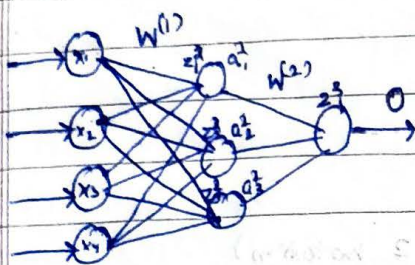


Main intuition for extra layer:

- The layer learns non-linear interactions between the input word vectors.

• Neural Networks: (Core Understanding)



Facts

- Neural network is not any real network, just a mathematical abstract idea.

→ In this we take an input vector multiply it with parameter vector and get output

- Input (is somewhat a hyperparameter); no. of hidden layers are hyper para.
- parameters or features are learned (W) which creates the initial gradient
- there is a loss function at the output of our choice
- Each hidden layer has basically 2 parts/function mixed together [adder, non-linearity].

Now we can define notation in 2 ways.

- either take $\overset{\text{this}}{\Rightarrow} \textcircled{S} \xrightarrow{W} \textcircled{E}$ as a row and $\textcircled{O} \xrightarrow{W} \textcircled{O} \xleftarrow{\text{this}}$ as column

where S denotes start of a function and E denotes end of a function

but then function will be written as

$$f(W^T x + b)$$

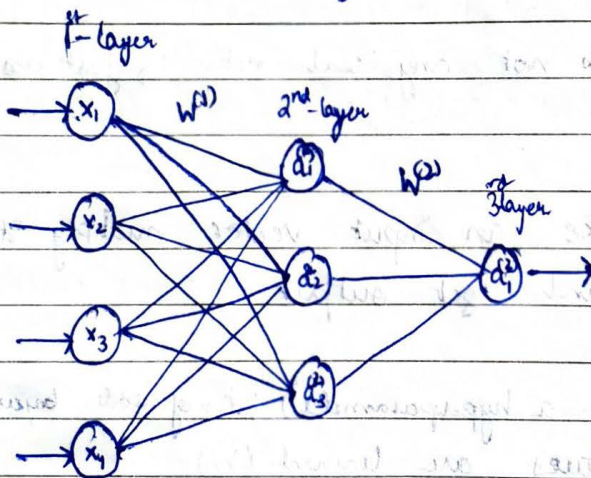
Another way, which is straight forward, thus used every ~~where~~ where is:-



where S denotes start of a connection & E denotes end of a connection
thus, now the functions have a nice look of

$$a = f(Wx + b)$$

Now, in matrix outlay (using 2 notation)



Here vector x is:

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

Here ~~matrix~~ dimension of parameters W are:-

* since end units of ~~1st~~ 1st hidden layer are 3, thus row is 3, and start is input layer thus column is 4. $\therefore R^{3 \times 4}$

$$W^{(1)} = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \end{bmatrix}$$

3x4

$$(d \times x \times W)$$

Similarly,

End of connection is $a^{(2)}$ hence row is 1, and start of connection is $a^{(1)}$ hence column is 3.

$$W^{(2)} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \end{bmatrix}_{1 \times 3}$$

Here we have 2 input vector (one is x that we give), other is $a^{(1)}$ that we form.

* Dry Implementation

$$Z^{(2)} = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \end{bmatrix}_{3 \times 4} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}_{4 \times 1}$$

$$Z^{(2)} = \begin{bmatrix} (w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4) \\ (w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + w_{24}x_4) \\ (w_{31}x_1 + w_{32}x_2 + w_{33}x_3 + w_{34}x_4) \end{bmatrix}_{3 \times 1}$$

{ where $Z^{(2)}$ is pre-activation / adder function }

$$a^{(2)} = \begin{bmatrix} f(w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4) \\ f(w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + w_{24}x_4) \\ f(w_{31}x_1 + w_{32}x_2 + w_{33}x_3 + w_{34}x_4) \end{bmatrix}_{3 \times 1}$$

{ where $f(x)$ is a non-linearity func }

Now, $a^{(1)}$ becomes input for (2) layer.

