

3/4/25

Gradient of Vector Valued Functions
or
Lee :- Jacobian matrix.

(1)

We know the defn of partial deriv. of a scalar valued function.

If $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a vector valued function where $n, m \geq 1$

& $x \in \mathbb{R}^n$, then $f(x)$ can be written as

$$f(x) = (f_1(x), f_2(x), \dots, f_m(x))$$

where ~~each~~ $f_i(x)$ denotes the i^{th} component of $f(x)$ for $1 \leq i \leq m$.

$$x = (x_1, x_2, \dots, x_n) = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = [x_1 \ x_2 \ \dots \ x_n]^T$$

we will use column vector notation for vectors, because we will deal with matrices. So $\mathbb{R}^n \cong M_{n \times 1}(\mathbb{R})$.

$$\text{Further } f(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_m(x) \end{bmatrix} = [f_1(x) \ \dots \ f_m(x)]^T \in \mathbb{R}^m.$$

$$\& f = [f_1 \ f_2 \ \dots \ f_m]^T$$

Defn (Partial deriv. of ~~f~~ \bar{f})

$$\frac{\partial f}{\partial x_i} = \left[\frac{\partial f_1}{\partial x_i} \ \frac{\partial f_2}{\partial x_i} \ \dots \ \frac{\partial f_m}{\partial x_i} \right]^T$$

Recall chain rule. If $w = f(x_1, x_2, x_3)$ is a scalar function in $D \subseteq \mathbb{R}^3$ which is $C^1(D)$ (ie all first ~~parti~~ order partials exist in D & are cont. in D) & x_1, x_2 & x_3 are themselves functions of u_1, u_2 , ie

$$x_1 = x_1(u_1, u_2)$$

$$x_2 = x_2(u_1, u_2)$$

$$x_3 = x_3(u_1, u_2)$$

each of which ~~satisfy~~ $x_i \in C^1(B)$ &, where B is in u_1, u_2 plane

where if $(u_1, u_2) \in B$, then $(x_1(u_1, u_2), x_2(u_1, u_2), x_3(u_1, u_2))$ $\in D$, then partial derivative of w wrt u_1 & u_2 can be written in terms of x_1, x_2 & x_3 as

$$\begin{aligned}\frac{\partial w}{\partial u_1} &= \frac{\partial w}{\partial x_1} \frac{\partial x_1}{\partial u_1} + \frac{\partial w}{\partial x_2} \frac{\partial x_2}{\partial u_1} + \frac{\partial w}{\partial x_3} \frac{\partial x_3}{\partial u_1}, \\ \text{&} \quad \frac{\partial w}{\partial u_2} &= \frac{\partial w}{\partial x_1} \frac{\partial x_1}{\partial u_2} + \dots + \frac{\partial w}{\partial x_3} \frac{\partial x_3}{\partial u_2} \\ \underbrace{\left[\begin{array}{cc} \frac{\partial w}{\partial u_1} & \frac{\partial w}{\partial u_2} \end{array} \right]}_{\nabla_u w} &= \underbrace{\left[\begin{array}{ccc} \frac{\partial w}{\partial x_1} & \frac{\partial w}{\partial x_2} & \frac{\partial w}{\partial x_3} \end{array} \right]}_{\nabla_x w} \underbrace{\left[\begin{array}{cc} \frac{\partial x_1}{\partial u_1} & \frac{\partial x_1}{\partial u_2} \\ \frac{\partial x_2}{\partial u_1} & \frac{\partial x_2}{\partial u_2} \\ \frac{\partial x_3}{\partial u_1} & \frac{\partial x_3}{\partial u_2} \end{array} \right]}_{\boxed{\text{Mat}}}\end{aligned}$$

where $h(u_1, u_2) = (x_1(u_1, u_2), x_2(u_1, u_2), x_3(u_1, u_2))$. The mat

To see chain rule as matrix multiplication, it is useful to see gradient as row vectors instead of column vectors.

