# Deep Learning

M.Tech. Data Science, Second Year, NMIMS

By,

Bilal Hungund, Data Scientist, Halliburton

Artificial Intelligence (Observing behaviour)
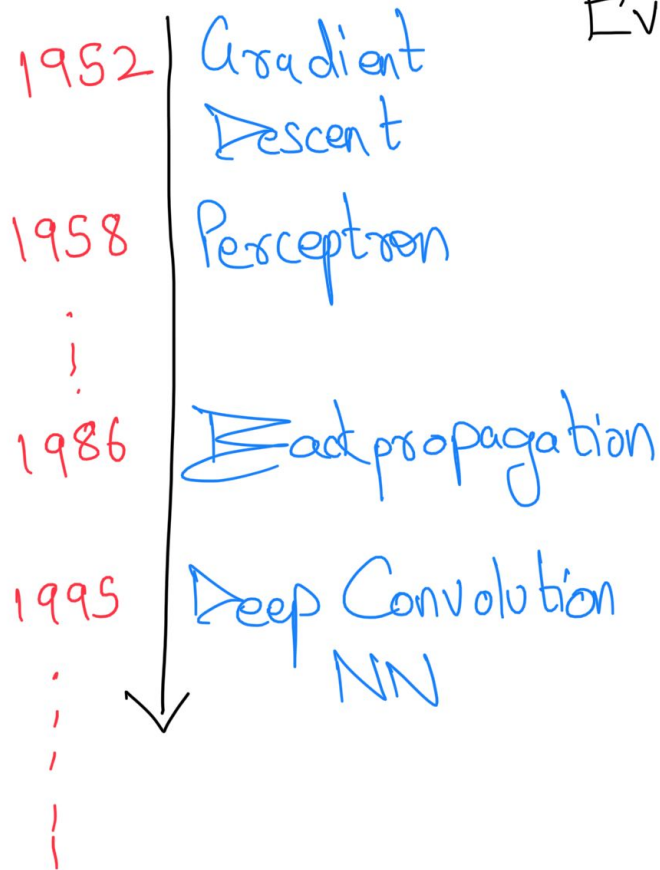
Machine Learning (Explicitly learn)

Deep Learning

(Extract Pattern from
Neural Network)

# Why Deep Learning

→ Time Consuming and brittle in unstructured data

→ Hard engineered features are not scalable

→ Ways to learn underlying features for unstructured data
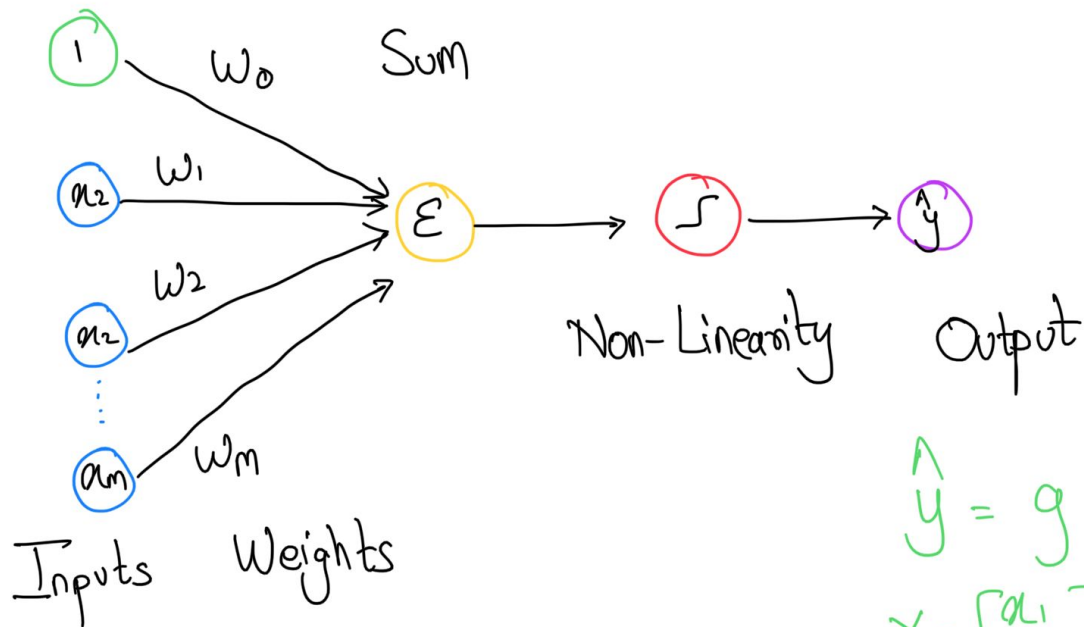
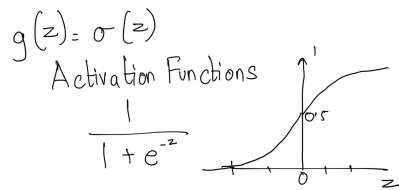# Evolution of Deep Learning

| | |
|---|---|
| 1952 | Gradient Descent |
| 1958 | Perceptron |
| ⋮ | |
| 1986 | Backpropagation |
| 1995 | Deep Convolution NN |
| ⋮ | |
| ⋮ | |

Why now?

