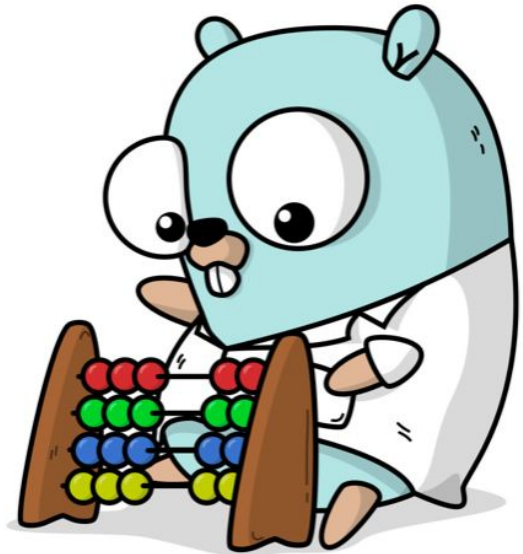# Deep Learning for Gophers

@iamrashminagpal

# Rise of Deep Learning

**Google AI Tool Identifies a Tumor's Mutations From an Image**

**Pit.ai puts a financial twist on reinforcement learning to outperform hedge funds**

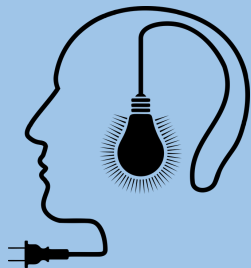Identifying artificial intelligence "blind spots"

Finding a good read among billions of choices

As natural language processing techniques improve, suggestions are getting speedier and more relevant.
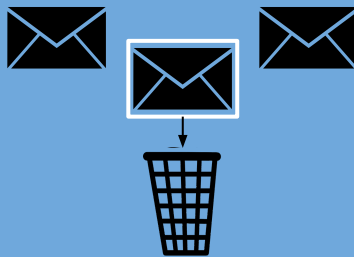
# What is Deep Learning?



**ARTIFICIAL INTELLIGENCE**

Any technique that enables computers to mimic human behavior

**MACHINE LEARNING**

Ability to learn without explicitly being programmed

**DEEP LEARNING**

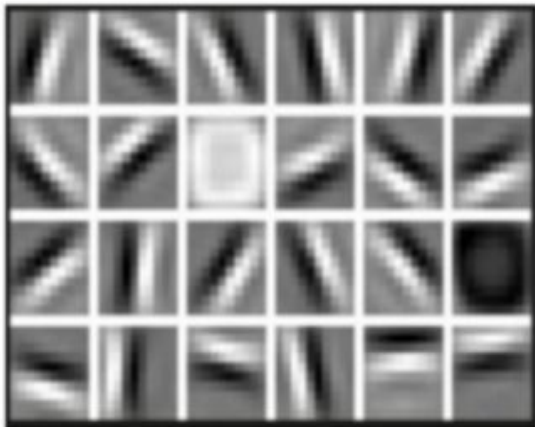Extract patterns from data using neural networks

3 1 3 5 6 7

1 4 5 9 2 3

# Why Deep Learning?

Hand engineered features are time consuming, brittle & not scalable in practise

Can we learn **underlying features** directly from data?
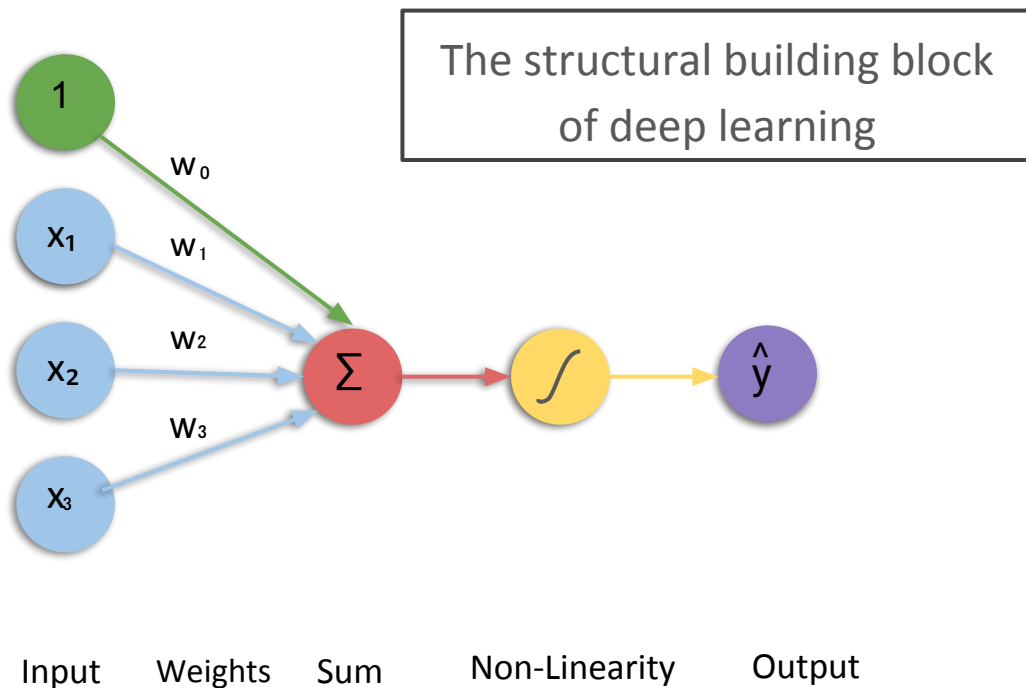
| Low Level Features | Mid Level Features | High Level Features |
| --- | --- | --- |

# The Perceptron : Feedforward Propagation

The structural building block of deep learning



Input    Weights    Sum    Non-Linearity    Output

Output

Linear combinations of input

$$\hat{y} = g\left(w_0 + \sum_{i=1}^{m} x_i\, w_i\right)$$

Non-linear activation function
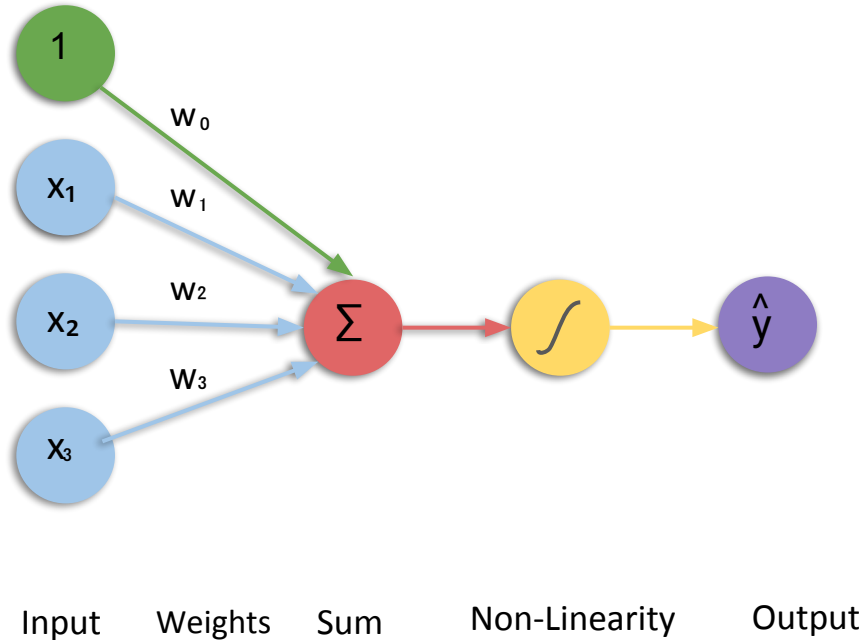
Bias

```go
package deep_learning_for_gophers

type Neuron struct {
    Act         ActivationFunc
    Input       []*Edge
    Output      []*Edge
    Value       float64
}

type Edge struct {
    Weight      float64
    Input       float64
    Output      float64
    IsBias      bool
}
```