NOTE : output of each layer is essentially a vector.

Now,

$$z^{(3)} = w^{(2)} a^{(2)}$$

$$z^{(3)} = \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} & w_{13}^{(2)} \\ \vdots & \vdots & \vdots \\ w_{31}^{(2)} & w_{32}^{(2)} & w_{33}^{(2)} \end{bmatrix}_{3 \times 3} \begin{bmatrix} f(w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3 + w_{14}^{(1)} x_4) \\ f(w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3 + w_{24}^{(1)} x_4) \\ f(w_{31}^{(1)} x_1 + w_{32}^{(1)} x_2 + w_{33}^{(1)} x_3 + w_{34}^{(1)} x_4) \end{bmatrix}_{3 \times 1}$$

$$z^{(3)} = \begin{bmatrix} w_{11}^{(2)} f(w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3 + w_{14}^{(1)} x_4) + w_{12}^{(2)} f(w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3 + w_{24}^{(1)} x_4) + w_{13}^{(2)} f(w_{31}^{(1)} x_1 + w_{32}^{(1)} x_2 + w_{33}^{(1)} x_3 + w_{34}^{(1)} x_4) \\ \vdots \\ w_{31}^{(2)} f(w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3 + w_{14}^{(1)} x_4) + w_{32}^{(2)} f(w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3 + w_{24}^{(1)} x_4) + w_{33}^{(2)} f(w_{31}^{(1)} x_1 + w_{32}^{(1)} x_2 + w_{33}^{(1)} x_3 + w_{34}^{(1)} x_4) \end{bmatrix}_{3 \times 1}$$

$$a^{(3)} = \begin{bmatrix} f(w_{11}^{(2)} f(w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3 + w_{14}^{(1)} x_4) + w_{12}^{(2)} f(w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3 + w_{24}^{(1)} x_4) + w_{13}^{(2)} f(w_{31}^{(1)} x_1 + w_{32}^{(1)} x_2 + w_{33}^{(1)} x_3 + w_{34}^{(1)} x_4)) \\ \vdots \\ f(w_{31}^{(2)} f(w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3 + w_{14}^{(1)} x_4) + w_{32}^{(2)} f(w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3 + w_{24}^{(1)} x_4) + w_{33}^{(2)} f(w_{31}^{(1)} x_1 + w_{32}^{(1)} x_2 + w_{33}^{(1)} x_3 + w_{34}^{(1)} x_4)) \end{bmatrix}$$

$$\text{Output } (o) = a^{(3)}$$

Notes

or dimension

- Size of parameter matrix : (no. of units in ending connection layer) (no. of unit in starting connection layer)

$$R^{(\text{unit in ending layer}) \times (\text{unit in starting layer})}$$

$$\rightarrow R^{(s_{j+1} \times s_j)}$$

Backpropagation in this network

* first of all we will design a cost function for the last output layer

→ let's take it as

$$L = \frac{1}{2} (\hat{y} - y)^2 \quad \text{where } \hat{y} \text{ is the output}$$

we know $\frac{\partial L}{\partial \hat{y}} = \hat{y} - y$

Now, find

$$\frac{\partial L}{\partial w_{ij}^{(2)}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_i^{(2)}} \cdot \left(\frac{\partial z_i^{(2)}}{\partial a_i^{(2)}} \cdot \frac{\partial a_i^{(2)}}{\partial z_i^{(2)}} \right) \cdot \frac{\partial z_i^{(2)}}{\partial w_{ij}^{(2)}}$$