1) Big Data

2) Hardware (GPUs, TPUs)

3) Software (Pytorch, Tensorflow)

# Perceptron

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

$g(z) = \sigma(z)$

Activation Functions

$$\frac{1}{1 + e^{-z}}$$

1 → $w_0$

Sum

$a_2$ → $w_1$

$a_2$ → $w_2$

$a_m$ → $w_m$

$\varepsilon$ → (non-linearity) → $\hat{y}$

Non-Linearity    Output

Inputs    Weights

Simplify to Matrix

Hyperbolic Tangent (tanh)

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

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

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

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \quad W = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$
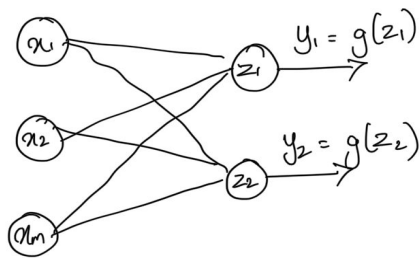
Rectified Linear Unit (ReLU)

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

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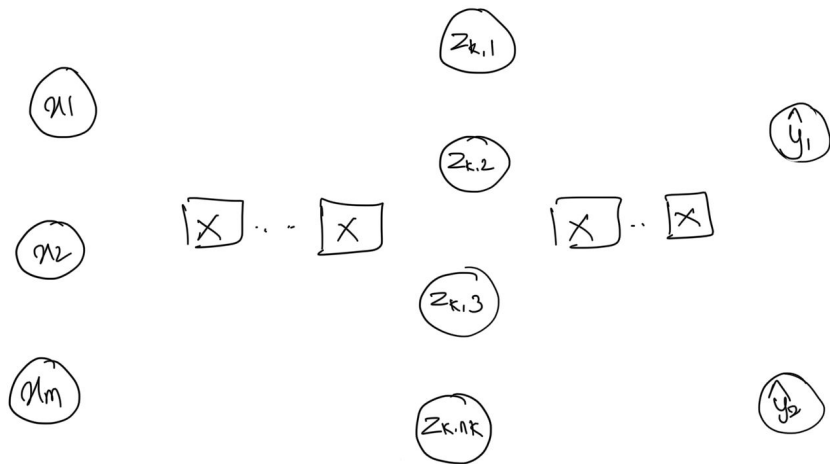
# Multi Output Perceptron

$x_1$ → $z_1$ → $y_1 = g(z_1)$

$x_2$

$z_2$ → $y_2 = g(z_2)$

$x_m$

$$z_i = w_{0,i} + \sum_{i=1}^{m} x_j w_{j,i}$$

# Single Layer Neural Network

$W^{(1)}$  $z_1$  $g(z_1)$  $W^{(2)}$

$x_1$

$z_2$  $g(z_2)$  $\hat{y}_1$

$x_2$

$g(z_3)$

$z_3$  $\hat{y}_2$

$x_m$

$z_{d_1}$  $g(z_{d_1})$

Inputs

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

Hidden
$$z_i = w_{0,i}^{(i)} + \sum_{j=1}^{m} x_j w_{j,i}^{(i)}$$

# Deep Neural Network

k # hidden layer

$z_{k,1}$

$x_1$

$z_{k,2}$  $\hat{y}_1$

$x_2$

$\boxed{\times} \cdots \boxed{\times}$    $\boxed{\times} \cdots \boxed{\times}$

$z_{k,3}$

$x_m$

$z_{k,n_k}$  $\hat{y}_2$

$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(z_{k-1,j}, j) \, w_{j,i}^{(k)}$$

Keep stacking hidden layers

# Neural Networks



Perseptron (P)

Feed Forward (FF)

Radial Basis Network (RBF)

Deep Feed Forward (DFF)

Recurrent Neural Network (RNN)

Long/Short Term Memory (LSTM)

Gated Recurrent Unit (GRU)

Auto Encoder (AE)

Variational AE (VAE)

Denoising AE (DAE)

Sparce AE (SAE)

Markov Chain (MC)

Hopfield Network (HN)

Boltzmann Machine (BM)

Restricted BM (RBM)

Deep Belief Network (DBN)

Deep Convolutional Network (DCN)

Deconvolutional Network (DN)

Deep Convolutional Inverse Graphics Network (DCIGN)

Generative Adversarial Network (GAN)

# Quantifying Loss

It measures the cost incurred from incorrect predictions

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

$$\underbrace{\phantom{\mathcal{L}\left(f(x^{(i)}; W), y^{(i)}\right)}}_{J(W) \to \text{Empirical Loss}}$$

Binary Cross Entropy Loss

$$J(w) = \frac{1}{n} \sum_{i=1}^{n} y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

Actuals    Predicted

Mean Squared Error Loss

$$J(w) = \frac{1}{n} \sum_{i=1}^{n} \left(y^{(i)} - \hat{y}^{(i)}\right)$$

Actuals   Predicted

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

$$W^* = \underset{W}{\mathrm{argmin}} J(W)$$

# Gradient Descent

**Algorithm**

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence:

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



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

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

# Optimization

- Learning Rate
- Regularization
- Dropout
- Early Checkpoint