

Father of Soft computing : Lotfi A. Zadeh sir.
Fuzzy logic

* Genetic Algorithm → John Holland

* Neural Network → Artificial NN : Geoffrey Hinton.
DL

Convolutional NN : Yann LeCun.

* Godfathers of AI → Yoshua Bengio

Geoffrey Hinton

Yann LeCun.

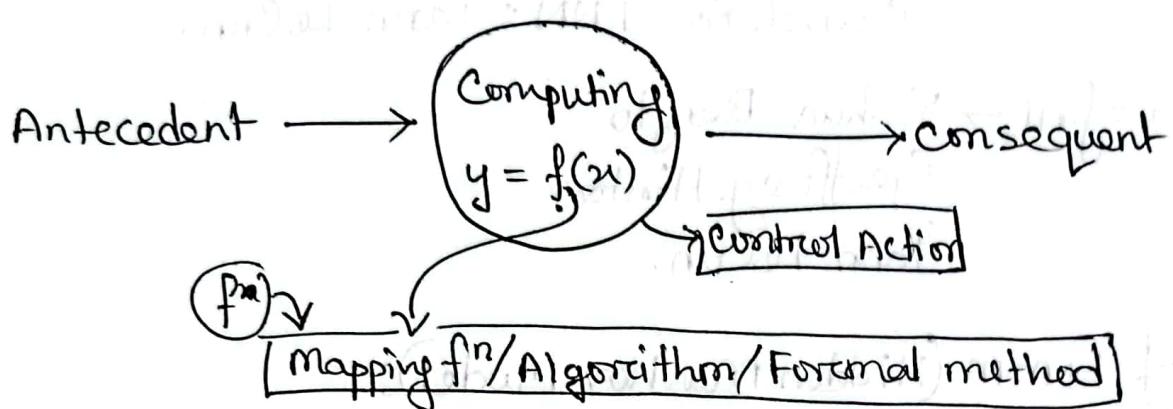
* Father of HMM (Hidden Markov Model) :

* Donald Knuth.

* Computation : Input- নিয়ে operators দিয়ে process করে outcome হয় এবং

$$f(x) \rightarrow y.$$

* Control Unit :



* Computing :

- process / act of calculation .
- action of mathematical calculation
- activity that uses Computers .
- development of hw and sw .
- critical , integral component

Accuracy - overall performance

Precision - কার্তুণ
precise - প্রতিবাধ
accuracy কর্মকে
না প্রাপ্ত করে ।

Brain vs Computer

1. 10^{10} neurons	10^8 transistors
2. Element size 10^{-6} m	10^{-6} m
3. Energy 30W	30W (CPU)
4. Speed 10^2 Hz	10^{12} Hz
5. Parallel computation Distributed "	Serial Centralized .
6. Efficiency 10^{-16} J.	10^{-6} J.
7. Fault Tolerant	Yes — No
8. Learns	Yes — a little

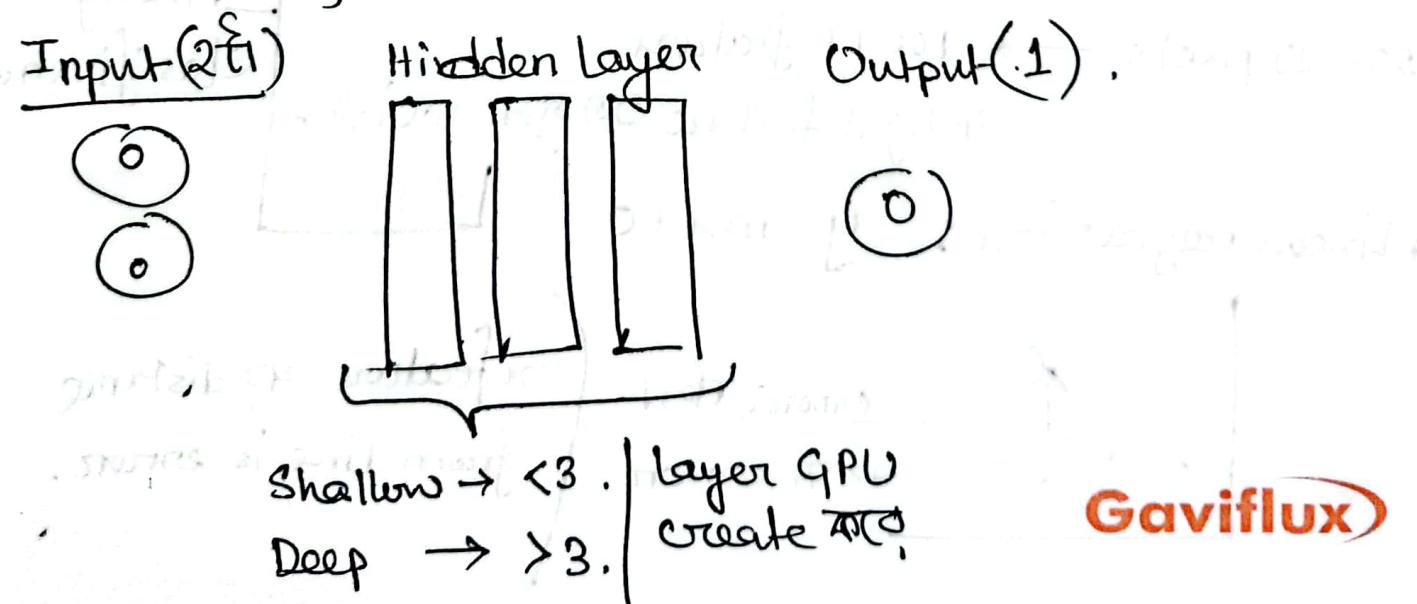
ML → data কর

DL → more data

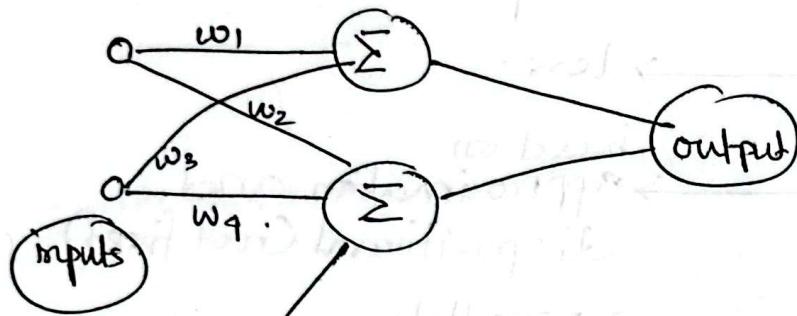
- * Hard Computing vs Soft computing (Computational Intelligence)
1. Precise result → imprecision, approximate, uncertain results
 2. Unambiguous → ambiguous control action
 3. Formally defined → adaptive
 4. Computation time more → less
 5. Crisp sw, binary logic, numerical sys → based on approximation and dispositional (not fixed)
 6. sequential computation → parallel
 7. works on exact data → on ambiguous & noisy data
 8. Uses two-valued logic → multivalued logic
 9. settled computing → incorporates randomness
 10. requires programs to be written → emerges its own programs
 11. deterministic nature → stochastic (random)

Neural Network (works great on at least - 10K data)

* Inspired by brain neurons



* আমাদের target Hardware-এ dot problem dot product -এ নিয়ে আমা as matrix, তুম্হা dot product -এ আনলে super fast হবে।



- * Hyper parameters : fixed [Batch size, Epoch size, iteration, learning rate]
- * Parameters (weight, bias) : change করা possible ; updates, or is learned

connection
between
input - hidden
and hidden-
output .

Controls
orientation
of decision
boundary .
(aka hyperplane)

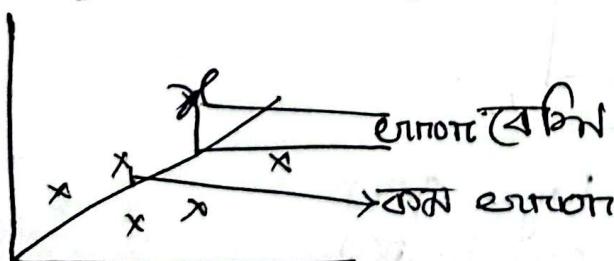
suppose phone বিষয়ে তাঁর
iphone নিয়ে না android
confused, যারে iphone wala
friend বললে iphone নিয়ে।
Iphone wala friend is bias.

* image এর pixels-এ feature এর রূপ :

20×20 pixels \rightarrow 400টি feature

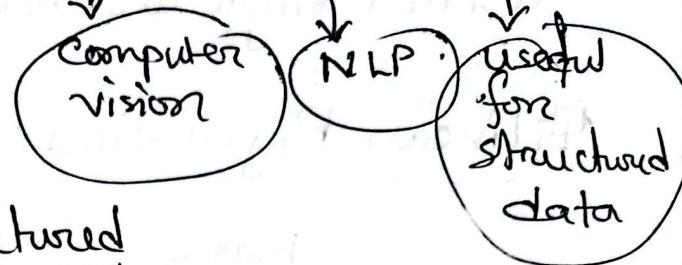
যতগুলো feature ততগুলো weight .

* Linear regression : $y = mx + c$



x feature - এর distance
from line is error .

④ Supervised learning - A CNN, RNN(recurrent), Standard NN
and Hybrid/custom NN we'll



* Structured vs Unstructured

1. Database, tables — Audio, text, images
2. More money
(companies big data prediction)

* Algorithms:

→ creative algorithms.

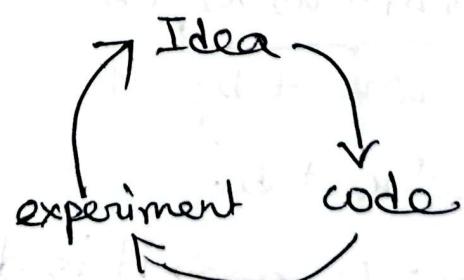
* ReLU, SIGMOID, TANH

↓
better than sigmoid.
(as helps with vanishing gradient problem)

* Why DL? → society के अनेकों data

- → यह algo generate IC को
- → GPU works faster.