$\delta_i^{(2)}$

$$\frac{\partial \hat{y}}{\partial z_i^{(2)}} = \frac{\partial (a_i^{(2)})}{\partial z_i^{(2)}} = f'(z_i^{(2)})$$

$$\frac{\partial z_i^{(2)}}{\partial a_i^{(2)}} = w_i^{(2)} \quad \left\{ \text{important} = \frac{\partial z_i^{(3)}}{\partial a_i^{(2)}} \cdot \frac{\partial a_i^{(2)}}{\partial z_i^{(2)}} = \delta_i^{(2)} \right\}$$

$$\frac{\partial a_i^{(2)}}{\partial z_i^{(2)}} = f'(z_i^{(2)}) \quad \left\{ \text{Here } f(z_i) = \sigma(z_i) \Rightarrow f(z_i) = f(z_i)(1 - f(z_i)) \right\}$$

$$\frac{\partial (z_i^{(2)})}{\partial w_{ij}^{(2)}} = \frac{\partial (w_{ij}^{(2)} x_j + b_i)}{\partial w_{ij}^{(2)}} = \frac{\partial (w_{ij}^{(2)} x_j + w_{i2}^{(2)} x_2 + w_{i3}^{(2)} x_3 + w_{i4}^{(2)} x_4 + b_i)}{\partial w_{ij}^{(2)}}$$

$$\frac{\partial (z_i^{(2)})}{\partial w_{ij}^{(2)}} = \frac{\partial (\sum_k w_{ki} x_k + b_i)}{\partial w_{ij}^{(2)}}$$

$$\frac{\partial (z_i^{(2)})}{\partial w_{ij}^{(2)}} = x_j$$

x_j is simpler, we choose it for w_{ij} choosing i in turn chooses z_i and choosing j takes out the specific element out of z_i

Substitute in eqn to get gradient

then update gradient wrt:-

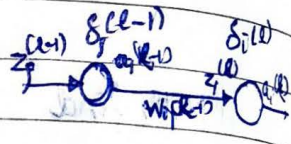
$$w_{ij} \leftarrow w_{ij} - \eta \frac{\partial L}{\partial w_{ij}} \quad (\text{sgd})$$

learning rate (hyperparameter)

So in general, from equation we can see

that, error received at $a_j^{(L-1)}$

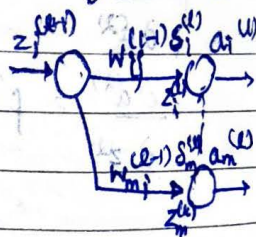
$$\delta_j^{(L-1)} = \delta_i^{(L)} w_{ij}^{(L)}$$



(derivative of j^{th} unit in $(L-1)$ layer) \equiv (derivative of i^{th} unit 'connecting' to $L-1$) \times (weight w_{ij} from connection)

If there are multiple ~~last~~ units in L layer connecting to j^{th} unit of $L-1$ layer then:

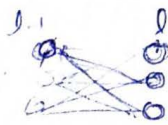
$$\delta_j^{(L-1)} = \sum_i \delta_i^{(L)} w_{ij}^{(L)}$$



Now propagating this $\delta_j^{(L-1)}$

$$\delta_j^{(L-1)} = f'(z_j^{(L-1)}) \sum_i \delta_i^{(L)} w_{ij}^{(L)}$$

Now this error goes back to other units
this is denoted as $\delta_i^{(L-1)}$



$$\delta_j^{(l-1)} = f'(z_j^{(l-1)}) \sum_i \delta_i^{(l)} w_{ij}^{(l-1)}$$

Obtain a matrix of $\delta_j^{(l-1)}$ that maps error δ unit to unit

$$\nabla W^{(l)} = \delta^{(l+1)} a^{(l)T}$$

$$\delta^{(l+1)} \rightarrow S_{l+1} \times 1$$

$$a^{(l)} \rightarrow 1 \times S_l$$

$$W^{(l)} \rightarrow S_{l+1} \times S_l$$

$$\delta^{(l)} = f'(z^{(l)}) \circ (W^{(l)T} \cdot \delta^{(l+1)})$$

$$W^{(l)} \rightarrow S_{l+1} \times S_l$$

$$\nabla W^{(l)} = \begin{bmatrix} \delta_1^{(l+1)} \\ \vdots \\ \delta_{S_{l+1}}^{(l+1)} \end{bmatrix} \begin{bmatrix} a_1^{(l)} & a_2^{(l)} & \dots & a_{S_l}^{(l)} \end{bmatrix}$$

