

Machine Learning Lect-2

Introduction to Neural Networks

DR. PREKSHA MUKHERJEE



Introduction

- Neural Networks

↳ Inspired by how human brain works & maps the inputs to logical outputs.

- Brain is the most complex part of human body & consists of 10 billion + neurons

- Neurons are highly interconnected & act as individual units.

- Characteristics of neurons:

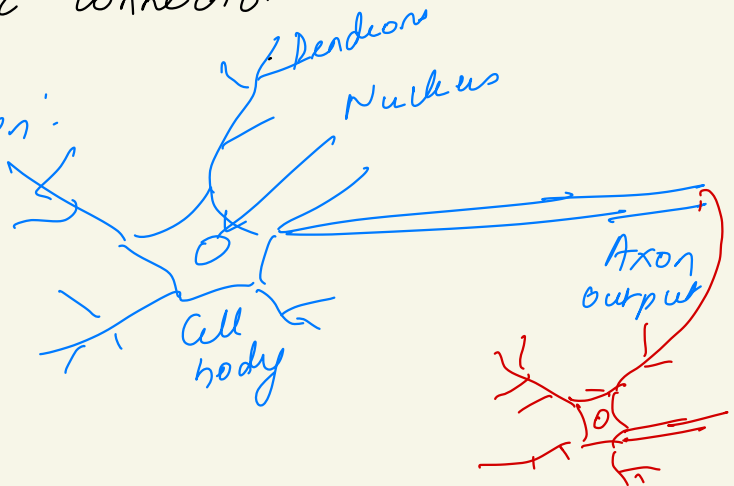
- Massive parallelism

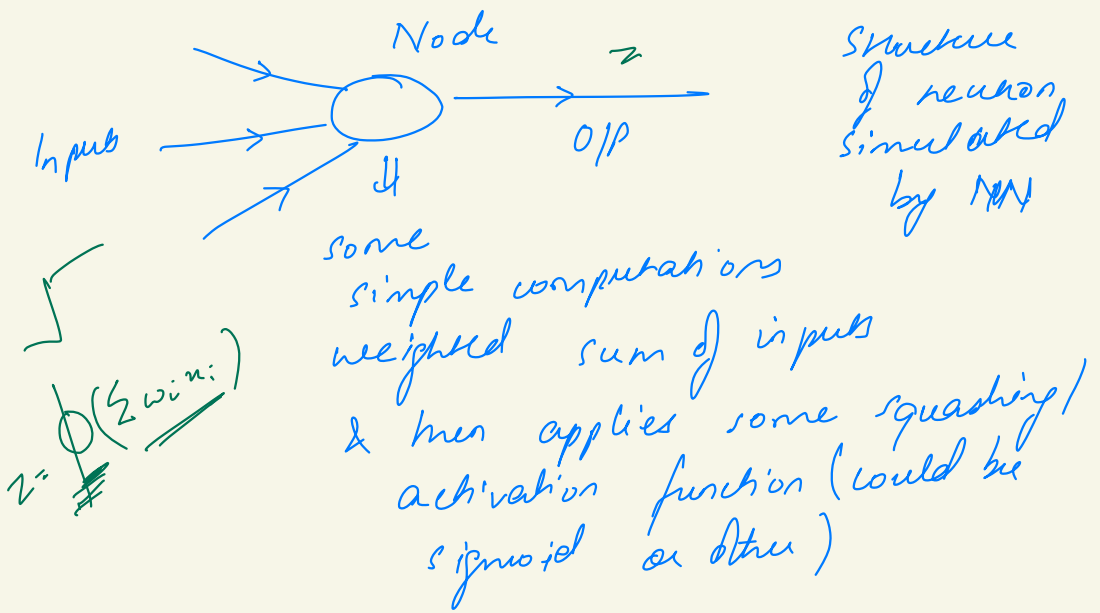
- Highly interconnected (solve the complex problems of your real world)

- Model distributive

associative memory through weights in the synaptic connections

Diagram of neuron:





ANNs through synaptic connections O/P meets another neuron input using electric impulses through the synaptic layers.

Incorporate the 2 fundamental components of the biological neural nets.

- Nodes - Neurons
- Weights - synapses

Basic unit in NN: Perception

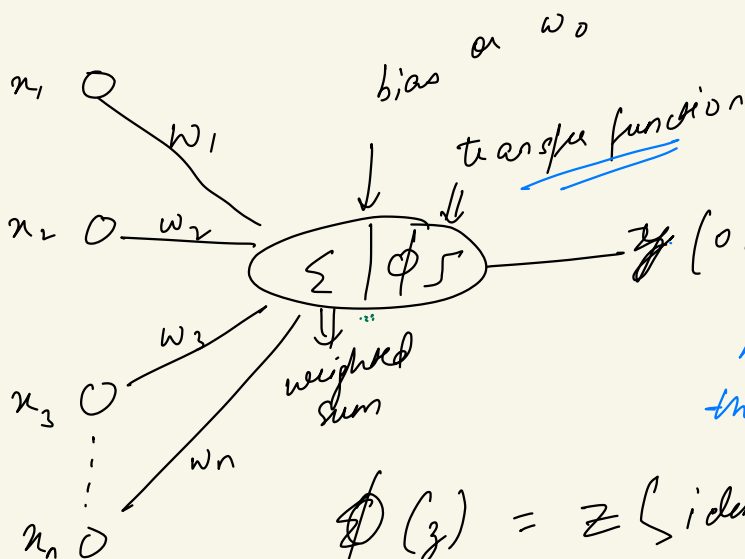
Perceptron

Inputs : n

Input
 $x_0 = 1$

Training Data

$\begin{matrix} \bar{x} & \bar{y} \\ \rightarrow \bar{x}_1 & \bar{y}_1 \\ & \bar{x}_2 & \bar{y}_2 \\ & \vdots & \vdots \\ & \bar{x}_m & \bar{y}_m \end{matrix}$