# The Perceptron : Feedforward Propagation



Input    Weights    Sum    Non-Linearity    Output

$$\hat{y} = g\left( w_0 + \sum_{i=1}^{m} x_i\, w_i \right)$$

$$\hat{y} = g\left( w_0 + X^T W \right)$$

where: $X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $W = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$
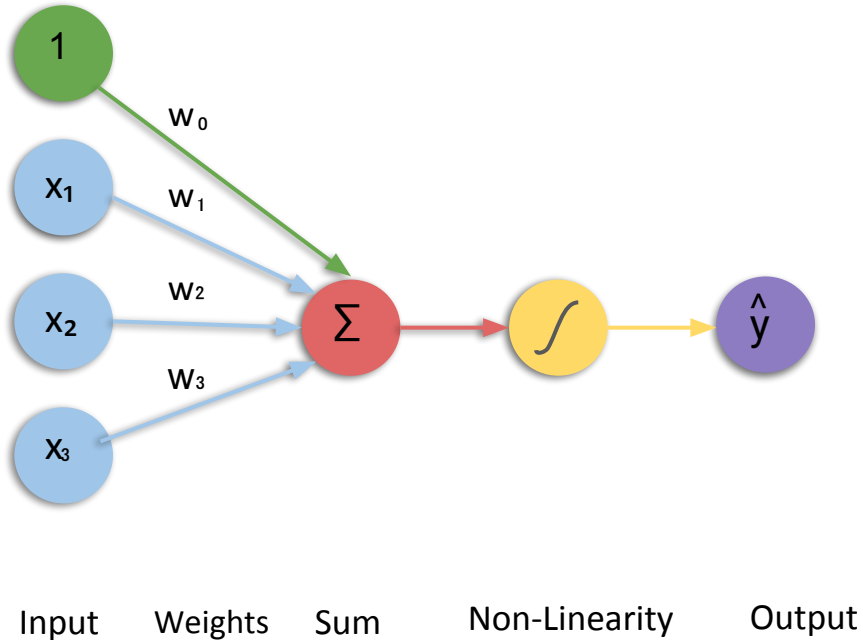
```go
package deep_learning_for_gophers

import ...

type initializeWeight func() float64

func NormalWeightInitialize(stdDev, mean float64) float64 {
    return rand.NormFloat64()*stdDev + mean
}

func UniformWeightInitialize(stdDev, mean float64) float64 {
    return (rand.Float64()-0.5)*stdDev + mean
}
```

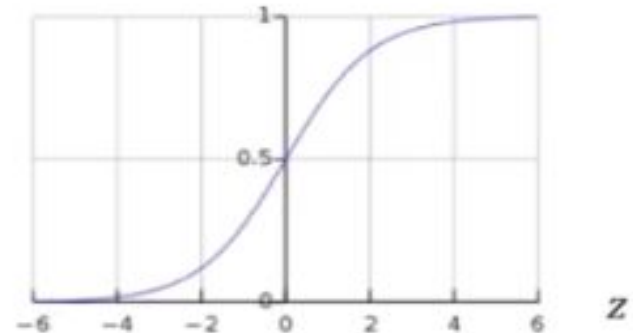# The Perceptron : Feedforward Propagation



Input    Weights    Sum    Non-Linearity    Output

**Activation Function**

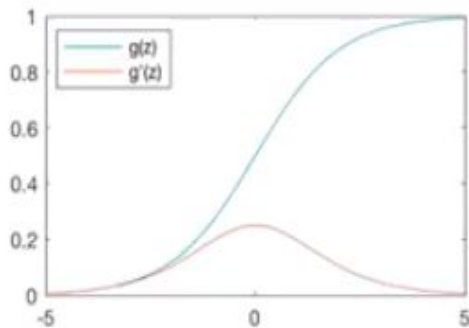$$\hat{y} = g( w_0 + X^T W )$$

- Example : Sigmoid function

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

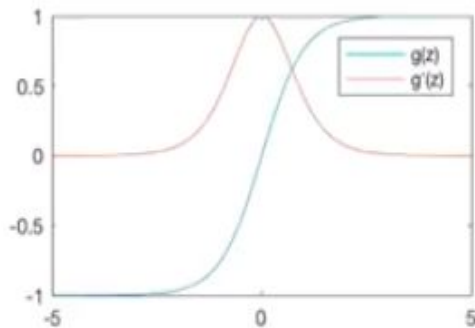# Common Activation Functions

### Sigmoid Function



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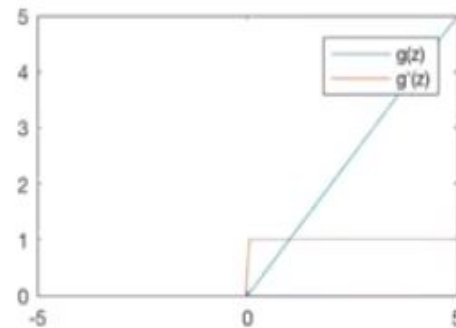
$$g'(z) = g(z)(1 - g(z))$$

### Hyperbolic Tangent



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

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

### Rectified Linear Unit (ReLU)



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

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

```go
package deep_learning_for_gophers

import (
    "math"
    _ "math"
)
type ActivationFunc int

const(
    NoActivation ActivationFunc      = 0
    SigmoidActivation ActivationFunc = 1
    TanhActivation ActivationFunc    = 3
    ReLuActivation ActivationFunc    = 4
    LinearActivation ActivationFunc  = 5
    SoftMaxActivation ActivationFunc = 6
)

func GetActivationFunc(act ActivationFunc) Differentiable{
    switch act {
    case SigmoidActivation:
        return Sigmoid{}
    case ReLuActivation:
        return ReLU{}
    case LinearActivation:
```

```go
func (s Sigmoid) Func(x float64) float64 {
    return Logistic(x, a: 1)
}

func (s Sigmoid) DFunc(y float64) float64{
    return y*(1-y)
}

type ReLU struct{}

func (a ReLU) Func(x float64) float64 { return math.Max(x, y: 0) }

func (a ReLU) DFunc(y float64) float64 {
    if y > 0 : 1 ↗
    return 0
}

type linear struct{}

func (l linear) Func(x float64) float64 {return x}

func (l linear) DFunc(x float64) float64 {return 1}
```
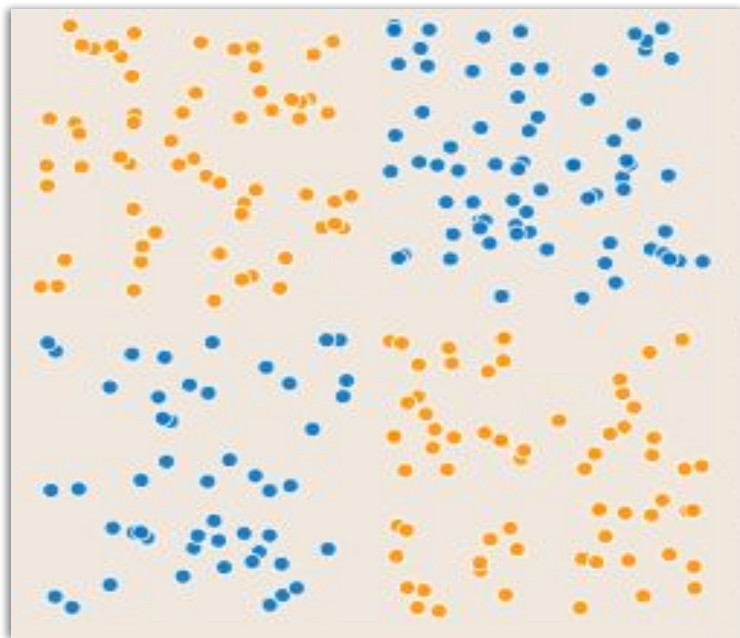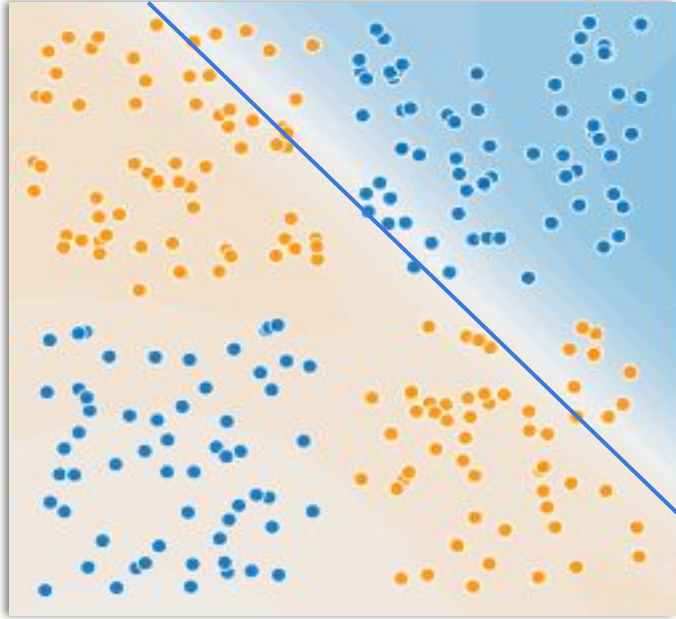
# Importance of Activation Functions

