



# Neural Networks

A basic and intuitive introduction to “shallow” fully-connected feed-forward networks

*Lion Krischer*



# Goal for Now + First Part of the Afternoon

- Fully understand and implement a simple feed-forward neural network from scratch.
- Little math (and probably slightly wrong notation).
- Go slow and step by step.

A bit tedious at times but not very hard.

Modern and deep networks are a bit more complicated but not fundamentally different.

# Machine Learning

---

—

Input  $\longrightarrow$   $f(x)$   $\longrightarrow$  Output

---

# Input



# $f(x)$



# Output

Source Parameters  
Earth Model

Physical Model

$F(\text{Earth Model, Source Parameters})$

Seismograms

---

**Input** —————→  **$f(x)$**  —————→ **Output**

Anything

Machine Learning Model

Anything

In a way we replace our knowledge of the physics or the logic of a problem with an arbitrary “learned” function.

---

---

**Input**

Anything



**$f(x)$**

Machine Learning Model



**Output**

Anything



---

**Input** —————→  **$f(x)$**  —————→ **Output**

Anything

Machine Learning Model

Anything

Data Driven

---

**Input** →  **$f(x)$**  → **Output**

Anything

Supervised Machine Learning Model  
Neural Network

Anything

Data Driven



# When is this a good idea/useful?

- If it is very hard to make a physical or logical model.
  - Image and speech recognition, self-driving cars, language translation, ...
  - Phase detection and recognition (see Thursday), signal classification (see Wednesday), ...
- If the physical model is very expensive to compute but the result is somehow “simple”.
- And: One has a LOT of data.

# Neural Networks

A photograph of a two-story wooden building, possibly a ski lodge or cabin, with a snow-covered roof. The building has a dark wood upper story with a balcony and a white lower story with several doors. It is situated in a snowy landscape with evergreen trees in the background. The text "Neural Networks" is overlaid in large white letters.



# Universal Function Approximators

- **Universal approximation theorem:** Feed-forward networks of finite size with a single hidden layer and non-linear activation functions can approximate any continuous function (under some mathematical restrictions).
- Newer results reduce the necessary width of the networks under some conditions.
- Only a statement that one could always create a neural network to perform a certain action. Not if it is actually learnable.

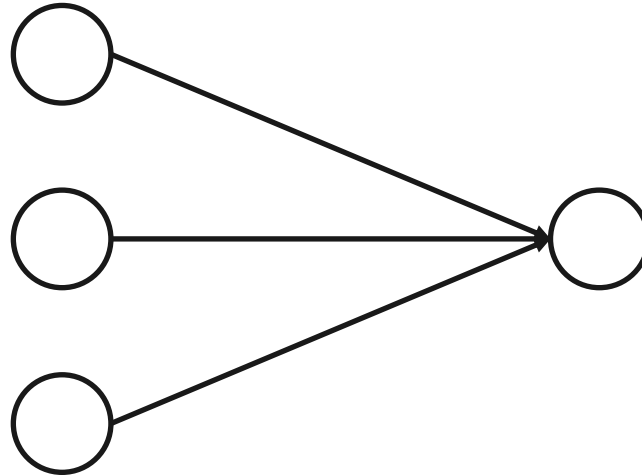
Thus the questions is not “Can a neural network do this?” but “Can I make it learn what I want it to do?”.

# Neural Networks Step by Step

---

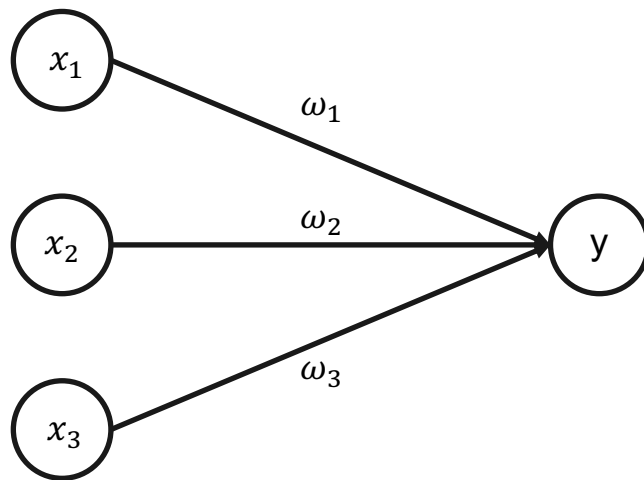


# Weighted Sum





# Weighted Sum

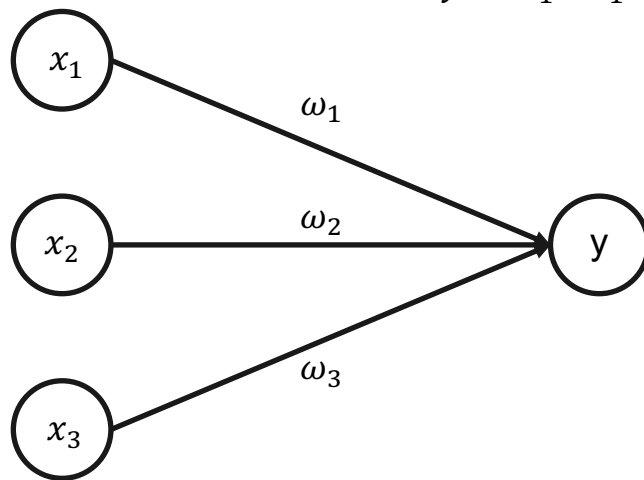






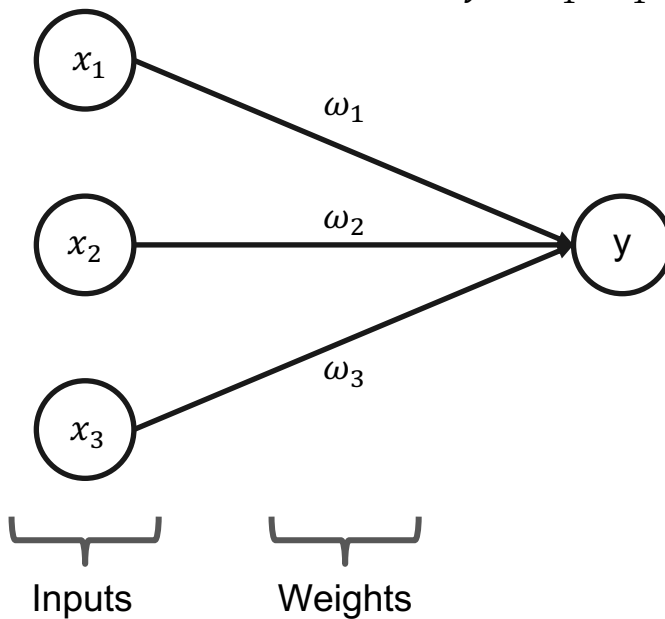
# Weighted Sum

$$y = \omega_1 * x_1 + \omega_2 * x_2 + \omega_3 * x_3 = \sum_i \omega_i x_i$$



# Weighted Sum

$$y = \omega_1 * x_1 + \omega_2 * x_2 + \omega_3 * x_3 = \sum_i \omega_i x_i$$

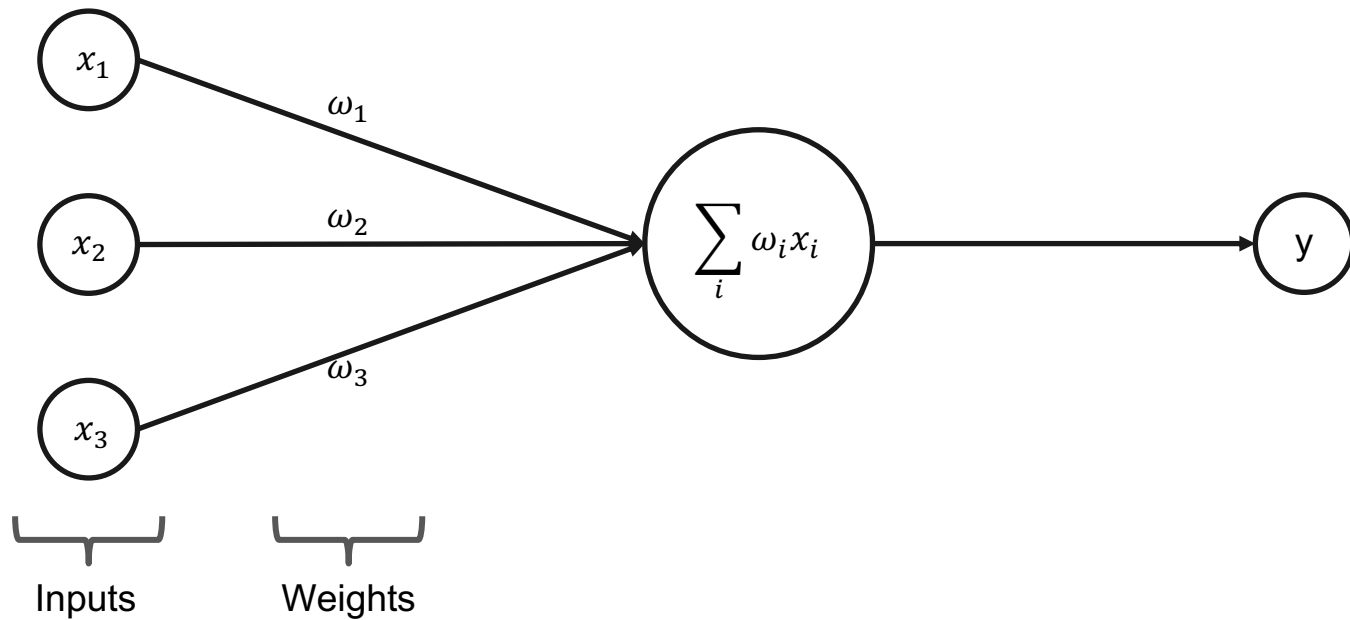




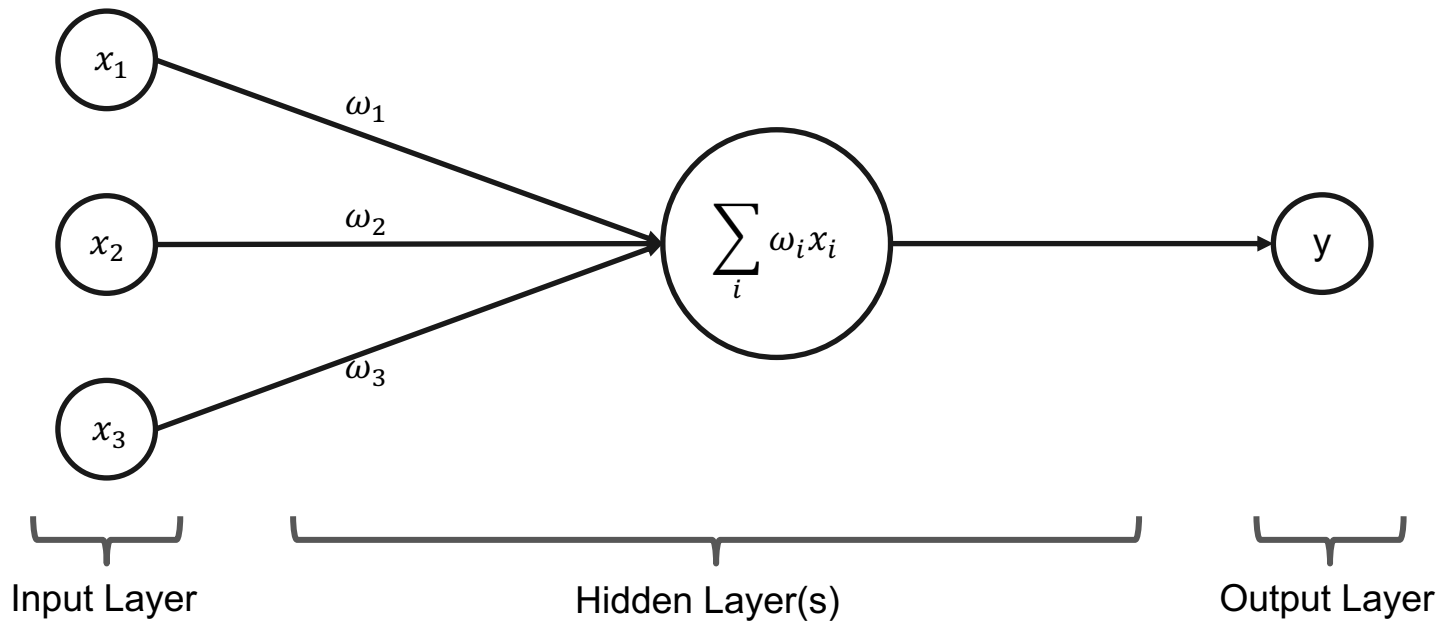
# Weighted Sum

Not too useful yet – basically a dot product

# Weighted Sum



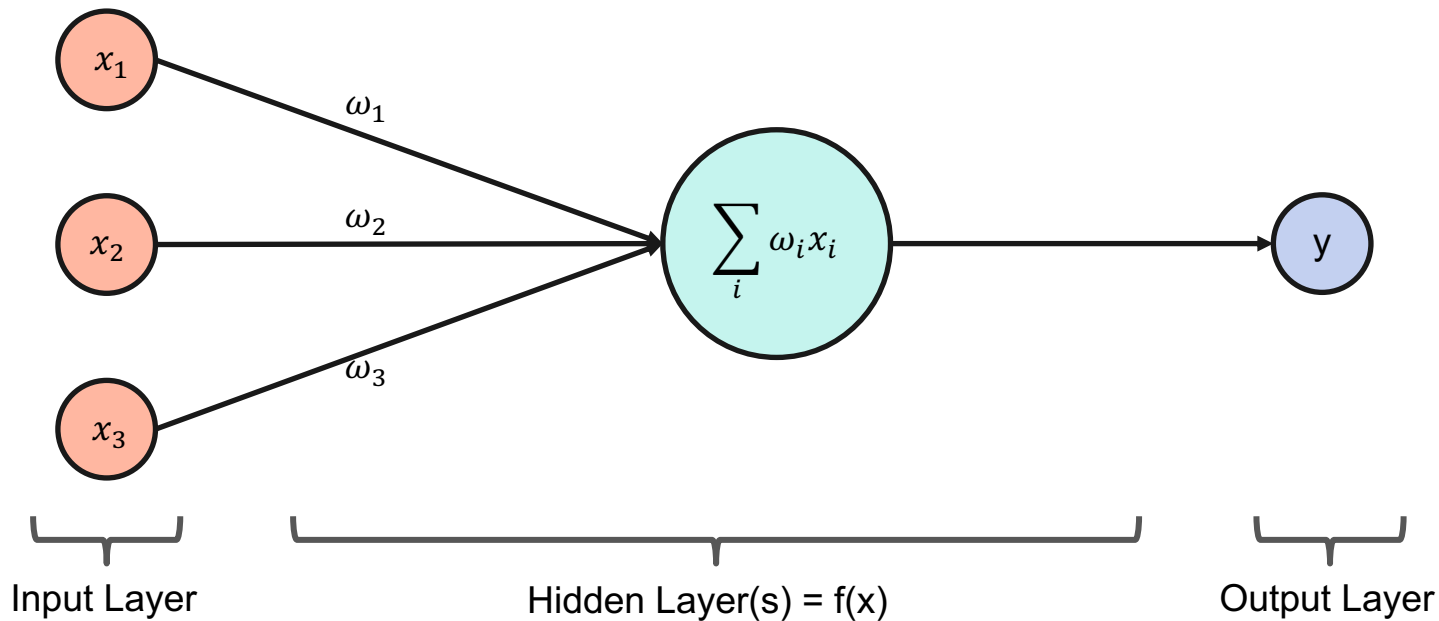
# Weighted Sum



—

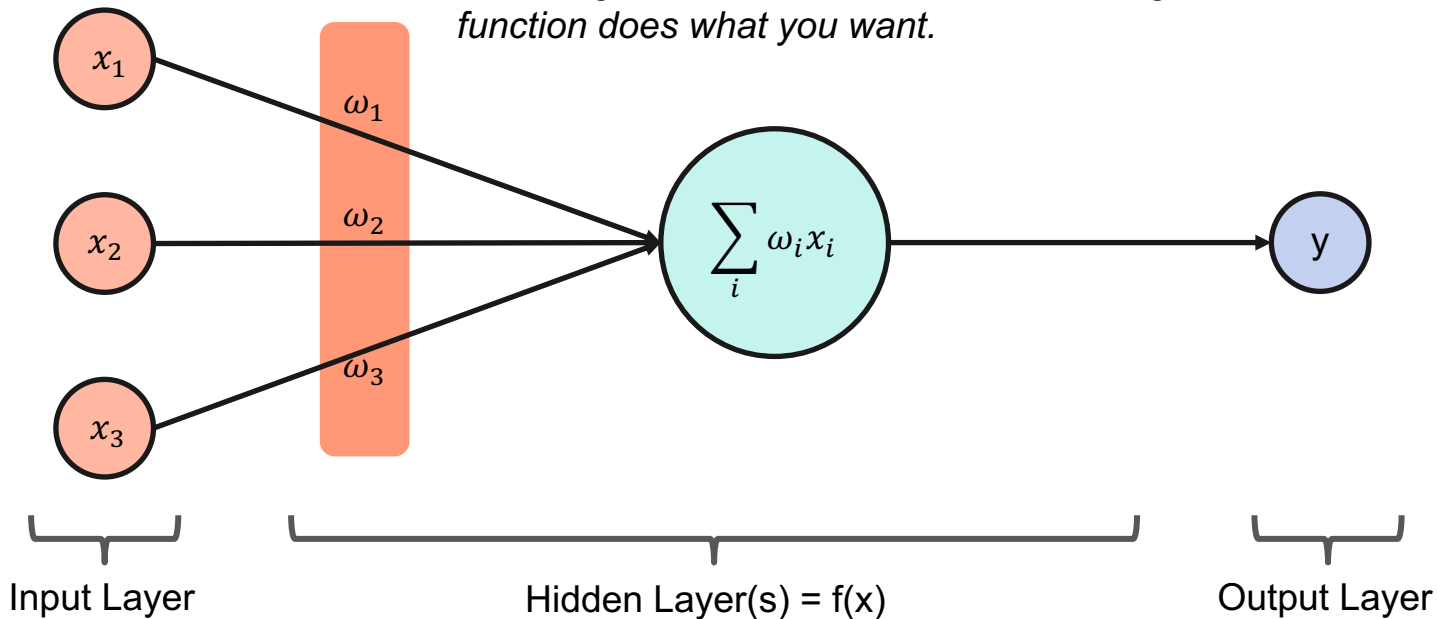
Input  $\longrightarrow$   $f(x)$   $\longrightarrow$  Output