Training the
n/w means learning
the weights to have
good fit for
training data.

$$\phi(z) = z \quad (\text{identity})$$

$$\phi(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (\text{unit step})$$

$$\phi(z) = \frac{1}{1 + e^{-z}} \quad (\text{sigmoid})$$

$$y = \sum_{i=1}^n w_i x_i + b \quad \text{or} \quad \left. \sum_{i=0}^n w_i x_i \right\} \text{ bias}$$

ϕ is identity function \nrightarrow same
passed as
output

or

$$\phi_2(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad \text{threshold function}$$

Perceptron training rules

- we feed 1 eg at a time & update the weights based on output.

if output (predicted) == output (target) \Rightarrow no update in weights required.

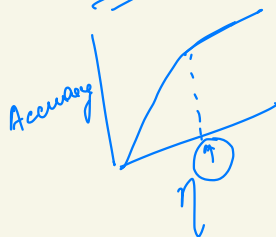
$$w_{i, \text{next}} = w_{i, \text{prev}} + \Delta w_i$$

$$\Delta w_i = \eta (y - \hat{y}) x_i$$

learning rate based on error from change weights

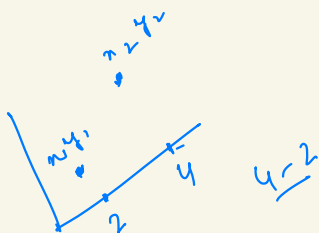
target predicted

Grid search



Stopping criteria \Rightarrow Perceptron learning convergence if D is linearly separable. (training data)

What if data is not linearly separable? \Rightarrow linear decision surface to separate + / -

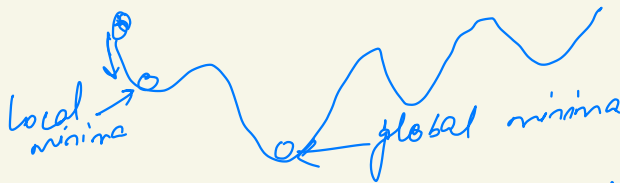


Solution : Gradient descent ↗

can be used when want to particular parameter set/weight values you can define the error of the n/w. Perform GD on error func. to find the optimal parameter set for which the error fn. is minimized

$$E = \frac{1}{2} \sum_{d \in D} (y_d - \hat{y}_d)^2 \Rightarrow \text{is like a surface \& we want to find minima of this surface}$$

\uparrow
so that differentiation is easy



Find partial derivative of error fn. with each & every weight and find the direction of the gradient & we go towards the -ve direction of gradient in order to go towards minima.

In certain cases, this error surface can be convex or quadratic & there will be single minima \Rightarrow thus gradient descent guarantees the minima.



Multi-layer ^{network error} surface will be ill-behaved (non-convex)
 there will be local minima
 ↓
 get stuck

$\phi_1(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$ step function \rightarrow is not differentiable \Rightarrow so gradient descent can't be done

$\phi_1(z) = z$ is differentiable \Rightarrow gradient descent can be done

Stochastic gradient descent

1 eg. from training set at a time

$$E = \frac{1}{2} (y_d - \hat{y}_d)^2$$

Training rule
 of G.D

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

$$\frac{\partial E}{\partial w_i} = x_{ij} (y_j - \hat{y}_j)$$

Single layer n/w capture linear decisions / linear separable data.

Capture non-linear functions \rightarrow multi-layer n/w's
 or complex
 \rightarrow different n/w's connected with each other
 \rightarrow if linear units are connected combination is again linear fn.

$$\phi_3(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Assignment
have it →

$$\phi'_3(z) = \phi_3(z) \cdot (1 - \phi_3(z)) \quad (*)$$

$$E = \frac{1}{2} \sum_{d \in D} \left(y_d - \underbrace{\sigma(w \cdot x_d)}_{\hat{y}_d} \right)^2$$

$$\frac{\partial E}{\partial w_i} = \frac{1}{2} \sum_d \frac{\partial E}{\partial \hat{y}_d} \cdot \frac{\partial \hat{y}_d}{\partial w_i}$$

$$\frac{\partial E}{\partial w_i} = \sum_d (y_d - \hat{y}_d) \cdot \sigma'(w \cdot x_d) \cdot x_{di}$$

$$\frac{1}{2} \sum_d (y_d - \hat{y}_d)$$

to represent complex functions we want non-linear units (which are differentiable)
 \Downarrow
 eg sigmoid / logistic transfer function
 $\sigma'(w \cdot x_d) = 1$

$$= \sum_d (y_d - \hat{y}_d) x_{di}$$

Using \odot in $\textcircled{1}$

$$\frac{\partial E}{\partial w_i} = \sum_d (y_d - \hat{y}_d) \hat{y}_d (1 - \hat{y}_d) n_{id}$$

$$\Delta w_i = -\eta \sum_d (y_d - \hat{y}_d) \hat{y}_d (1 - \hat{y}_d) n_{id}$$

training
rule
of sigmoid
unit

↑
single layered
NN

→ for linearly
separable function.