The purpose of activation functions is to **introduce non-linearities** into the network
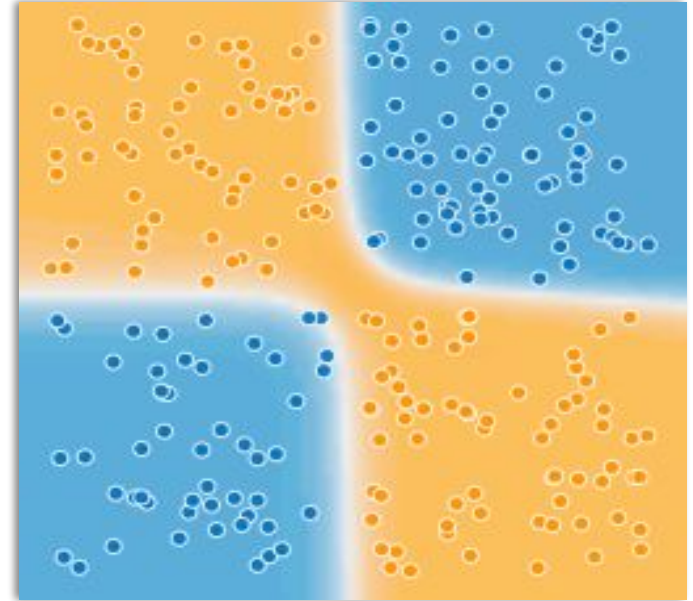


What is we want to build a Neural Network to distinguish between
blue vs orange points?
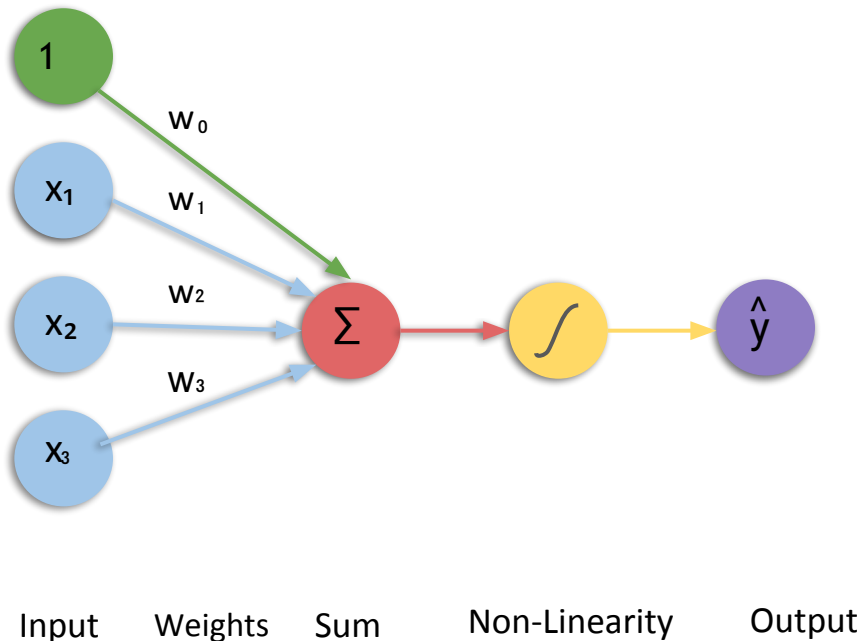
# Importance of Activation Functions



Linear activations produce linear decisions
no matter what the network size

Non-linearities allow us to approximate
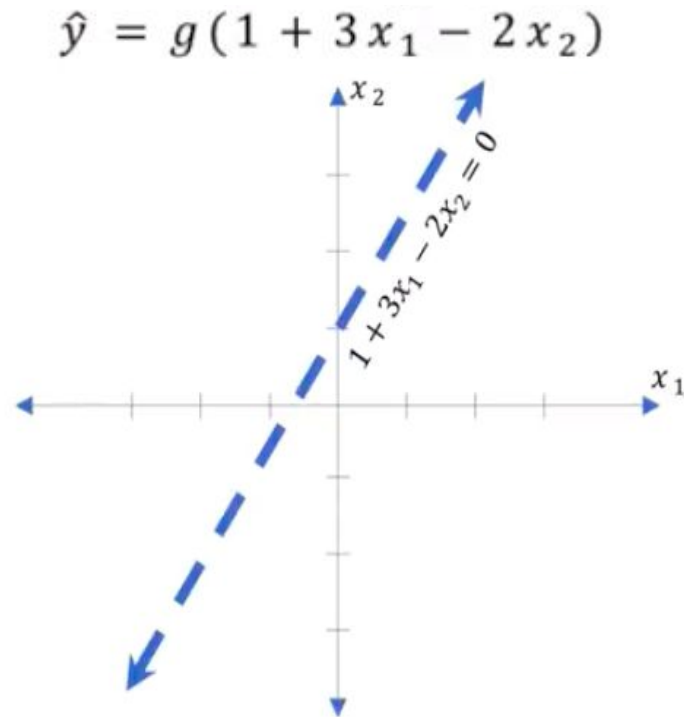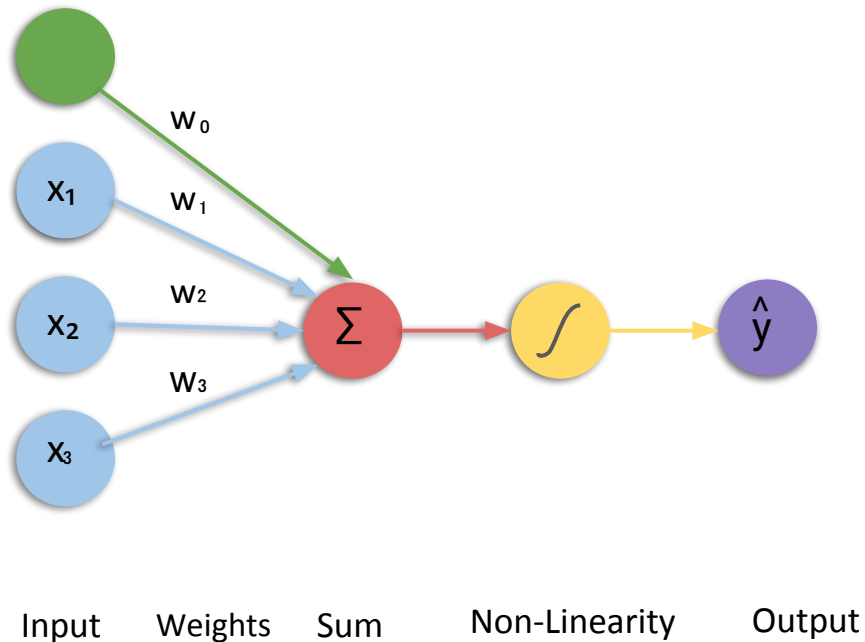arbitrarily complex functions