Defn (Gradient of vector valued fns)

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ & $x \in \mathbb{R}^n$. Then, define

$$\nabla_x f = \left[\begin{array}{ccc} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \frac{\partial f}{\partial x_n} \end{array} \right]$$

From defn of pd for vector valued f , we rewrite this as

$$= \left[\begin{array}{ccc} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & & \frac{\partial f_m}{\partial x_n} \end{array} \right]$$

This matrix is called the gradient of f or Jacobian of f .

$$\text{for } a \in \mathbb{R}^n, \nabla_x f(a) = \left[\frac{\partial f_i}{\partial x_j}(a) \right]_{i,j}$$

Eg 1) a) find the gradient of $f(x) = Ax$ where (3)

$$f: \mathbb{R}^N \rightarrow \mathbb{R}^M \quad \& \quad A \in M_{M \times N}(\mathbb{R}).$$

Soln:- To compute $\nabla_x f$, we need to know $\frac{\partial f_i}{\partial x_j}$ where $f = [f_1 \dots f_M]^T$
 $\& x = [x_1 \ x_2 \ \dots \ x_N]^T$.

$$f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_M(x) \end{bmatrix} = Ax = \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1N} \\ A_{21} & A_{22} & \dots & A_{2N} \\ \vdots & \vdots & & \vdots \\ A_{M1} & A_{M2} & \dots & A_{MN} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

$$\Rightarrow f_i(x) = A_{1i}x_1 + \dots + A_{Ni}x_N = \sum_{j=1}^N A_{ij}x_j$$

$$\text{More generally, } f_i(x) = \sum_{j=1}^N A_{ij}x_j$$

$$\therefore \frac{\partial f_i}{\partial x_j} = A_{ij} \quad \# \quad 1 \leq i \leq M \quad \& \quad 1 \leq j \leq N.$$

$$\therefore \nabla_x f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_N} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \frac{\partial f_M}{\partial x_2} & \dots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = A$$

$\& b \text{ (constant)} \in \mathbb{R}^M$

Ques 9/4/25 ie, $\nabla_x f$ evaluated at any $x \in \mathbb{R}^N$ is A itself. $\nabla_x f = A$

Lee 32 Ex. If $f(x) = Ax + b$ where $x \in \mathbb{R}^N$; $A \in M_{M \times n}(\mathbb{R})$, show that ~~if~~ $\nabla_x f = A$

(Supervised Learning ~~task~~) - Applying Machine Learning $\hat{y} = \underline{\Phi}\underline{\theta}$ to predict y

Eg 2) Let us consider the linear model $\hat{y} = \underline{\Phi}\underline{\theta}$ where $\underline{\theta}$ is a parameter vector, $\underline{\theta} \in \mathbb{R}^D \cong M_{D \times 1}(\mathbb{R})$ & $\underline{\Phi} \in M_{N \times D}(\mathbb{R})$ are input features

Let $e(\underline{\theta}) = \hat{y} - \underline{\Phi}\underline{\theta}$ & $L(e) = \|e\|^2$. Find $\nabla_{\underline{\theta}} L$ wrt $\underline{\theta}$.

Soln:- $\hat{y} \in \mathbb{R}^N$ & hence $e(\underline{\theta}) \in \mathbb{R}^N$. So, $R^D \xrightarrow{e} R^N$
 $\underline{\theta} \rightarrow \hat{y} - \underline{\Phi}\underline{\theta}$.

$$L(e) = \|e\|^2 \Rightarrow \mathbb{R}^N \xrightarrow{L} \mathbb{R}$$

$$x \rightarrow \|x\|^2.$$
(4)

$$\mathbb{R}^D \xrightarrow{e} \mathbb{R}^N \xrightarrow{L} \mathbb{R}$$

$$\theta \rightarrow y - \frac{1}{2}\theta \rightarrow \|y - \frac{1}{2}\theta\|^2$$

$$\nabla_{\theta} L = \nabla_e^T L \times \nabla_{\theta} e \quad \begin{matrix} \text{Gradient of vector valued } e \\ \text{or Jacobian} \end{matrix}$$

$$\nabla_{\theta} e = \left[\begin{array}{ccc} \frac{\partial e_1}{\partial \theta_1} & \dots & \frac{\partial e_1}{\partial \theta_D} \\ \vdots & & \vdots \\ \frac{\partial e_N}{\partial \theta_1} & \dots & \frac{\partial e_N}{\partial \theta_D} \end{array} \right]$$

~~$e(\theta) = y$~~ $e(\theta) = y - \frac{1}{2}\theta$ is a linear fn & hence
by prev. problem $(\nabla_{\theta} e)(a) = -\frac{1}{2}$.

