# Foundations of NN

Perceptrons → Gradient → non-threshold
(single neuron)   Learning    activation functions

↓                     sigmoid ReLu   tah   softmax
primarily for                                SBAF
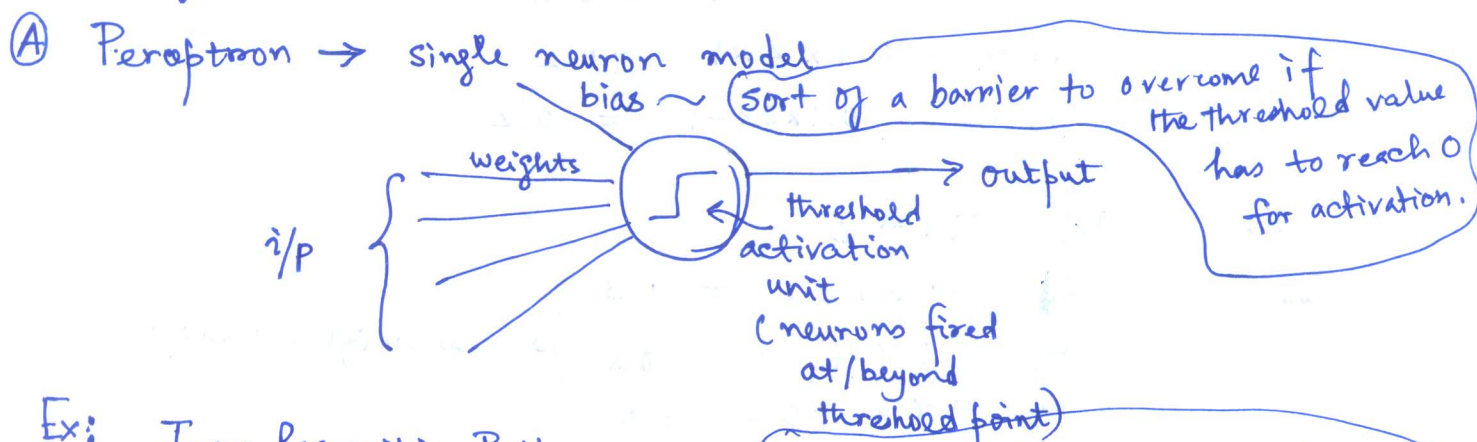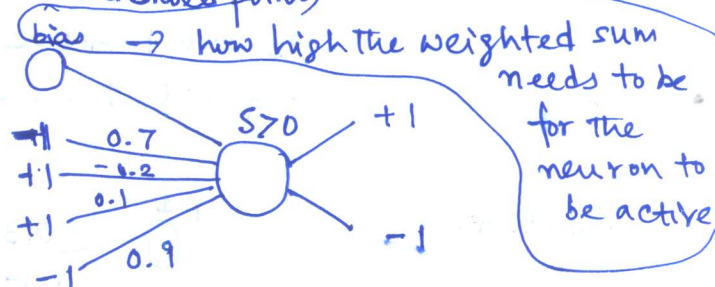object recognition
(visual cortix)

ANN mimics human brains (only partially) for pattern recognition jobs — basic functions for living animals. — ANN's don't have the ability to recall an answer to a search problem at a later point.

## Building blocks:

Ⓐ Perceptron → single neuron model
bias ~ (sort of a barrier to overcome if the threshold value has to reach 0 for activation.)

weights           → output

i/p {      threshold activation unit (neurons fired at/beyond threshold point)

bias → how high the weighted sum needs to be for the neuron to be active

Ex: Image Recognition Problem

→ dark (-1)
→ bright (+1)

target label: bright image
(+1)

-1   0.7     S>0    +1
+1   -0.2
    0.1
+1
-1   0.9        -1

Goal: to classify image as bright or dark.