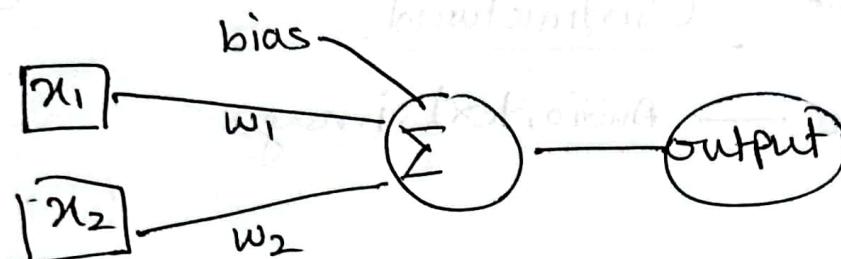
* Iterative training :



Logistic Regression: nonlinear data तिए ढंग करते

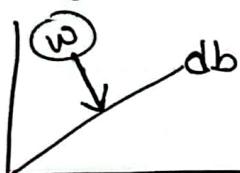
एक single unit neuron network (एक node वाला)

Linear Regression: works with linear data.



* linear sum equation $\rightarrow w_1x_1 + w_2x_2 + \text{bias} = \Sigma$.

* weight is perpendicular to db; $w \perp db$.



Decision boundary's equation:

(*) In vector $\rightarrow [w_1, w_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

so, db - का equation $= w_1x_1 + w_2x_2$

$$g(x) = \underbrace{w_1x_1 + w_2x_2}_{} + b.$$

$$= w^T x + b.$$

$$= w^T x + b.$$

if $w^T x + b = 0 \rightsquigarrow$ it is decision boundary

" $w^T x + b > 0 \dots \dots x \in \text{Class 1}$

" $w^T x + b < 0 \dots \dots x \in \text{Class 2}$

threshold value.

Facts of logistic Regression.

- ↳ supervised.
- ↳ output between 0-1.
- ↳ goal is to minimize error b/w prediction & actual data.

(VS) linear regression.

- ↳ data linearly separated
- but ↳ log R-1 nonlinearly separated.

Maths of Logistic Regression:

Let y = actual data

\hat{y} = predicted data

$$\hat{y} = P(y=1|x)$$

↳ probability of x feature of being class 1.

* Model always returns probability.

Mathematical parameters:

① Input $\rightarrow x \in \mathbb{R}_{\geq 0}^{n_x}$ ($n_x \in \mathbb{N}$ feature)
 ↳ real no.

② Training Label $\rightarrow y \in \{0, 1\}$
 ↳ labels.
 ↳ actual value

Feature যতগুলো weight ততগুলা \rightarrow so $w \in \mathbb{R}^{n \times 1}$

Bias $\rightarrow b \in \mathbb{R}$

So, input, weight, bias সব Real numbers.

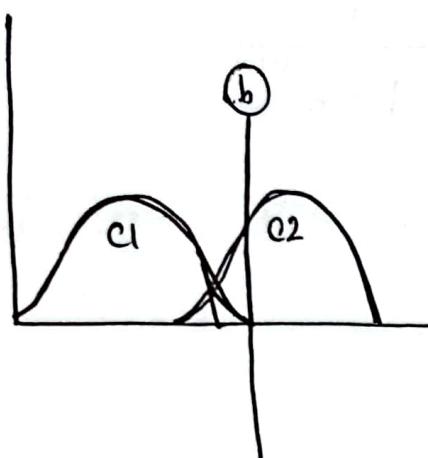
Output $\rightarrow \hat{y} = \sigma \cdot (w^T u + b)$

Sigmoid function: $s = \sigma(w^T u + b) = \sigma(z) = \frac{1}{1+e^{-z}}$

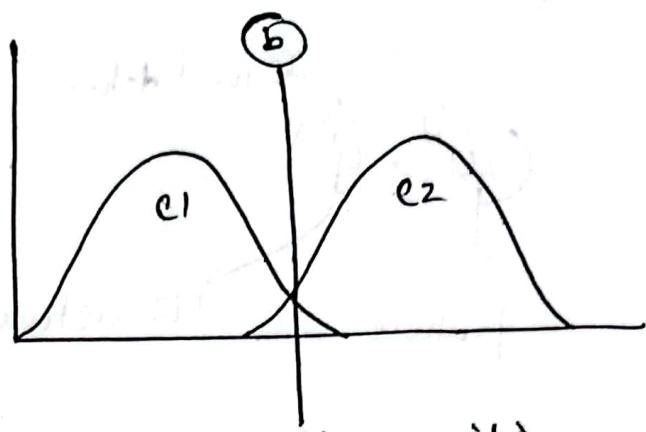
$w_i x_i + b$; b বাড়লে \hat{y} value বাঢ়ব

*

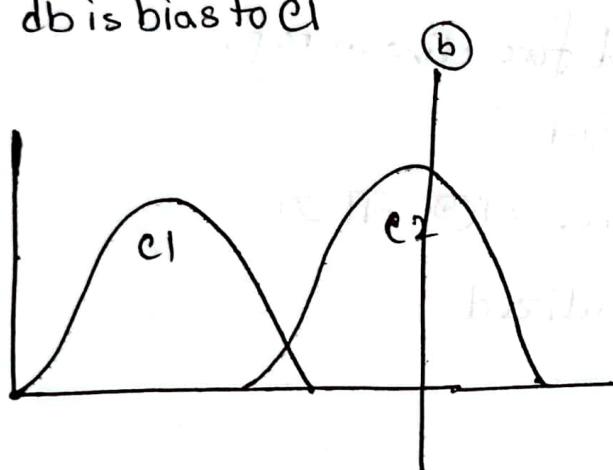
Bias positions in graph meanings:



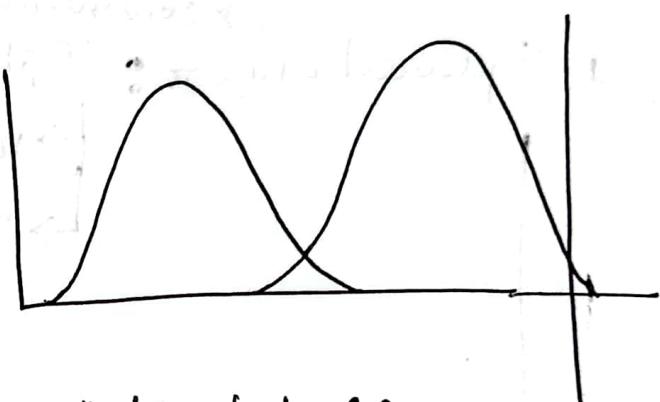
b is bias to c_1



better position



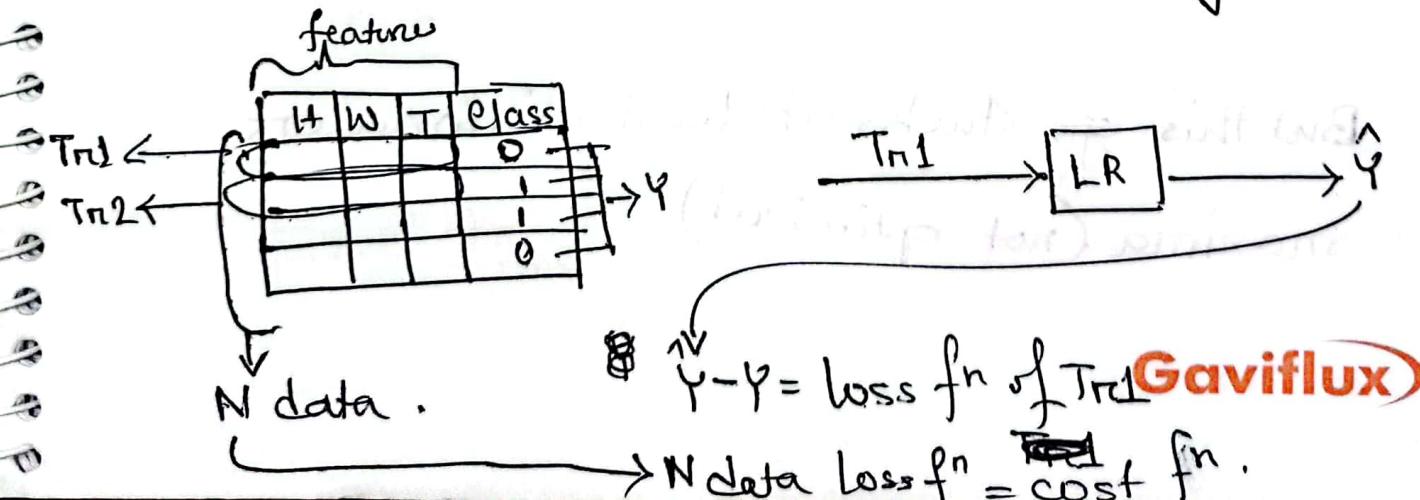
b is bias to c_2



b biased to c_2

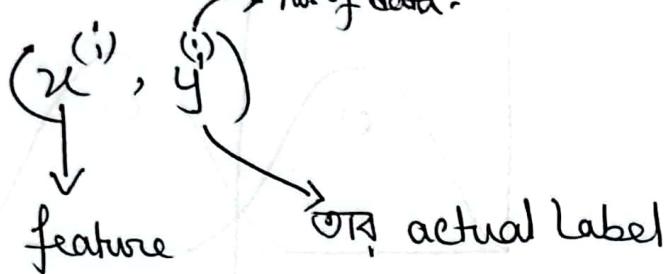
loss f_n : $(\hat{y} - y) = \text{loss } f_n \rightarrow$ gap between predicted and actual data (single training example - \hat{y} , error)

Cost f_n : full dataset \rightarrow loss \rightarrow arg. fn.