$1 \times S_{l+1}$

$$\nabla W^{(l)} = \begin{bmatrix} \delta_1^{(l+1)} a_1^{(l)} & \delta_1^{(l+1)} a_2^{(l)} & \dots & \delta_1^{(l+1)} a_{S_l}^{(l)} \\ \delta_2^{(l+1)} a_1^{(l)} & \delta_2^{(l+1)} a_2^{(l)} & \dots & \delta_2^{(l+1)} a_{S_l}^{(l)} \\ \vdots & \vdots & \ddots & \vdots \\ \delta_{S_{l+1}}^{(l+1)} a_1^{(l)} & \delta_{S_{l+1}}^{(l+1)} a_2^{(l)} & \dots & \delta_{S_{l+1}}^{(l+1)} a_{S_l}^{(l)} \end{bmatrix}$$

$S_{l+1} \times S_{l+1}$

where S_{l+1} is size of layer $(l+1)$

& S_l is size of layer (l)

$$W \leftarrow W - \alpha (\nabla W)$$

$$\left\{ W_{S_{l+1} \times S_l} \leftarrow W_{S_{l+1} \times S_l} - \alpha (\nabla W_{S_{l+1} \times S_l}) \right\}$$

$$\delta^{(l)} = f'(z^{(l)}) \circ (W^{(l)T} \cdot \delta^{(l+1)})$$

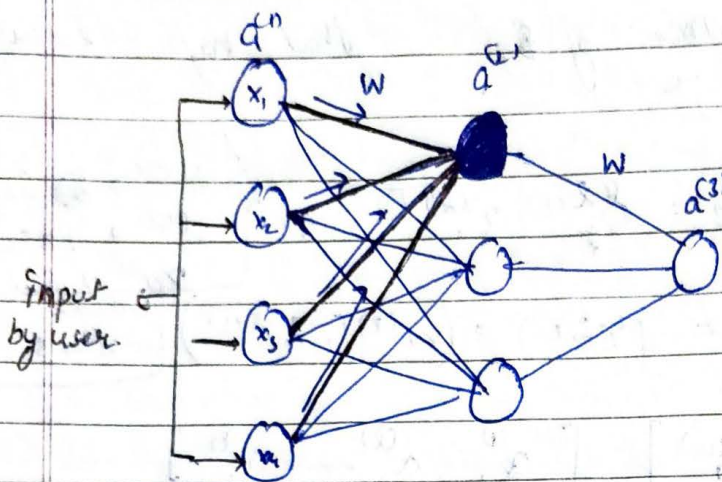
dimension

$$W^{(l)T} \rightarrow S_l \times S_{l+1}$$

$$\delta^{(l+1)} \rightarrow S_{l+1} \times 1$$

$$\delta^{(l)} \rightarrow S_l \times 1$$

Visualization of feed-forward propagation



In forward propagation, we select one unit of forward layer,

here $a_1^{(2)}$ (coloured)

take all the weights connecting to that layer from previous layer.

$$w_{ij}^{(1)}$$

and calculate 'pre-activation' at that layer by:

$$z_1^{(2)} = \sum_j w_{1j} x_j$$

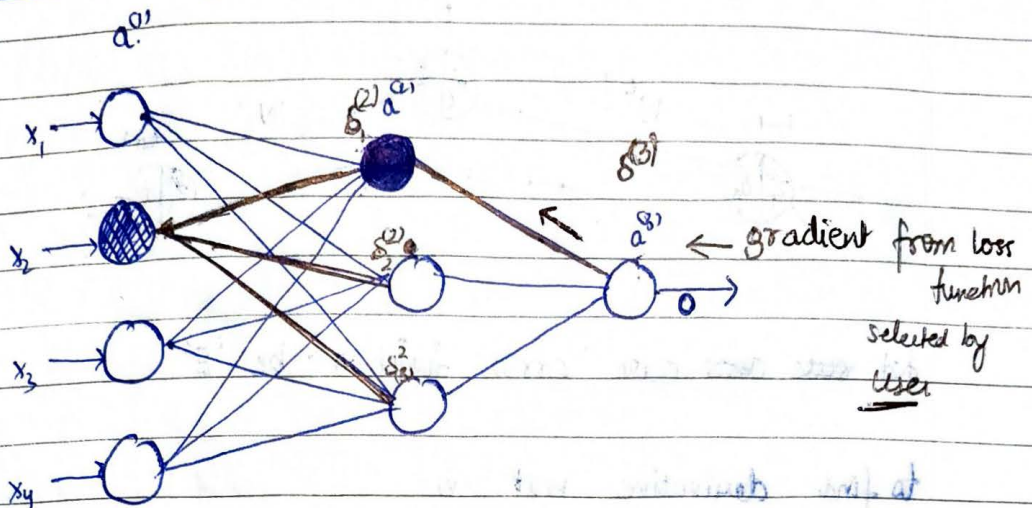
and $a_1^{(2)} = \sigma(\sum_j w_{1j} x_j)$ {where $\sigma(x)$ is non-linearity}

$$(1 \times 1) \cdot (1 \times 4) \cdot (4 \times 3) = 1 \times 3$$

neuron

$$(1 \times 1) \cdot (1 \times 4) \cdot (4 \times 3) = 1 \times 3$$

Visualization of back-prop



In backward propagation, we go from output to input layer propagating 's' (gradients)

$$\delta^{(l)} = (W^{(l+1)^T} \delta^{(l+1)}) \odot f'(z^{(l)}) \quad l \rightarrow L-1 \text{ to } 1$$

where, $\delta^{(last)}$ is calculated before by the help of cost function

here, the selected, $a^{(2)}_1$

$$\delta^{(2)}_1 = \delta^{(3)} \cdot w^{(2)}_{11} \cdot f'(z^{(2)}_1)$$

$$\nabla W^{(l)} = \delta^{(l+1)} \cdot a^{(l)T} \quad \left\{ \begin{array}{l} \text{outer product: as dimension} \\ \text{is expanding} \\ (5 \times 1) \times (1 \times 5) = 5 \times 5 \end{array} \right.$$

Now, propagating ahead let's consider another unit $a^{(1)}_2$

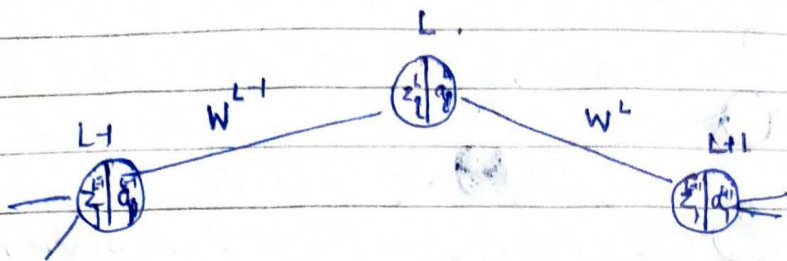
$$\delta^{(1)}_2 = \sum_i \delta^{(2)}_i \cdot w^{(1)}_{i2} \cdot f'(z^{(1)}_2)$$

$$\nabla W^{(1)} = \delta^{(2)} \cdot a^{(1)T}$$

$$\nabla W_{12} = \delta^{(2)}_1 \cdot x_2$$

Just for w_{12}

- How to get generalized gradients?



let our error function be 'E'.