## Steps:

1.   $S = \sum w x = w^T x$ ; weighted combination of inputs

2.   provide S to the activation function: $0 = f(s)$

3.   Compute Error: $E = t - 0$

4.   Calculate $\Delta w_i = \alpha E x_i$   ($\alpha$ → learning rate)

5.   update weight, $w_i = w_i + \Delta w_i$

Iterate till Error is minimized

Perceptron Learning Rule

**Iteration 1**

$$[1\ -1\ 1\ 1\ -1]\begin{bmatrix} -0.5 \\ 0.7 \\ -0.2 \\ 0.1 \\ 0.9 \end{bmatrix} \implies W^T X \; ; \quad \begin{array}{l} W \rightarrow \text{weight vector} \\ X \rightarrow \text{input vector} \end{array}$$

1. $S = -0.5 - 0.7 - 0.2 + 0.1 - 0.9 = -2.2 < 0$

2. $0 = -1, \quad t = +1;$

3. $E = t - 0 = 2$

4. $\Delta w_i = \alpha E x_i \; ; \quad \text{set } \alpha = 0.1$

$\Delta w_0 = 0.1 \times 2 \times 1 = 0.2 \qquad \rightarrow \text{for bias element}$

$\Delta w_1 = 0.1 \times 2 \times -1 = -0.2$

$\Delta w_2 = 0.1 \times 2 \times 1 = 0.2$

$\Delta w_3 = 0.1 \times 2 \times 1 = 0.2$

$\Delta w_4 = 0.1 \times 2 \times -1 = -0.2$

5. $w_i = w_i + \Delta w_i;$

$\left.\begin{array}{l} w_0 = -0.5 + 0.2 = -0.3 \\ w_1 = 0.7 + 0.2 = 0.5 \\ w_2 = -0.2 + 0.2 = 0 \\ w_3 = 0.1 + 0.2 = 0.3 \\ w_4 = 0.9 - 0.2 = 0.7 \end{array}\right\} \text{new weight vector}$

**Iteration 2**



1. $S = 1 \times (-0.3) + 1 * (0.5) + 1 (0) + 1 * (0.3) - 1 * (0.7)$

$= -1.2 < 0$

2. $0 = -1, \quad t = +1$

3. $E = t - 0 = 2$

4. $\Delta w_i = \alpha E x_i$

$\Delta w_0 = 0.1 \times 2 \times 1 = 0.2$

$\Delta w_1 \rightarrow \Delta w_4 \text{ (remain the same)}$

updated weight vector

$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix} = \begin{bmatrix} -0.3 + 0.2 \\ 0.5 + 0.2 \\ 0 + 0.2 \\ 0.3 + 0.2 \\ 0.7 + 0.2 \end{bmatrix} = \begin{bmatrix} -0.1 \\ 0.7 \\ 0.2 \\ 0.5 \\ 0.9 \end{bmatrix}$$

Continue iterating....

### 3rd iteration

$$\begin{bmatrix} 0.1 \\ 0.1 \\ 0.4 \\ 0.7 \\ 0.3 \end{bmatrix} \longrightarrow$$

### 4th iteration

$$[1 - 1 + 1 + 1 - 1] \begin{bmatrix} 0.1 \\ 0.1 \\ 0.4 \\ 0.7 \\ 0.3 \end{bmatrix}$$

1. $S = 0.8 > 0$
2. $0 = +1$
3. $t = +1$
4. $E = t - 0 = 0$

(STOP)

Converged weights

---

(B) Building Block 2 → Gradient Learning

How machines learn? → Calculus based tools → comes down
(ANN)          to finding minimum of some function

Initially, ANN ~~perform~~ performs horribly → random initialization
of bias & weight vectors ⟹ combination of neurons to another
neuron — bad computing! ⟹ There's a cost ⟹ define a
loss function!!

Cost of a single training example = difference b/w what we
want (target) & output produced by activation (trash)

What's the cost of this difference? Squared difference.

Consider the average cost → determines how good/bad the
network is!

