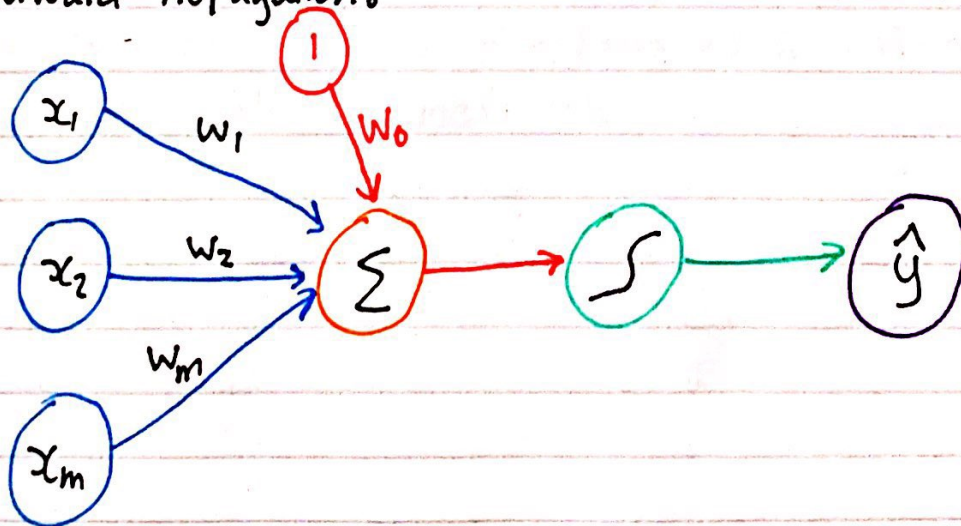# Intro to Deep Learning

$(W^T X + W)_p = \beta$

**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

**Perception/Neuron** : the Structural building block of deep learning

⤷ **Forward Propagation:**



| Inputs | Weights | Sum | Non-linearity | Output |

Output

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

bias term

Non-linearity activation function

linear combination of inputs

$$\rightarrow \hat{y} = g(W_0 + X^T W) \quad \text{where} \quad X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}, W = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$