Now,

$$\nabla_e^T L = \left[\begin{array}{ccc} \frac{\partial L}{\partial e_1} & \dots & \frac{\partial L}{\partial e_N} \end{array} \right] \quad \begin{matrix} L(e_1, \dots, e_N) \\ = e_1^2 + e_2^2 + \dots + e_N^2 \end{matrix}$$

$$= [2e_1 \ 2e_2 \ \dots \ 2e_N]$$

$$= 2e^T = 2(y - \frac{1}{2}\theta)^T = 2(y^T - \theta^T \frac{1}{2}^T)$$

$$\Rightarrow \nabla_{\theta}^T L = -2 \underbrace{\bullet (y^T - \theta^T \frac{1}{2}^T)}_{1 \times N} \underbrace{\frac{1}{2}}_{N \times D}$$

Supervised learning problem

Suppose you have m many training examples say (x^i, y^i) $1 \leq i \leq m$ where $x^i \in \mathbb{R}^n$. Say for example

| | Living area (x_1) | No of bedrooms (x_2) | ... | n^{th} feature (x_n) | Price of the house (y) |
|-----|--------------------------|-----------------------------|-----|-------------------------------|-------------------------------|
| 1 | x_1^1 | x_2^1 | .. | x_n^1 | y^1 |
| 2 | x_1^2 | x_2^2 | .. | x_n^2 | y^2 |
| .. | .. | .. | .. | .. | .. |
| m | x_1^m | x_2^m | .. | x_n^m | y^m |

We want to find a nice function $h: \mathbb{R}^n \rightarrow \mathbb{R}$ such that $h(\vec{x})$ (where $\vec{x} = (x_1, x_2, \dots, x_n)$) is a good predictor for the corresponding value of \tilde{y} .

These problems are called regression problems. We will consider the simplest model among those called linear regression problems where we approximate y using a linear function of \vec{x} .

$$\text{i.e., } h(\vec{x}) = h(x_1, x_2, \dots, x_n) = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

where θ_i -s are called parameters or weights. Note that

$$h\left(\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}\right) = \vec{x}^T \theta \quad \begin{array}{l} \text{see } \vec{x} \text{ as column vectors} \\ (\text{A} \theta) \end{array}$$

~~Let~~ x^i for each $1 \leq i \leq m$ is also a column vector

$$\Rightarrow h\left(\begin{bmatrix} x_1^i \\ x_2^i \\ \vdots \\ x_n^i \end{bmatrix}\right) = (x^i)^T \theta$$

(4)

$$\left[\begin{array}{c} h(x_1) \\ h(x_2) \\ \vdots \\ h(x_m) \end{array} \right] = \left[\begin{array}{c} (x_1)^T \theta \\ (x_2)^T \theta \\ \vdots \\ (x_m)^T \theta \end{array} \right]$$

Here $x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}$

$$= \Phi \in n \times m \text{ matrix}$$

$y^i - \frac{h(x^i)}{\theta}$ will be diff b/w actual value & predicted value, i.e., the error. Here $\theta \in \mathbb{R}^n$ & $x^i \in \mathbb{R}^n$

$$e(\theta) = \begin{bmatrix} y^1 - \frac{h(x^1)}{\theta} \\ y^2 - \frac{h(x^2)}{\theta} \\ \vdots \\ y^m - \frac{h(x^m)}{\theta} \end{bmatrix} = \begin{bmatrix} y^1 - (x^1)^T \theta \\ y^2 - (x^2)^T \theta \\ \vdots \\ y^m - (x^m)^T \theta \end{bmatrix}$$