<u>Cost function:</u>   I/P : thousands of weights/biases ⎤ ↓
                    O/P : 1 number (cost)            ⎥ change those
                    Parameters : Many training examples. ⎦ weights/biases
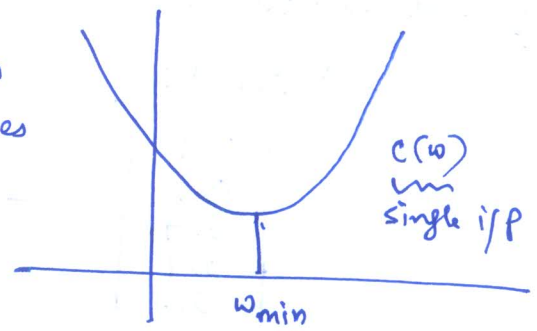                                                          so that the
                                                          performance
                                                          improved!

Define the Cost function: $C(w_1, -- w_n)$

Imagine: Single input single o/p function
— how do we find a number that minimizes
the o/p?

$$\frac{dc}{dw} = 0 \Rightarrow \text{easy for a}$$
simple function

but very complicated for a function w/ thousands of i/p's

Ball downhill problem $\Rightarrow \nabla C(x,y)$

$C(w)$ in single i/p

$w_{min}$

weights & biases

$$\vec{w} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \qquad - \nabla C(\vec{w}) = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \longrightarrow \text{Minimize Cost}$$

↓
direction
of steepest
descent
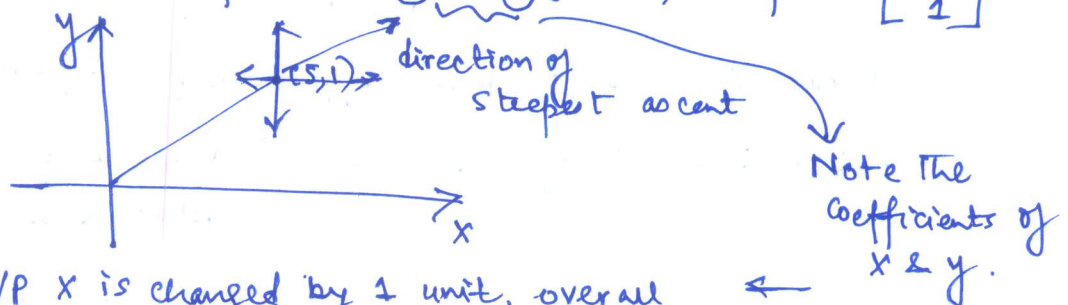
Local (valley)
minima

Column vectors

What does this
tell us?

relative magnitude of
the components of $\nabla C(\vec{w})$ $\Rightarrow$ how to make changes in
the i/p matrix (since
a huge difference b/w i/p & o/p matrix will make the
cost look bad & the task of classification won't benefit.

Imagine a function: $f(x,y) = \left(\frac{5}{2}\right)x^2 + \left(\frac{1}{2}\right)y^2$ ; $\nabla f = \begin{bmatrix} 5 \\ 1 \end{bmatrix}$

$y$

(5,1) direction of
steepest ascent

$x$

Note the
coefficients of
$x$ & $y$.

if the i/p $x$ is changed by 1 unit, overall
change in o/p $(f)$ is more as compared to change in $y$.
This tells us, gradients can help us decide which components
in i/p should be changed & by how much so that the
difference b/w i/p & o/p is squashed.

# More on Gradient descent:

Gradient, $\nabla E(\vec{w}) = \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \cdots \frac{\partial E}{\partial w_n} \right]$

Training Rule: $\Delta \vec{w} = -\eta \nabla E(\vec{w})$

i.e. $\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$ ; for each weight component

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 = \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2$$

$$= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d)$$

$$= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d)$$

$$= \frac{1}{2} \sum_d 2(t_d - o_d)(-\vec{x}_d)$$

## GD (training examples, $\eta$):

Each training example is a pair $\langle \vec{x}, t \rangle$ where $\vec{x}$ is the vector of i/p values & $t$ is the target o/p value. $\eta$ is the learning rate.