# Weighted Sum



# Weighted Sum

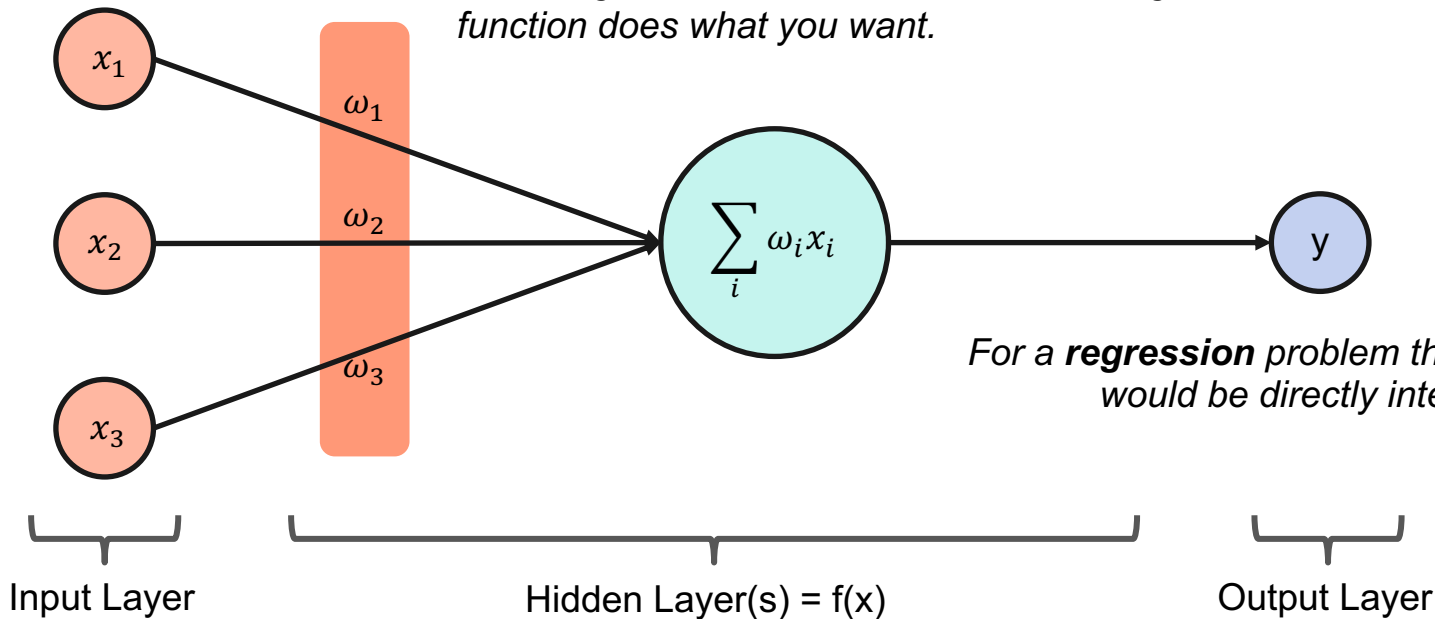
*Change weights to change the behavior of the function.  
"Learning" would thus be to find better weights so the  
function does what you want.*





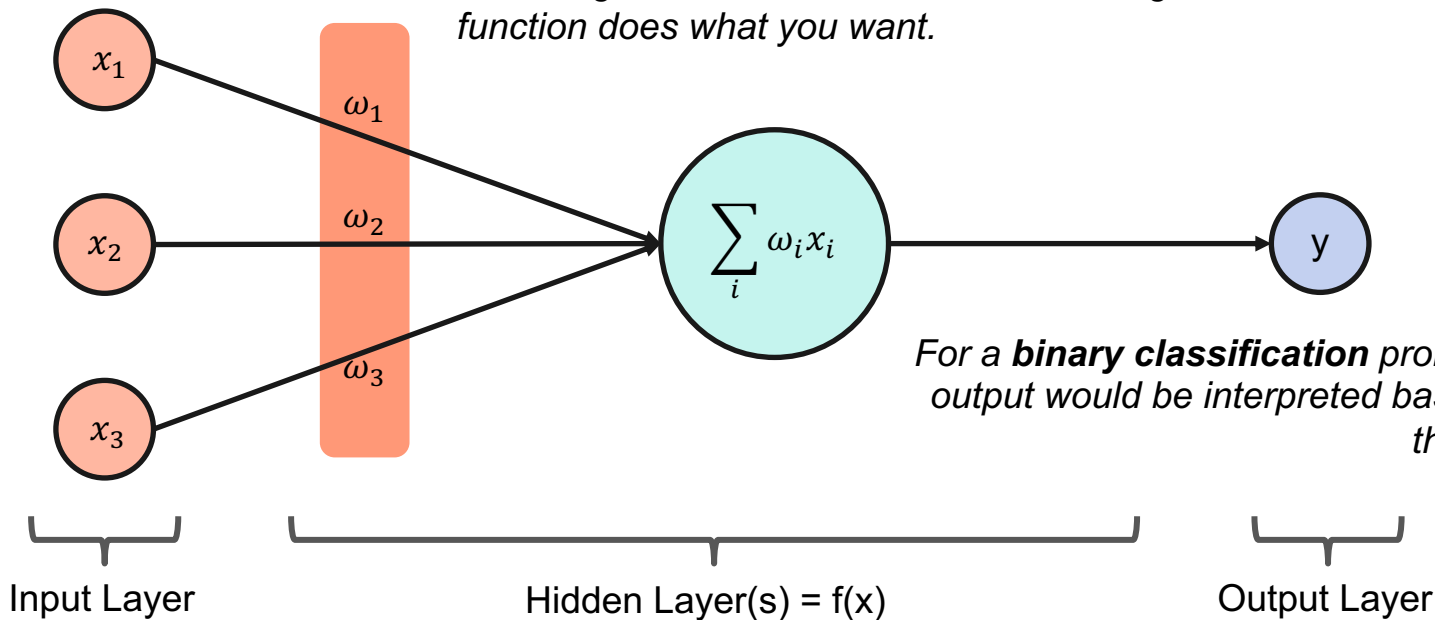
# Weighted Sum

*Change weights to change the behavior of the function.  
"Learning" would thus be to find better weights so the  
function does what you want.*



# Weighted Sum

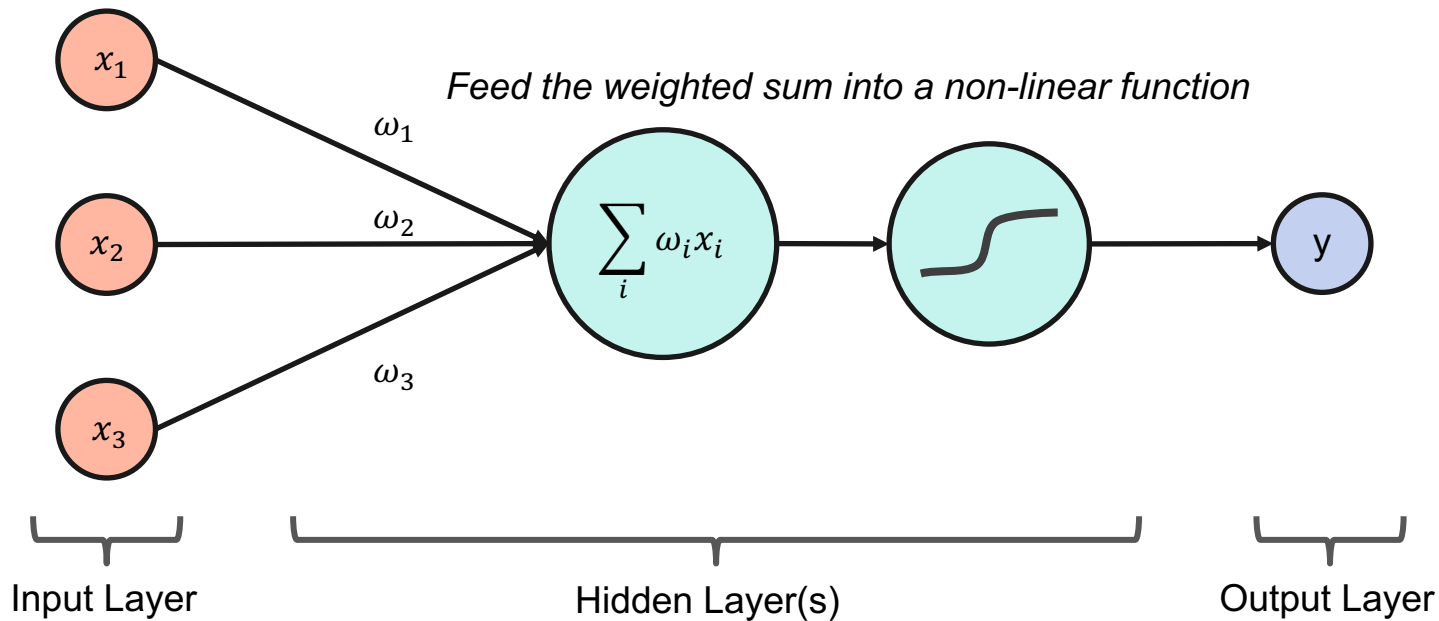
*Change weights to change the behavior of the function.  
"Learning" would thus be to find better weights so the  
function does what you want.*



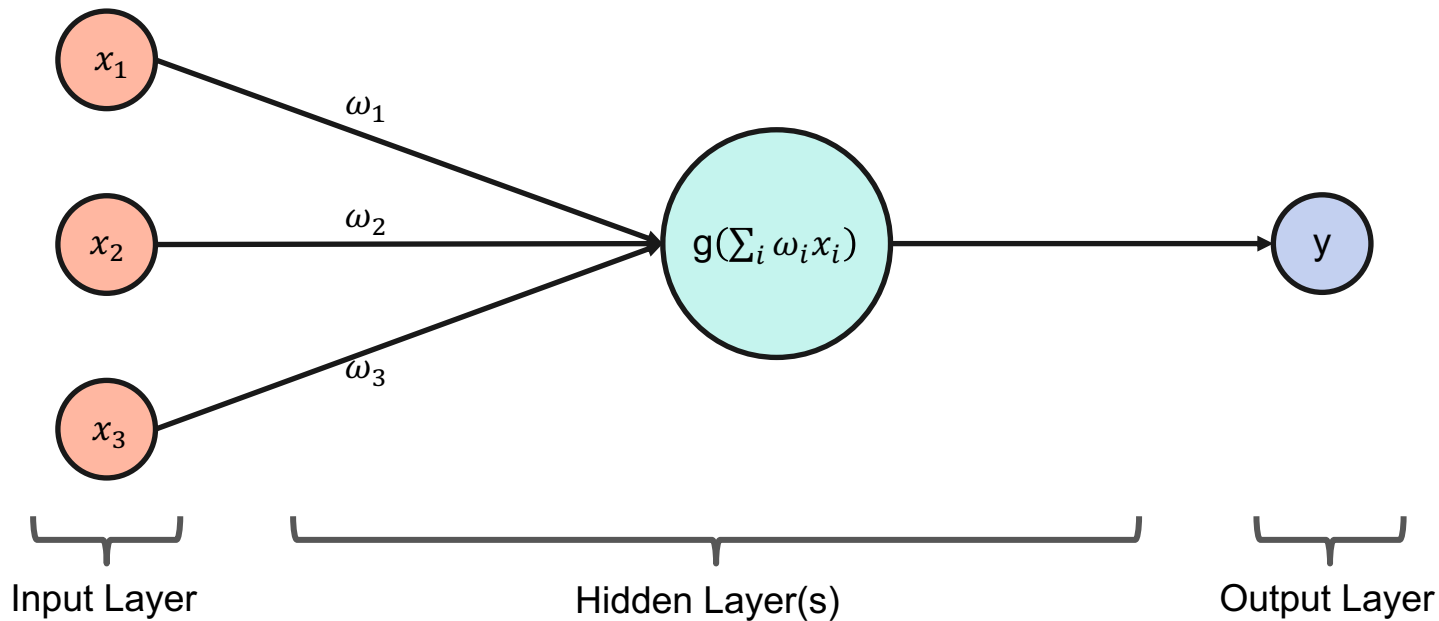


Computing something non-linear might be useful.

# Activation Function



# Activation Function





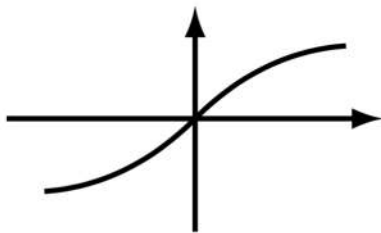
# Why do we need an activation function?

- Some non-linearity is needed to compute anything non-linear.
- More complex relationships between inputs and outputs cannot be represented by something purely linear.
- (I think) they are called activation functions because the first one used was a Heaviside function and thus either “activated” a neuron or not for a given set of inputs.

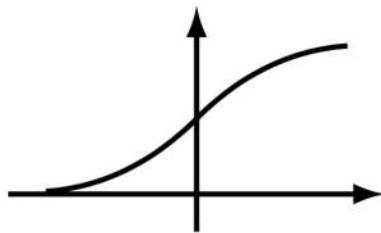
# Popular Choices

- Features people are looking for in activation function:
  - At least piecewise differentiable for reasons we'll elaborate upon later
  - Monotonous
  - Helpful if they bound the output (but not all modern choices do that)
  - Computationally cheap

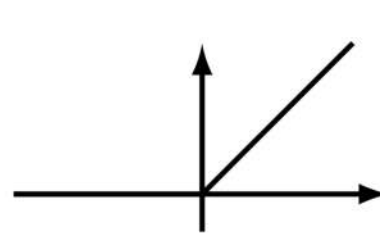
$\tanh(x)$



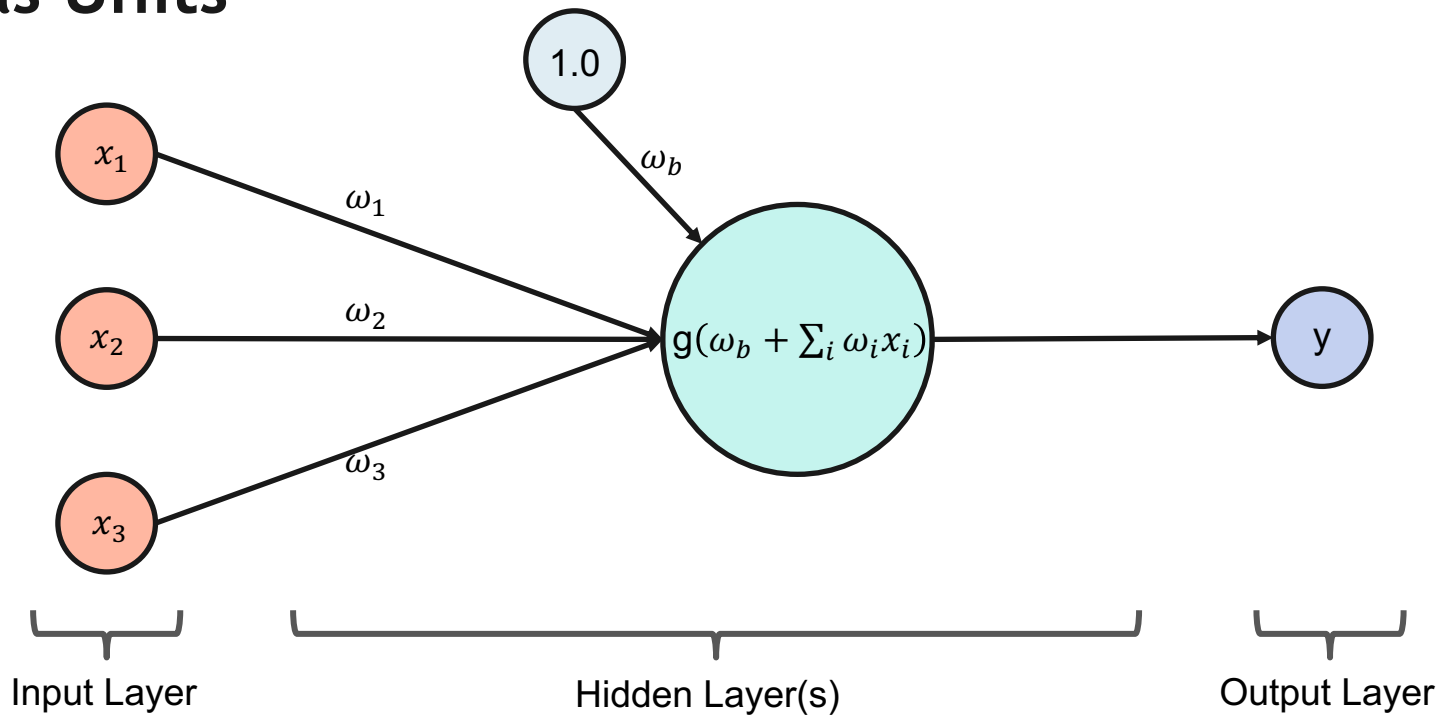
$\text{sigmoid}(x) = \frac{1}{1+e^{-x}}$