# The Perceptron : Example



Input     Weights     Sum     Non-Linearity     Output

We have: $w_0 = 1$ and $W = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

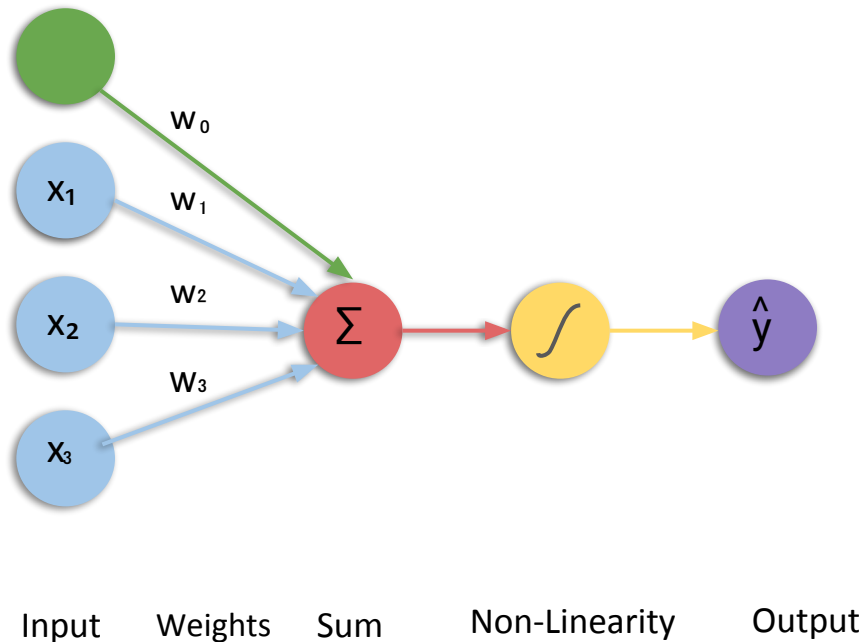$$\hat{y} = g(w_0 + X^T W)$$
$$= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right)$$
$$\hat{y} = g(1 + 3x_1 - 2x_2)$$
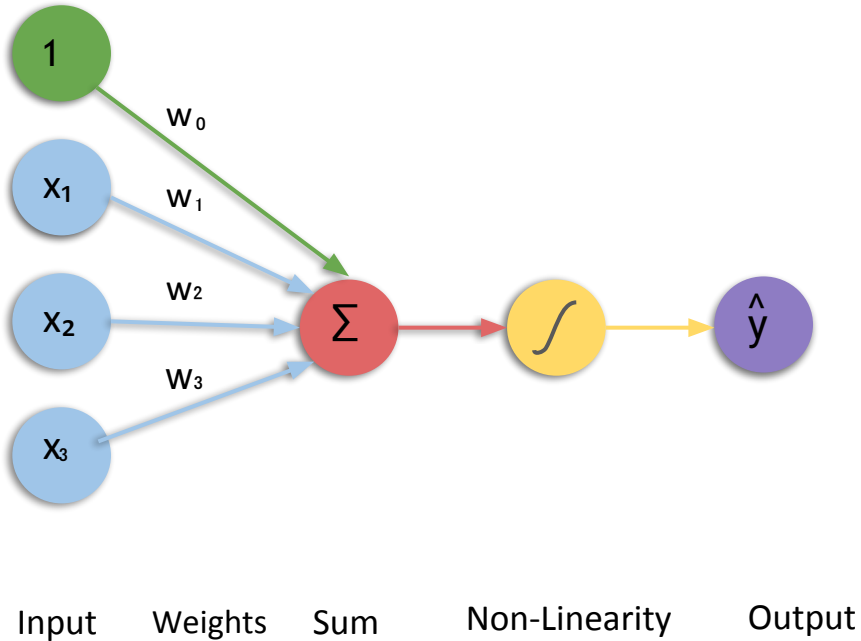
This is just a line in 2D!

# The Perceptron : Example



Input    Weights    Sum    Non-Linearity    Output

$$\hat{y} = g(1 + 3x_1 - 2x_2)$$

$1 + 3x_1 - 2x_2 = 0$

# The Perceptron : Example



Input   Weights   Sum   Non-Linearity   Output

$$\hat{y} = g(1 + 3x_1 - 2x_2)$$

$z < 0$
$y < 0.5$

$1 + 3x_1 - 2x_2 = 0$

$z > 0$
$y > 0.5$

# Building Neural Networks with Perceptrons

# The Perceptron : Simplified



Input    Weights    Sum    Non-Linearity    Output

# The Perceptron : Simplified



$$z = w_0 + \sum_{j=1}^{m} x_j\, w_j$$

# Single Layer Neural Network



W$^1$    W$^2$

$x_1$   $x_2$   $x_m$   $z_1$   $z_2$   $z_3$   $z_d$   $\hat{y}_1$   $\hat{y}_2$
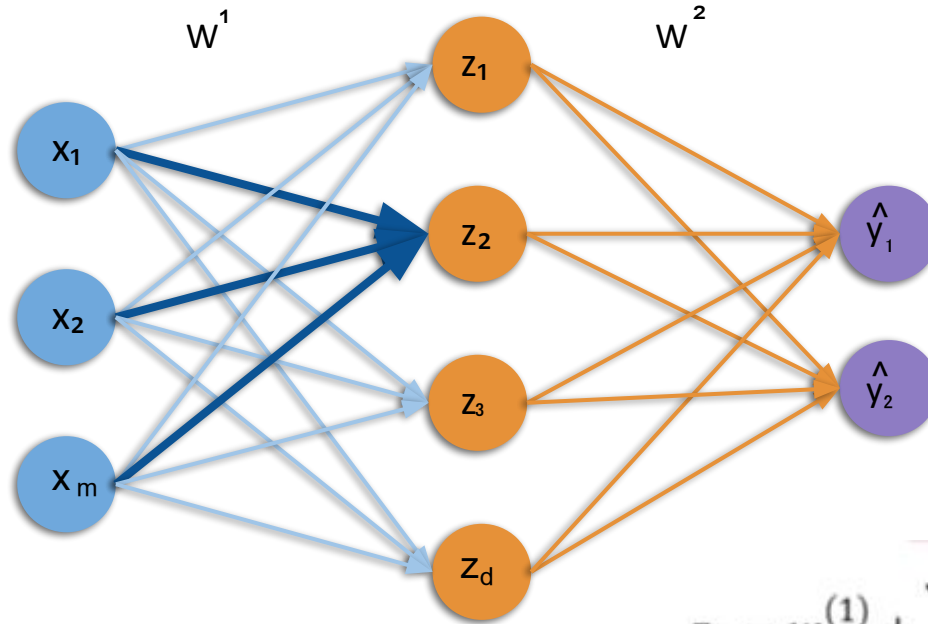
Inputs                Hidden Layer                Final Output

$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^{m} x_j \, w_{j,i}^{(1)} \qquad \hat{y}_i = g\left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j \, w_{j,i}^{(2)}\right)$$

