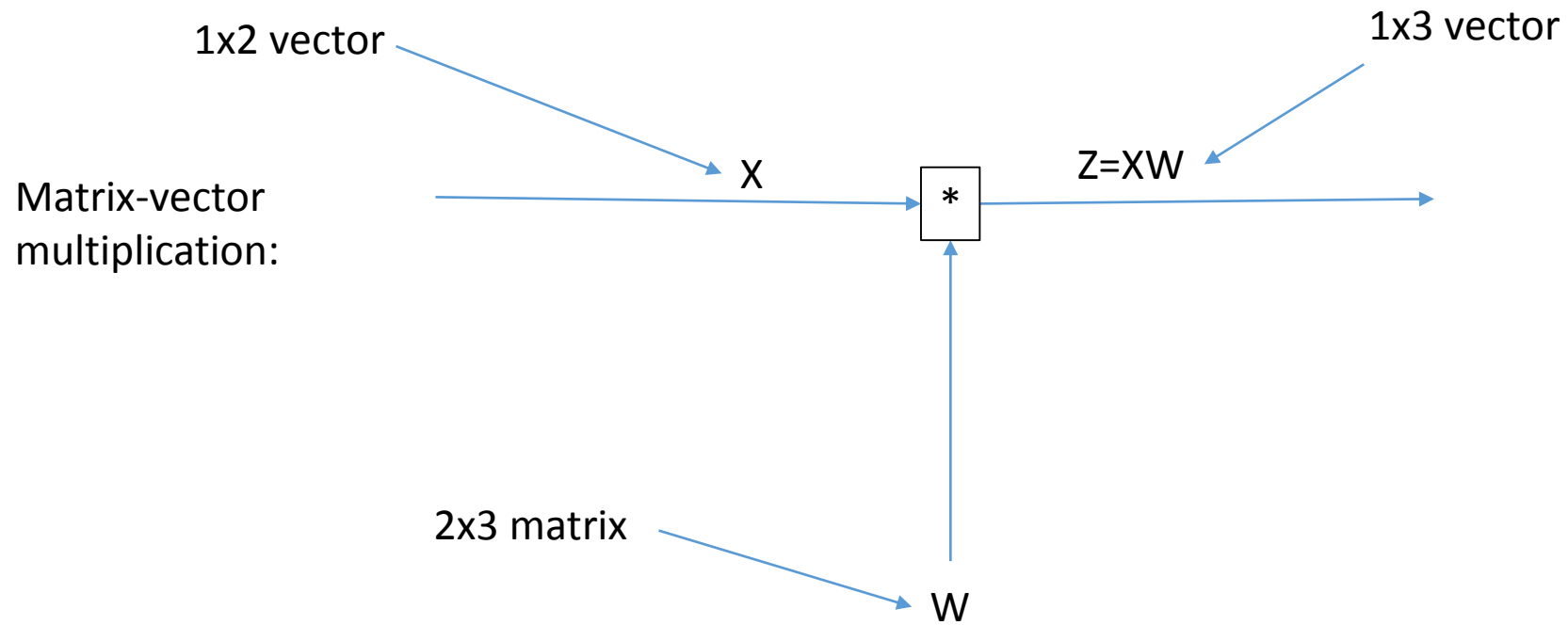
"*" refers to matrix vector multiplication

$$\nabla_W(\text{Loss}) = \left(\frac{\partial Z}{\partial W}\right) * \nabla_Z(\text{Loss})$$

# Example 1

1x2 vector

scalar

X=[x1, x2]

Z =XW= x1w1+x2w2

$*$

2x1 vector

W=[w1, w2]$^\mathsf{T}$

Chain rules: $\nabla_X(\text{Loss}) = W^T \dfrac{\partial(\text{Loss})}{\partial Z} = \begin{bmatrix} w1 & w2 \end{bmatrix} \dfrac{\partial(\text{Loss})}{\partial Z}$

Why?