$\text{relu}(x) = \max(0, x)$

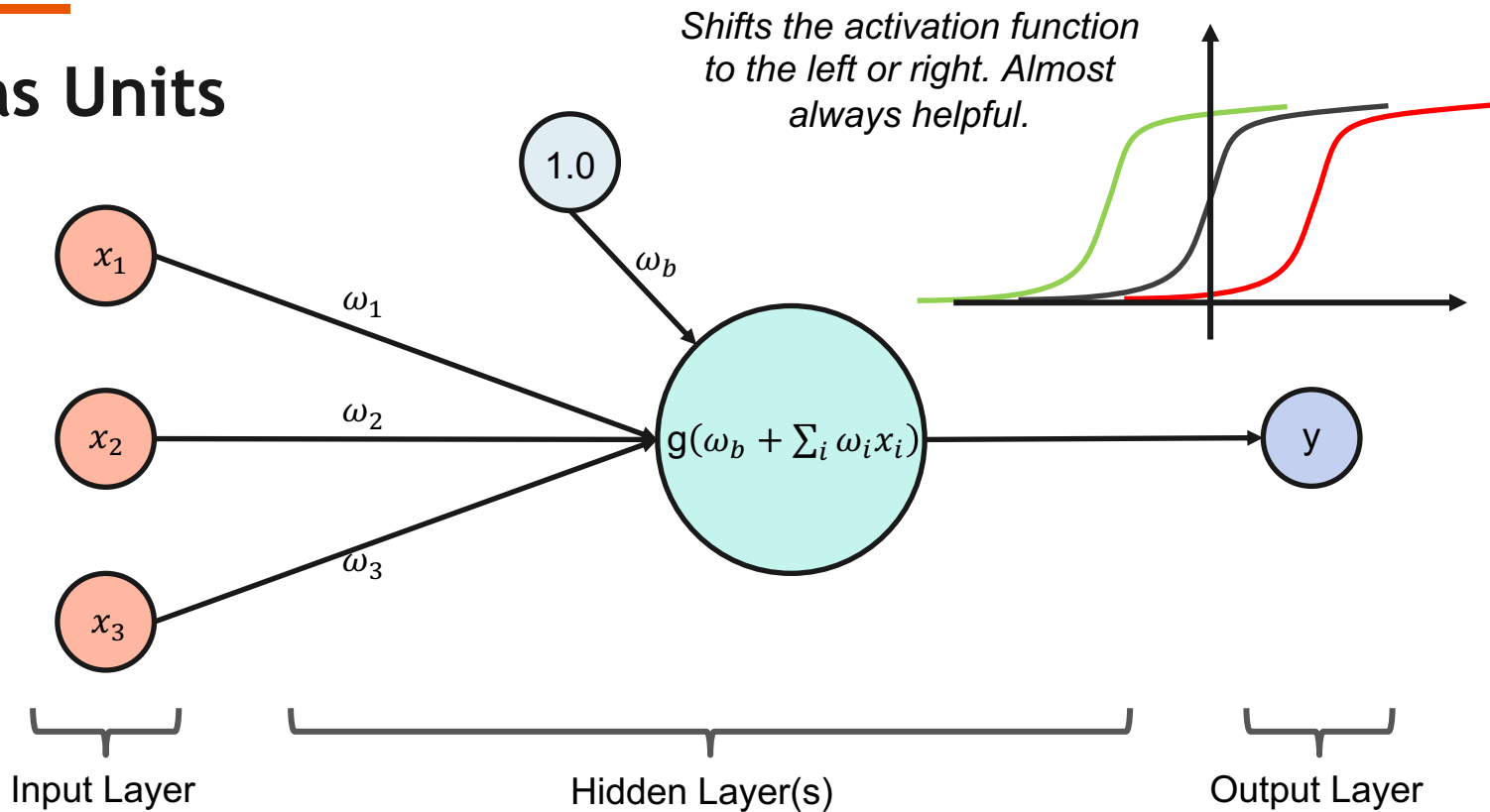


## Bias Units





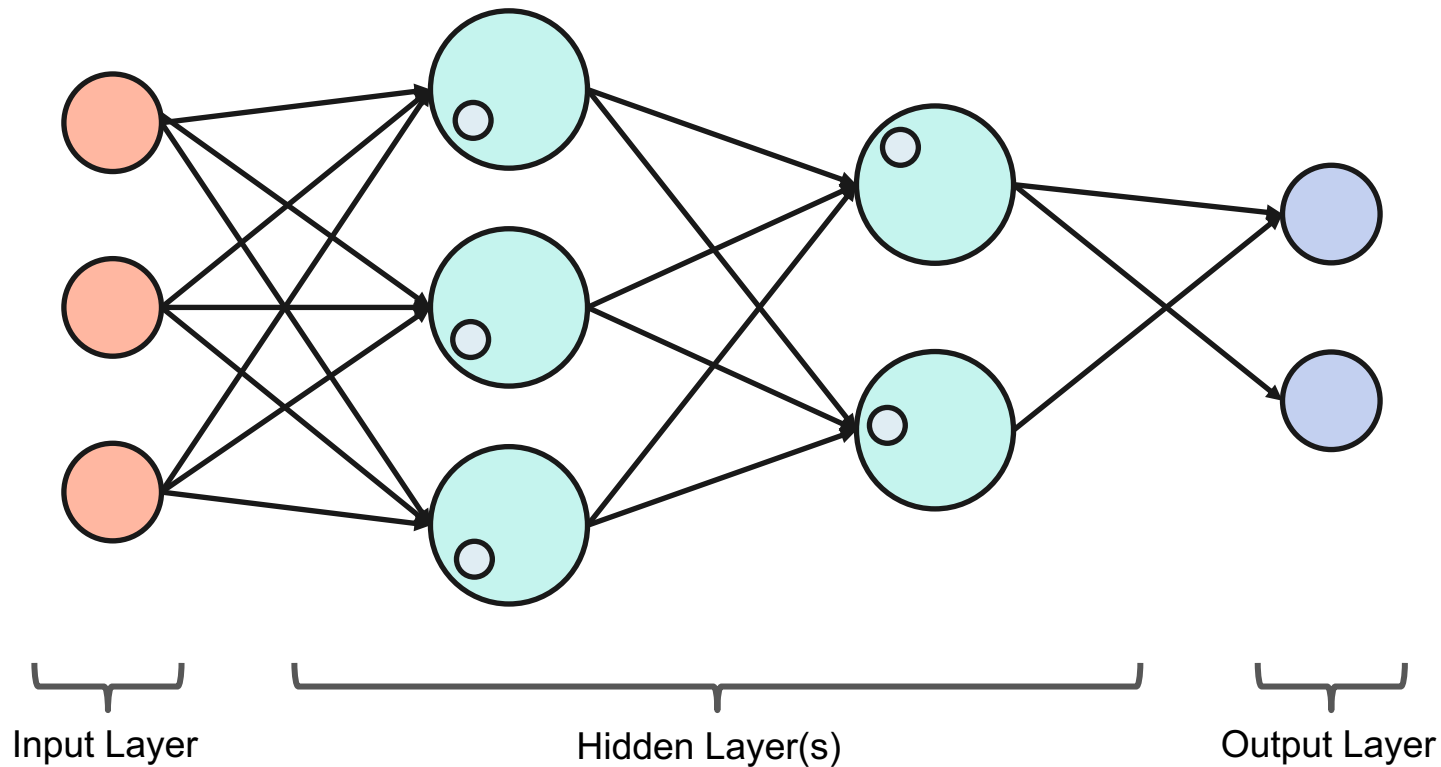
## Bias Units





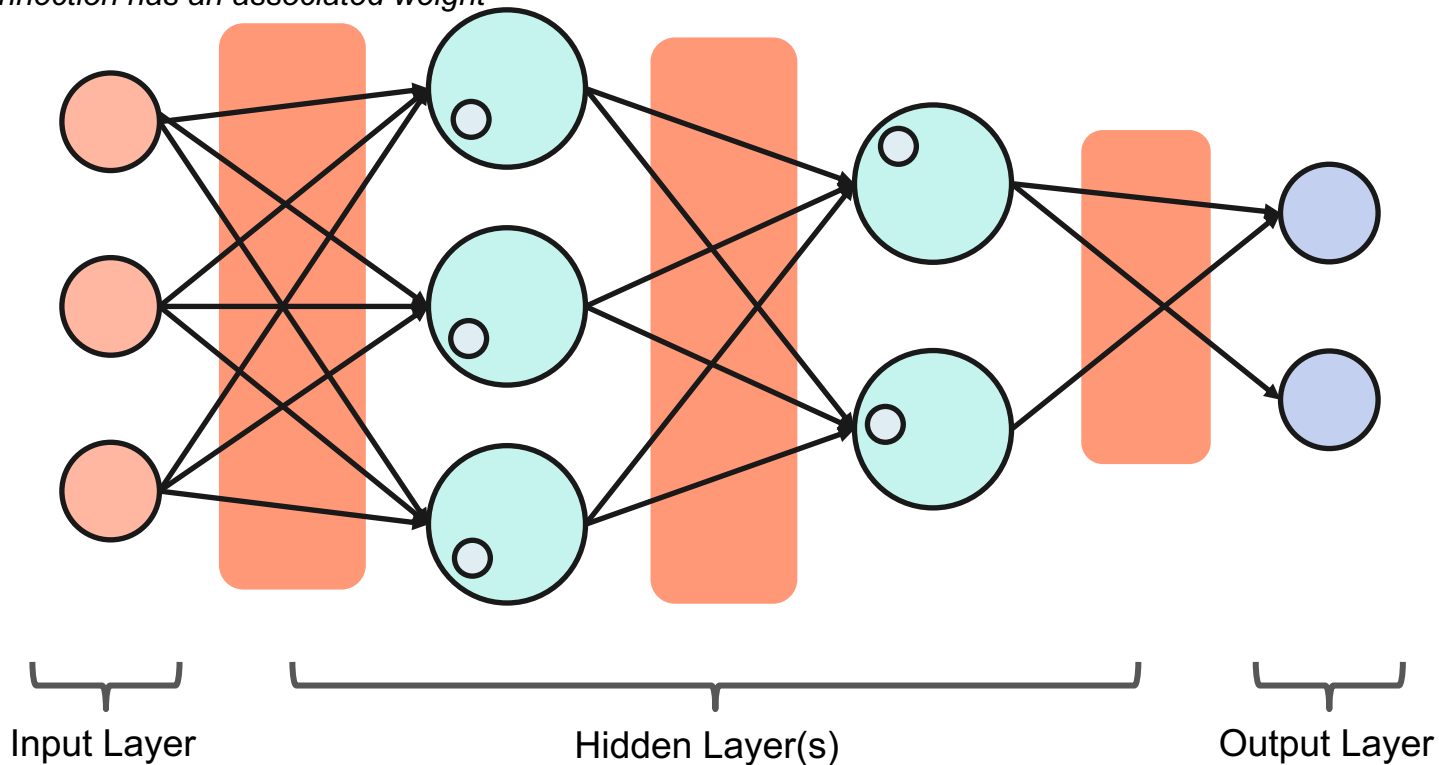
Now let's start to build a proper network.

# Fully-connected Feed-forward Neural Network



# Fully-connected Feed-forward Neural Network

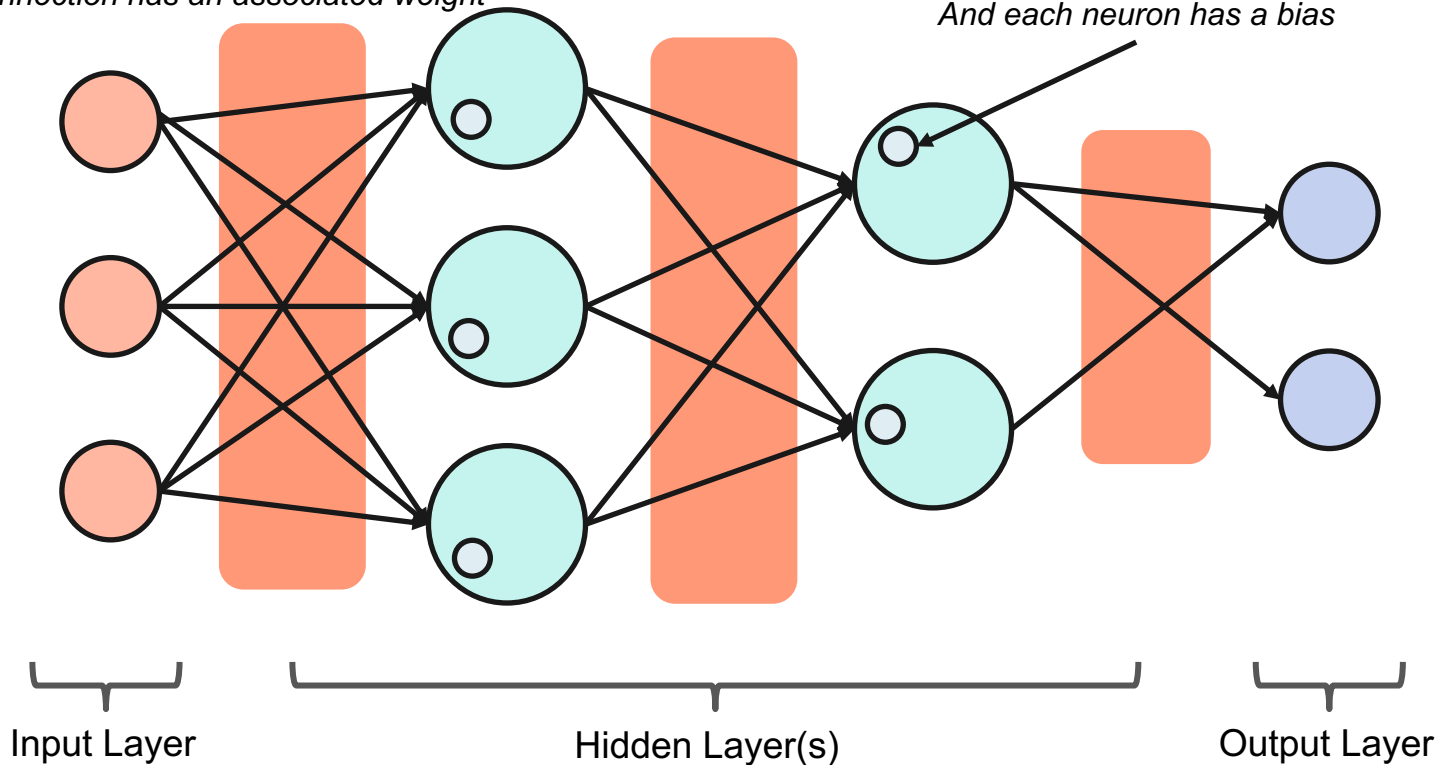
*Every connection has an associated weight*