# Single Layer Neural Network



$$z_2 = w_{0,2}^{(1)} + \sum_{j=1}^{m} x_j \, w_{j,2}^{(1)}$$

$$= w_{0,2}^{(1)} + x_1 \, w_{1,2}^{(1)} + x_2 \, w_{2,2}^{(1)} + x_m \, w_{m,2}^{(1)}$$
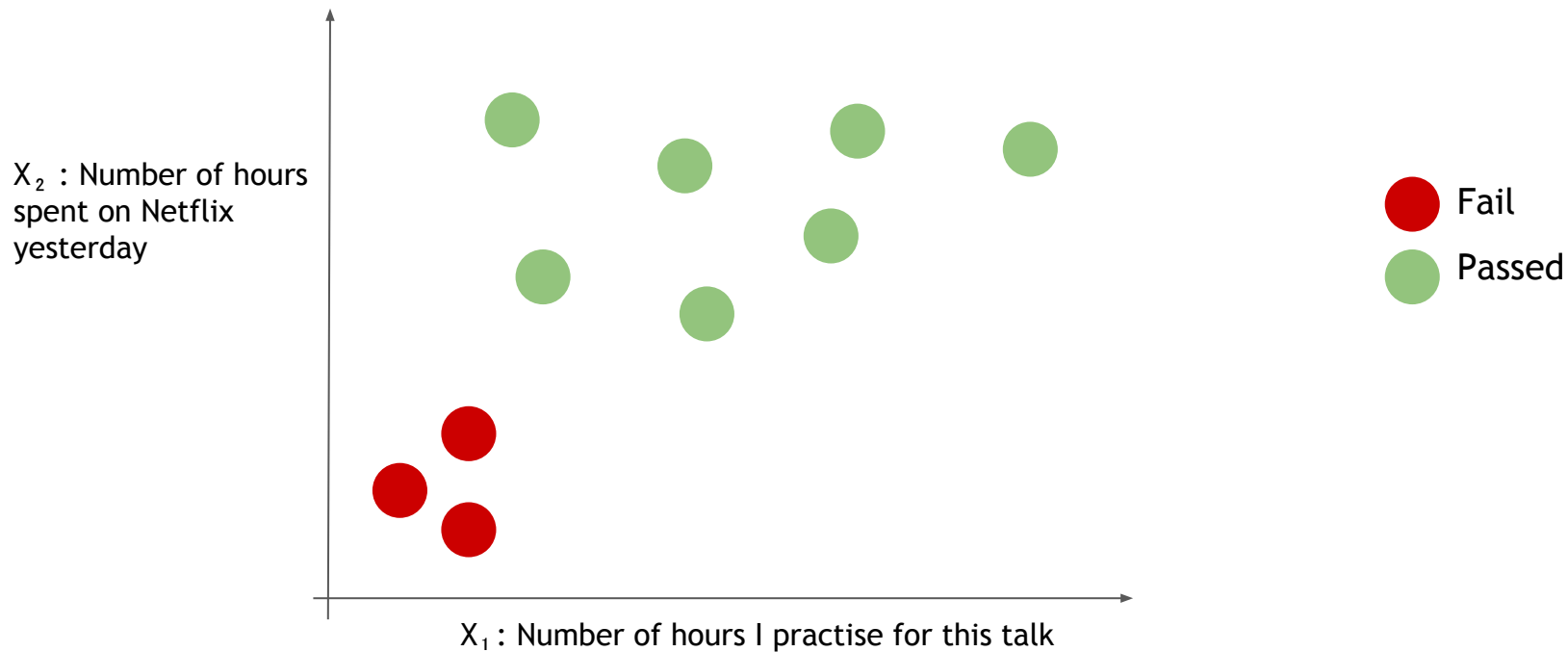
# Applying Neural Networks

# Example Problem

Will I be able to do justice to this topic?

Let's start with simple two feature model

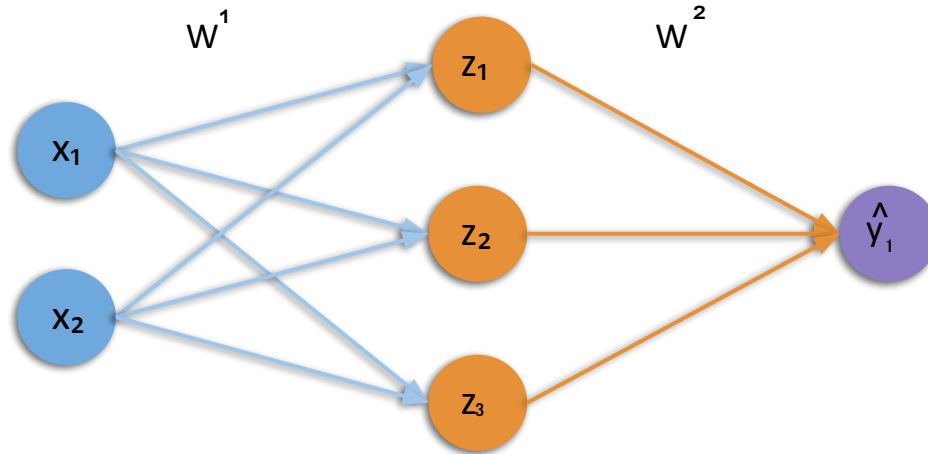X1 : Number of hours I practise for this talk
X2 : Number of hours spent on Netflix yesterday
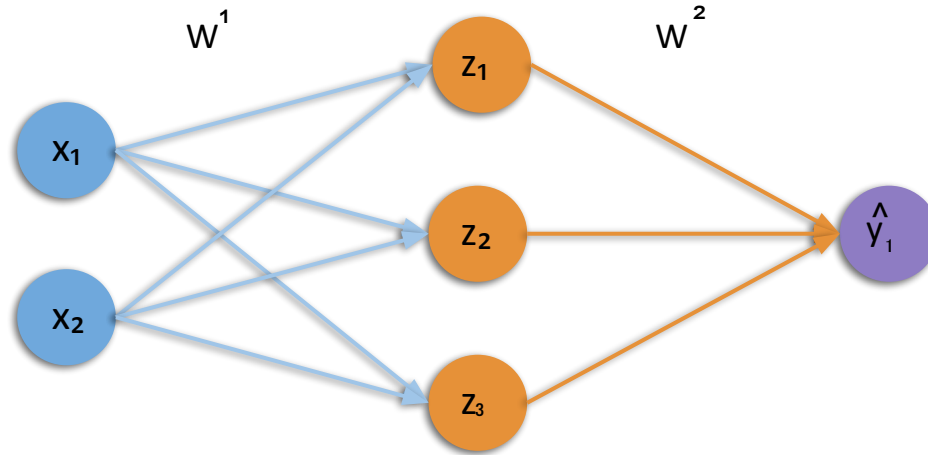
# Will I be able to do justice to this topic?



$X_2$ : Number of hours spent on Netflix yesterday

$X_1$ : Number of hours I practise for this talk

Fail

Passed

# Single Layer Neural Network



$W^1$

$W^2$

$x_1$

$x_2$

$z_1$

$z_2$

$z_3$

$\hat{y}_1$

Predicted : 0.1
Actual : 1

# Quantifying Loss



$W^1$

$W^2$

$x_1$

$x_2$

$z_1$

$z_2$

$z_3$

$\hat{y}_1$

Predicted : 0.1
Actual : 1

# Training Neural Networks

# Loss Optimization

We want to find the network weights that **achieve the lowest loss**

$$W^* = \underset{W}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}\left(f\left(x^{(i)}; W\right), y^{(i)}\right)$$