$\nabla_W(\text{Loss}) = X^T \dfrac{\partial(\text{Loss})}{\partial Z} = \begin{bmatrix} x1 \\ x2 \end{bmatrix} \dfrac{\partial(\text{Loss})}{\partial Z}$
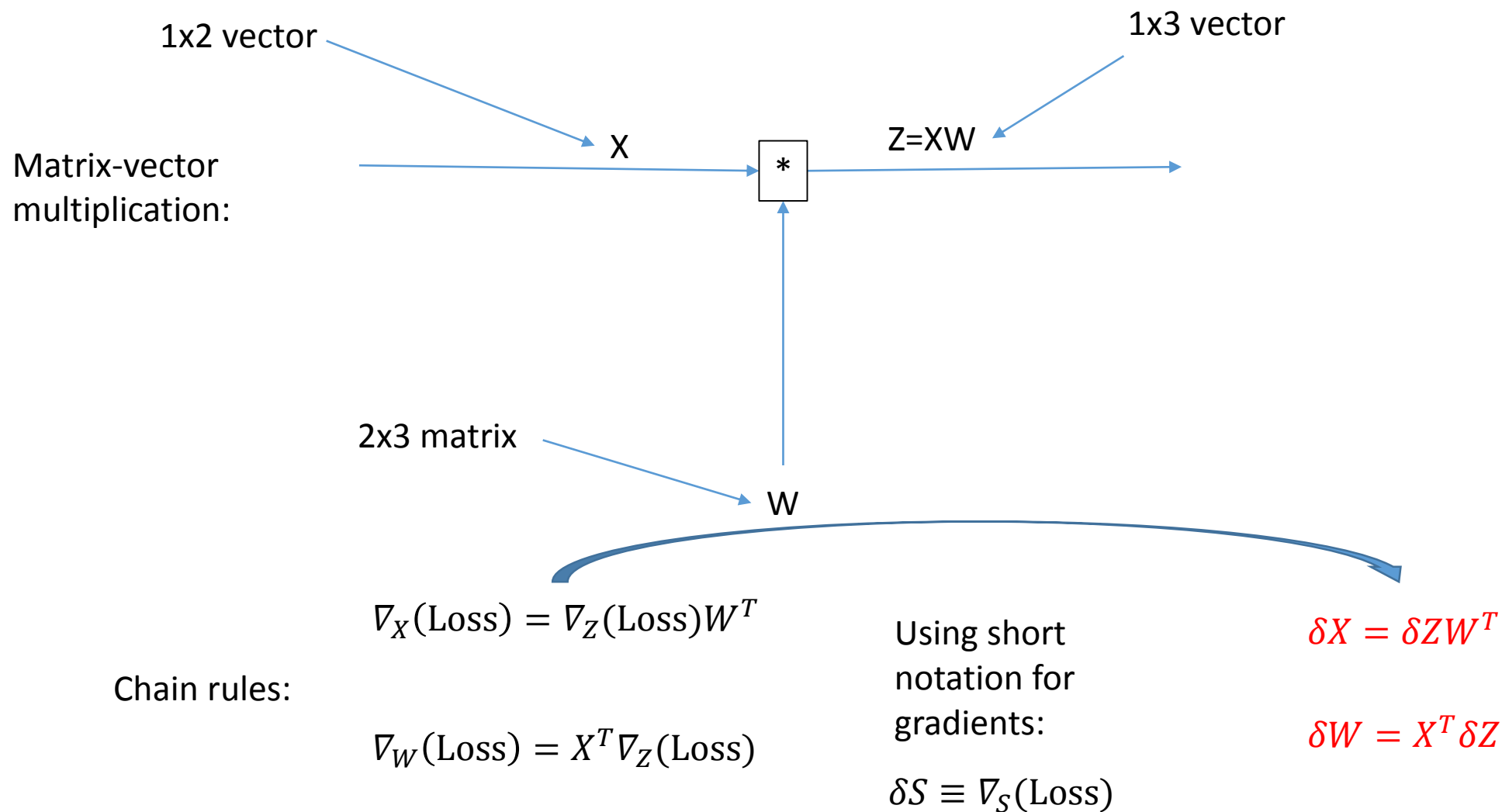
# Example 2

1x2 vector

1x3 vector

Matrix-vector multiplication:

X

Z=XW

*

2x3 matrix

W

$$\nabla_X(\text{Loss}) = \nabla_Z(\text{Loss})W^T$$

Chain rules:

Why?

$$\nabla_W(\text{Loss}) = X^T\nabla_Z(\text{Loss})$$

# Backprop derivation

1x2 vector

1x3 vector

Matrix-vector multiplication:

X

$Z=XW$

*

2x3 matrix

W

Chain rules:

$$\nabla_X(\text{Loss}) = \nabla_Z(\text{Loss})W^T$$

$$\nabla_W(\text{Loss}) = X^T\nabla_Z(\text{Loss})$$

Using short notation for gradients:

$$\delta S \equiv \nabla_S(\text{Loss})$$

$$\delta X = \delta Z W^T$$

$$\delta W = X^T \delta Z$$

# Backprop derivation...

$$\delta X_i = \sum_k \frac{\partial Z_k}{\partial X_i} \frac{\partial (\text{Loss})}{\partial Z_k} = \sum_k \frac{\partial}{\partial X_i} \left[ \sum_j X_j W_{jk} \right] \delta Z_k = \sum_k W_{ik} \delta Z_k \quad \Longrightarrow \quad \textcolor{red}{\delta X = \delta Z W^T}$$

$i^{\text{th}}$ component of $\delta X$ vector

Chain rule of derivative

Substitute $Z_k$

$k^{\text{th}}$ component of $\delta Z$ vector

Because,

$$\frac{\partial}{\partial X_i} \left[ \sum_j X_j W_{jk} \right] = W_{ik}$$

Writing in matrix-vector multiplication form

# Backprop derivation...

$$\delta W_{ij} = \sum_k \frac{\partial Z_k}{\partial W_{ij}} \frac{\partial(\text{Loss})}{\partial Z_k} = \sum_k \frac{\partial}{\partial W_{ij}} \left[ \sum_m X_m W_{mk} \right] \delta Z_k = X_i \delta Z_j$$

$$\Longrightarrow \quad \delta W = X^T \delta Z$$

$(i,j)$th component of $\delta W$ matrix

Chain rule of derivative

Substitute $Z_k$

Because,

$k$th component of $\delta Z$ vector

Writing in matrix-vector multiplication form

$$\frac{\partial}{\partial W_{ij}} \left[ \sum_m X_m W_{mk} \right] = \begin{cases} X_i, & \text{if } i = m \text{ and } j = k, \\ 0, & \text{otherwise.} \end{cases}$$
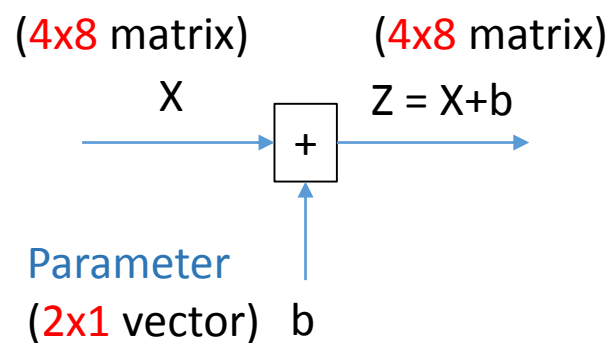
# Backprop derivation…

"Broadcast" addition:

(4x8 matrix)

X

(4x8 matrix)

Z = X+b

$+$

Parameter

(2x1 vector)  b

$$\delta X_{i,j} = \sum_k \sum_l \frac{\partial Z_{k,l}}{\partial X_{i,j}} \delta Z_{k,l} = \sum_k \sum_l \frac{\partial}{\partial X_{i,j}} [X_{k,l} + b_k] \delta Z_{k,l} = \delta Z_{i,j} \implies \delta X = \delta Z$$
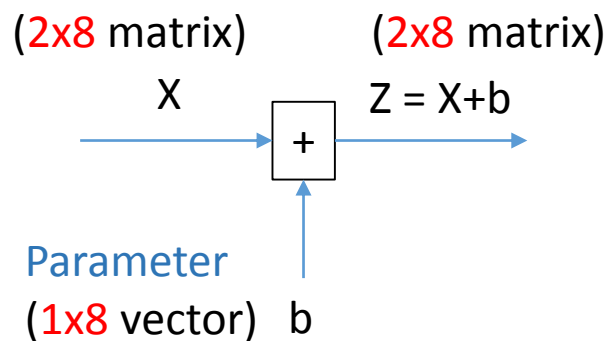
Chain rule