to find derivative w.r.t. w^{L-1}

$$\frac{\partial E}{\partial w_{ij}^{L-1}} = \frac{\partial E}{\partial a_i^L} \frac{\partial a_i^L}{\partial w_{ij}^{L-1}}$$

$$= \left(\frac{\partial E}{\partial a_i^L} \frac{\partial a_i^L}{\partial z_i^L} \right) \left(\frac{\partial z_i^L}{\partial w_{ij}^{L-1}} \right)$$

δ^L because every term on L

We know,

$$a_i^L = f(z_i^L) \quad \text{if } f(x) \text{ is a non-linearly function}$$

$$\Rightarrow \frac{\partial a_i^L}{\partial z_i^L} = f'(z_i^L)$$

$$\frac{\partial z_i^L}{\partial w_{ij}^{L-1}} = \frac{\partial \left(\sum_j w_{ij}^{L-1} a_j^{L-1} \right)}{\partial w_{ij}^{L-1}} = a_j^{L-1}$$

$$\frac{\partial E}{\partial w_{ij}^{L-1}} = \frac{\partial E}{\partial a_i^L} \cdot f'(z_i^L) \cdot a_j^{L-1}$$

$$\frac{\partial E}{\partial w_{ij}^{L-1}} = \underbrace{\frac{\partial E}{\partial a_i^L} \cdot f'(z_i^L)}_{\delta^L} \cdot a_j^{L-1} = \delta^L a_j^{L-1}$$

$$\delta_i^L = \frac{\partial E}{\partial a_i^L} \cdot f'(z_i^L)$$

$$\frac{\partial E}{\partial a_i^L} = \frac{\partial E}{\partial a_i^{L+1}} \frac{\partial a_i^{L+1}}{\partial a_i^L}$$

$$\delta_i^L = \delta_i^{L+1} \frac{\partial a_i^{L+1}}{\partial a_i^L}$$

$$\frac{\partial E}{\partial a_i^L} = \underbrace{\left(\frac{\partial E}{\partial a_i^{L+1}} \frac{\partial a_i^{L+1}}{\partial z_i^{L+1}} \right)}_{\delta_i^{L+1}} \left(\frac{\partial z_i^{L+1}}{\partial a_i^L} \right)$$

$$\delta_i^L = \delta_i^{L+1} \left(\frac{\partial z_i^{L+1}}{\partial a_i^L} \right) \left[f'(z_i^L) \right]$$

$$\delta_i^L = \delta_i^{L+1} \left(\frac{\partial (W_{ij}^{L+1} a_j^{L+1})}{\partial a_i^L} \right) f'(z_i^L)$$

$$\delta_j^L = \sum_i \delta_i^{L+1} W_{ij}^{L+1} \cdot f'(z_i^L)$$

$$\frac{\partial E}{\partial W_{ij}^{(L)}} = \delta_j^{(L+1)} a_i^L$$

~~Backpropagation~~

where

$$\delta_j^L = \sum_i \delta_i^{L+1} W_{ij}^{L+1} f'(z_i^L)$$

→ can be computed before with forward prop

For vectorization:

$$\nabla_{W^{(L)}} \rightarrow \delta^{(L+1)} a^{(L)T} + \underbrace{\lambda W^{(L)}}_{\text{regularization term}} \left[S^{(L+1)} \times S^{(L)} \rightarrow (S^{(L+1)} \times I) \times (I \times S^{(L)}) \right]$$

$$\delta^L = (W^{(L)T} \delta^{(L+1)}) f'(z^L) S^{(L)} \rightarrow \frac{\delta^{(L+1)} S^{(L+1)T}}{(S^{(L)} \times S^{(L+1)}) \times (S^{(L+1)} \times I)}$$

• Why go through derivations!

