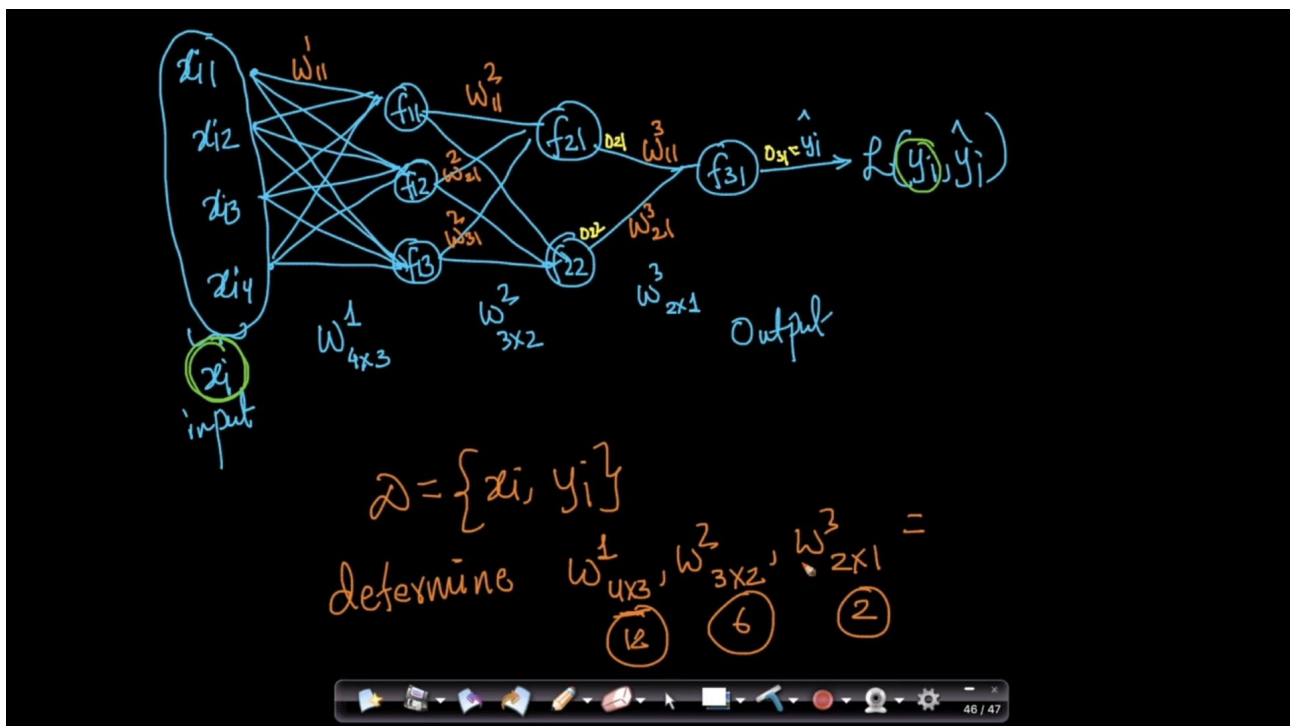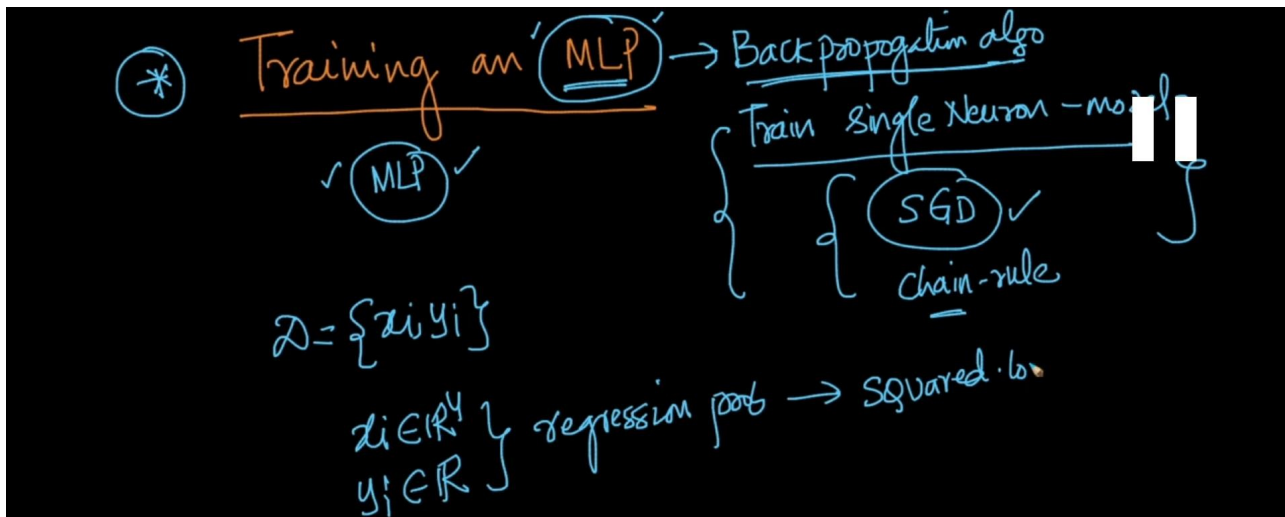We're using squared loss for our regression problem





we have to determine the weights above, w1 has 12 weights, w2 has 6, w3 has 2 weights

total of 20 weights

training an MLP means computing all these weights

the loss function is given below, we can choose any of the given regularisers, we ignore for it for now for simplicity :

the text in green in above slide summarises it at once

2nd step is we do GD or SGD

in SGD we initialise variables 1st randomly, there are techniques todo it in a better way

at iteration 0, $w_{ij}^{k}$ will have a random value



a very imp thing is we compute this partial derivative, which is circled above

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{11}^3}$$

above is chain rule



$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{21}^3}$$



$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w_{11}^2} \quad \bigg\} \text{ -sanity check}$$

$$\boxed{W^2} \qquad \frac{\partial L}{\partial W_{11}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W_{11}^2}$$

$$\frac{\partial L}{\partial W_{21}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W_{21}^2}$$