$$\begin{bmatrix} (x^1)^T \theta \\ (x^2)^T \theta \\ \vdots \\ (x^m)^T \theta \end{bmatrix} = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_n^1 \\ x_1^2 & x_2^2 & \dots & x_n^2 \\ \vdots & & & \\ x_1^m & x_2^m & \dots & x_n^m \end{bmatrix}_{m \times n} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}_{n \times 1}$$

$$\therefore \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^m \end{bmatrix} - \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_n^1 \\ x_1^2 & x_2^2 & \dots & x_n^2 \\ \vdots & & & \\ x_1^m & x_2^m & \dots & x_n^m \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

or

$$e(\theta) = y - \underbrace{\Phi \theta}_{\substack{\mathbb{R}^m \\ M_{m \times n}(\mathbb{R})}} \rightarrow \mathbb{R}^n$$

input features

In such problems, we are interested in minimizing the error fn $e(\theta)$

i.e., $\min \|e(\theta)\|^2$ for a suitable θ . ATM: finding θ such that $J(\theta)$ is min

Gradients of Matrices

Suppose we have ~~two~~ ^{two} matrices $A \in M_{m \times n}(R)$ & $B \in M_{p \times q}(R)$

& wish to compute the gradient of A w.r.t B , (think of B as a ~~real~~ variable, but matrix) & A is dependent on B . ~~i.e., $A = f(B)$~~ the resulting Jacobian ~~is also a matrix~~

Then, $A_{11}, \dots, A_{1n}, | A_{21}, \dots, A_{2n} | \dots, | A_m, \dots, A_{mn}$ are ~~the~~ component functions of A & each of them ~~were~~ is a function of pq many variables. Its gradient would be the row vector

$$\nabla_B A_{11} = \left[\begin{array}{cccc|ccccc} \frac{\partial A_{11}}{\partial B_{11}} & \frac{\partial A_{11}}{\partial B_{12}} & \dots & \frac{\partial A_{11}}{\partial B_{1q}}, & \frac{\partial A_{11}}{\partial B_{21}} & \dots & \frac{\partial A_{11}}{\partial B_{2q}}, & \dots, & \frac{\partial A_{11}}{\partial B_{p1}} & \dots & \frac{\partial A_{11}}{\partial B_{pq}} \end{array} \right]$$

Sometimes it is also written

If we wish to write it as matrix

$$\left[\begin{array}{cccc} \frac{\partial A_{11}}{\partial B_{11}} & \frac{\partial A_{11}}{\partial B_{12}} & \dots & \frac{\partial A_{11}}{\partial B_{1q}}, \\ \frac{\partial A_{11}}{\partial B_{21}} & \dots & \frac{\partial A_{11}}{\partial B_{2q}}, \\ \vdots & & & \\ \frac{\partial A_{11}}{\partial B_{p1}} & & \frac{\partial A_{11}}{\partial B_{pq}} & \end{array} \right]$$

Each component's gradient is a pq matrix & since there are mn component functions, the Jacobian would be $m \times n \times pq$ matrix, i.e. a four-dimensional tensor] whose entries are given as $J_{ijkl} = \frac{\partial A_{ij}}{\partial B_{kl}}$.

This form of representation of gradient of matrices is highly inconvenient & since $M_{m \times n}(R) \cong R^{\frac{mn}{1}}$, we see that f is equivalently a function as below :- $R^{\frac{pq}{1}} \xrightarrow{f} R^{mn}$

i.e. Matrix is reshaped into vector
($\tilde{B} \in R^{pq}$)

Eg 1) (Gradient of vectors with respect to matrices) (6)

Consider the map $A \xrightarrow{f} Ax \quad (Ax \in \mathbb{R}^M)$

where $A \in \mathbb{R}^{M \times N}$ matrix & $x \in \mathbb{R}^N$

Ax - vector
 A - matrix

1) What is the dimension of $\nabla_A f$?

2) Compute the gradient of f w.r.t A .