④ Cost fn $\rightarrow \hat{y}^{(i)} = \sigma(w^\top x^{(i)} + b)$, where

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



■ Loss fn:

1. Squared Error: \rightarrow commonly used for linear reg.

- negative মাত্রে না হয়
- To penalized

$$L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

\rightarrow square করে negative value avoid করতে.

$\rightarrow 2$ করে, square-এর পর value-কে extra 2.

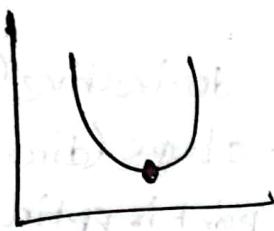
কাটাফাটি করে ফলাব জন্য.

But this ~~gives~~ stuck at local minima or maxima (not optimized)

Optimization \rightarrow getting the global maxima or minima; maximum/minimum f^n .

* Convex graph is better ~~as~~ as minima or maxima.

একটি



* Non-convex \rightarrow local -> তাঁগোনা মুসলিম পাঠ



2. Cross Entropy loss :-

$$L(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)}))$$

If $y^{(i)} = 1$; $L(\hat{y}^{(i)}, y^{(i)}) = -\log(\hat{y}^{(i)})$ where $\hat{y}^{(i)}$ and $\log(\hat{y}^{(i)})$ near to 1.

আব 0 - হলে 0 -এ কাঠে

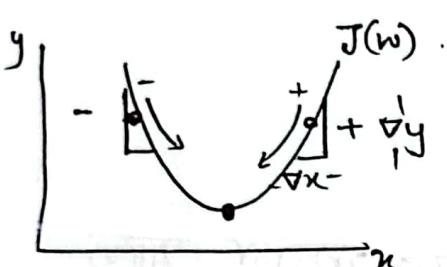
Cost fn \rightarrow avg of lossfn of the entire training set.

$$\begin{aligned} J(w, b) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})] \end{aligned}$$

We want to find, w and b to minimize cost fn $J(w, b)$

Gradient Descent :-

- cost fn is convex.
- initialize w and b with random numbers
(0,0 → in case of logistic reg) + जारी करने के लिए IMPROVE.
- यद्युपि fn convex - initialize अचूक मिला या ग्लोबल optimal or near to it जारी।



Gradient - derivative (slope of a fn at a point)
Descent - slope (direction of ↓)
if $f^n = 0$; point is optimal.

For update we use, $w(t+1) = w(t) - \alpha \frac{\partial J(w, b)}{\partial w}$

$$b = b - \alpha \frac{\partial J(w, b)}{\partial b}$$

Partial derivative

Necessity of GD:

solve this problem.

optimal solution $f_{\min}(x)$

"Premature optimization is the root of all evil" - Donald Knuth.

① → learning rate.
(how bigger step we choose at each ite).

② Computational Graph: Why we need it? → to represent full NN + explains why it is organized that way

③ Forward pass: Computes output

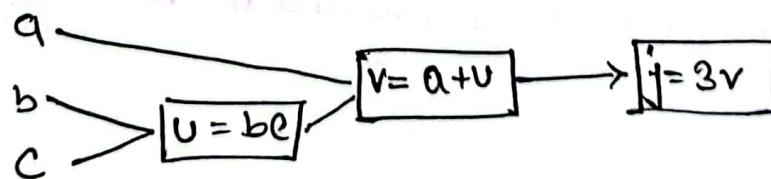
④ Backward pass: " derivatives + gradient to update w and b.

$$J(a, b, c) = 3(a + bc)$$

$$\textcircled{1} \quad U = bc$$

$$\textcircled{2} \quad V = a + bU$$

$$\textcircled{3} \quad j = 3V$$

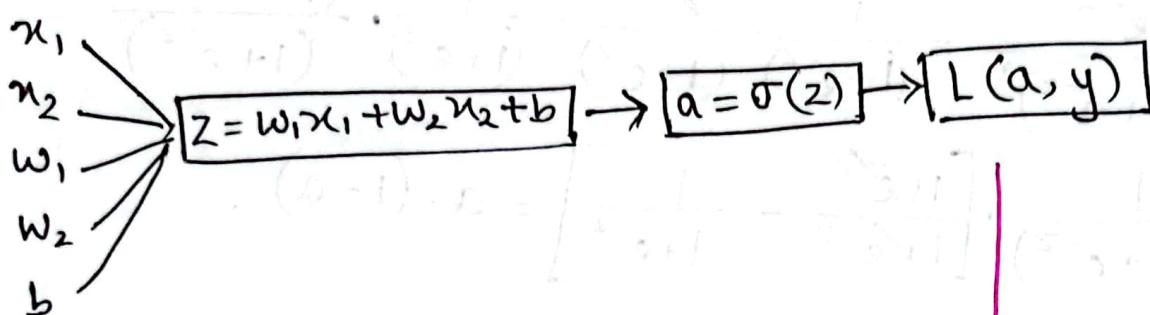


Forward:

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$L(a, y) = -(y \log(a) + (1-y) \log(1-a))$$



Backward:

$$dz = \frac{\partial L(a, y)}{\partial z} = \frac{\partial L}{\partial z}$$

deriv. of L with respect to z.

$$da = \frac{\partial L(a, y)}{\partial a}$$

↳ deriv. of Loss with respect to a.

$$\text{for, } \frac{\partial L(a, y)}{\partial a} = \frac{d}{da} [y \log(a) + (1-y) \log(1-a)] \quad [\because \text{without -sign}]$$

$$\begin{aligned} &= y \cancel{\frac{d}{da} [\ln a]} + (1-y) \frac{d}{da} [\ln(1-a)] \\ &= y \cdot \frac{1}{a} \frac{da}{da} + (1-y) \cdot \frac{1}{(1-a)} \frac{da}{da} \\ &= y \cdot \frac{1}{a} + \frac{(1-y)}{(1-a)} \left(\frac{da}{da} - \frac{da}{da} \right) \\ &= \frac{y}{a} + \frac{(1-y)}{(1-a)} \times (0-1) \\ &= \frac{y}{a} + \frac{(1-y)}{(1-a)} \times (-1) \\ &= \frac{y}{a} + \frac{(y-1)}{(1-a)} \end{aligned}$$

$$\text{adding '-' sign} \rightarrow \boxed{-\frac{y}{a} + \frac{(1-y)}{(1-a)}}$$

a - কে sigmoid-এর পরিমাণ 2-তে সাফেলভি desire করা,

$$\begin{aligned}\frac{\partial a}{\partial z} &= \frac{\partial \sigma(z)}{\partial z} = \frac{\partial}{\partial z} \left[\frac{1}{1+e^{-z}} \right] \\ &= \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{(1+e^{-z})} \cdot \frac{e^{-z}}{(1+e^{-z})} = \frac{1}{(1+e^{-z})} \cdot \frac{(1+e^{-z})-1}{(1+e^{-z})} \\ &= \frac{1}{(1+e^{-z})} \cdot \left[\frac{1+e^{-z}}{1+e^{-z}} - \frac{1}{1+e^{-z}} \right] = a \cdot (1-a).\end{aligned}$$

* Applying chain Rule $\rightarrow \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = \frac{\partial L}{\partial z}$

$$\cancel{\frac{\partial L}{\partial a}} / \cancel{\frac{\partial w_1}{\partial w_1}} \cancel{\frac{\partial L}{\partial w_1}} \cancel{\frac{\partial w_1}{\partial z}} \quad \downarrow \\ dz = (a - y)$$

$$\frac{\partial w_1}{\partial z} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$$= dz \cdot \frac{\partial \sigma(w_1x_1 + w_2x_2 + b)}{\partial w_1}$$

$$= dz \cdot x_1$$

$$\text{So, } \boxed{\frac{\partial w_1}{\partial z} = dz \cdot x_1}$$

$$\boxed{\frac{\partial w_2}{\partial z} = dz \cdot x_2}$$

$$\boxed{\frac{\partial b}{\partial z} = dz \cdot 1}$$

$$\text{update: } w_1 = w_1 - \alpha * dw_1$$

$$w_2 = w_2 - \alpha * dw_2$$

$$b = b - \alpha * db$$

④ Logistic regression QD on m examples.

$J=0; dw1=0; dw2=0; db=0;$ (initial values) $\xrightarrow{\text{modified}} \text{(1 for loop)}$,
 $w1=0; w2=0; b=0;$

for i = 1 to m

$\rightarrow dw = np.zeros((n_x, 1))$

Forward

$$z(i) = w1 * x1(i) + w2 * x2(i) + b,$$

$$a(i) = \text{sigmoid}(z(i)).$$

$$J += (y(i) * \log(a(i)) + (1 - y(i)) \log(1 - a(i))).$$

Backward.

$$dz(i) = a(i) - y(i)$$

$$dw1 += dz(i) * x1(i)$$

$$dw2 += dz(i) * x2(i)$$

$$db += dz(i)$$

$$\rightarrow dw += z(i) * dz(i)$$

$$J /= m$$

$$dw1 /= m$$

$$dw2 /= m$$

$$db /= m$$

$$\rightarrow dw /= m$$

Gradient.

$$w1 = w1 - \alpha * dw1$$

$$w2 = w2 - \alpha * dw2$$

$$b = b - \alpha * db.$$

problem:

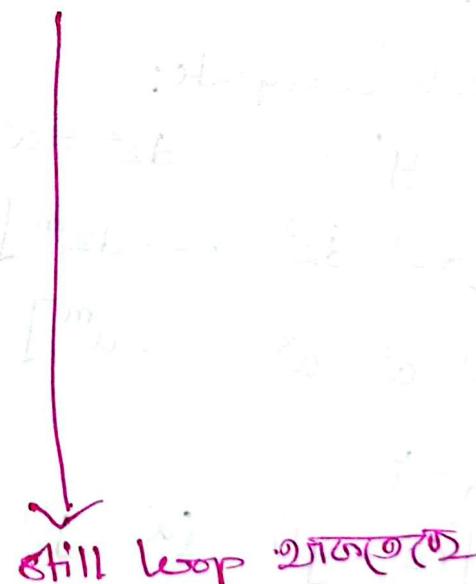
① 2 for loops. for 1 step of QD.