# Fully-connected Feed-forward Neural Network

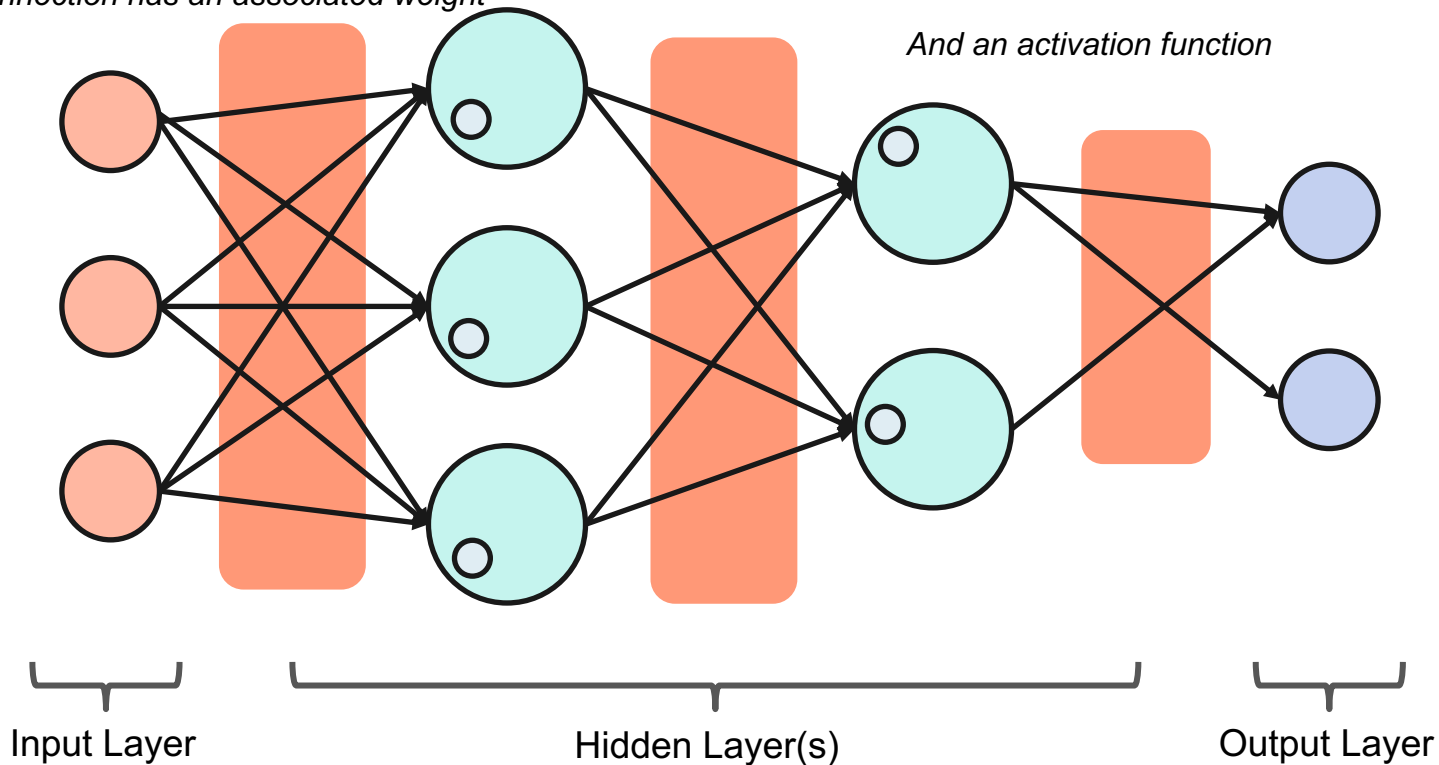
*Every connection has an associated weight*

*And each neuron has a bias*

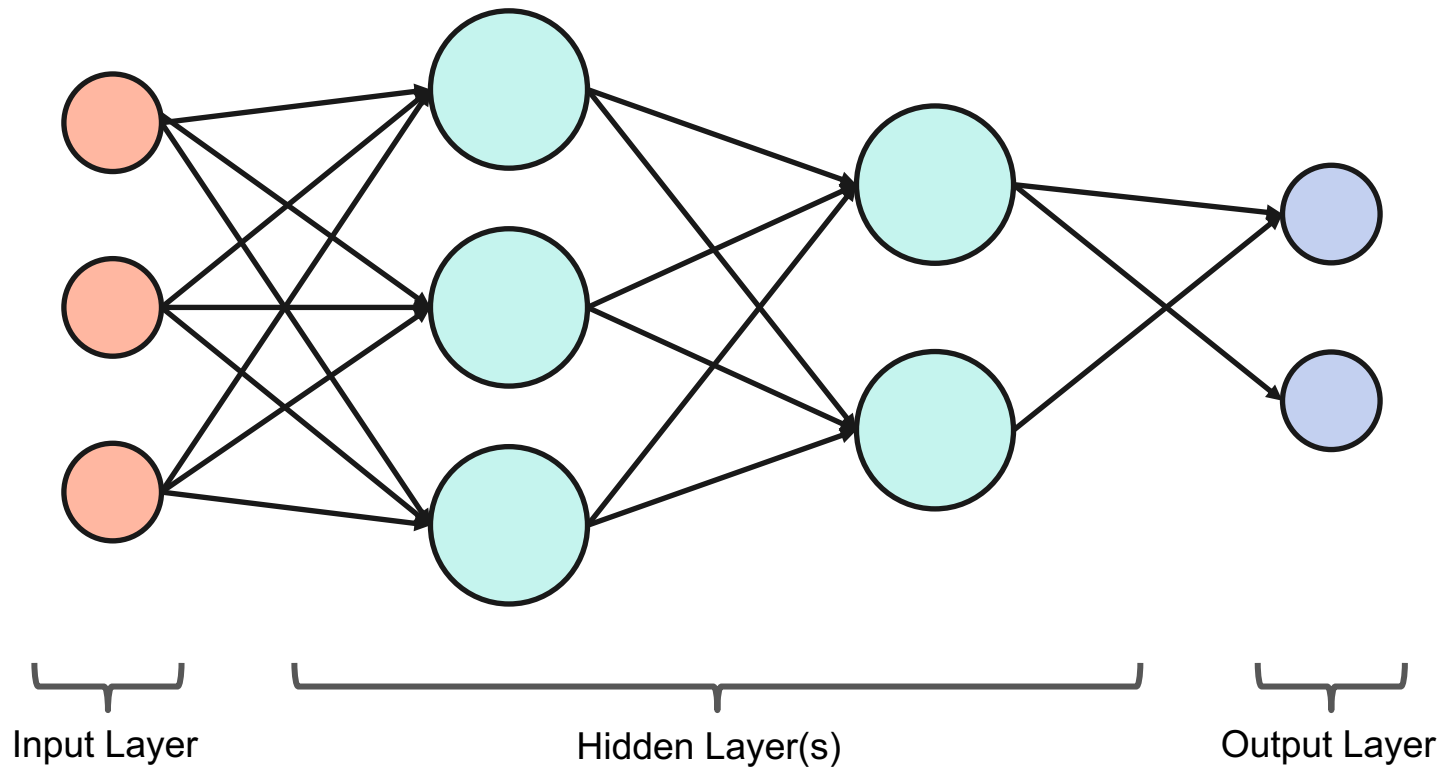


# Fully-connected Feed-forward Neural Network

*Every connection has an associated weight*



# Fully-connected Feed-forward Neural Network





# What is gained by adding more “depth”?

- Remember that everything can also be done with a single layer but this layer might get very large.
- If done properly the hope is to decompose a problem in many separate processing or computation steps.
- More layers thus means more steps which allows more complicated functions.
- This generalizes better and also requires fewer parameters.
- Deep networks are discussed in depth on Thursday.

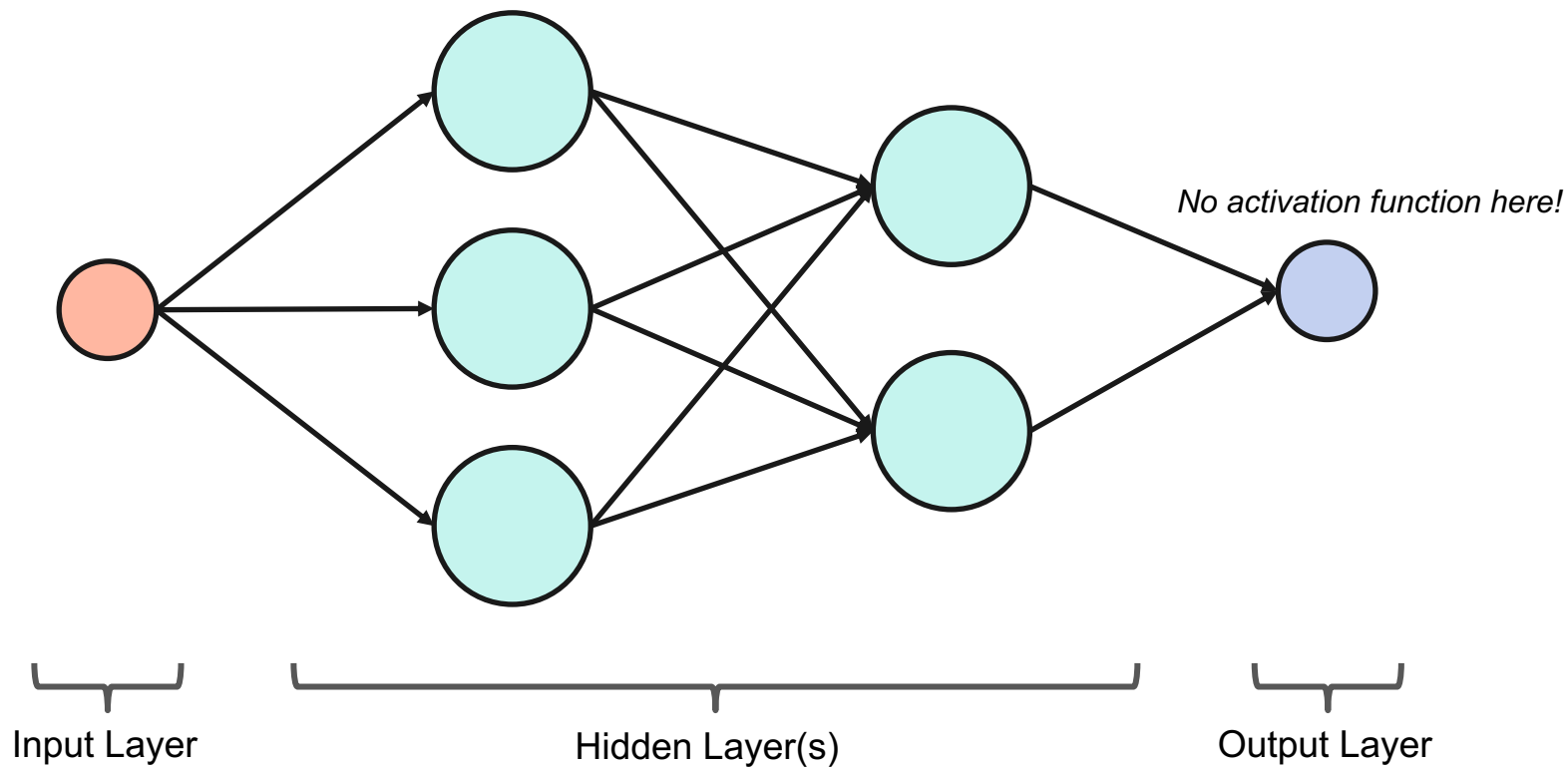


Now we know how a neural network computes something.

Let's implement this!

---

## Exercise: Forward Pass of this Network



**Open the Neural Networks from Scratch  
Notebook.**

**How to make it do what we  
want?**

—



So far we just multiplying numbers together...

—

Input  $\longrightarrow$   $f(x)$   $\longrightarrow$  Output

---

Input  $\longrightarrow$   $f(x)$   $\longrightarrow$  Output

Known input  
data

Known output  
data

---

Input  $\longrightarrow$   $f(x)$   $\longrightarrow$  Output

Known input  
data

Optimize for the  
parameters of  $f(x)$

Known output  
data



\_\_\_\_\_

Known input  
data

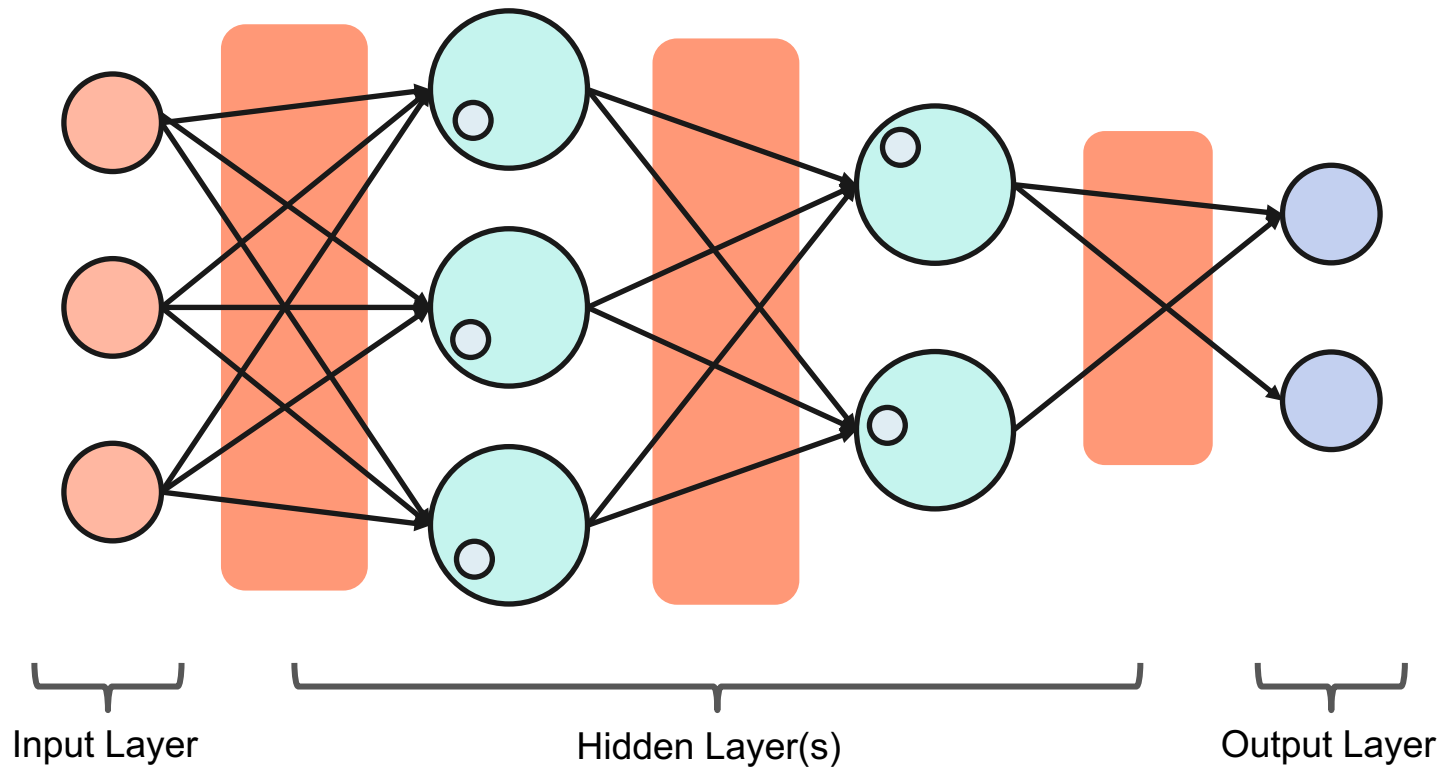
Optimize for the  
parameters of  $f(x)$

Known output  
data

**Objective for neural networks: Find a set of weights so that it can convert the inputs to output with as little error as possible.**

**The actual network layout is fixed.**

# Fully-connected Feed-forward Neural Network





# Numerical Optimization Problem

- These networks can have a lot of parameters => Global search methods infeasible
- Derivatives can be computed analytically
  - => Local gradient based optimization

We require two ingredients:

- A way to compute derivatives (back-propagation)
- A way to perform the actual optimization (here we use simple gradient descent)

# Back-Propagation with Pen & Paper

---



# What is Back-Propagation?

- Returns the analytic gradients of some objective function/misfit measurement/loss function with respect to the weights of a neural network
- Thus tells us how we have to change the network parameters to increase that loss function. We of course go the other way.
- It is really just the repeated application of the chain rule.
- Works just the same as auto-differentiation does.



# Simple Example

$$f(x, y, z) = ((x + y)z)^2$$



# Simple Example

$$f(x, y, z) = ((x + y)z)^2$$

Goal:

$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$



## Simple Example

$$f(x, y, z) = ((x + y)z)^2$$

Goal:

$$\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

Chain rule:

$$\frac{\partial f}{\partial z} = 2((x + y)z)(x + y)$$

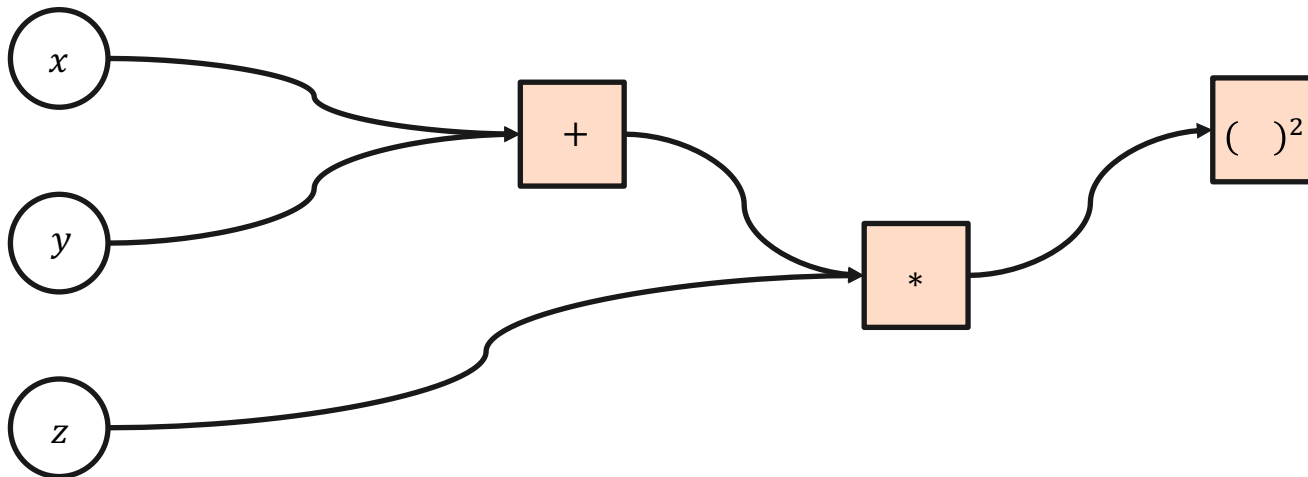
of course works, but tedious for deeply nested computations.