Soln: Since $Ax \in M_{m \times 1}(\mathbb{R})$, $F: \mathbb{R}^{MN} \rightarrow \mathbb{R}^M$

~~Each component function f_i~~ & hence $\nabla_A F \in M_{M \times (M \times N)}$

2) Let $f_1(A), f_2(A), \dots, f_M(A)$ be its component functions.

What is $\frac{\partial f_i}{\partial A_{kj}}$? where $1 \leq i \leq M, 1 \leq k \leq M$
 $1 \leq j \leq N$

Let $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$.

$Ax \in \mathbb{R}^M$. i^{th} entry of Ax , ie, $(Ax)_i$ is given by

$$Ax = \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1N} \\ A_{i1} & A_{i2} & \dots & A_{iN} \\ \vdots & & & \\ A_{m1} & A_{m2} & \dots & A_{mN} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

$$(Ax)_i = [A_{i1} \ A_{i2} \ \dots \ A_{iN}] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \sum_{k=1}^N A_{ik} x_k$$

$$\therefore \frac{\partial f_i}{\partial A_{kj}} = \frac{\partial (Ax)_i}{\partial A_{kj}} = \begin{cases} x_j & \text{if } k=i \\ 0 & \text{if } k \neq i \end{cases}$$

$$\therefore \nabla_A f = \begin{bmatrix} \frac{\partial f_1}{\partial A_{11}} & \frac{\partial f_1}{\partial A_{12}} & \dots & \frac{\partial f_1}{\partial A_{1N}} & \frac{\partial f_1}{\partial A_{21}} & \dots & \frac{\partial f_1}{\partial A_{2N}} & \dots & \frac{\partial f_1}{\partial A_{M1}} & \dots & \frac{\partial f_1}{\partial A_{MN}} \\ \vdots & & & & & & & & & & & \\ \frac{\partial f_M}{\partial A_{11}} & \dots & \frac{\partial f_M}{\partial A_{1N}} & \frac{\partial f_M}{\partial A_{21}} & \dots & \frac{\partial f_M}{\partial A_{2N}} & \dots & \frac{\partial f_M}{\partial A_{M1}} & \dots & \frac{\partial f_M}{\partial A_{MN}} \end{bmatrix}_{M \times (M \times N)}$$

$$= \left[\begin{array}{cccc|cc|cc|cc} x_1 & x_2 & \dots & x_N & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & & & & x_1 & x_2 & \dots & x_N & 0 & \\ & & & & & & & & \ddots & \\ & & & & & & & & & x_1 & x_2 & \dots & x_N \end{array} \right]$$

Follow (Diezengoff, Faisal, Cheng Soon Ong) : Mathematics for Machine Learning.

Lee 33 :- Eg 2) (Gradient of f. Matrices w.r.t matrices)

10/4/25 Let $\overset{A}{\underset{\in M_{M \times N}(R)}{\bullet}}$ & $f(A) = A^T A = K$

i.e., $f : M_{M \times N}(R) \rightarrow M_{N \times N}(R)$.

What is $\nabla_A K$?

Soln :- Let A_i denote i^{th} row of A & A_{ij} denote i, j^{th} entry of A .

Then ~~$A_{ij} = A$~~ $K_{pq} = (A^T A)_{pq}$

$$= (p^{\text{th}} \text{ row of } A^T) \cdot (q^{\text{th}} \text{ col. of } A) = (p^{\text{th}} \text{ col. of } A) \cdot (q^{\text{th}} \text{ col. of } A)$$

$$= A_p \cdot A_q = \sum_{m=1}^M A_{pm} A_{mq}$$

$$A = \left(\begin{array}{c|ccc} A_{1p} & & & \\ A_{2p} & & & \\ \vdots & & & \\ A_{mp} & & & \\ \hline & A_{1N} & \dots & A_{mN} \end{array} \right) \quad q=p \Rightarrow K_{pp} = \sum_{m=1}^M (A_{mp})^2$$

$$\therefore \frac{\partial K_{pq}}{\partial A_{ij}} = \begin{cases} 2A_{ip} & p=q \text{ & } j=p \\ -A_{iq} & p \neq q \text{ but } j=p \\ A_{ip} & p \neq q \text{ but } j=q \\ 0 & \text{elsewhere} \end{cases}$$

$$\nabla_A f \in R^{(N \times N) \times (M \times N)}$$

(8)

Suppose 1
Here, we will follow the convention that

$$f_1 = K_{11}$$

$$f_2 = K_{21}$$

:

$$f_N = K_{N1}$$

:

$$f_{N^2} = K_{NN}$$

i.e., component functions are ordered along the columns first & then rows.