② less efficient.

solt!

① Vectorization.

② do not use for loop



Vectorization :- (LR)

$$z^{(1)} = w^T x^{(1)} + b \quad \dots \quad z^{(m)} = w^T x^{(m)} + b$$

$$a^{(1)} = \sigma(z^{(1)}) \quad \dots \quad a^{(m)} = \sigma(z^{(m)})$$

input matrix

$$X = \begin{bmatrix} & & \\ \vdots & \vdots & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \vdots & \vdots & & \vdots \end{bmatrix} \text{ while } R^{n_x \times m}$$

$$\begin{bmatrix} w^T \end{bmatrix} \begin{bmatrix} & & \\ \vdots & \vdots & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \vdots & \vdots & & \vdots \end{bmatrix}$$

① So, $\begin{bmatrix} z^{(1)} & z^{(2)} & z^{(3)} & \dots & z^{(m)} \end{bmatrix} = w^T X + \begin{bmatrix} b & b & b & \dots & b \end{bmatrix}$ broadcasting

$$= \begin{bmatrix} w^T x^{(1)} + b & w^T x^{(2)} + b & \dots & w^T x^{(m)} + b \end{bmatrix}$$

② $A = \begin{bmatrix} a^{(1)} & a^{(2)} & a^{(3)} & \dots & a^{(m)} \end{bmatrix} = \sigma(Z)$

③ Gradient compute:

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)} \quad \dots \quad dz^{(m)} = a^{(m)} - y^{(m)}$$

$$dZ = [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}]$$

$$A = [a^{(1)} \ a^{(2)} \ a^{(3)} \ \dots \ a^{(m)}] \quad Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$dZ = A - Y$$

$$= [(a^{(1)} - y^{(1)}) \ (a^{(2)} - y^{(2)}) \ \dots \ (a^{(m)} - y^{(m)})]$$

Vectorization implementation:

$$Z = w^T X + b.$$

$$Z = \text{np.dot}(w.T, X) + b.$$

$$A = \sigma(Z).$$

$$dZ = A - Y.$$

$$dw = \frac{1}{m} X \cdot dZ^T.$$

$$db = \frac{1}{m} dZ(i)$$

$$db = \frac{1}{m} \text{np.sum}(dZ)$$

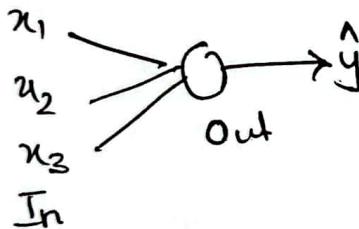
update

$$w = w - \alpha * dw$$

$$b = b - \alpha * db$$

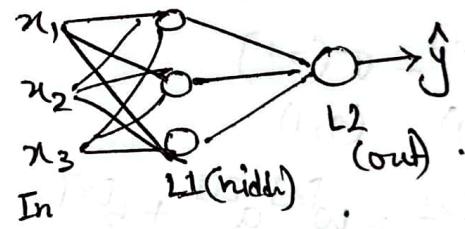
Composition of brain and NN.

logistic single node.

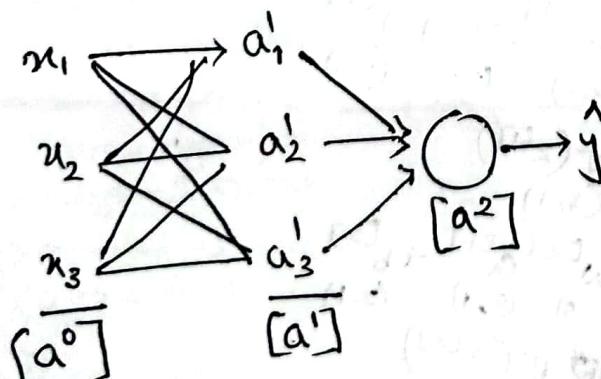


let a [layer]
node no.

NN.



Total connection = input x hidden layer.

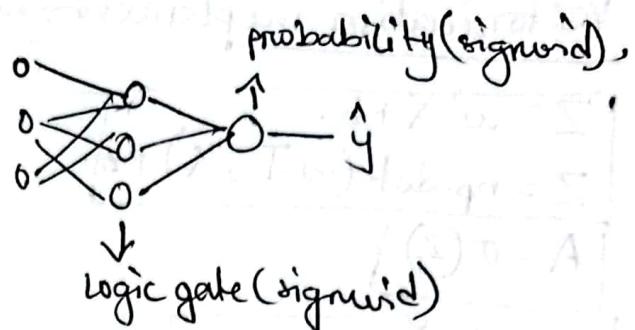
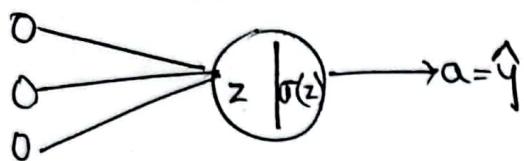


do not count input layer

→ 2 layered neural network.

Gaviflux

Neural Net



Neural Network Representation

$$\begin{matrix} x_1 & a_1 \\ x_2 & a_2 \\ x_3 & a_3 \\ & a_4 \end{matrix} \rightarrow \hat{y}$$

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ \vdots \\ z_n^{[1]} \end{bmatrix} = \begin{bmatrix} \cdots w_1^T \cdots \\ \vdots \\ \cdots w_n^T \cdots \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

$x \in \mathbb{R}^{3 \times 1}, w \in \mathbb{R}^{4 \times 3}, b \in \mathbb{R}^{4 \times 1}$

$$\text{So, } z^{[1]} = w^{[1]} x + b^{[1]}$$

$(4,1) \quad (4 \times 3) \quad (3 \times 1) \quad (4 \times 1)$

$$a^{[1]} = \sigma(z^{[1]})$$

$(4,1) \quad (4,1)$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$(1,1) \quad (1,4) \quad (4,1) \quad (1,1)$

$$a^{[2]} = \sigma(z^{[2]})$$

$(1,1) \quad (1,1)$

Suppose

$$\begin{matrix} x_1 & 0 & 0 \\ x_2 & 0 & 0 \\ x_3 & 0 & 0 \end{matrix} \rightarrow \hat{y}$$

$\frac{0}{a^1} \quad \frac{0}{a^2} \quad \frac{0}{a^3}$

$$z^{[1]} = w^{[1]} \cdot x + b^{[1]}$$

$(5,1) \quad (5,3) \quad (3,1) \quad (5,1)$

$$a^{[1]} = \sigma(z^{[1]})$$

$(5,1) \quad (5,1)$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$(3,1) \quad (3,5) \quad (5,1) \quad (3,1)$

$$a^{[2]} = \sigma(z^{[2]})$$

$(3,1) \quad (3,1)$

$$z^{[3]} = w^{[3]} a^{[2]} + b^{[3]}$$

$(1,1) \quad (1,3) \quad (3,1) \quad (1,1)$

$$a^{[3]} = \sigma(z^{[3]})$$

$(1,1) \quad (1,1)$

Justification for vectorized implementation.

$$z^{[0]}(1) = w^{[0]}x^{(1)} + b^{[0]}$$

$$w^{[0]} = \begin{bmatrix} - & - \\ - & - \\ - & - \end{bmatrix}$$

$$w^{[0]}x^{(1)} = \begin{bmatrix} \square \\ \square \\ \square \end{bmatrix}$$

$$z^{[1]}(2) = w^{[1]}x^{(2)} + b^{[1]}$$

$$w^{[1]}x^{(2)} = \begin{bmatrix} \Delta \\ \Delta \\ \Delta \\ \Delta \end{bmatrix}$$

$$w^{[1]}x^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$w^{[1]} \times \begin{bmatrix} x^1 & x^2 & x^3 \end{bmatrix} = \begin{bmatrix} \square & \Delta & 0 \\ \square & \Delta & 0 \\ \square & \Delta & 0 \end{bmatrix} = \begin{bmatrix} z^{1} & z^{[1](2)} & z^{[1](3)} \end{bmatrix} = z^{[1]}$$

$$z^{[1]} = w^{[1]}x + b^{[1]}$$

Activation function: ~~in~~ a NN without an activation fn or with a linear act. fn

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_i} \quad (L, a, z, w)$$

is just a linear regression model.

$$z = w^T x + b$$

$$a = g(z)$$

activation fn (may be linear or non linear).

① Binary Step Fn :- {0, 1}

* Threshold based.

* above or below input

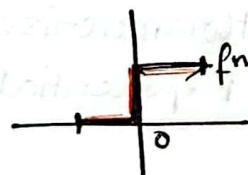
एते neuron active

एते same signal

next layer-2 मात्रा

* $y > 0 \rightarrow \text{out} = 1$

$y < 0 \rightarrow \text{out} = 0$.



→ symbol.

Dis:

① multivalue outputs - not allowed

② derivative - 0 एवं of constant

does not represent derivative

LR - so hidden layer-2 use

करने की बात

② Linear Act. Fn: (-∞, ∞)

* Allows multivalue outputs

But

1. No backpropagation,
 w, b no update.

2. Linear- Σ , combination is
always linear - so NN will
turn into one LAYER.

Gaviflux

3. Sigmoid: $\sigma = \frac{1}{1+e^{-z}}$ (0,1)

- * Non linear
- * Smooth gradient
- * [0 → 1] range
↓
output.
- * prevents jumps.
* back propagation allowed.