Substitute $Z_{k,l}$

Because,

$$\frac{\partial}{\partial X_{i,j}} [X_{k,l} + b_k] = \begin{cases} 1, \text{if } i = k \text{ and } j = l, \\ 0, \text{otherwise.} \end{cases}$$

# Backprop derivation for broadcast addition

"Broadcast" addition:

(2x8 matrix)        (2x8 matrix)

X      +      Z = X+b

Parameter
(1x8 vector)   b

$$\delta b_i = \sum_k \sum_l \frac{\partial Z_{k,l}}{\partial b_i} \delta Z_{k,l} = \sum_k \sum_l \frac{\partial}{\partial b_i} [X_{k,l} + b_l] \delta Z_{k,l} = \sum_k \delta Z_{k,i}$$

$$\delta b = \sum_k \delta Z_{k,:}$$

Chain rule

Substitute $Z_{k,l}$
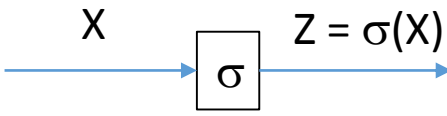
Because,

$$\frac{\partial}{\partial b_i} [X_{k,l} + b_l] = \begin{cases} 1, \text{if } i = l, \\ 0, \text{otherwise.} \end{cases}$$

# Backprop derivation for activation function

Non-linear function: (applied pointwise)

X → $\boxed{\sigma}$ → Z = σ(X)

Using chain rule:   $\delta X_{i,j} = \dfrac{dZ_{i,j}}{dX_{i,j}} \delta Z_{i,j} = \dfrac{d\sigma(X_{i,j})}{dX_{i,j}} \delta Z_{i,j}$

If the non-linear function is sigmoid,   $\sigma(a) = \dfrac{1}{1 + \exp(-a)}$

$$\frac{d\sigma}{da} = \frac{\exp(-a)}{(1 + \exp(-a))^2} = \frac{1}{1 + \exp(-a)}\left(1 - \frac{1}{1 + \exp(-a)}\right) = \sigma(a)(1 - \sigma(a))$$

$$\delta X_{i,j} = \sigma(X_{i,j})(1 - \sigma(X_{i,j}))\delta Z_{i,j}$$
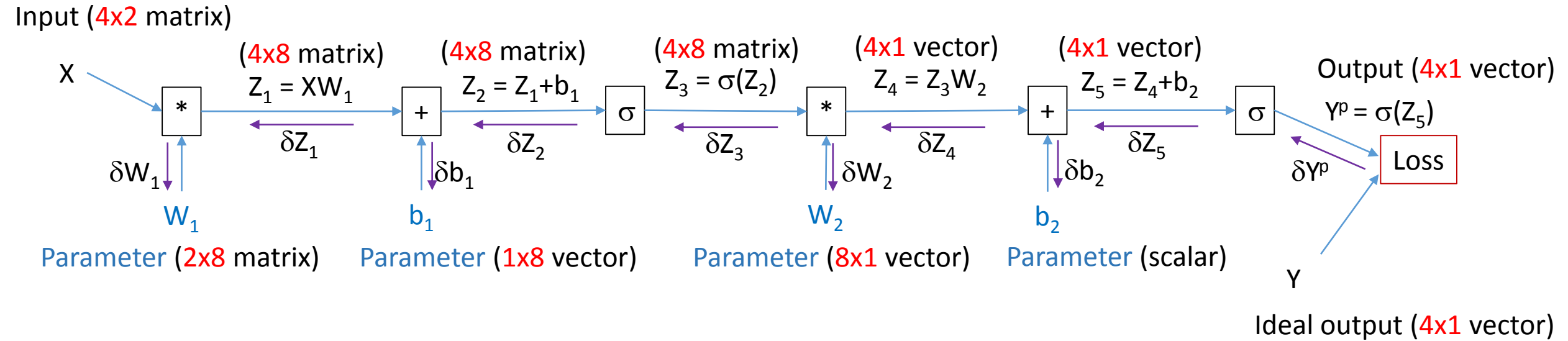
# Backprop derivation for loss function