And inefficient, assuming that we already calculated for example  $(x + y)$  if we evaluated the function once.



$$f(x, y, z) = ((x + y)z)^2$$

## Simple Example

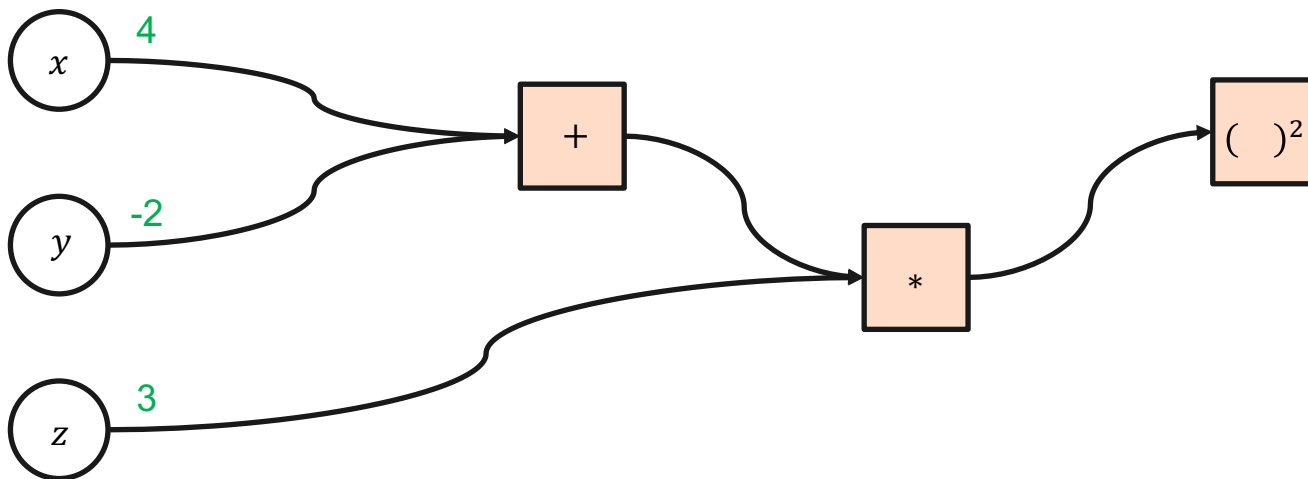




## Simple Example

$$f(x, y, z) = ((x + y)z)^2$$

Forward Pass

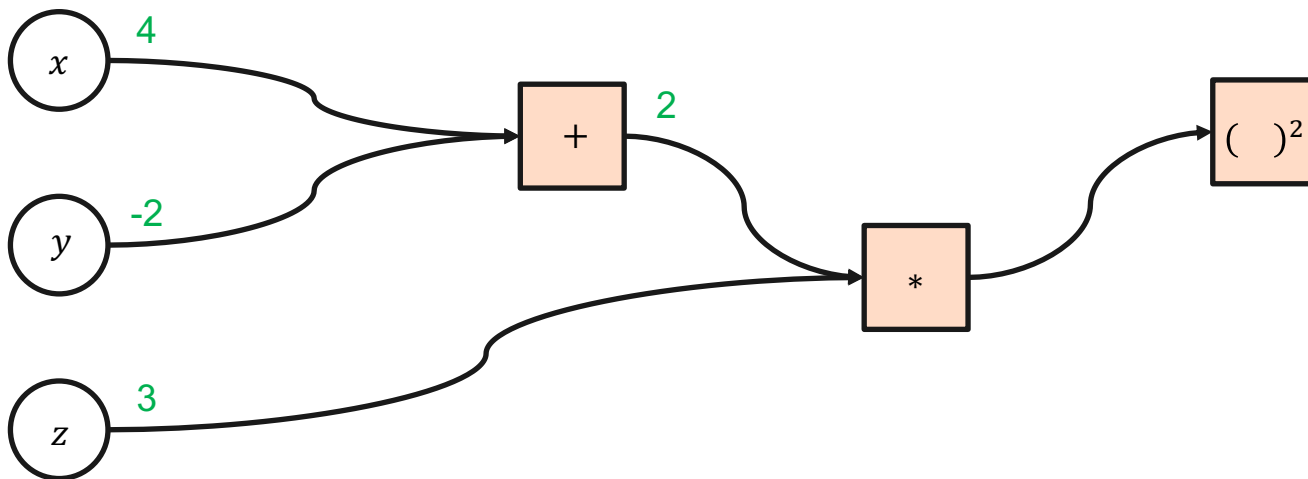




## Simple Example

$$f(x, y, z) = ((x + y)z)^2$$

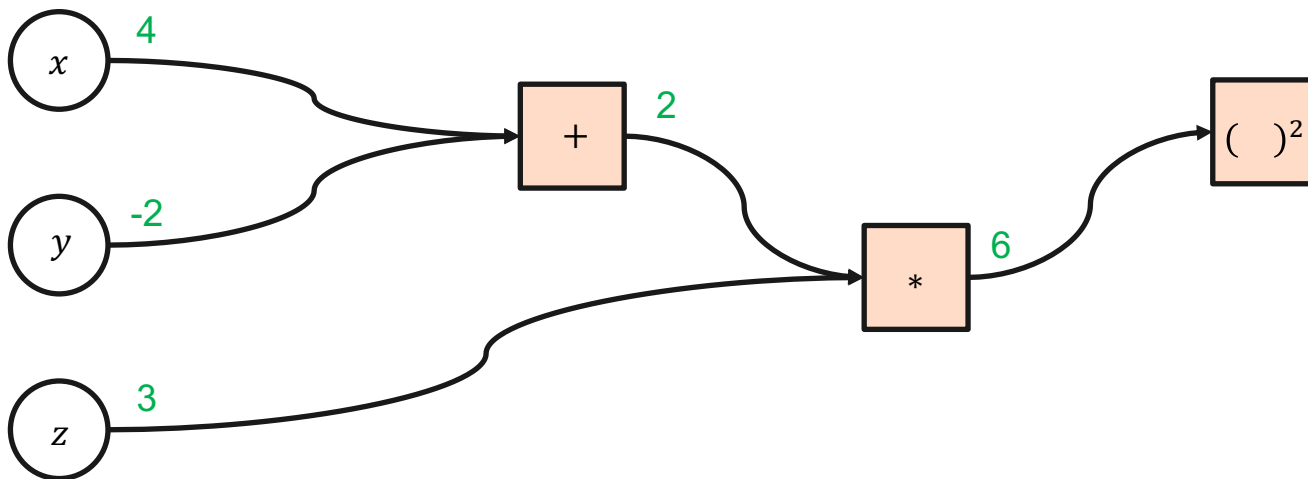
Forward Pass



# Simple Example

$$f(x, y, z) = ((x + y)z)^2$$

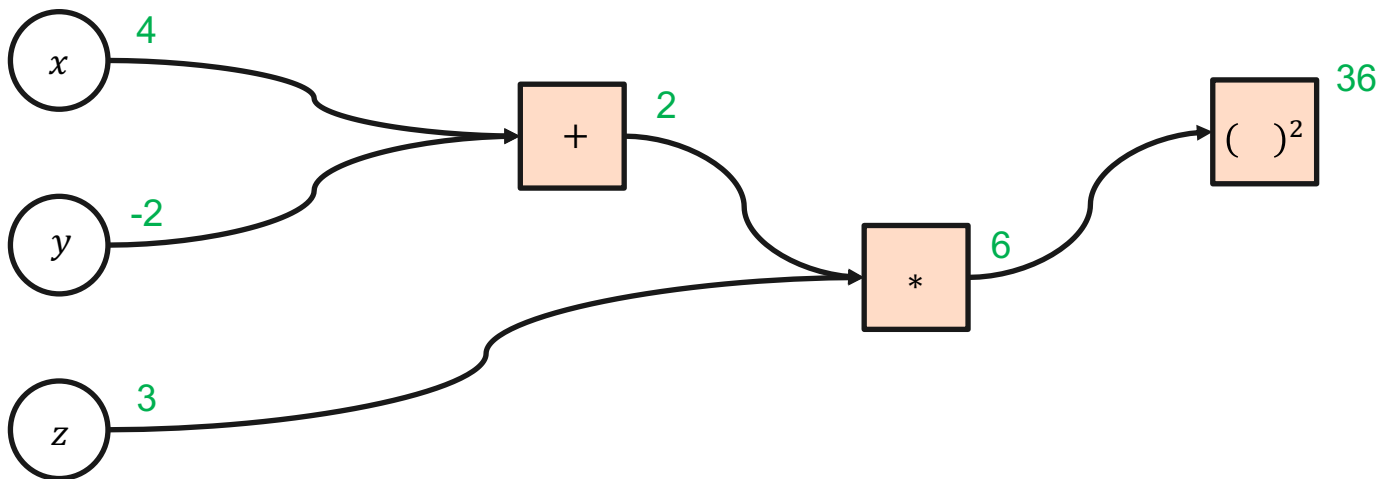
Forward Pass



$$f(x, y, z) = ((x + y)z)^2$$

## Simple Example

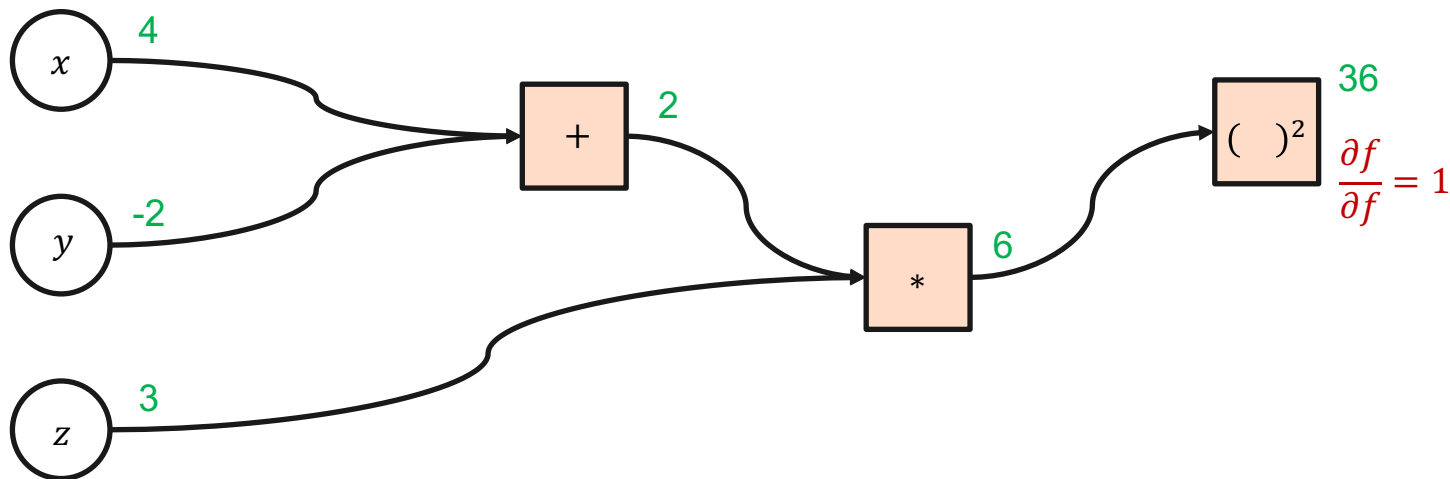
Forward Pass



$$f(x, y, z) = ((x + y)z)^2$$

## Simple Example

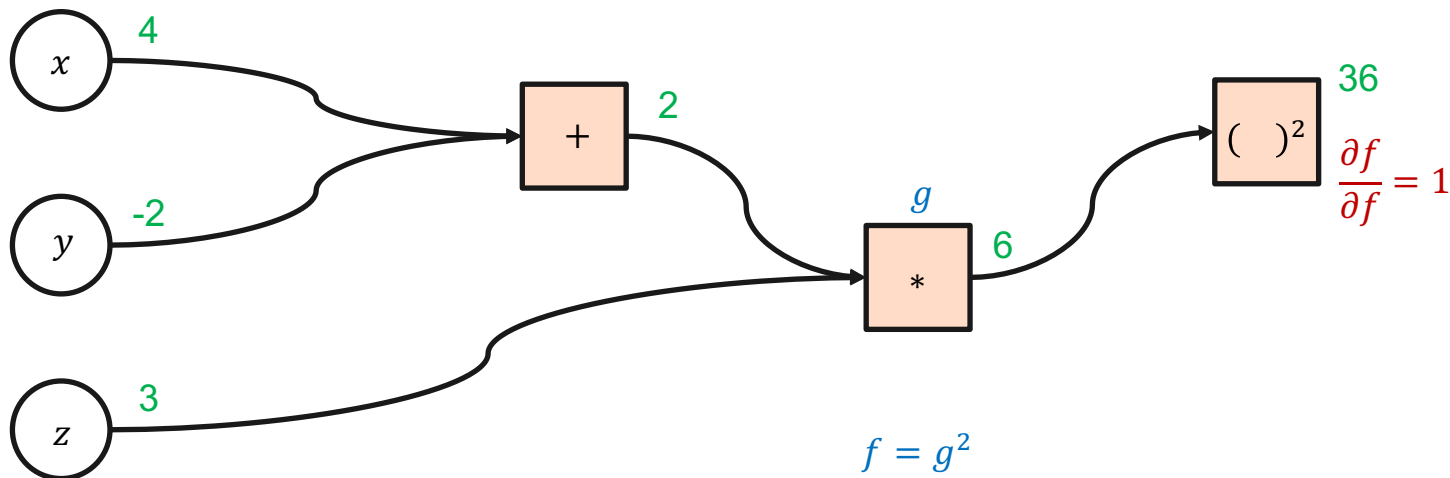
Forward Pass  
Backward Pass



$$f(x, y, z) = ((x + y)z)^2$$

## Simple Example

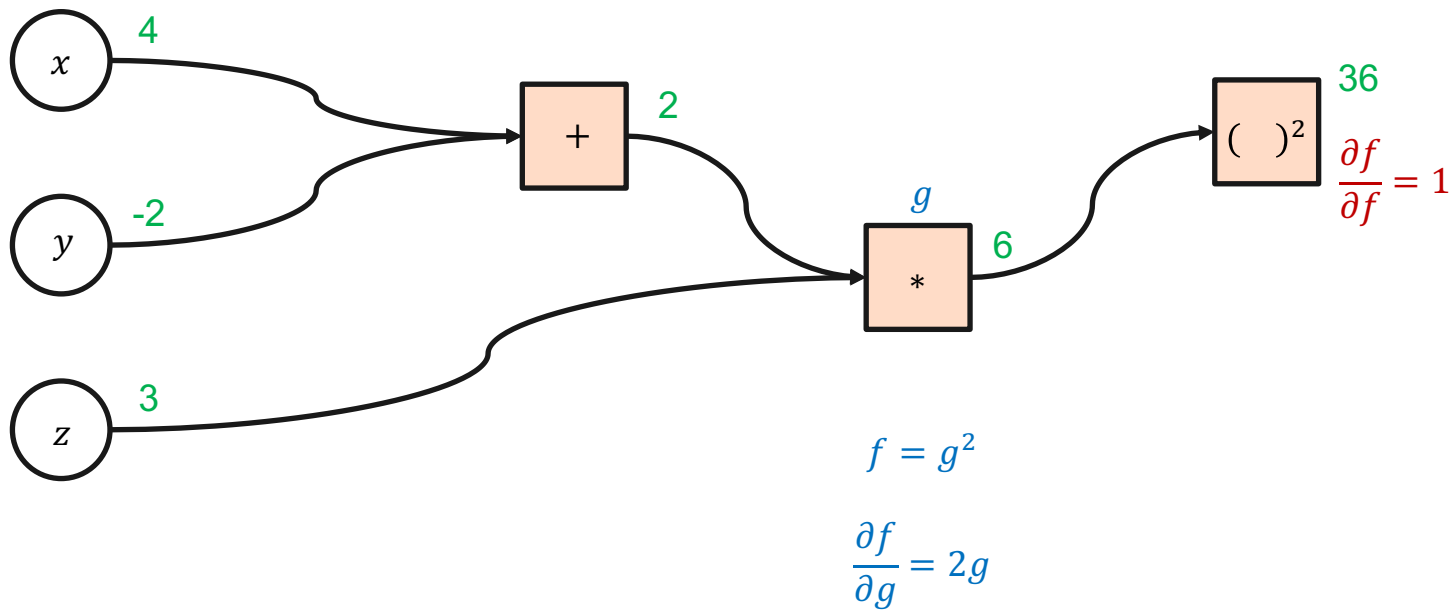
Forward Pass  
Backward Pass



# Simple Example

$$f(x, y, z) = ((x + y)z)^2$$

Forward Pass  
Backward Pass

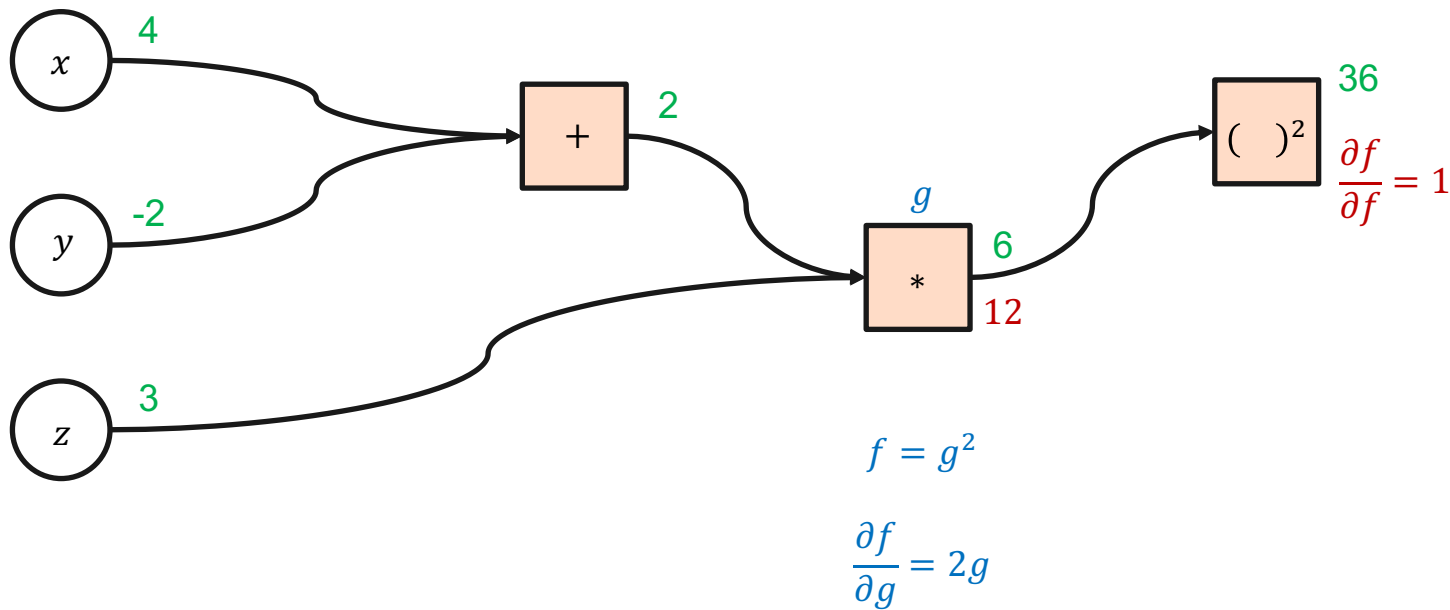




# Simple Example

$$f(x, y, z) = ((x + y)z)^2$$

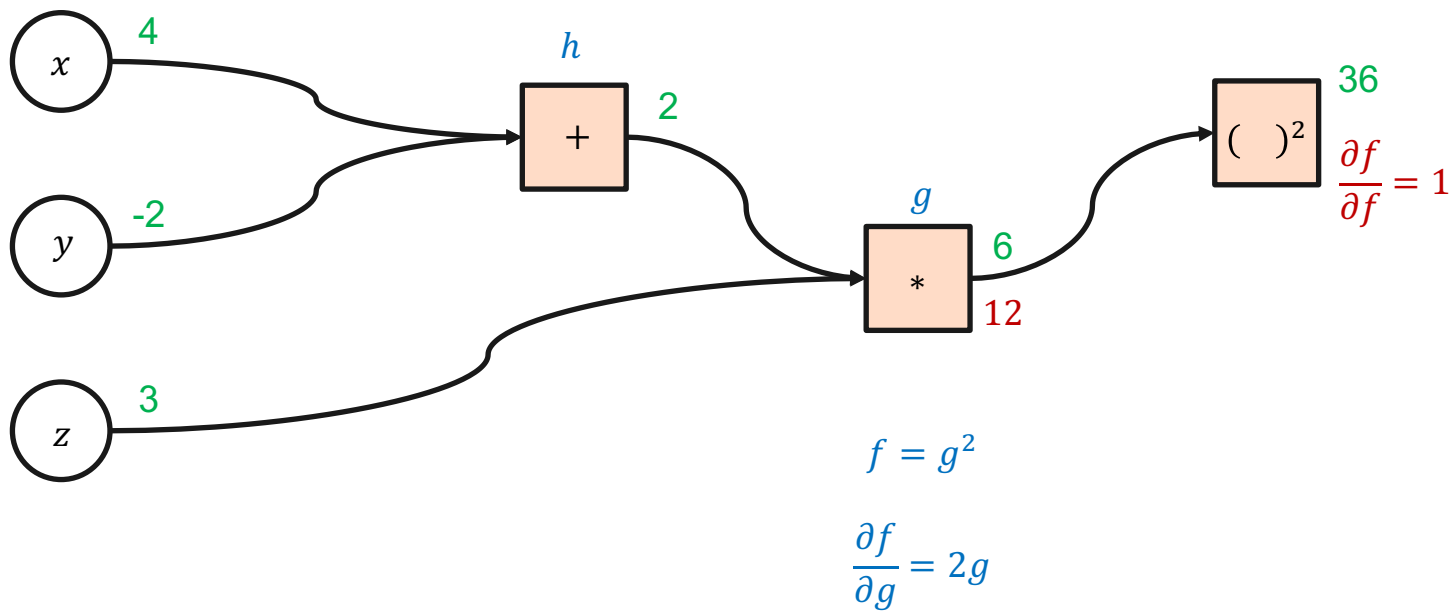
Forward Pass  
Backward Pass



# Simple Example

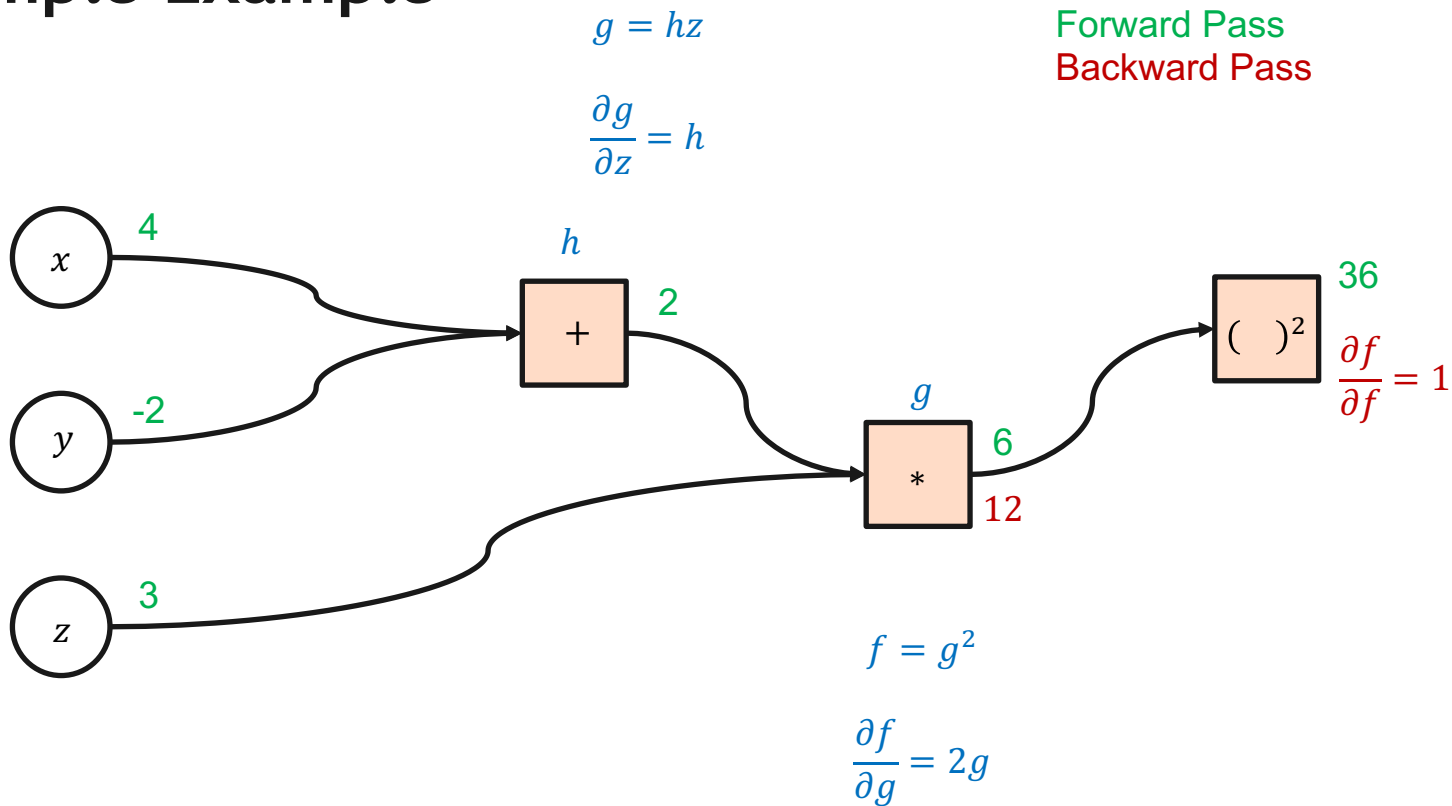
$$f(x, y, z) = ((x + y)z)^2$$

Forward Pass  
Backward Pass



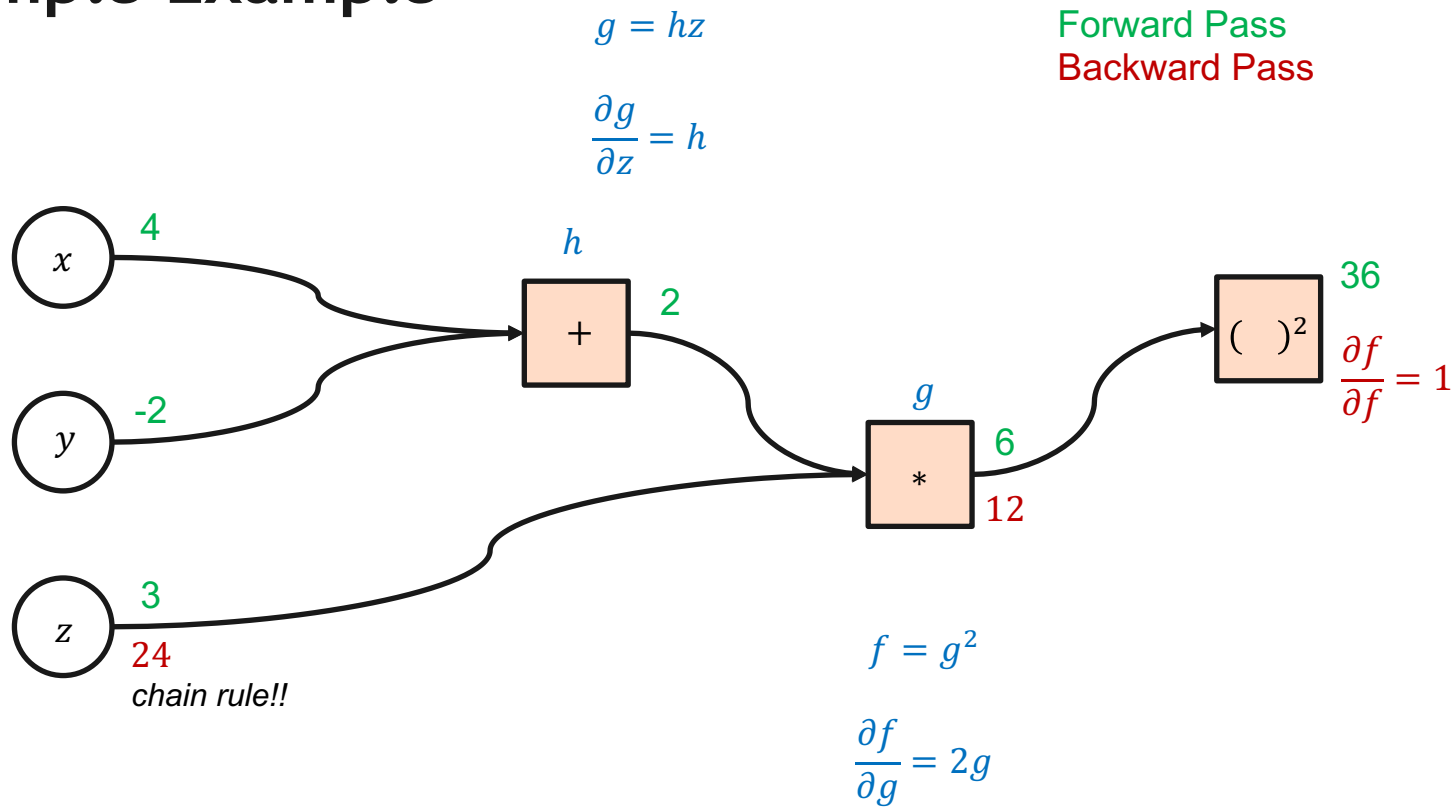
# Simple Example

$$f(x, y, z) = ((x + y)z)^2$$



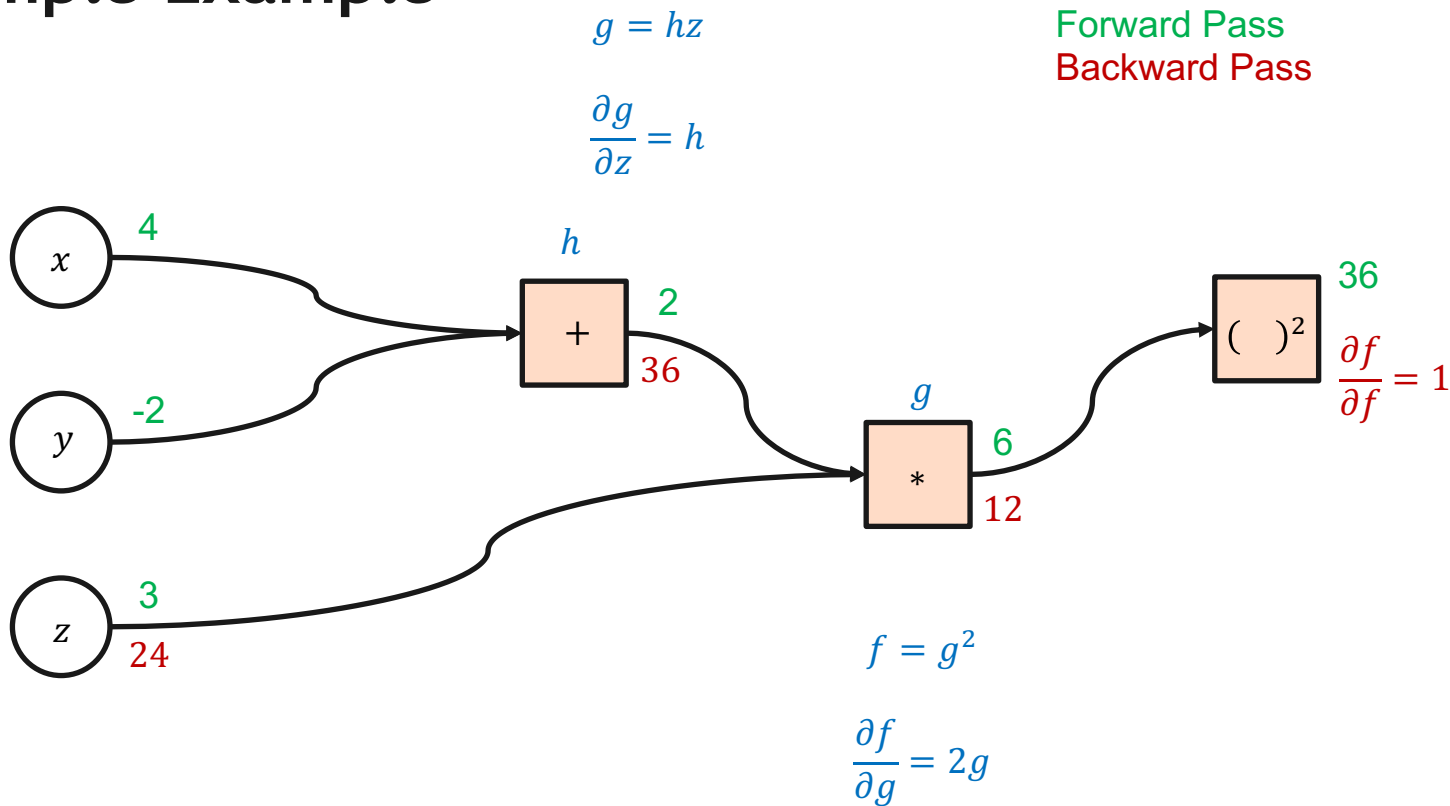
# Simple Example

$$f(x, y, z) = ((x + y)z)^2$$