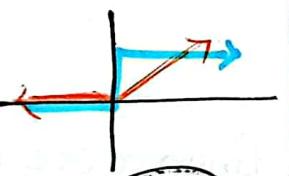
Dis:

- ① Vanishing Gradient: For High or very low value - no change, deriv converges to 0.
- ② Certain point → flat.
- ③ Refuses to learn further + too slow.
- ④ Stuck at local minima.
- ⑤ Output not ZERO CENTERED.
→ so sign same as data, (but real life data neg 3/4)
- ⑥ Computationally Expensive.

5. Rectified Linear Unit (ReLU) [0, ∞)

- ⑦ Solves VG.

↳ by piecewise linear.



$\max(0, x)$
orange.

- * converges quickly
- * computationally efficient
- * Back prop allowed.

Dis:

- ① Dropout problem: Rel Derivative 0 ना होती तो off एवं याद नहीं।
- ② Dying ReLU "": when input reaches 0 or neg, grad f' = 0 होती है → back propagation off होती है thus.
- ③ Regularization problem: क्वांटिनेशन off होती है but in my hand.

But dropout - नहीं करने का calculation speed ज्यादा, Load भी।

4. Tanh Hyperbolic Tangent: $\tanh\left(\frac{e^{-x} - e^x}{e^{-x} + e^x}\right)$

- * Sigmoid - इसका विकल्प होता है → range of output is now neg - pos; as zero centered.
- * Derivative is steeper.
- * Wider range + fast learning.
- * Back propagation allowed

Dis:

- ① Vanishing Gradient (VG)
- ② Computationally Expensive
- ③ No back propagation.

6. Leaky ReLU: (-∞, ∞)

- * Solves VG + Dying ReLU problem.

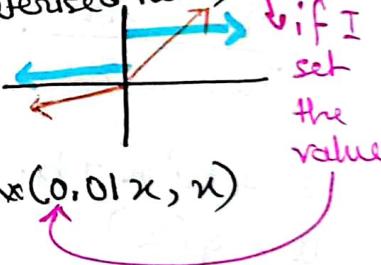
by small pos. slope in the neg direction (Parameterised ReLU)

- * enables back prop.

* value is $\max(0.01x, x)$

Dis:

- ① Results inconsistent for neg input values.
- ② Parameterised ReLU
- ③ Exponential ReLU.



Softmax: Combination of multiple sigmoid.

→ sigmoid আলাদা আলাদা output-ই use করব - it always gives 0 অকৃত - যিন্মধ্যে probability এর a data point belonging to a class —— but এমন না হয় সবগুলো output হ্রাস করবলে I summation পাব।

But probability সব যাগ করে তা ১ পাঞ্চাব করণ।

Thus sigmoid binary classification - 2 use 2 out of 3

multiclass classification - softmax use

$$\text{Softmax-equation} \rightarrow a_i = \frac{e^{z_i}}{\sum_{k=1}^c e^{z_k}}$$

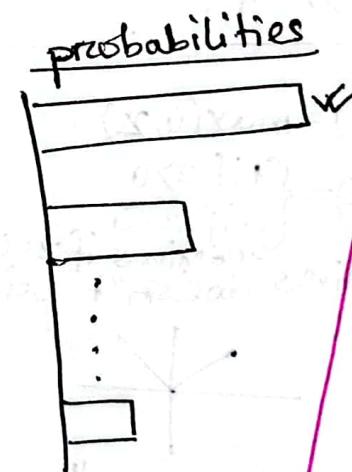
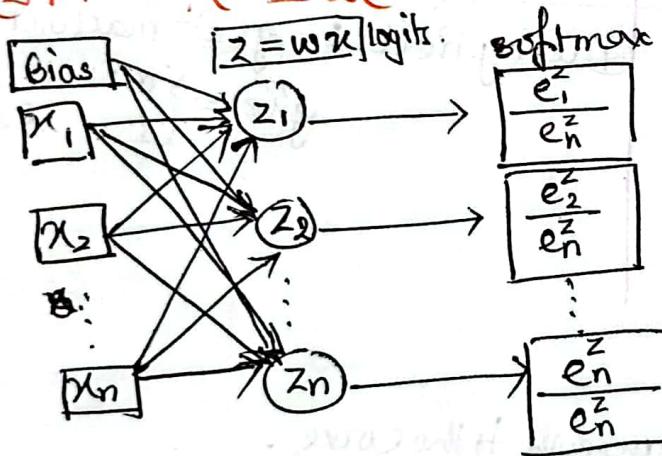
i → individual c
s → sum of output-
c → no. of class .

Suppose 3 newtons gave 3 output

such that $1.2, 0.9, 0.75$ — but summation is more than 1

softmax - 3类 probability $f(x) \rightarrow 0.42, 0.31, 0.27$

ପ୍ରାଚୀ କବଳୀ । ୧୮



e याएँ
यहाँ यहाँ
monotonic
function.

small value
କେ ମାତ୍ରେ small,
large value
କେ ମାତ୍ରେ large
କାହାରେ ମାତ୍ର

Q] Which activation fn is preferred?

④ Which activation fn is preferred?
 ⑤ For CLASSIFIERS, (especially Binary) → sigmoid, tanh (but vq problem)

④ For CLASSIFIERS, operating on large NN → ReLU (speedy)

④ For Deep + high complexity
(mostly used for INTERMEDIATE LAYERS)

④ To avoid dying ReLU → use leaky ReLU

Start with R&U - then change into others.

- softmax is used in (OUTPUT layers)

* Use of softmax :-

Distinguish blurry and
crisp image.

* all layer of NN collapse into 1
linear layer if linear act. fn is used
in each layer. - as it gives linear
output

* HL necessity :- to learn pattern of
input and feature

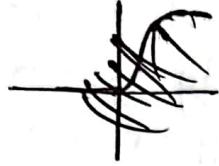
* Why more node in HL is better

- to learn better and give
better output

Derivatives :-

① Sigmoid :-

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



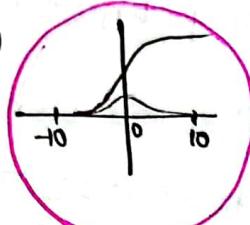
$$g'(z) = \frac{d}{dz} g(z) \rightarrow \text{slope of } g(z) \text{ at } z$$

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

Suppose $z = 10$.

$$g(z) = 1$$

$$g'(z) = 1(1-1) = 0$$



or, $z = -10$.

$$g(z) = 0$$

$$g'(z) = 0(1-0) = 0$$

or, $z = 0$.

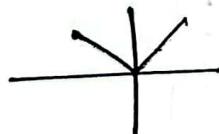
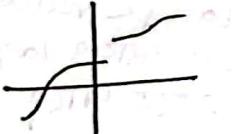
$$g(z) = \frac{1}{2}$$

$$g'(z) = \frac{1}{2}(1-\frac{1}{2}) = \frac{1}{4}$$

② ReLU :- $g(z) = \max(0, z)$

$$g'(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z < 0 \\ \text{undefined if } z = 0 \end{cases}$$

when derivatives doesn't exist?



① discontinuous

② more than 1 tangent to the curve.

at $z=0$ - (more than) one tangent

Derivative exists :-

① continuous curve

② Only one tangent to the curve.

② Tanh :-

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

$$\frac{d}{dz} g(z) = \frac{(e^z + e^{-z})(e^z - e^{-z}) - (e^z - e^{-z})(e^z + e^{-z})}{(e^z + e^{-z})^2}$$

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

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

$$= 1 - \tanh(z)^2$$



$$z=10; \tanh(z) = 1; g'(z) = 1-1^2 = 0$$

$$z=-10; \tanh(-10) = -1; g'(z) = 1 - (-1)^2 = 0$$

$$z=0; \tanh(0) = 0; g'(z) = 1-0^2 = 1$$

④ Leaky ReLU :- $g(z) = \max(0.01z, z)$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \end{cases}$$

Entropy :- Indiscipline - the uncertainty inherent in the a random variable x 's possible outcome.

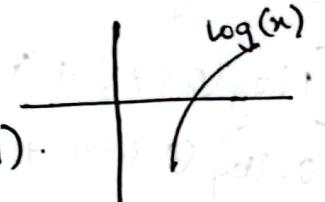
$$H(x) = \begin{cases} -\int_x P(x) \log P(x) & , x \text{ is continuous} \\ -\sum_x p(x) \log p(x) & , x \text{ is discrete} \end{cases}$$

neg sign why?

(0,1) - ৰঠ মাত্ৰ প্ৰিব থাকে
so - log of (0,1) - ৰঠ মাত্ৰ
ক্লিনা value - এ neg value
আসে - so neg-neg pos. কাণ্ডে
neg value থাকে

$p(x)$ সুন ক্ষয়তি cause
logarithm ক্ষয়লৈ value
(2) আস

$\log(p(x)) < 0$ for all $p(x)$ in (0,1).



Entropy vs Information:

provides a measure of avg amt of info needed to represent an event drawn from a probability dist.

a way to quantify the amount of surprise for an event measured in bits.

Example : (একটি)

$$\Delta = 26 \quad 0 = 4$$

①

$$\Delta = 19 \quad 0 = 16$$

②

$$\Delta = 1 \quad 0 = 29$$

③

uncertainty
হ্যাম

uncertainty
high-entropy

uncertainty হ্যাম

Cross Entropy / Logarithm loss fn / log loss / Logistic Loss fn.