1. Initialize each $w_i$ to some small random value

2.    Until the termination condition is met, do

   — initialize each $\Delta w_i$ to zero
   — for each $\langle \vec{x}, t \rangle$, do
      * input the instance $\vec{x}$ to the activation unit & compute the o/p, $o_d = f(\vec{x}_d)$
      * for each linear unit. $w_i$ do
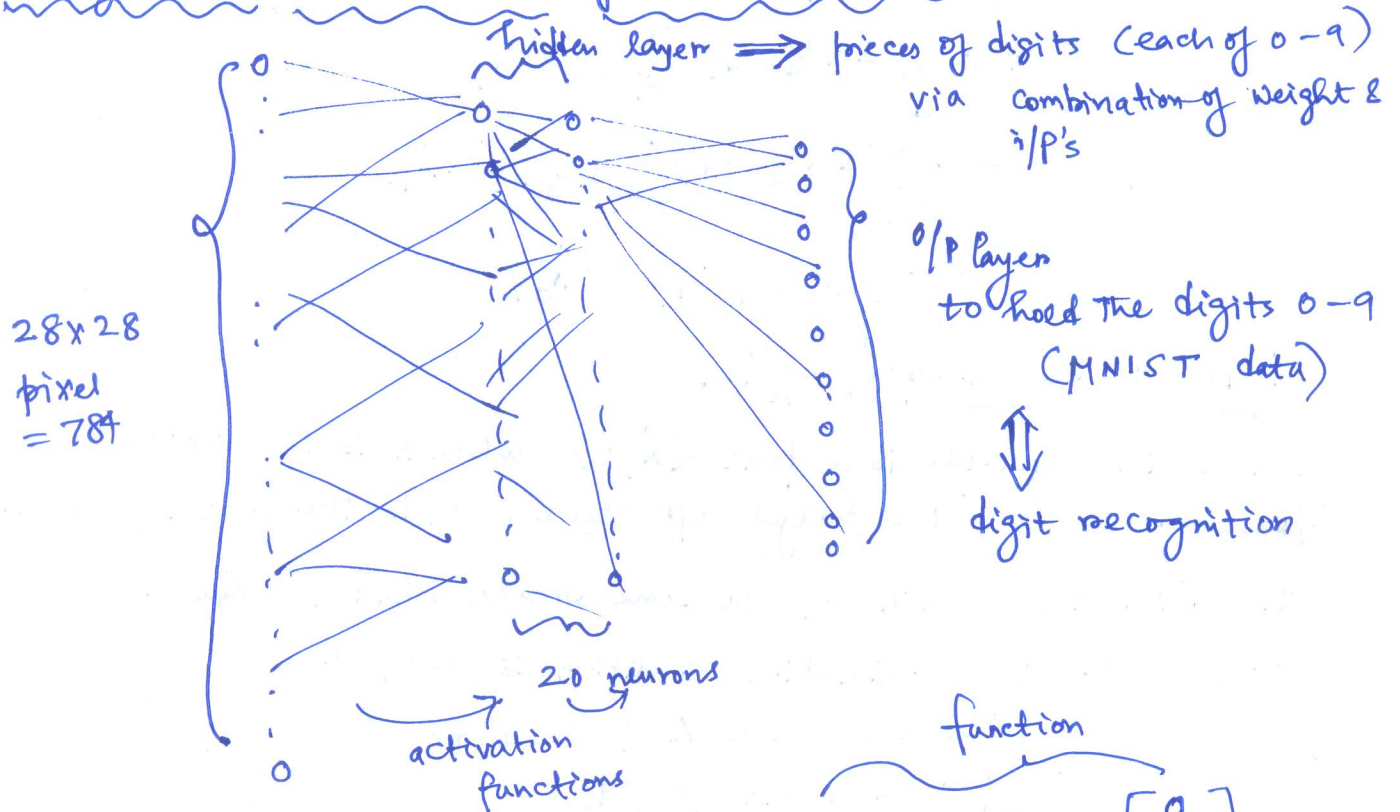      $$\Delta w_i \leftarrow \Delta w_i + \eta (t_d - o_d) x_d$$

   — for each linear unit weight $w_i$, do
      $$w_i \leftarrow w_i + \Delta w_i$$