## Common Activation Functions

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

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

**Hyperbolic Tangent:**
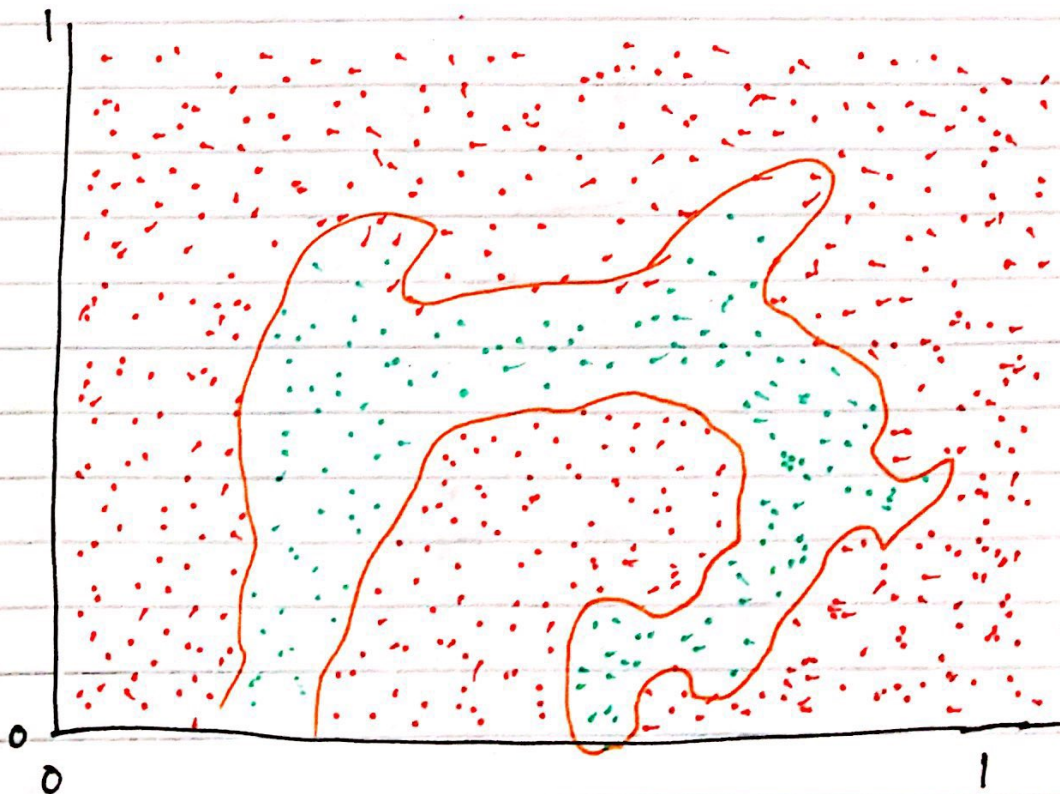$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

**Rectified Linear Unit (ReLU):**
$$g(z) = max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & o/w \end{cases}$$
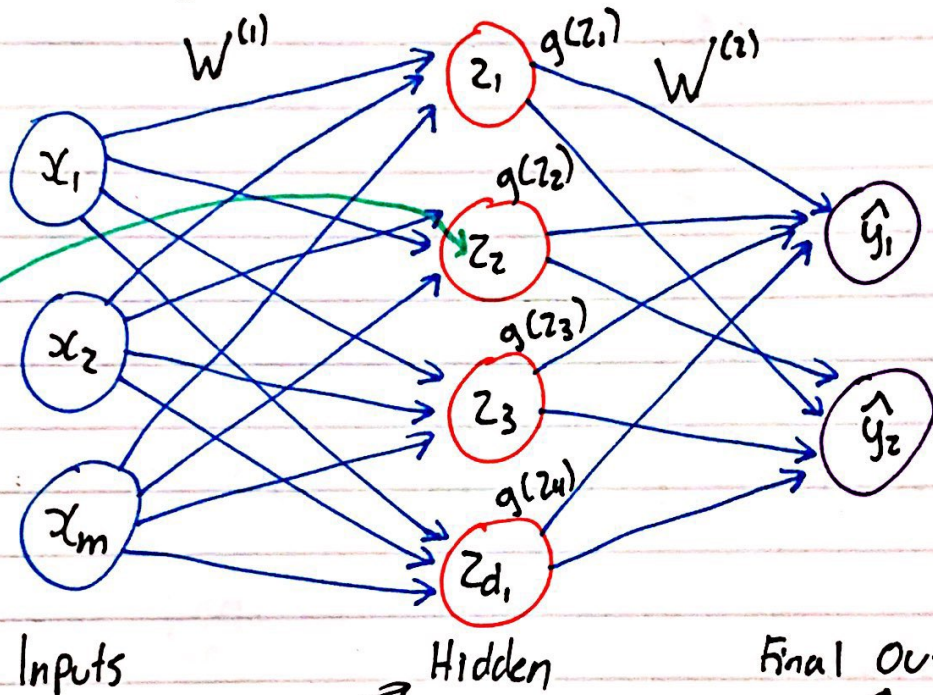
**Importance of Activation Functions:**