$$W^* = \underset{W}{\operatorname{argmin}} J(W)$$
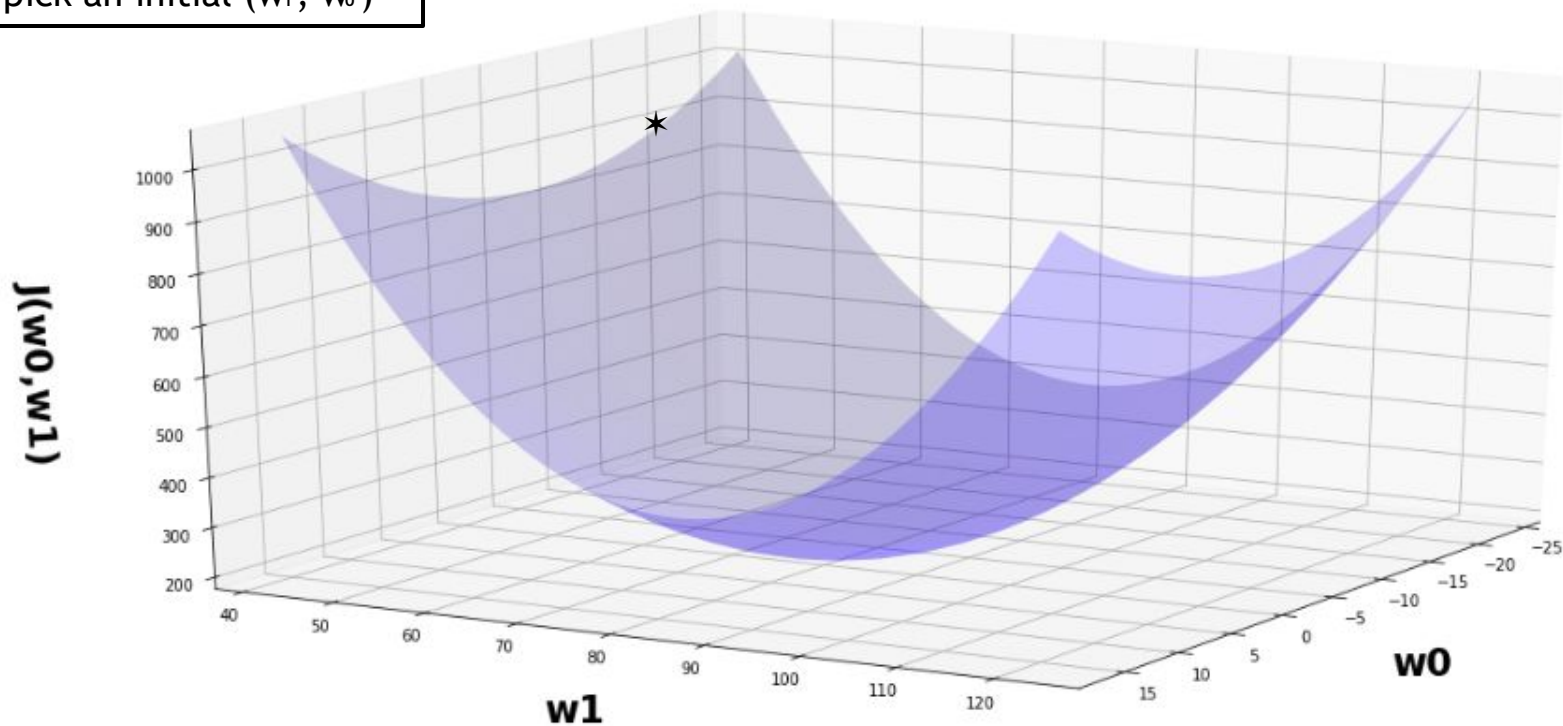
Remember:
$$W = \left\{W^{(0)}, W^{(1)}, \cdots\right\}$$