## Additional Reads:
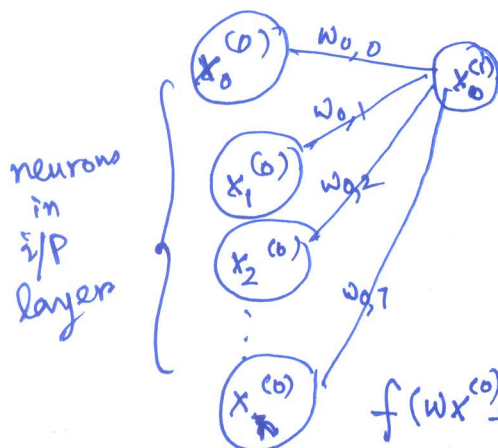
1. Hinton's slides
2. Goodfellow's text on Deep Learning.

: A high level abstraction of a Neural Network:

hidden layer $\Longrightarrow$ pieces of digits (each of 0-9) via combination of weight & i/p's

O/P layer to hold the digits 0-9 (MNIST data)

$\Updownarrow$

digit recognition

28×28 pixel = 784

20 neurons

activation functions

function

1) Each neuron holds a # i.e $f(x_0, \ldots x_{783}) = \begin{bmatrix} O_0 \\ \vdots \\ O_9 \end{bmatrix}$

2) Entire network is just a function : takes in 784 units as i/P & produces 10 units as O/P.

$$784 \times 20x + 20 \times 20 + 20 \times 10 \quad \text{—— weights}$$

$$20 + 20 + 10 \quad \text{—— biases}$$

$$\left. \begin{array}{} \end{array} \right\} \quad \boxed{\begin{array}{l} 15740 + 600 + 50 \\ = 16390 \text{ i/P's} \end{array}}$$
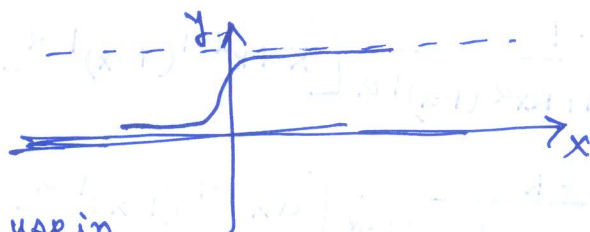
neurons in i/P layer

$$X_0^{(1)} = f\left( W_{0,0} X_0^{(0)} + W_{0,1} X_1^{(0)} + \cdots W_{0,n} X_n^{(0)} + b_0 \right)$$

$$f\left( \begin{bmatrix} W_{0,0} & W_{0,1} & -- & W_{0,n} \\ W_{1,0} & W_{1,1} & -- & W_{1,n} \\ \vdots & & & \\ \vdots & & & \\ W_{K,0} & W_{K,n} & -- & W_{K,n} \end{bmatrix} \begin{bmatrix} X_0^{(0)} \\ X_1^{(0)} \\ \vdots \\ X_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix} \right)$$

$f(wx^{(0)} + b) \Longleftrightarrow$

Communicates b/w layers

Gradient descent regression: ⇒ printed notes:

Building block 3 : (Activation function)

(i) Sigmoid:



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

Nice Properties: (to use in back propagation)

$$\frac{d\sigma}{dx}$$

$$= \frac{-1}{(1+e^{-x})^2}(+e^{-x}) = \frac{1+e^{-x}-1}{(1+e^{-x})^2}$$
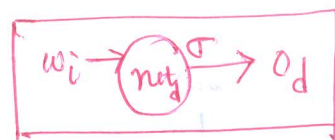
$$= \frac{1}{1+e^{-x}} - \frac{1}{(1+e^{-x})^2}$$

O/P of the activation function
⇒ how positive the relevant weighted sum is )

$$= \sigma(x) - \sigma(x)^2 = \sigma(x)(1-\sigma(x))$$

Error Gradient for a sigmoid unit:

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D}(t_d - O_d)^2$$

$$= \frac{1}{2} \sum_{d \in D}\frac{\partial}{\partial w_i}(t_d - O_d)^2$$

$$= \frac{1}{2} \sum_{d \in D} 2(t_d - O_d)\frac{\partial}{\partial w_i}(t_d - O_d)$$

$$= \frac{1}{2} \sum_{d \in D}(t_d - O_d)\left(-\frac{\partial O_d}{\partial w_i}\right)$$

$$= -\sum_{d \in D}(t_d - O_d)\frac{\partial O_d}{\partial net_d} \cdot \frac{\partial net_d}{\partial w_i}$$



$$w_i \to \boxed{net_d} \xrightarrow{\sigma} O_d$$

Note,

$$\frac{\partial O_d}{\partial net_d} = \frac{\partial \sigma(net_d)}{\partial net_d} \quad ; \quad \text{since } O_d \to \sigma(net_d)$$

network data passed through the activation function $\sigma$

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

$$\& \quad \frac{\partial net_d}{\partial w_i} = \frac{\partial(\vec{w}.\vec{x}_d)}{\partial w_i} = x_{i,d}$$

See, how helpful the property
$$\frac{d\sigma}{dx} = \sigma(1-\sigma)$$

So, $$\frac{\partial E}{\partial w_i} = -\sum_{d \in D}\underbrace{(t_d - O_d)}_{\text{loss}}\underbrace{O_d(1-O_d)}_{\text{slope of Loss function}}\underbrace{x_{i,d}}_{\text{i/p data}}$$

(ii) SBAF: $\quad y = \dfrac{1}{1+Kx^\alpha(1-x)^{1-\alpha}}$

$\ln y = \ln 1 - \ln\left(1+Kx^\alpha(1-x)^{1-\alpha}\right) = -\ln\left(1+Kx^\alpha(1-x)^{1-\alpha}\right)$

$\Rightarrow \dfrac{1}{y}\dfrac{dy}{dx} = \dfrac{-1}{1+Kx^\alpha(1-x)^{1-\alpha}}\left[K\alpha x^{\alpha-1}(1-x)^{1-\alpha} - Kx^\alpha(1-\alpha)(1-x)^{1-\alpha-1}\right]$

$= \dfrac{-K}{1+Kx^\alpha(1-x)^{1-\alpha}}\left[\alpha x^{\alpha-1}(1-x)^{1-\alpha} - (1-\alpha)x^\alpha(1-x)^{1-\alpha-y}\right]$

$= -\dfrac{Kx^\alpha(1-x)^{1-\alpha}}{1+Kx^\alpha(1-x)^{1-\alpha}}\left[\dfrac{\alpha}{x} - \dfrac{(1-\alpha)}{1-x}\right]$

$= -\dfrac{Kx^\alpha(1-x)^{1-\alpha}}{1+Kx^\alpha(1-x)^{1-\alpha}}\left[\dfrac{\alpha(1-x) - (1-\alpha)x}{x(1-x)}\right]$

$= -\dfrac{Kx^\alpha(1-x)^{1-\alpha}}{1+Kx^\alpha(1-x)^{1-\alpha}}\left(\dfrac{\alpha-x}{x(1-x)}\right)$

Note,

$1 + Kx^\alpha(1-x)^{1-\alpha} = \dfrac{1}{y}$

$\Rightarrow Kx^\alpha(1-x)^{1-\alpha} = \dfrac{1}{y} - 1 = \dfrac{1-y}{y}$

$\therefore \dfrac{1}{y}\dfrac{dy}{dx} = -\dfrac{1-y}{y}\cdot y\dfrac{\alpha-x}{x(1-x)}$

$\Rightarrow \dfrac{dy}{dx} = \dfrac{(y-1)y}{x(1-x)}(\alpha-x)$

(note) Exploding Gradient: $\quad x = 0, 1;$ $\qquad\qquad$ input needs to be normalized

Vanishing condition: $\quad \alpha = x$ (critical point)

$\alpha > x,\ \dfrac{dy}{dx} < 0;\qquad \alpha < x,\ \dfrac{dy}{dx} > 0$