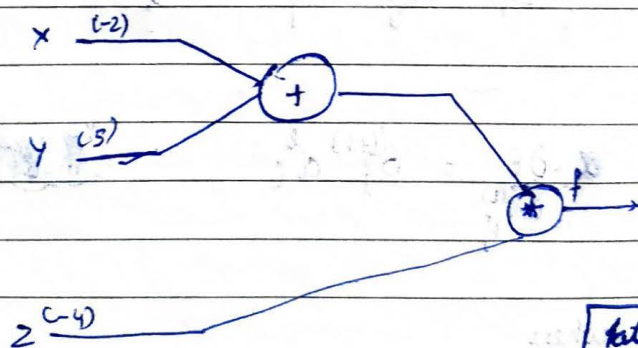
1. ~~Act~~ Actual understanding of math behind most of deep learning
2. Backprop can be an imperfect abstraction e.g. — issues such as vanishing gradient
3. Enables you to design models, think of and implement completely new models that aren't yet supported by any framework.

(Giving general idea)

• Explanation of Backprop: Circuit

Function as circuits

$$f(x, y, z) = (x + y)z \quad ; \quad x = -2, y = 5, z = -4$$



$q = x + y$	$\frac{\partial q}{\partial x} = 1$	$\frac{\partial q}{\partial y} = 1$
$f = qz$	$\frac{\partial f}{\partial q} = z$	$\frac{\partial f}{\partial z} = q$

Let's trace the backward flow

add gate: gradient distributed

mul gate: gradient sent

multiply: gradient sent

mul gate = min(x, y)

suppose x > y
Then we will route the grad to y

Hint: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

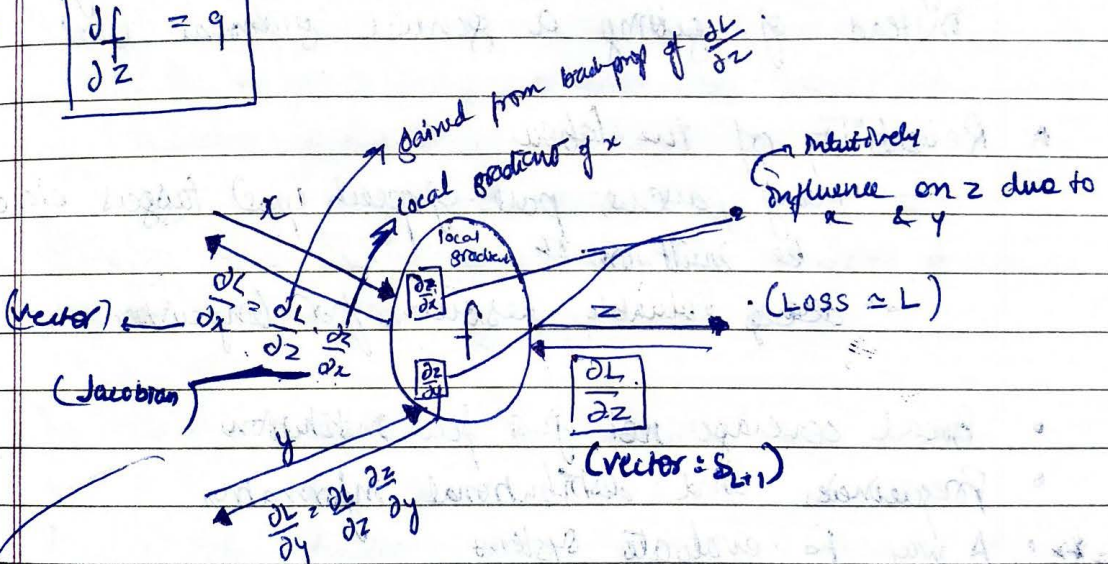
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x}$$

$$\frac{\partial f}{\partial x} = z$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial y}$$

$$\frac{\partial f}{\partial y} = z$$

$$\frac{\partial f}{\partial z} = g$$



• Another example { In lecture 5 ; similar to circuit ex }

- I (compute local gradient/derivative)
- II Calculate the derivative with the value obtained in forw prop
- III To completely compute the gradient of function :
(II) * (gradient of the successive layer)