# Loss Optimization
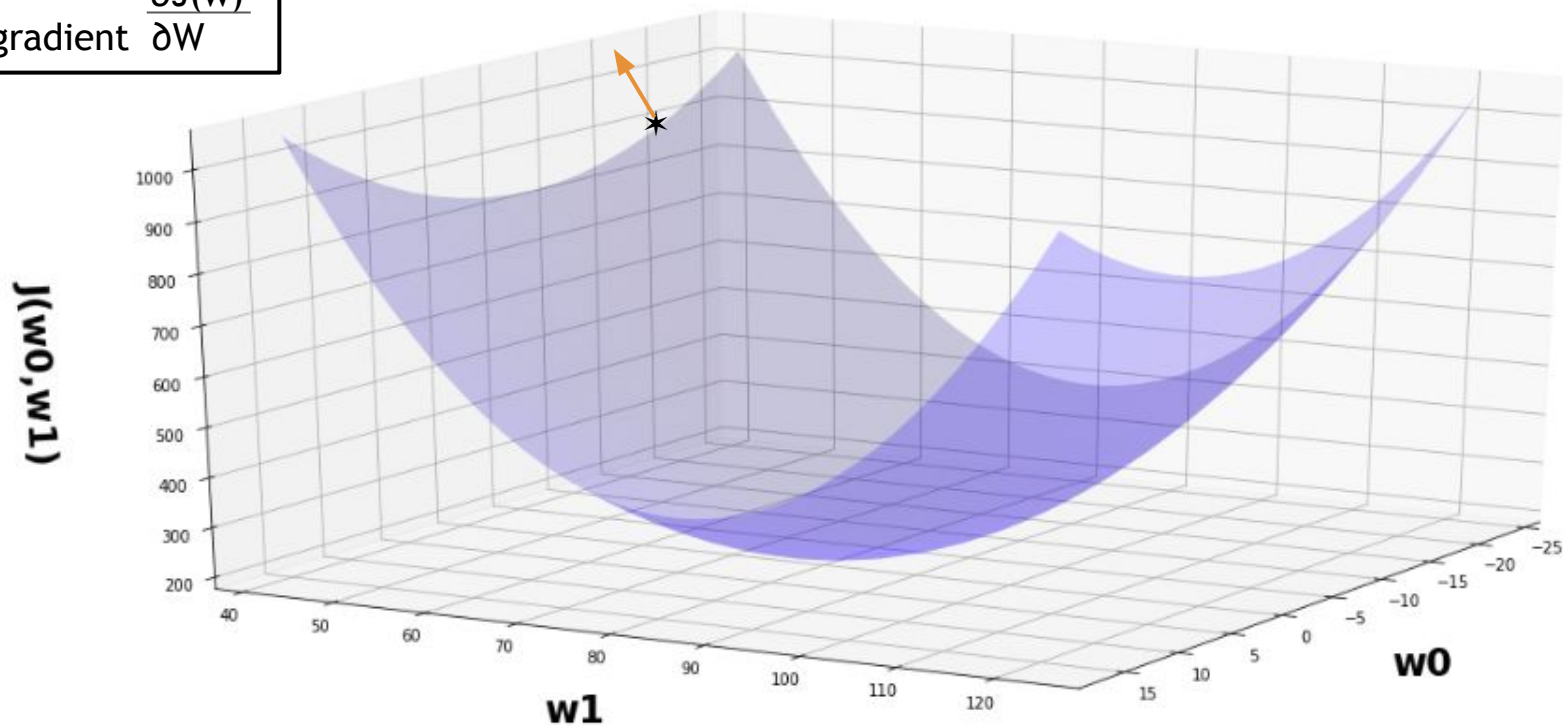
$$W^* = \underset{W}{\arg\min}\ J(w)$$

# Loss Optimization
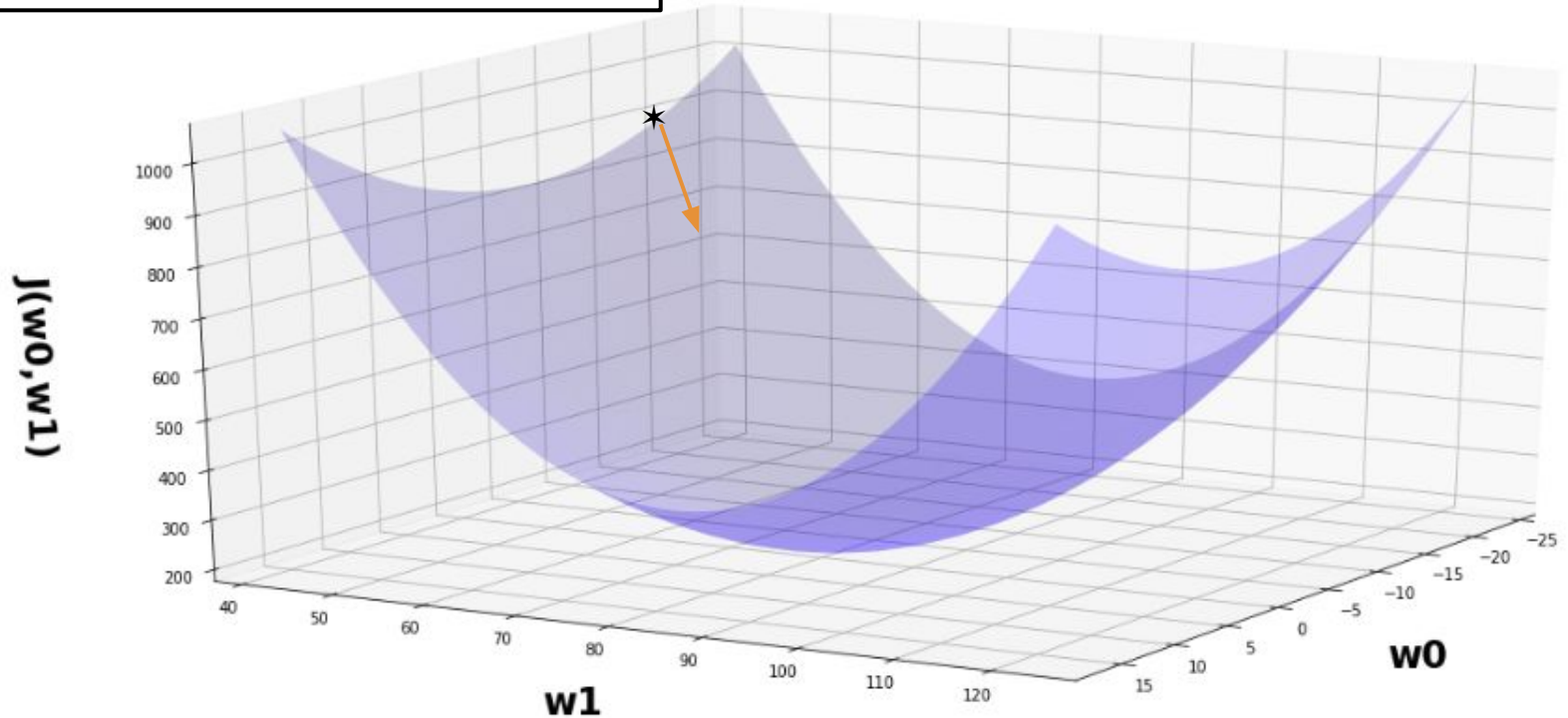
Randomly pick an initial ($w_1$, $w_0$)

# Loss Optimization

Compute gradient $\dfrac{\partial J(w)}{\partial W}$
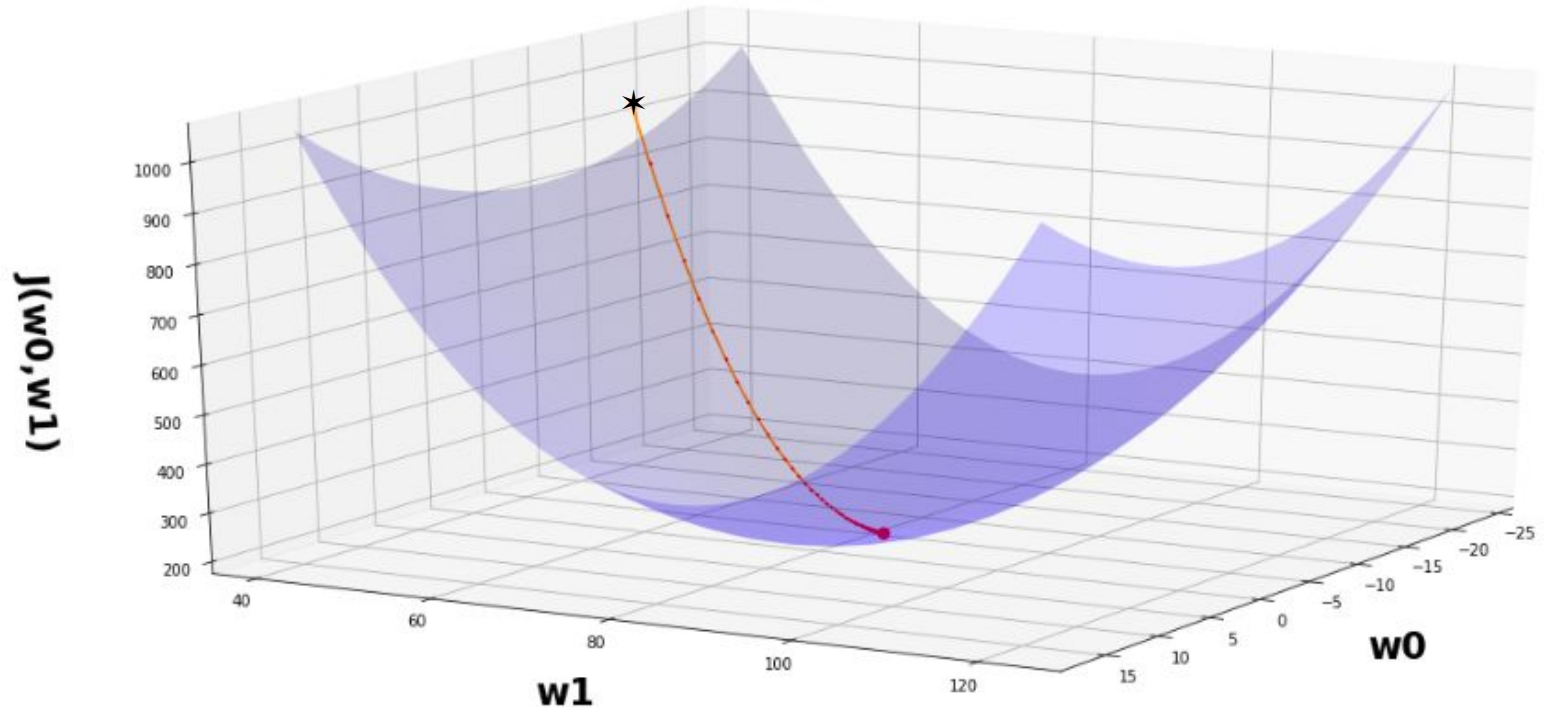
# Loss Optimization

Take direction in opposite direction of gradient

# Gradient Descent

Repeat until convergence

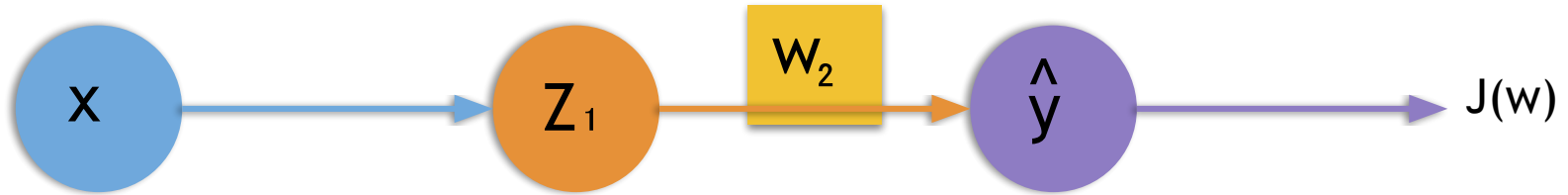# Gradient Descent

Algorithm

1.  Initialize weights randomly

2.  Loop until convergence

3.  Compute gradient $\boxed{\dfrac{\partial J(w)}{\partial W}}$

4.  Update weights $W \leftarrow W - \blacklozenge\blacklozenge \dfrac{\partial J(w)}{\partial W}$
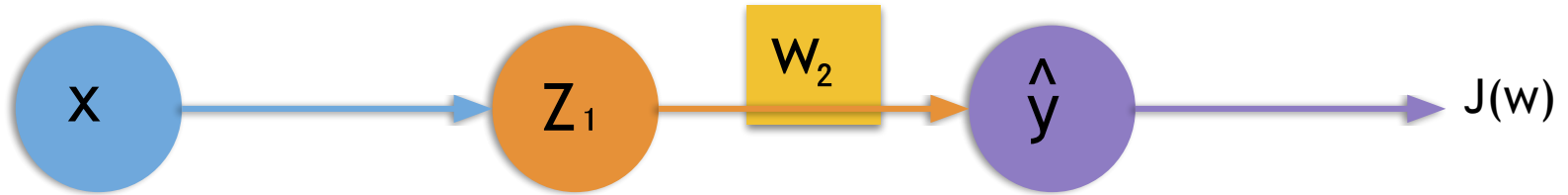
5.  Return weights

# Computing Gradients : Backpropagation



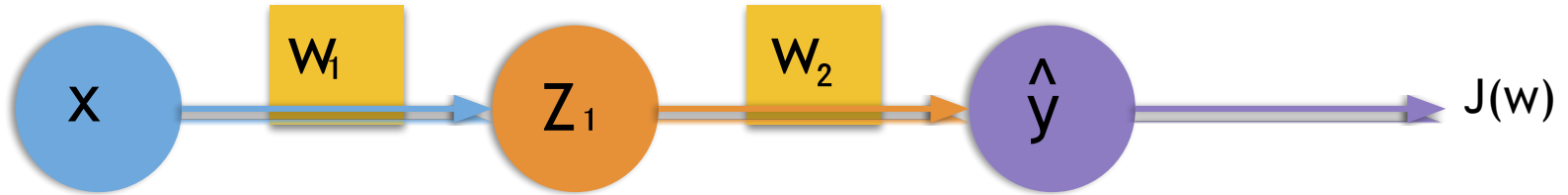How does small change in one weight (eg $W_2$) affect the final loss **J(W)**?