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

Non-linearities allow us to approximate arbitrarily complex functions

## Building Neural Networks with Perceptrons
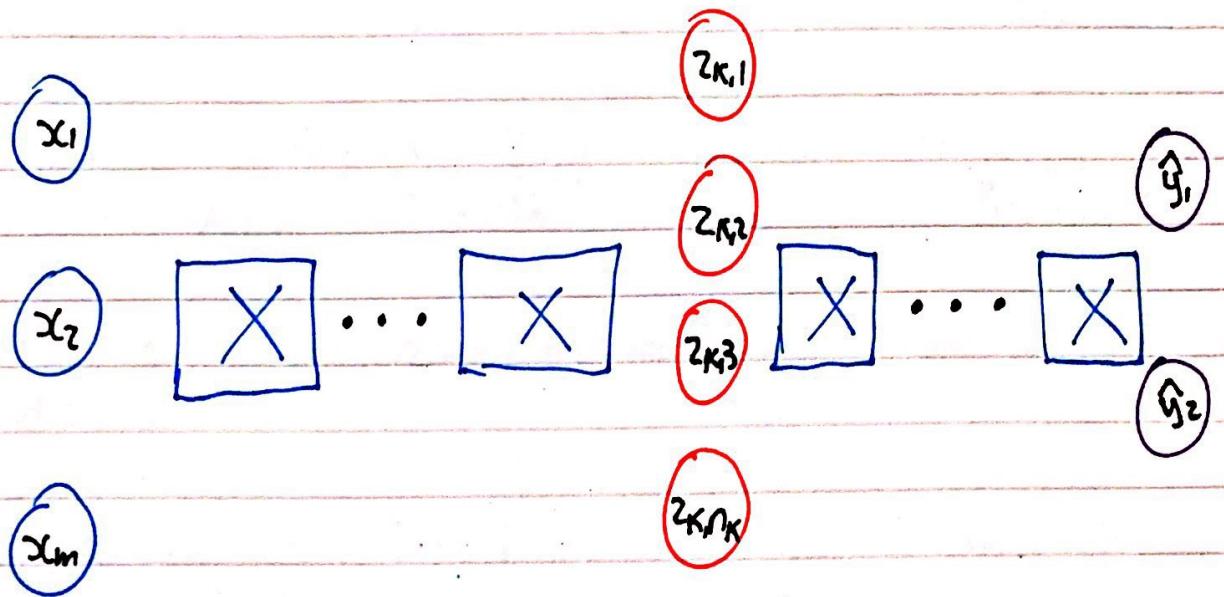
Single Layer Neural Network:



Inputs          Hidden          Final Output

$$Z_i = W_{0,i}^{(1)} + \sum_{j=1}^{m} x_j W_{j,i}^{(1)}$$

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

$$Z_2 = W_{0,2}^{(1)} + \sum_{j=1}^{m} x_j W_{j,2}^{(1)}$$

$$= W_{0,2}^{(1)} + x_1 W_{1,2}^{(1)} + x_2 W_{2,2}^{(1)} + x_m W_{m,2}^{(1)}$$

# Deep Neural Networks:



| Inputs | Hidden | Output |

Hidden Layer: $Z_{k,i} = W_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(Z_{k-1,j}) w_{j,i}^{(k)}$

## Quantifying Loss:

The loss of our network measures the cost incurred from incorrect predictions

$$L(\underbrace{f(x^{(i)}; W)}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$

Empirical Loss: measures the total loss over our entire dataset

$$J(W) = \frac{1}{n} \sum_{i=1}^{n} L(\underbrace{f(x^{(i)}; W)}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$

Note: It is the average loss amongst all individual losses

## Binary Cross Entropy Loss: can be used with models that output a probability between 0 and 1

$$J(W) = \frac{1}{n} \sum_{i=1}^{n} y^{(i)} \log(f(x^{(i)}; W)) + (1-y^{(i)}) \log(1 - f(x^{(i)}; W))$$

Actual    Predicted    Actual    Predicted

↳ This is the difference between actual and predicted distributions

## Mean Squared Error Loss: can be used with regression models that output continuous real numbers

$$J(W) = \frac{1}{n} \sum_{i=1}^{n} (y^{(i)} - f(x^{(i)}; W))^2$$

Actual    Predicted

## Training Neural Networks

Goal: we want to find the network weights that achieve the lowest loss

$$W^* = \underset{W}{\arg\min} \frac{1}{n} \sum_{i=1}^{n} L(f(x^{(i)}; W), y^{(i)})$$

$$= \underset{W}{\arg\min} J(W) \longrightarrow W = \{W^{(0)}, W^{(1)}, \ldots\}$$

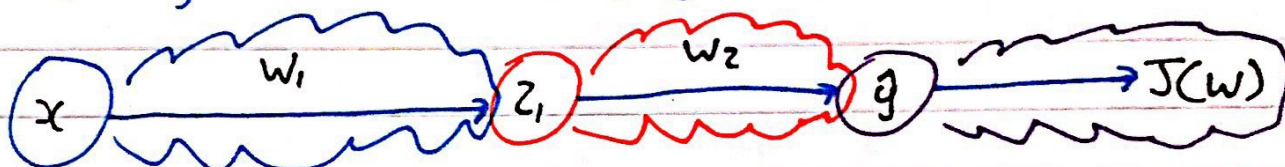Total Loss Function

⟹ Try to find weights that minimize $J(W)$

# Gradient Descent

1. Initialize weights randomly $\sim N(0, \sigma^2)$
2. Loop until convergence (local minimum):
3.      Compute gradient $\frac{\partial J(w)}{\partial w}$
4.      Update weights, $w \leftarrow w - \eta \frac{\partial J(w)}{\partial w}$
5. Return weights

$\eta$ ↳ Learning rate

Computing Gradients: Back propagation



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

Note: ML languages can do this for you

Goal: Pick a learning rate that isn't too small and get stuck at a false local minima but also don't pick an overly aggressive learning rate that causes divergence

↳ Pick a learning rate that adapts to the landscape

Examples: SGD, Adam, Adadelta, Adagrad, RMSProp

# Neural Networks in Practice: Mini-batchies

Note: Calculate gradient descent for each data point within a dataset is computationally intensive
↳ Instead we can choose a single point and compute gradient descent relative to that point

However, this is very noisy (Stochastic)

Middle ground: Pick a batch of data points

Modified Step 4: $\dfrac{\partial J(w)}{\partial w} = \dfrac{1}{B} \sum_{k=1}^{B} \dfrac{\partial J_k(w)}{\partial w}$

This allows us to converge at a more optimal rate

# Neural Networks in Practice: Overfitting

We want a middle ground between underfitting and overfitting
↳ Our model continues to generalize when it sees new data

Regularization: technique that constrains our optimization problem to discourage complex models

↳ Dropout: During training, randomly set some activations to 0. Typically drop 50% of activation layer. Forces network not to rely on any one node (in hidden layers)

→ Early Stopping: Stop training before we have a chance to overfit



Loss (y-axis)

Training Iterations (x-axis)

Stop training here

Possibly started to memorize training data