Euclidean loss function:
$$Loss(Y^p, Y) = \frac{1}{2}\|Y^p - Y\|^2 = \frac{1}{2}\sum_i (Y_i^p - Y_i)^2$$

$i^{th}$ component of $\delta Y^p$ vector:
$$\delta Y_i^p = \frac{\partial}{\partial Y_i^p} Loss(Y^p, Y) = \frac{\partial}{\partial Y_i^p}\frac{1}{2}\sum_i (Y_i^p - Y_i)^2 = Y_i^p - Y_i$$

Using vector notation:
$$\delta Y^p = Y^p - Y$$

# Apply chain rule to XOR neural network

Input (4x2 matrix)

(4x8 matrix)
$Z_1 = XW_1$

(4x8 matrix)
$Z_2 = Z_1+b_1$

(4x8 matrix)
$Z_3 = \sigma(Z_2)$

(4x1 vector)
$Z_4 = Z_3W_2$

(4x1 vector)
$Z_5 = Z_4+b_2$

Output (4x1 vector)
$Y^p = \sigma(Z_5)$

X

$*$ $\delta Z_1$ $+$ $\delta Z_2$ $\sigma$ $\delta Z_3$ $*$ $\delta Z_4$ $+$ $\delta Z_5$ $\sigma$ $\delta Y^p$ Loss

$\delta W_1$ $\delta b_1$ $\delta W_2$ $\delta b_2$

$W_1$

$b_1$

$W_2$

$b_2$

Parameter (2x8 matrix)

Parameter (1x8 vector)

Parameter (8x1 vector)

Parameter (scalar)

Y

Ideal output (4x1 vector)

Chain rule of derivatives:

$$\delta Y^p = Y^p - Y$$

$$\delta Z_5 = \sigma(Z_5)(1-\sigma(Z_5))\delta Y^p$$

$$\delta Z_4 = \delta Z_5$$

$$\delta Z_3 = \delta Z_4 W_2^T$$

$$\delta Z_2 = \sigma(Z_2)(1-\sigma(Z_2))\delta Z_3$$

$$\delta Z_1 = \delta Z_2$$

New notation:
$$\delta S \equiv \nabla_S(\text{Loss})$$

Gradient of "Loss" with respect to input signals

Propagates backward

$$\delta W_2 = Z_3^T \delta Z_4$$

$$\delta b_2 = \sum_k (\delta Z_5)_k$$

$$\delta W_1 = X^T \delta Z_1$$

$$\delta b_1 = \sum_k (\delta Z_2)_{k,:}$$

Gradient of "Loss" with respect to parameters

# Backprop to train a neural net

Initialize all parameters of the neural network

Initialize learning rate variable *lr*

Iterate:

If loading the whole training data, do it once outside the "Iterate" loop, to be efficient

    (Load Data): Get training data batch

    (Forward pass): Compute $Z_1, Z_2, \ldots, Y^p$

    (Backward pass): Compute gradients $\delta Y^p, \delta Z_5, \ldots, \delta Z_1, \delta W_2, \delta W_1, \delta b_2, \delta b_1$

    (Gradient descent to update parameters): $\quad W_2 \leftarrow W_2 - lr * \delta W_2, \quad b_2 \leftarrow b_2 - lr * \delta b_2, \ldots,$

    (Diagnostics): Compute "Loss" from time to time to check if it is decreasing

"Learn_XOR_manualBP.ipnyb" implements this learning algorithm

# PyTorch magic!

- Fortunately, PyTorch can automatically compute derivatives by chain rules!

- Also, it has several optimizers that can use these derivatives in the gradient descent optimization method.

- Look at "Learn_XOR.ipnyb" and "Learn_XOR_with_LBFGS.ipnyb"

# Universal function approximation

- A neural network with a single hidden layer can approximate "any" function!
  - Wikipedia has a clear statement:
    https://en.wikipedia.org/wiki/Universal_approximation_theorem
- A non-technical explanation of universality theorem:
  - http://neuralnetworksanddeeplearning.com/chap4.html

Why do we then even need multiple layers and why even deep nets?