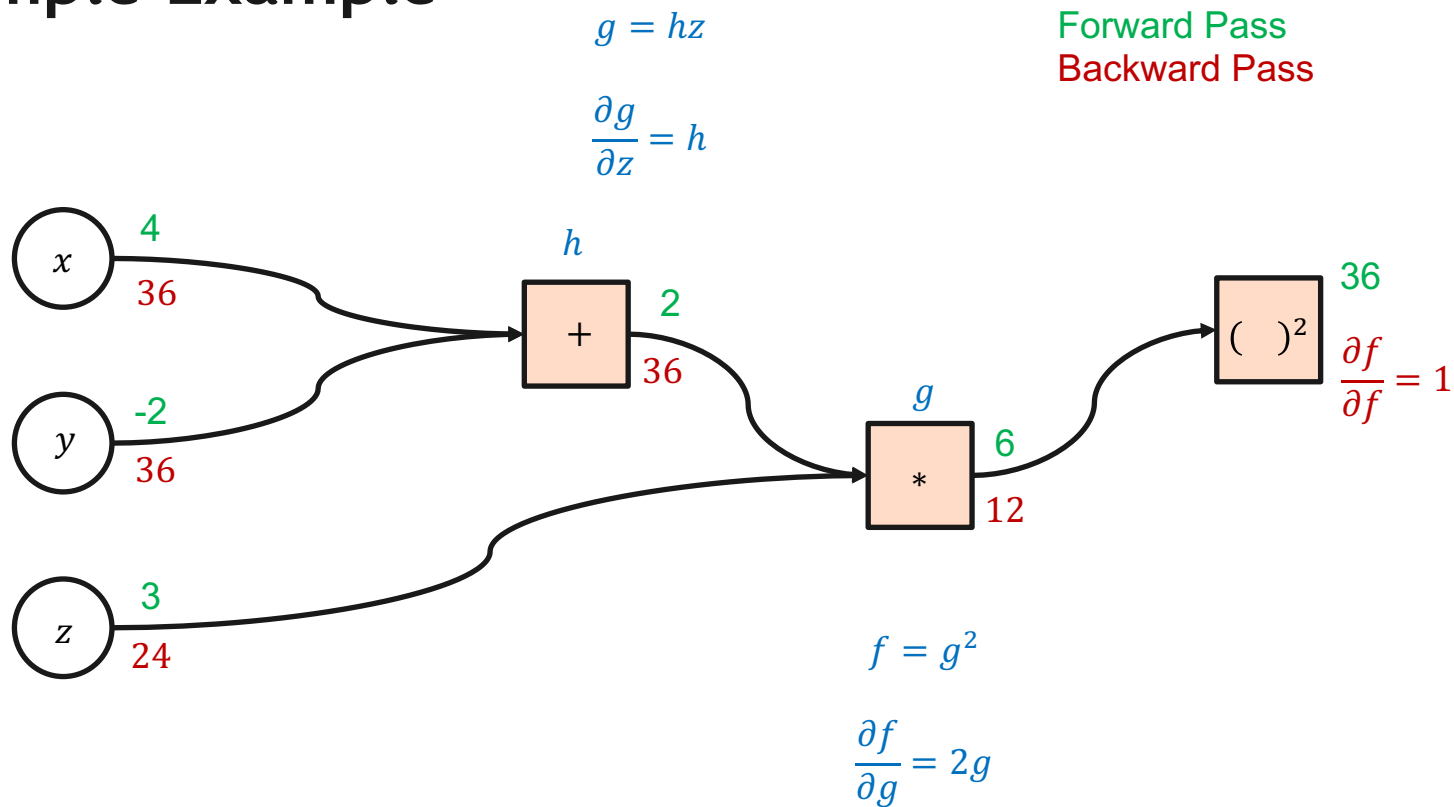
# Simple Example

$$f(x, y, z) = ((x + y)z)^2$$



# Simple Example

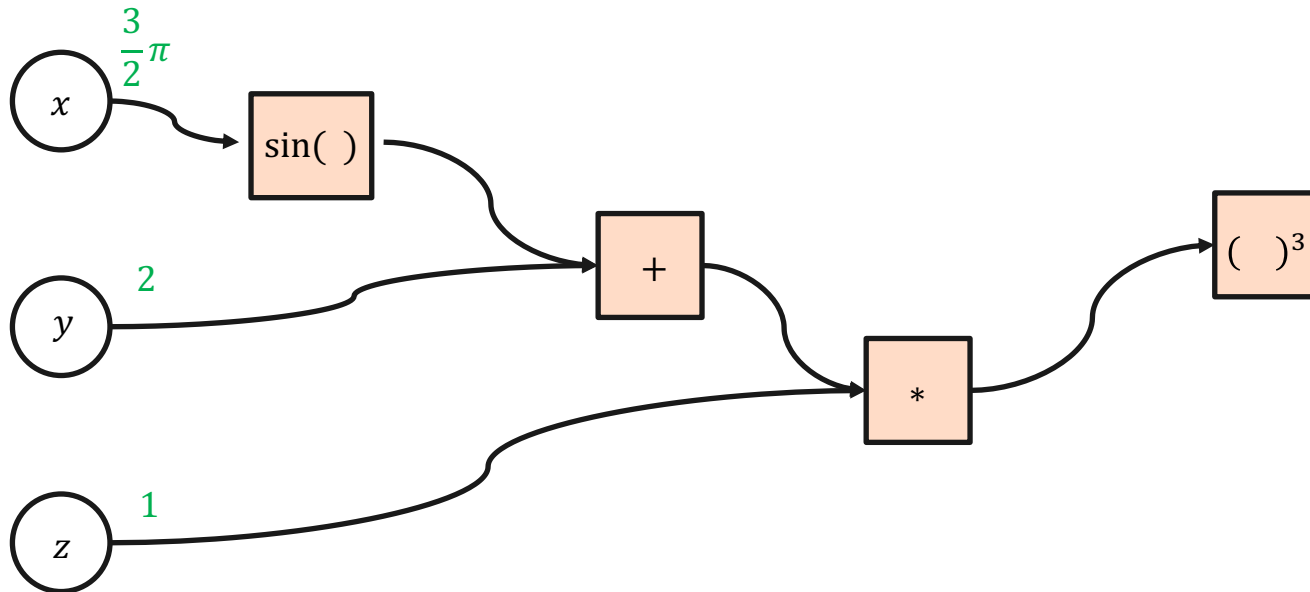
$$f(x, y, z) = ((x + y)z)^2$$



$$f(x, y, z) = ((\sin(x) + y)z)^3$$

## Exercise

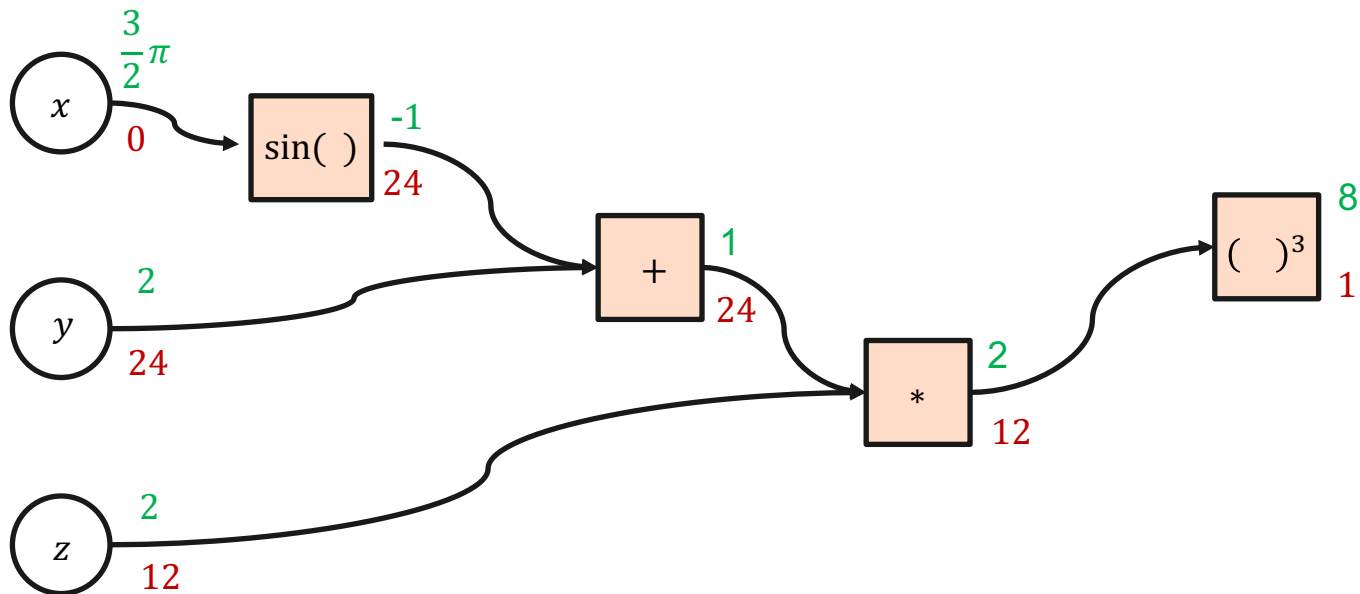
Forward Pass



$$f(x, y, z) = ((\sin(x) + y)z)^3$$

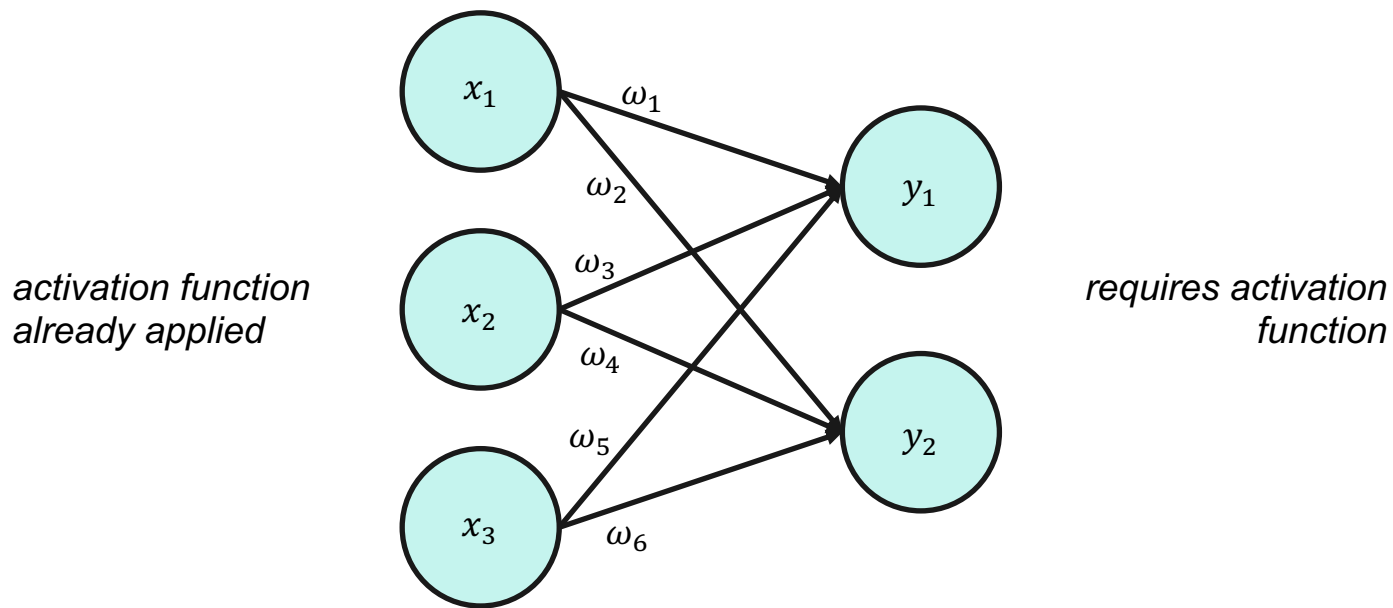
## Exercise

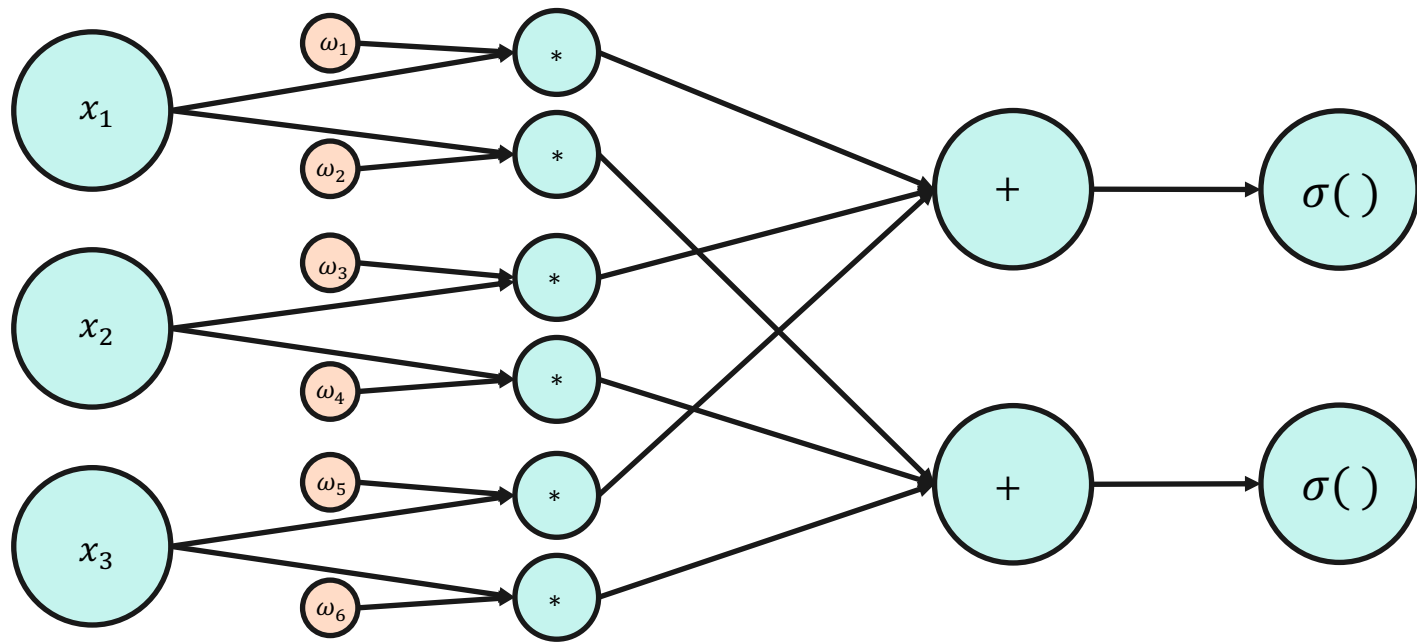
Forward Pass  
Backward Pass

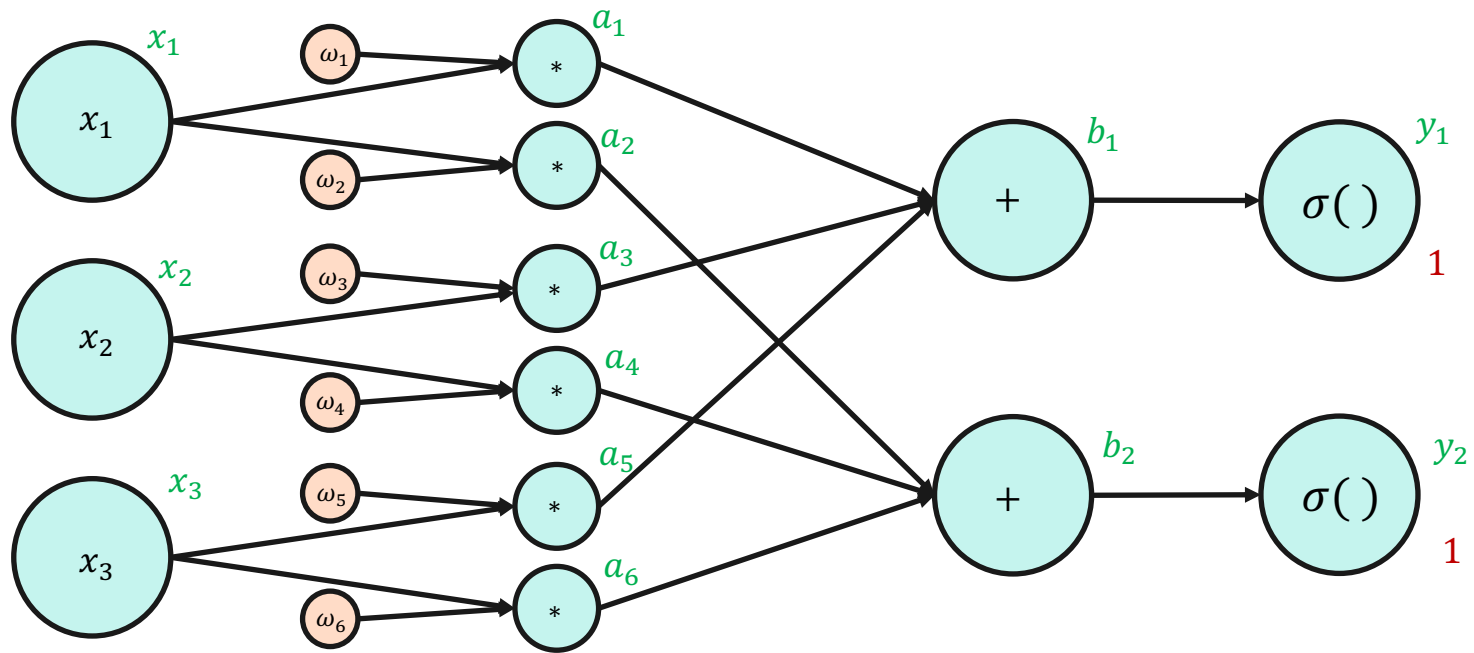




## How to do Backpropagation Through This Layer?

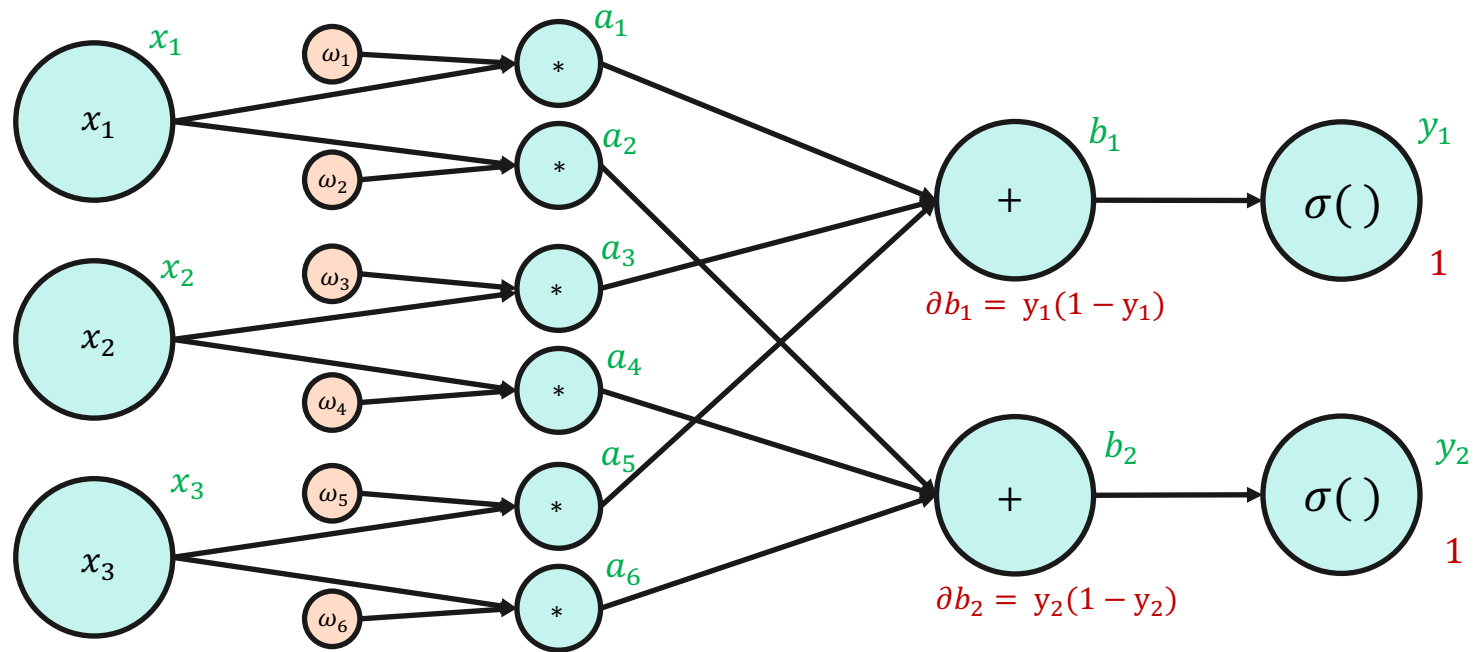






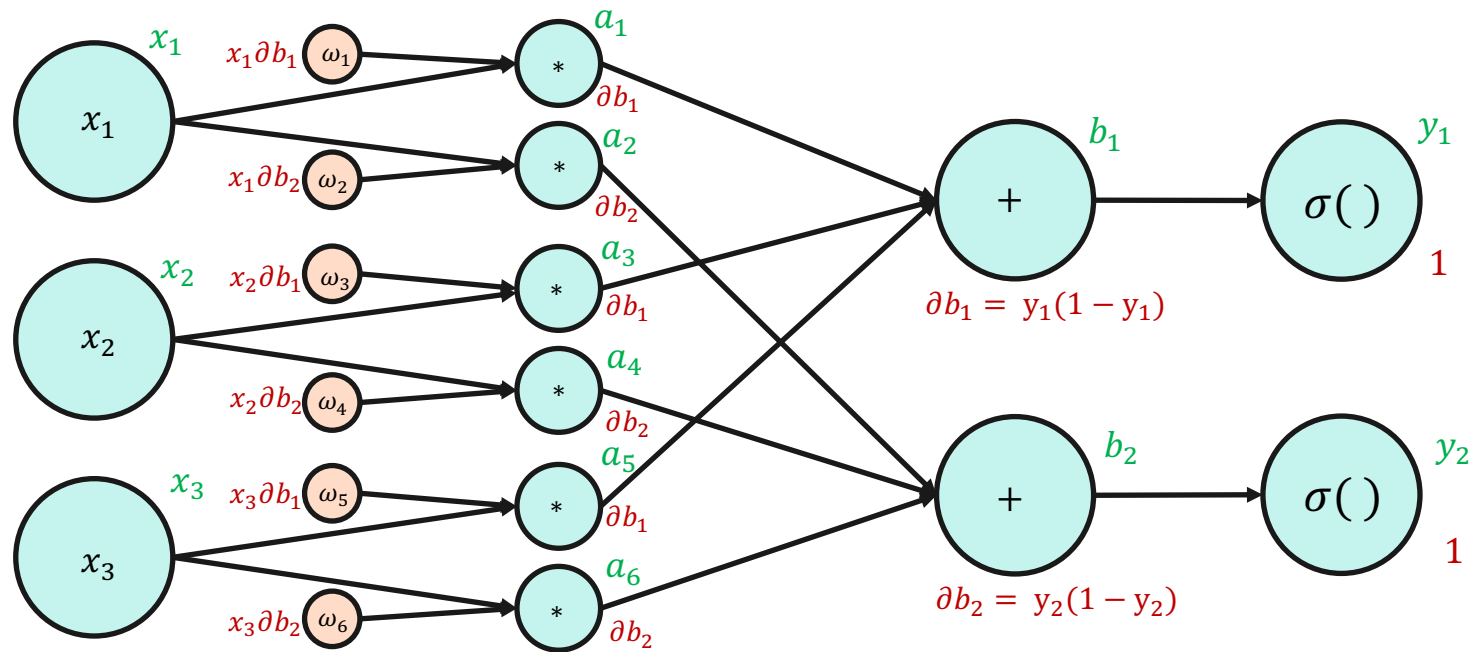
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{\partial \sigma}{\partial x} = \sigma(x)(1 - \sigma(x))$$



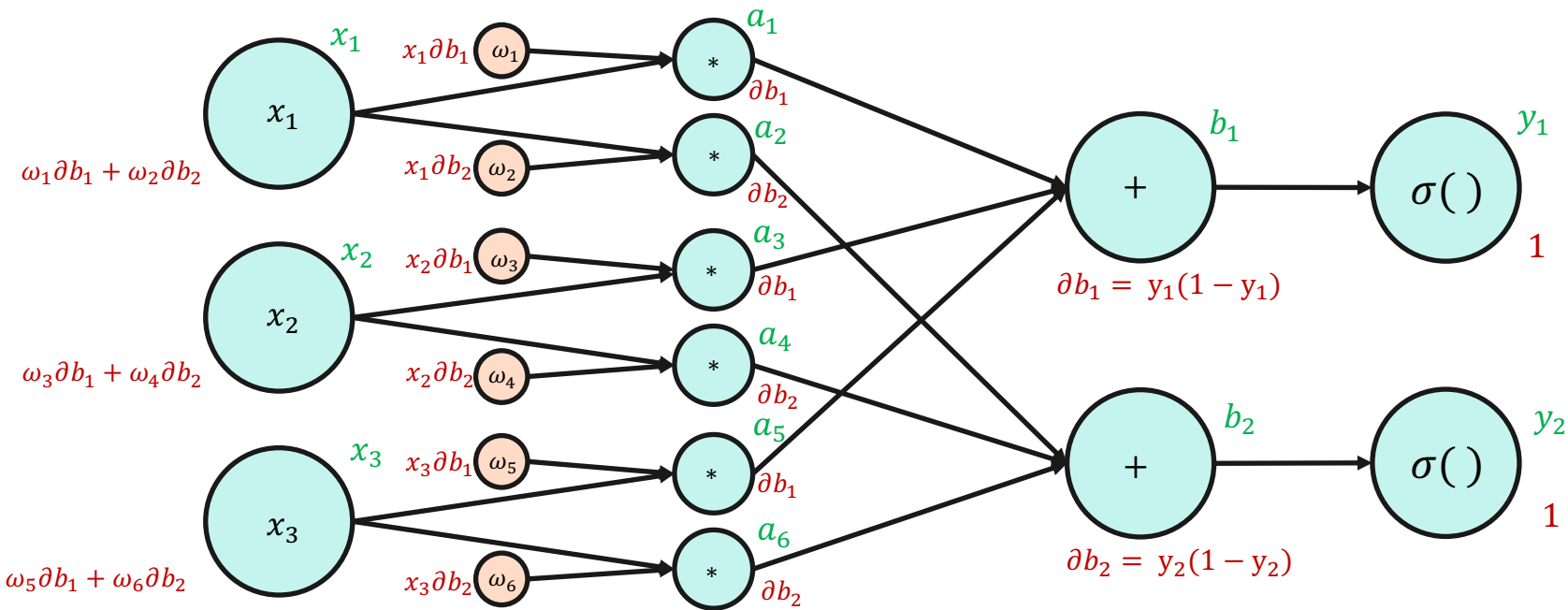
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{\partial \sigma}{\partial x} = \sigma(x)(1 - \sigma(x))$$



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

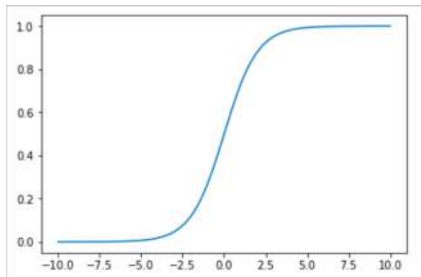
$$\frac{\partial \sigma}{\partial x} = \sigma(x)(1 - \sigma(x))$$



Everything can again be expressed using matrix multiplications.

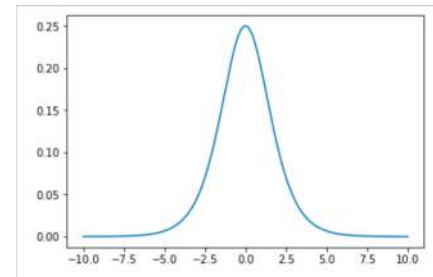
Also note that we require all  $x$  and  $y$  values from the forwards pass!

# Sigmoid Activation Function



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

$$\frac{\partial \sigma}{\partial x} = \sigma(x)(1 - \sigma(x))$$



- Remember that back-propagation always multiplies with the previous gradients