Eg. If $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \\ A_{31} & A_{32} \end{bmatrix}_{3 \times 2}$

then $K = A^T A = \begin{bmatrix} (A_{11})^2 + (A_{21})^2 + (A_{31})^2 = f_1 \\ A_{11} A_{12} + A_{21} A_{22} + A_{31} A_{32} = f_3 \\ A_{12} A_{11} + A_{22} A_{21} + A_{32} A_{31} = f_2 \\ (A_{12})^2 + (A_{22})^2 + (A_{32})^2 = f_4 \end{bmatrix}$

$$\nabla_A f_1 = \begin{bmatrix} \frac{\partial f_1}{\partial A_{11}} & \dots & \frac{\partial f_1}{\partial A_{32}} \end{bmatrix} = \begin{bmatrix} 2A_{11} & 0 & 2A_{21} & 0 & 2A_{31} & 0 \end{bmatrix}$$

$$\nabla_A f_2 = \nabla_A f_3 = \begin{bmatrix} A_{12} & A_{11} & A_{22} & A_{21} & A_{32} & A_{31} \end{bmatrix}$$

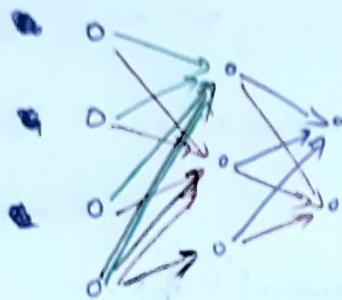
:

$$\nabla_A f_4 = [0 \ 2A_{12} \ 0 \ 2A_{22} \ 0 \ 2A_{32}]$$

~~$\nabla_A f \in R^{(N \times N) \times (M \times N)}$~~

$$\nabla_A f = \begin{bmatrix} 2A_{11} & 0 & 2A_{21} & 0 & 2A_{31} & 0 \\ A_{12} & A_{11} & A_{22} & A_{21} & A_{32} & A_{31} \\ A_{12} & A_{11} & A_{22} & A_{21} & A_{32} & A_{31} \\ 0 & 2A_{12} & 0 & 2A_{22} & 0 & 2A_{32} \end{bmatrix}_{(2 \times 2) \times (3 \times 2)}$$

Anyone familiar with neural networks would have seen the following diagram.



Here, circles denote neurons & lines correspond to the network of connection between them.

Neural network is fundamentally a mathematical function that takes an input (which could be vector) X & gives an output Y (vector)

Case 1 :-

$$a^{(0)} \rightarrow a^{(1)}$$

Define $a^{(1)} = \sigma(wa^{(0)} + b)$ where $b, w \in \mathbb{R}$

σ is a function called activation function.

a activity

w weight

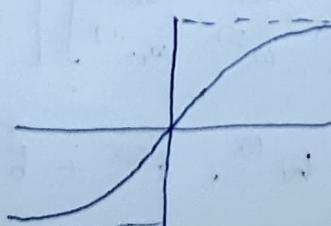
b bias

When the stimulation of one or more ~~should~~ neurons goes beyond a threshold, the neuron on which ~~it~~ interacts gets "activated".