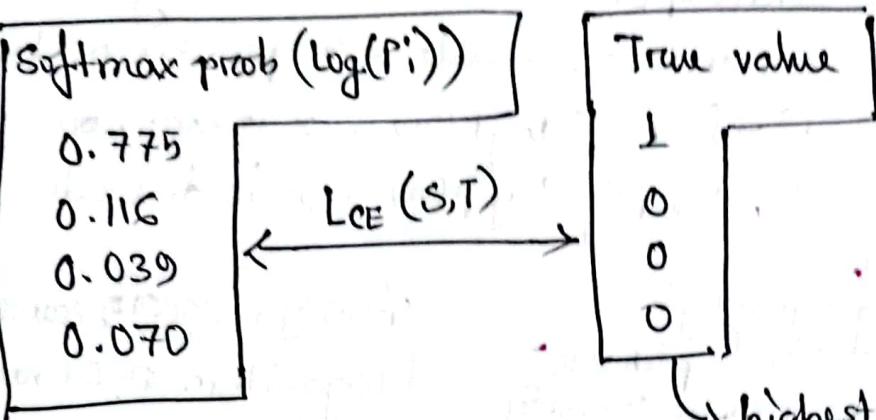
$$L_{CE} = - \sum_{i=1}^n t_i \log(p_i) ; \text{ for } n \text{ class.}$$

true → predicted.

$$* * \boxed{- \sum t_i \log(p_i)}$$

Gaviflux

Example:



highest value 1

বাসি 0 রয়ে দুওয়ারে

One hot encoding vector

$$\text{So, } L_{CE} = - \sum T_i \log(S_i)$$

$$= - \left[1 \log_2(0.775) + 0 \log_2(0.116) + 0 \log_2(0.039) + 0 \log_2(0.070) \right]$$
$$= - \log_2(0.775) = 0.3677$$

For Binary Cross Entropy :

$$L = - \sum_{i=1}^2 t_i \log(P_i)$$
$$= - [t \log(P) + (1-t) \log(1-P)]$$

* Categorical Cross Entropy \rightarrow one hot encoded.

* Sparse " " \rightarrow integer encoded.

Digit

Viterbi Algorithm

$$\pi = \begin{bmatrix} n & y \end{bmatrix}$$

transition matrix \rightarrow

		(to)	
		x	y
(from)	n	.	.
	y	.	.

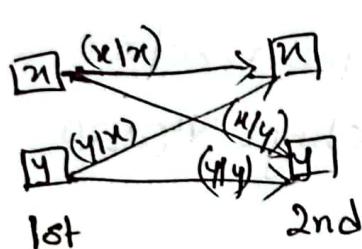
① sequence. ফরওয়া

② স্টেট স্টার্ট

③ for each situation
x, y - গতি probability

Observation/
emission probability

	a	b	c	d
n
y



1st

2nd

3rd

For 1st,

$$P(1st, n) = P(1st | x) P(x)$$

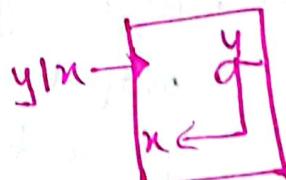
$$P(1st, y) = P(1st | y) P(y)$$

For 2nd,

$$P(2nd, x) = P(2nd | x) P(x | x)$$

$$P(2nd, y) = P(2nd | y) P(y | x)$$

given that, y



[LR] \rightarrow find dw_1, dw_2, db .

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$\frac{\partial L(a, y)}{\partial a} = \frac{\partial}{\partial a} (-y \log a + (1-y) \log(1-a))$$

$$\frac{\partial a}{\partial z} = \frac{\partial}{\partial z} \text{ (activation fn)}$$

$$\frac{\partial z}{\partial w} = \frac{\partial (wz + b)}{\partial w}, \frac{\partial (wz + b)}{\partial b}$$

Final $dw_1 = \frac{\partial L}{\partial w} \Rightarrow \text{II আসলো তা}$

Gaviflux

CNN Architecture:

- * CNN output shape : $\left\lfloor \frac{\text{input shape} - \text{kernel} + 2 \times \text{padding}}{\text{stride}} \right\rfloor + 1$, out-features.
parameter : $((\text{Kernel H} \times \text{W}) + 1) \times \text{out-features}$.
- * Maxpool : $\left\lfloor \frac{\text{input shape} - \text{kernel}}{\text{stride}} \right\rfloor + 1$.
parameter : 0. (no weights involved)
- * Minpool : $\left\lfloor \frac{\text{input shape} - \text{kernel}}{\text{stride}} \right\rfloor + 1$
parameter: 0
- * Dropout : (drops network) output \rightarrow same as before
parameters $\Rightarrow 0$, (use π , $1-\pi$, change π)
- * Flatten : (multiD \rightarrow 1D conversion)
output $\rightarrow (\text{H} \times \text{W} \times \text{out-features})$ [shape-product]
parameter $\rightarrow 0$.
- * Linear : logits / final prediction.
output \rightarrow (out-features).
parameters $\rightarrow (\text{input} \times \text{out})$ if Bias=F
 $(\text{input} \times \text{out}) + \text{out}$ if Bias=T
- * Act-fn : output \rightarrow out-features - অন্তর্ভুক্ত
parameter $\rightarrow 0$

	output	parameter
CNN	$\left\lfloor \frac{\text{input} - \text{kernel} + (\text{f} \times \text{pad})}{\text{stride}} \right\rfloor + 1, \text{out_f}$	$((\text{kernel H} \times \text{W}) + 1) \times \text{out_feats}$
Max/minpool	$\left\lfloor \frac{\text{in} - \text{kernel}}{\text{stride}} \right\rfloor + 1$	0
Dropout	(output - shape)	0
Flatten	(shape product)	0
Activation.fn	(output shape)	0
Linear	(output shape)	$(\text{in} \times \text{out}) + \text{out}$ $(B=T)$ $(\text{in} \times \text{out})$ $(B=F)$

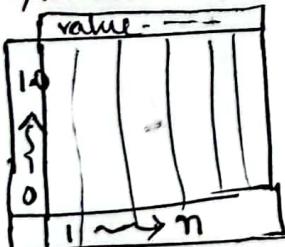
Fuzzy Extension

A \ B	1 ~ n
1	
n	

- गणना box- $\min\{A, B\}$ द्वारा।
- For n , $(n-i)$ diagonally check करके \max value नियुत करावी।

Fuzzy Addition (α -cut)

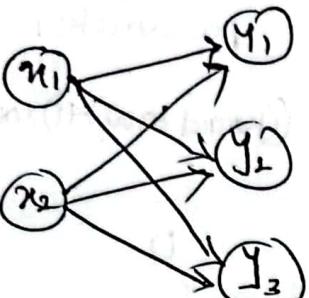
- A/B box गणना



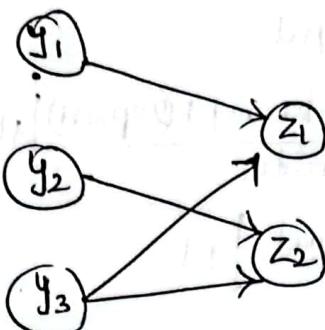
$$C_{\alpha} = \underline{\alpha} A + \overline{\alpha} B.$$

- $C_{\alpha=0.1} = [A_{0.1}, \underline{\alpha} A + \overline{\alpha} B, A_{0.1}, \underline{\alpha} A + \overline{\alpha} B, \dots]$
- $C_{\alpha} [\text{lowest, highest}] - \text{value wise}$ (dup allowed), box गणना करके value द्वारा (except 1).

Fuzzy Relation:



R1



R2

$y \rightarrow z \rightarrow R3$

$$H_{R1}(x, y) \cdot \underset{\text{OP}}{\cancel{H_{R2}(y, z)}} \quad H_{R1R2}(x_1, z_1) = \frac{(x_1 y_1) + (x_1 y_2) + (x_1 y_3) + (x_2 y_1) + (x_2 y_2) + (x_2 y_3)}{6}$$

$$\text{Max-Min: } H_{R1R2}(x_1, z_1) = ((x_1 y_1) \wedge (y_1 z_1)) \vee ((x_1 y_2) \wedge (y_2 z_1)) \vee \dots$$

$$\text{Min-Max: } H_{R1R2}(x_1, z_1) = ((x_1 y_1) \vee (y_1 z_1)) \wedge ((x_1 y_2) \vee (y_2 z_1)) \wedge \dots$$

$$\text{Max-Prod: } H_{R1R2}(x_1, z_1) = ((x_1 y_1) \times (y_1 z_1)) \vee ((x_1 y_2) \times (y_2 z_1)) \vee \dots$$

$$\text{Max-Avg} \rightarrow \frac{(+ \vee + \vee +)}{3}$$

$$\text{Max Min} \rightarrow (\wedge) \vee (\wedge) \vee (\wedge) \dots$$

$$\text{Min Max} \rightarrow (\vee) \wedge (\vee) \wedge (\vee) \dots$$

$$\text{Max Prod} \rightarrow (x) \vee (x) \vee (x) \dots$$

Fuzzy membership function:

$$\text{pos}(\checkmark) \rightarrow \frac{x-a}{b-a}; a \leq x \leq b$$

$$\text{neg}(\times) \rightarrow \frac{a-x}{a-b}; a \leq x \leq b$$

$$\text{cons}(-) \rightarrow 1; -a \leq x \leq b$$

Analogy reasoning

Suppose; given, Man \rightarrow Woman then King \rightarrow ?

features	Woman	Man	King	Queen	Apple
f_1					
f_2					
f_3					

① Analogy = $(f_{\text{woman}} - f_{\text{man}}) + f_{\text{king}}$

let, $f_{\text{random}} = f_{\text{queen}}$

② $(f_{\text{random}}, \text{Analogy})$

③ Compare \rightarrow Euclidean distance এবং কোণ f_R আব্দি Analogy-র

④ $f_R = f_{\text{man}} = f_{\text{woman}} = f_{\text{man}}$ --- JT DT table-এ
মানে distance এবং কোণ
মানে distance এবং কোণ
কোণ কোণ কোণ কোণ
মানে কোণ কোণ কোণ