- The maximum slope of the sigmoid function is 0.25  $\Rightarrow$  very slow training for the first layers in deep network if sigmoid or similar activation functions are used.
- Gradients also very small for large absolute values
- Vanishing gradients
- Other activation functions are used for deep learning (see Thursday)

**Let's look at some code.**

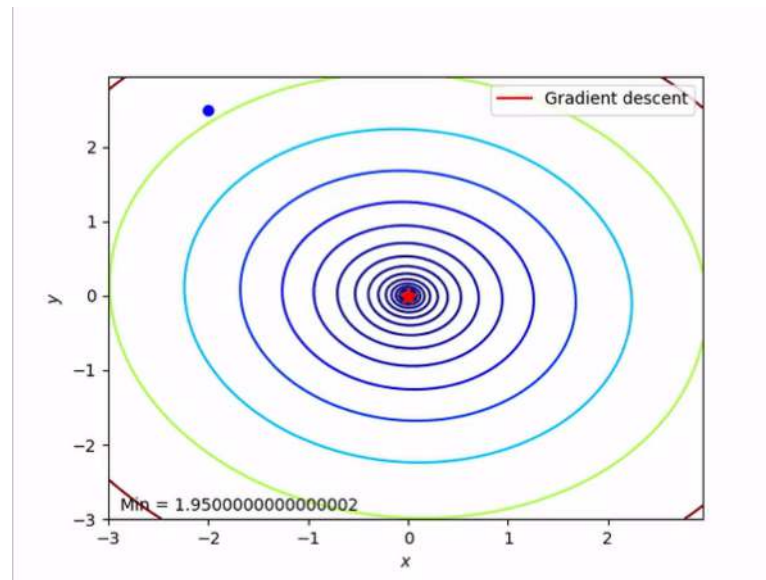
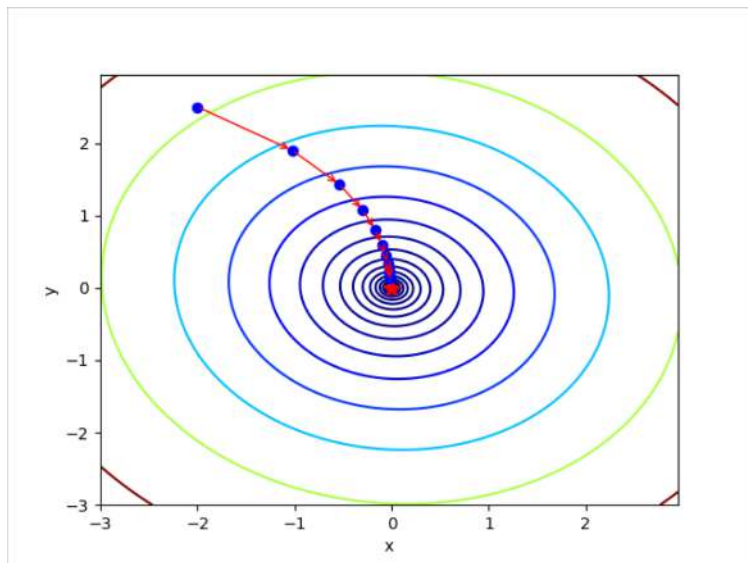




# Gradient Descent

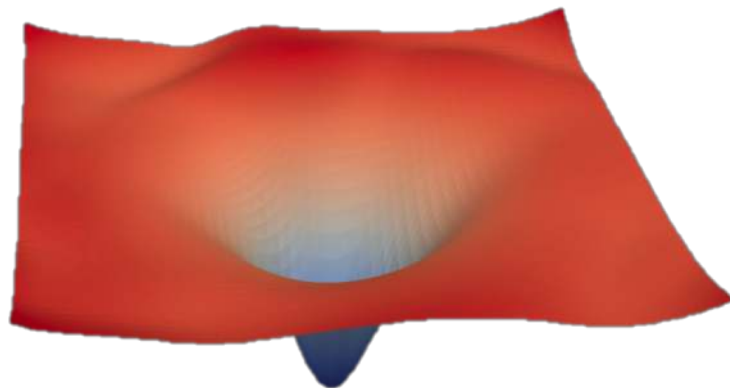
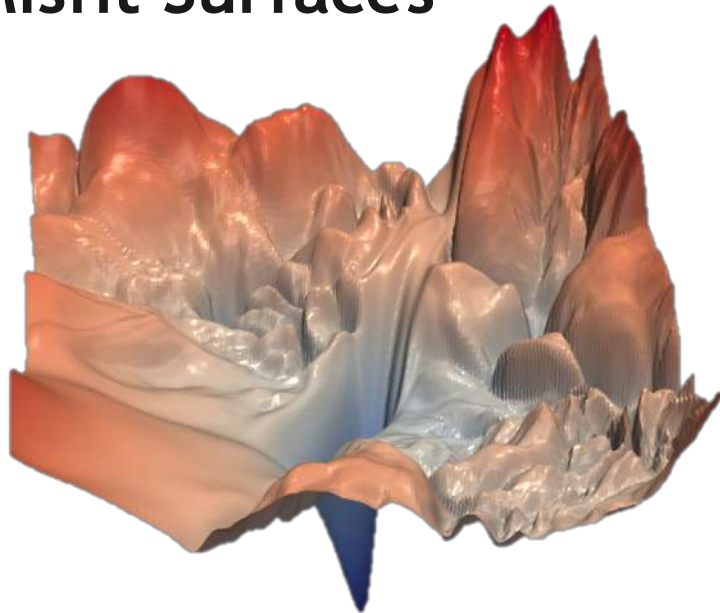
- Recall that we ruled out global optimization as there are just too many parameters in the network.
- And we have an analytic way to compute exact derivatives to all the parameters in our network.
- The most basic gradient based optimization algorithm is gradient descent:
  - Compute the gradient of the objective/loss function at the current point
  - Choose a step-length and walk in the direction of the negative gradient
  - Rinse and repeat

# Gradient Descent Illustrations





# Misfit Surfaces



## Exercise: Implement Gradient Descent



# How to Learn a Neural Network

- Design the network architecture (e.g number of inputs and outputs, number of layers, neurons in each layer, activation functions, ...).
- Loop until convergence:
  - Take a data sample, run a forward pass
  - Compute the loss of the result by comparing to the training data
  - Compute the gradient of the loss function with respect to the weights of the neural network via backpropagation
  - Update the weights with a chosen gradient based optimization algorithm

**Have another look at the code.**



# Practicalities

- Straightforward gradient descent is rarely used. People instead use some form of stochastic gradient descent
- Higher-order methods like L-BFGS are also rarely used.
- The vanishing gradient problem needs to be dealt with.
- Regularization is important and all kinds of fancy tricks are done.
- Inputs and outputs are usually normalized in some fashion. Thus to actually use a neural network the outputs then have to be “denormalized” again.



# Stochastic Gradient Descent

- Instead of computing the cumulative gradient of the whole training data set at once, compute it in so-called mini batches and update after each batch.
- This is called mini-batch gradient descent (stochastic gradient descent would be to update every time).
- Some momentum term is usually added to the optimization algorithm to deal with outliers in the gradients.
- Ends up being very fast (more efficient on GPUs to run many in parallel).
- Results in faster updates for highly redundant data.
- As the gradient is slightly different for each batch it can potentially get out of local minima again.
- State-of-the-art are adaptive mini-batched gradient descent methods like Adam, AdaGrad, RMSProb, ...

<http://ruder.io/optimizing-gradient-descent>





# Thank you.