An example of a function that has this thresholding property is the hyperbolic

tangent $(\text{Tanh})(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

Note : If $x \rightarrow \infty$ $\text{Tanh} x = 1$ & If $x \rightarrow -\infty$ $\text{Tanh} x = -1$



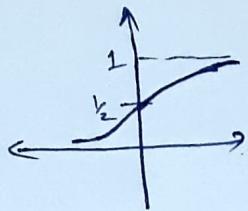
which have this thresholding property
Class of such functions ~~in general~~ are called sigmoids

(10)

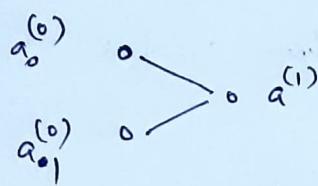
E.g. $f(x) = \frac{1}{1+e^{-x}}$

$$\lim_{x \rightarrow \infty} f(x) = 1$$

$$\lim_{x \rightarrow -\infty} f(x) = 0$$



Case 2 :- (Input) 2 neurons activating 1 neuron (Output)

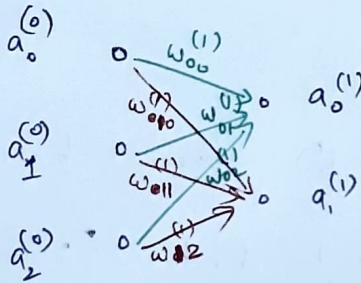


b^1 - bias of neuron 1

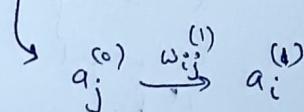
$$a^{(1)} = \sigma(\omega_0 a_0^{(0)} + \omega_1 a_1^{(0)} + b^1)$$

Layer 0 Layer 1

Case 3 :- Input 3 neurons & Output 2 neurons



$w_{ij}^{(1)}$ is the weight / line connecting neuron $a_i^{(0)}$ of 1st layer (i.e., $a_i^{(1)}$) & $a_j^{(0)}$ of 0th layer (i.e., $a_j^{(0)}$).



$b_0^{(1)}$ is the bias of neuron $a_0^{(1)}$ in 1st layer.

$$b_1^{(1)} \leftarrow " " \rightarrow a_1^{(1)} \leftarrow " " \rightarrow$$

$$a_0^{(1)} = \sigma(w_{00}^{(1)} a_0^{(0)} + w_{01}^{(1)} a_1^{(0)} + w_{02}^{(1)} a_2^{(0)} + b_0^{(1)})$$

$$a_1^{(1)} = \sigma(w_{10}^{(1)} a_0^{(0)} + w_{11}^{(1)} a_1^{(0)} + w_{12}^{(1)} a_2^{(0)} + b_1^{(1)})$$

$$\begin{aligned} \vec{a}^{(1)} &= \begin{bmatrix} a_0^{(1)} \\ a_1^{(1)} \end{bmatrix} = \sigma \left(\begin{bmatrix} w_{00}^{(1)} & w_{01}^{(1)} & w_{02}^{(1)} \\ w_{10}^{(1)} & w_{11}^{(1)} & w_{12}^{(1)} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ a_2^{(0)} \end{bmatrix} + \begin{bmatrix} b_0^{(1)} \\ b_1^{(1)} \end{bmatrix} \right) \\ &= \sigma \left(W^{(1)} \cdot \vec{a}^{(0)} + \vec{b}^{(1)} \right) \end{aligned}$$

This is an example of a single layered network with 3 inputs & 2 output & can be represented as:

$$a^{(1)} = \sigma(W^{(1)} \cdot a^{(0)} + b^{(1)}).$$

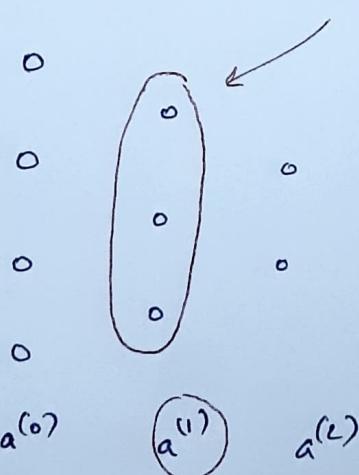
More generally, when there are n inputs & m outputs, we can represent it as

$$a^{(1)} = \sigma(W^{(1)} \cdot a^{(0)} + b^{(1)}) \text{ where}$$

$$W^{(1)} = \begin{bmatrix} w_{0,0}^{(1)} & w_{0,1}^{(1)} & \dots & w_{0,n-1}^{(1)} \\ w_{1,0}^{(1)} & w_{1,1}^{(1)} & \dots & w_{1,n-1}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m-1,0}^{(1)} & w_{m-1,1}^{(1)} & \dots & w_{m-1,n-1}^{(1)} \end{bmatrix}$$