TF-IDF

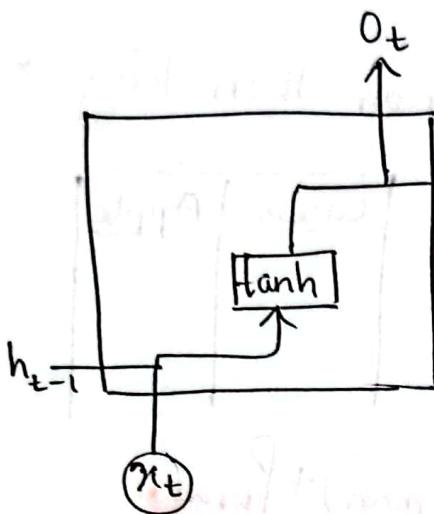
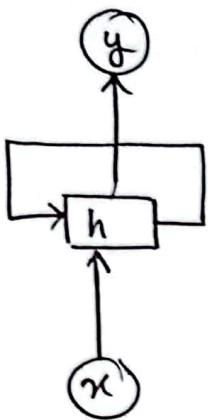
TF (Term f_q) = $\frac{\text{no. of times this term appears in the sentence (Doc)}}{\text{no. of total terms in the sentence}}$

IDF (Inverse Doc f_q) = $\frac{\text{no. of total corpus}}{\text{no. of sentence this word is in}}$

term	DOC	TF	IDF	TF-IDF (product)

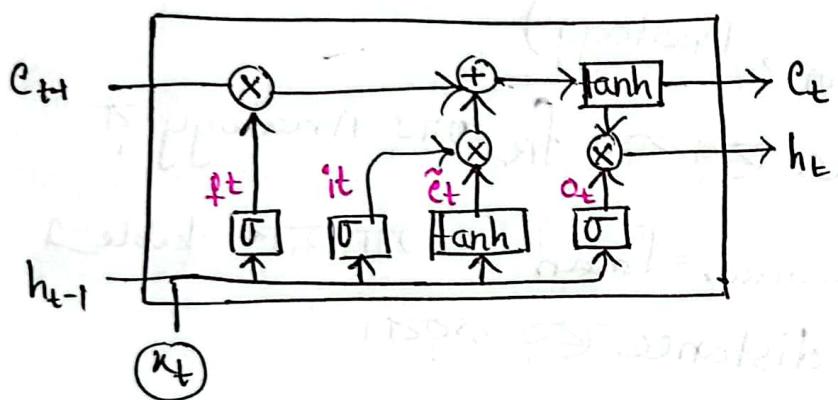
Gaviflux

RNN:



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

LSTM:



Gates:

$$\text{forget: } f_t = \sigma(W_f[h_{t-1}, h_t, x_t] + b_f)$$

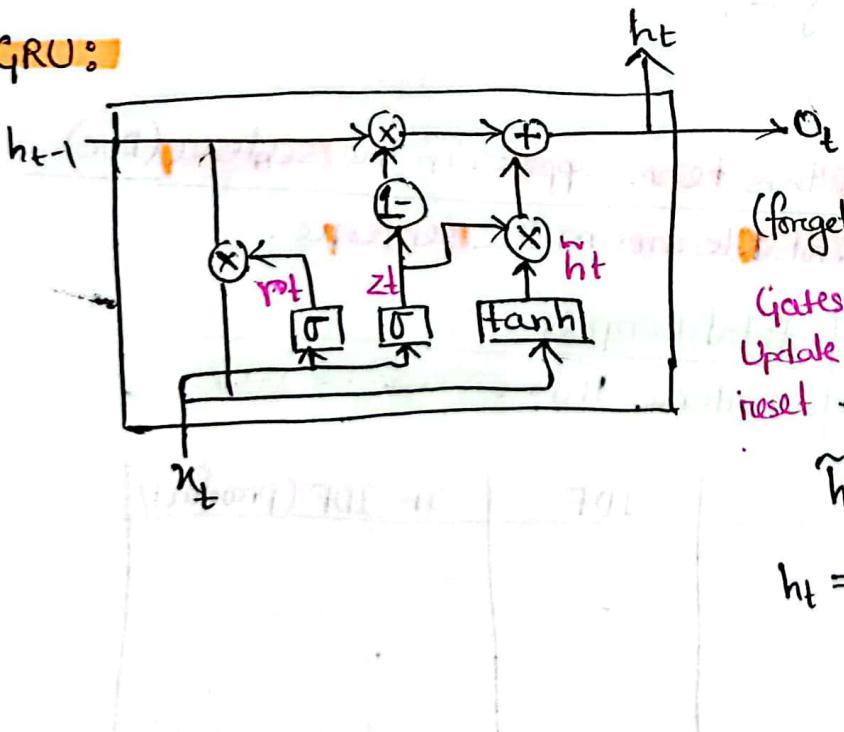
$$\text{input: } i_t = \sigma(W_i[h_t, x_t, c_{t-1}] + b_i)$$

$$\text{Candidate: } c_t^* = \tanh(W_c[h_{t-1}, c_{t-1}, x_t] + b_c)$$

$$\text{mem: } c_t = f_t \cdot c_{t-1} + i_t \cdot c_t^*$$

$$\text{Output: } o_t = \sigma(W_o[c_t, h_{t-1}, x_t] + b_o)$$

GRU:



(forget+input) \rightarrow update gate

Gates

$$\text{Update: } z_t = \sigma(W_z[h_{t-1}, x_t])$$

$$\text{reset: } r_t = \sigma(W_r[h_{t-1}, x_t])$$

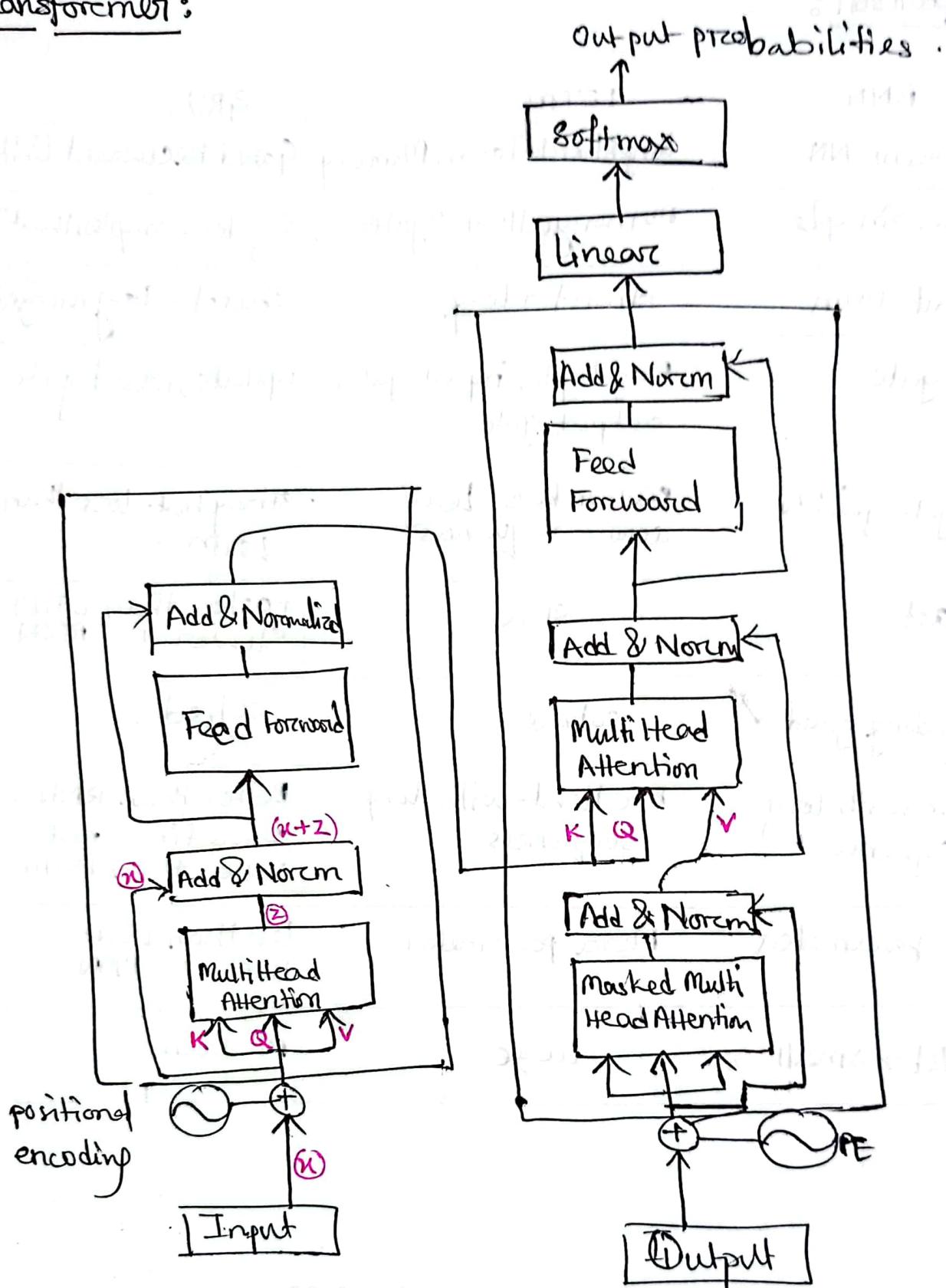
$$\tilde{h}_t = \tanh(W_r[r_t * h_{t-1} + W_z[z_t * h_t]] + b_t)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Comparison:

RNN	LSTM	GRU.
Recurrent NN	Long Short-Term Memory	Gated Recurrent Unit
Basic, simple	Memory cells + 3 gates	2 gates, simpler than LSTM
short term	short + long	short - long merged
No gate	Forget gate, input gate, update, reset gate output gate	
Forgets quickly	remembers long across sequences	Strong but less than LSTM
Fast	Slow	Faster than LSTM Slower " RNN
Vanishing grad ✓	solved	solved .
Poor with long sequences	Excellent with long sequences	Better than RNN, sometimes not better than LSTM
Few parameters	More parameters	less than LSTM more " RNN
Model → small	large	Medium

Transformer:

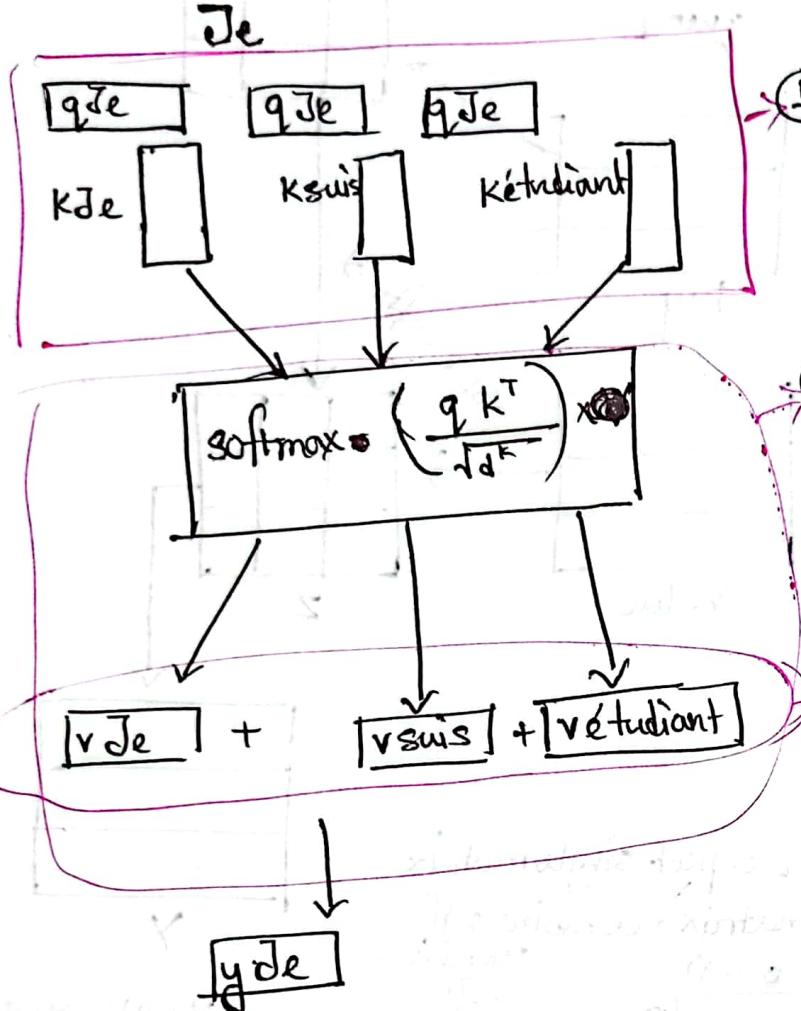


Je suis étudiant

I am a Student

Concept of Self Attention:

There are three vectors of each word \rightarrow Query, key, Value:



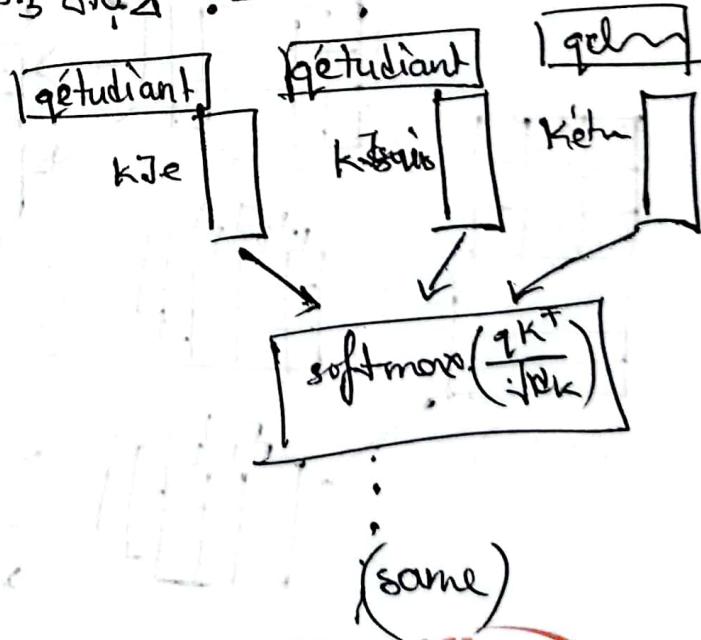
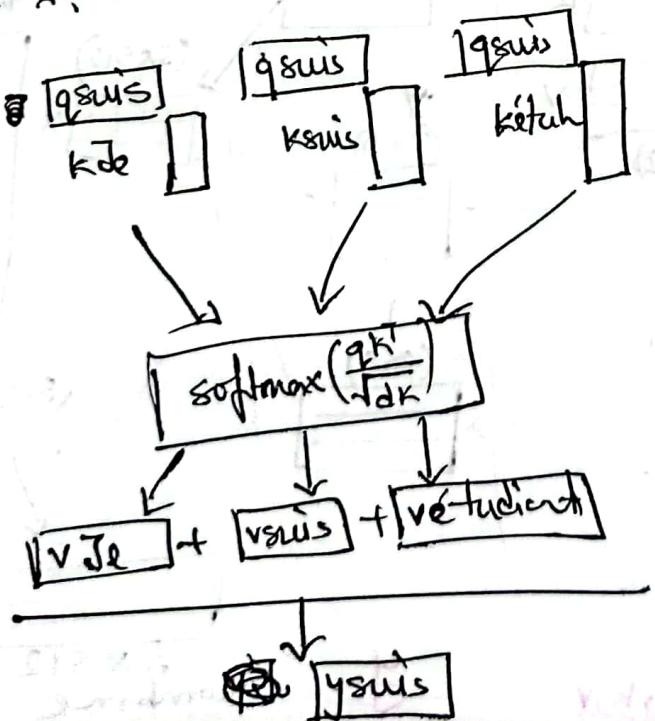
Suppose, Je

① query vector of Je will get the dot product with each key vectors of three words "Je", "suis", "étudiant"

② will apply softmax on query and key vector's dot product then multiply with each value vectors to get z . and then

③ Sum up embedding with z and proceed to feed forward.

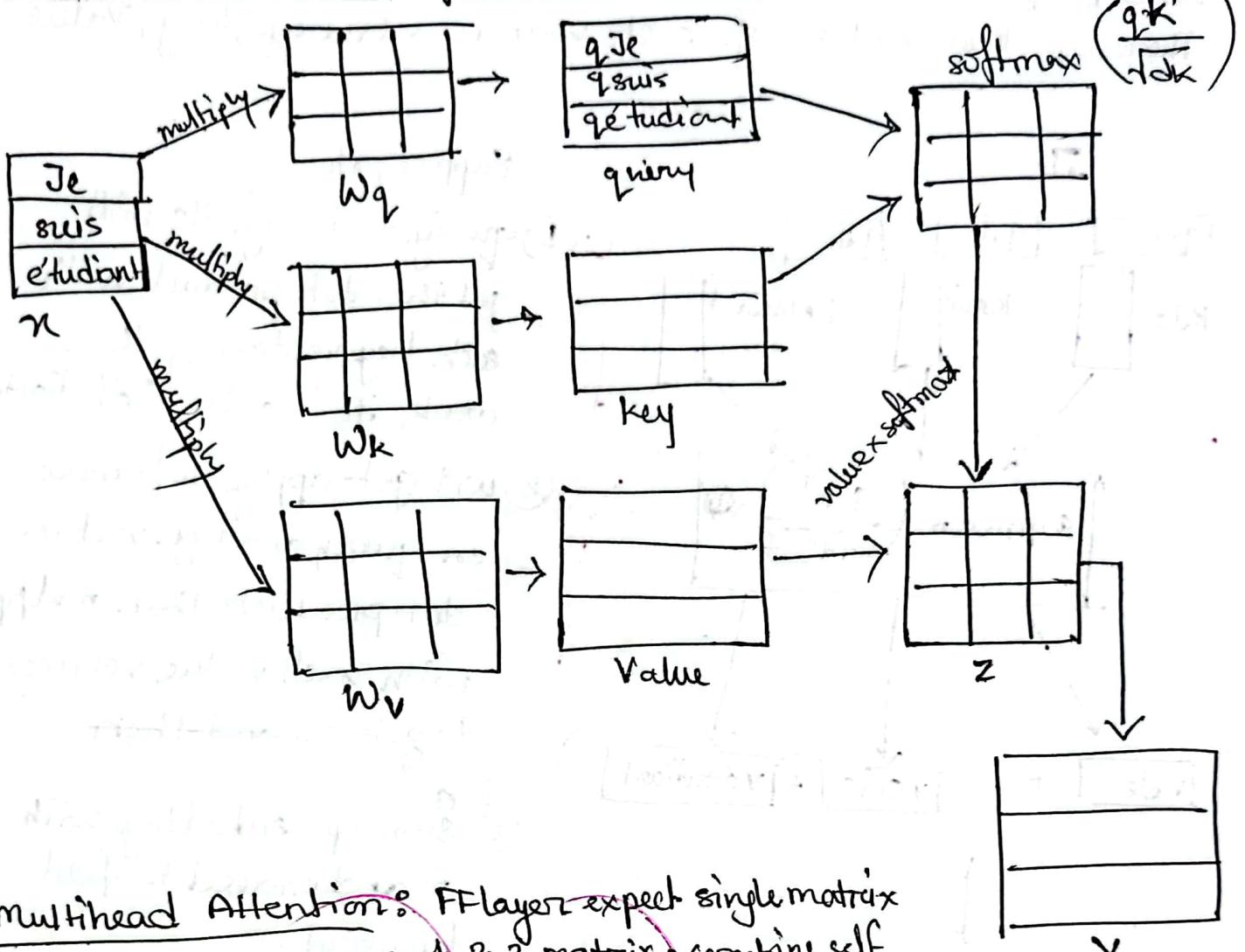
Je - २० तत्त्व suis का étudiant - २०३ का :-



(same)

Gaviflux

Matrix Calculation of self Attention:



multihead Attention: FFLayer expect single matrix
not 3×3 matrix; combine self attention