# Computing Gradients : Backpropagation



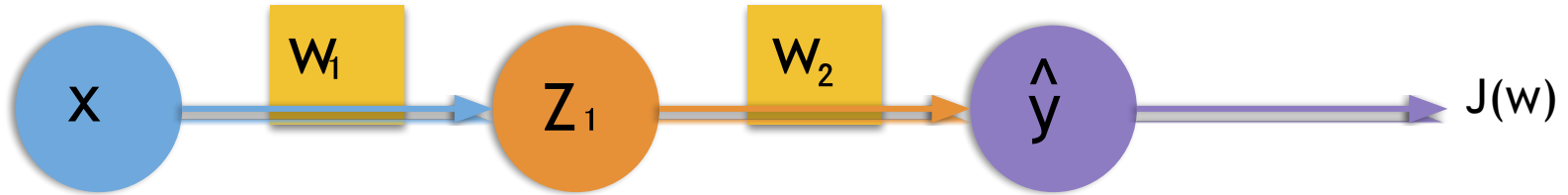$$\frac{\partial J(W)}{\partial w_2}$$

Using chain rule

# Computing Gradients : Backpropagation



$$\frac{\partial J(\boldsymbol{W})}{\partial w_2} = \frac{\partial J(\boldsymbol{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_2}$$
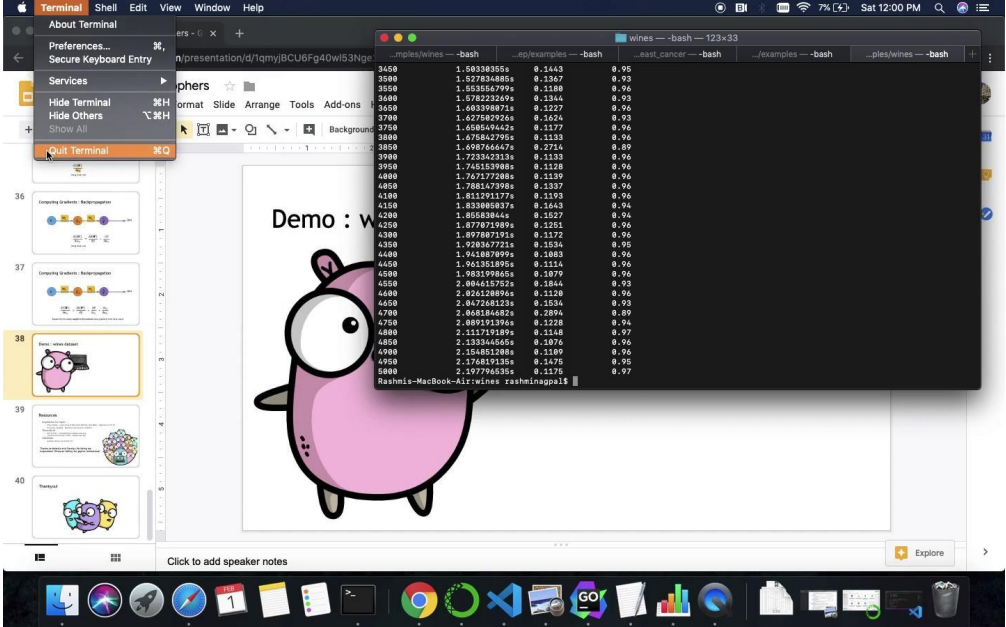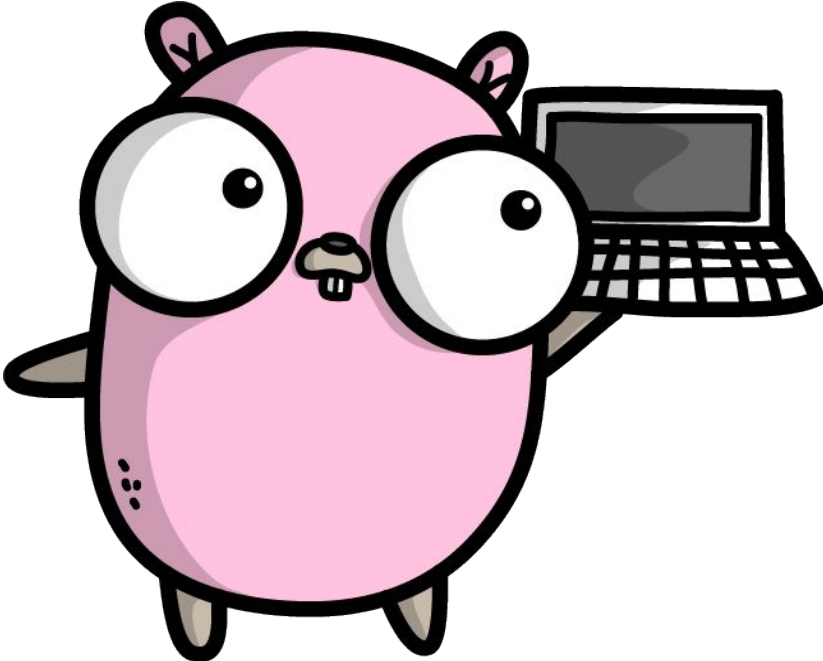
Using chain rule

# Computing Gradients : Backpropagation



$$\frac{\partial J(\boldsymbol{W})}{\partial w_1} = \frac{\partial J(\boldsymbol{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

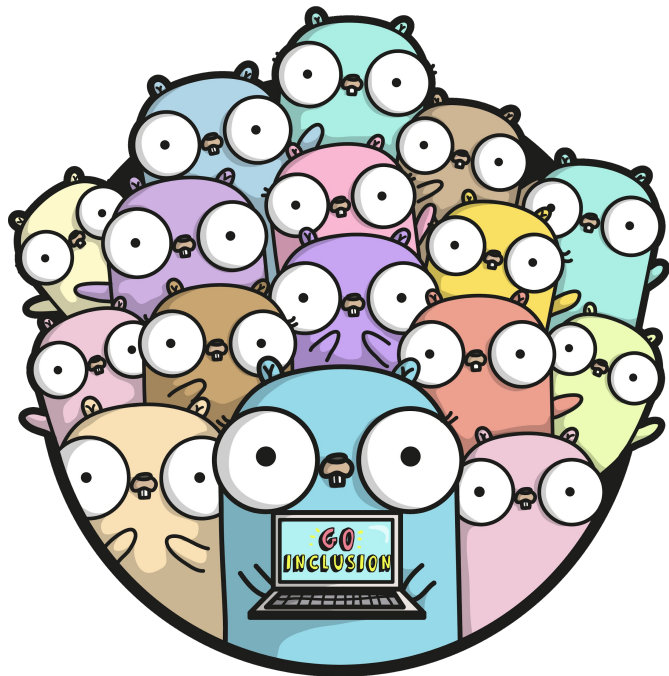Repeat this for **every weight in the network** using gradients from later layers

# Demo

# Resources

- Inspiration for topic :
    - Ellen Kobes - Learn Neural Networks WithGo, Not Math : GopherCon EU'19
    - Francesc Campoy - Machine Learning for Gophers
- Theoretical
    - MIT 6.S191 - Introduction to Deep Learning
    - Stanford University CS230 - Deep Learning
- Codebase
    - godeep library as starter kit

Thanks to Maartje and Carolyn for being my inspiration! Ofcourse Ashley,for gopher animations!

Thankyou!