$$\vec{a}^{(0)} = \begin{bmatrix} a_0^{(0)} \\ \vdots \\ a_{n-1}^{(0)} \end{bmatrix} \quad \& \quad \vec{b}^{(1)} = \begin{bmatrix} b_0^{(1)} \\ b_1^{(1)} \\ \vdots \\ b_{m-1}^{(1)} \end{bmatrix}.$$

Neural networks may have ~~more than one~~^{one or more} layers of neurons between the input & the output. These are called hidden layers. (The outputs of ~~the~~ a ~~become~~ ~~hidden layer here to the~~ next ~~layer~~ layer)



Neurons in layer (1) becomes inputs for neurons in layer (2).

$$\vec{a}^{(1)} = \sigma(W^{(1)} \cdot \vec{a}^{(0)} + \vec{b}^{(1)})$$

$$\vec{a}^{(2)} = \sigma(W^{(2)} \cdot \vec{a}^{(1)} + \vec{b}^{(2)}).$$

More generally,

$$a^{(L)} = \sigma(w^{(L)} \cdot a^{(L-1)} + b^{(L)})$$

Such neural networks are sometimes called "feed forward" neural networks.

Persuading a neural network to do tasks (like image recognition etc) now becomes just a job of ~~teaching it~~ the right weights & biases.

Eg 1)

To train the network to give a not function

(ie, if you input 1, it returns 0

if you input 0, it returns 1.)

Q. Let us use $\sigma(x) = \tanh x$.

$$a^{(1)} = \sigma(a^{(0)} w^{(1)} + b^{(1)})$$

among which of the following weight & bias gives the best result?

a) $w^{(1)} = 10, b^{(1)} = 0$

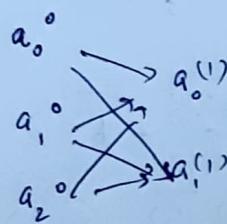
b) $w^{(1)} = -3, b^{(1)} = 0$

c) $w^{(1)} = 3, b^{(1)} = 1$

d) $w^{(1)} = -5, b^{(1)} = 5$

e) $w^{(1)} = 0, b^{(1)} = 5$

Eg 2) ~~at~~ Consider



$$\text{Let } W^{(1)} = \begin{bmatrix} -2 & 4 & -1 \\ 6 & 0 & -3 \end{bmatrix} \quad \& \quad \vec{b}^{(1)} = \begin{bmatrix} 0.1 \\ 0.25 \end{bmatrix} \quad (13)$$

$$\text{in } \vec{a}^{(1)} = \sigma(W^{(0)} \cdot \vec{a}^{(0)} + \vec{b}^{(0)}) \quad \text{with}$$

$$\& \quad \vec{a}^{(0)} = \begin{bmatrix} 0.3 \\ 0.4 \\ 0.1 \end{bmatrix}$$

Find $\vec{a}^{(1)}$.

$$\text{Soln: } W^{(0)} \cdot \vec{a}^{(0)} = \begin{bmatrix} -2 & 4 & -1 \\ 6 & 0 & -3 \end{bmatrix} \begin{bmatrix} 0.3 \\ 0.4 \\ 0.1 \end{bmatrix} = \begin{bmatrix} -0.6 + 1.6 - 0.1 \\ 1.8 - 0.3 \end{bmatrix} = \begin{bmatrix} 0.9 \\ 1.5 \end{bmatrix}$$

$$\therefore W^{(0)} \cdot \vec{a}^{(0)} + \vec{b}^{(0)} = \begin{bmatrix} 1.0 \\ -1.0 \end{bmatrix}$$

$$\Rightarrow \vec{a}^{(1)} = \begin{bmatrix} \sigma(1.0) \\ \sigma(-1.0) \end{bmatrix} = \begin{bmatrix} 0.76 \\ -0.76 \end{bmatrix}.$$