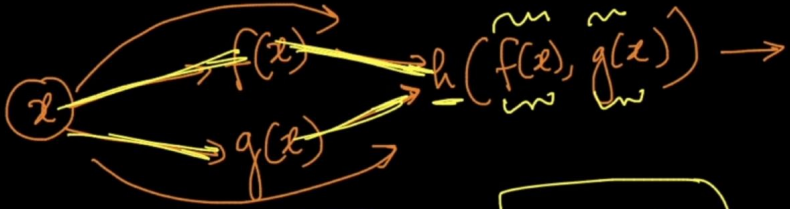$$\frac{\partial L}{\partial W_{31}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W^2.}$$

$$\boxed{W^3} \qquad \frac{\partial L}{\partial W_{11}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial W_{11}^3} \qquad \longleftarrow \text{chain rule}$$

$$\frac{\partial L}{\partial W_{21}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial W_{21}^3} \qquad \longleftarrow \text{chain rule}$$

next case is more interesting since O11 is taking 2 paths



$$\frac{\partial h}{\partial x} = \boxed{\frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x}} \enspace \boxed{+} \enspace \boxed{\frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}} \qquad \longleftarrow \text{chain rule}$$
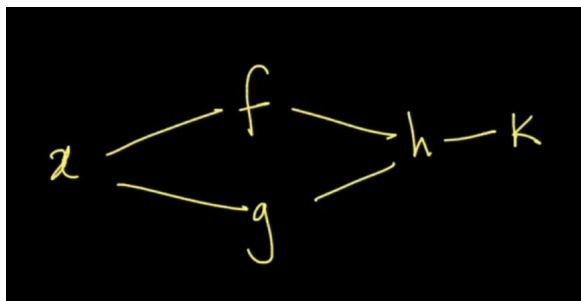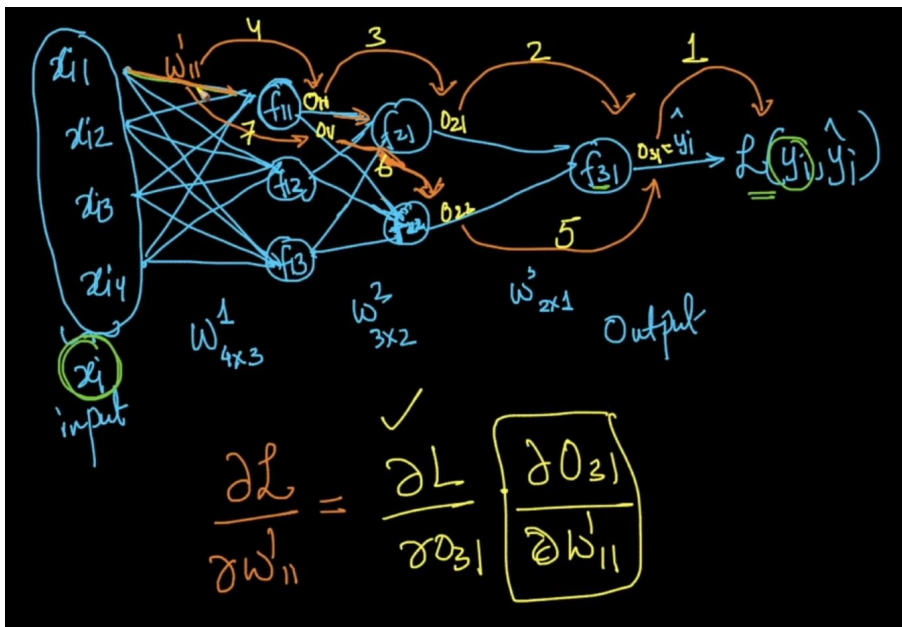
$$\frac{\partial L}{\partial W_{11}} = \frac{\partial L}{\partial O_{31}} \cdot \boxed{\frac{\partial O_{31}}{\partial W_{11}'}}$$

Neural network diagram with nodes $x_{i1}, x_{i2}, x_{i3}, x_{i4}$ (input), weights $W_{11}$, $f_{11}$ ($O_{11}$, $O_{1v}$), $f_{12}$, $f_{13}$, $f_{13}$, $f_{21}$ ($O_{21}$), $f_{22}$ ($O_{22}$), $f_{31}$ ($O_4 = \hat{y}_i$), $L(y_i, \hat{y}_i)$. Weight matrices $W^1_{4\times3}$, $W^2_{3\times2}$, $W^3_{2\times1}$, Output.

---



this is the structure that we have

---



$$\frac{\partial k}{\partial x} = \boxed{\frac{\partial k}{\partial h}} \cdot \frac{\partial h}{\partial x}$$

$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}$$

$$\frac{\partial k}{\partial x} = \frac{\partial k}{\partial h} * \left\{ \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x} \right\}$$

---



$$\frac{\partial k}{\partial x} = \boxed{\frac{\partial k}{\partial h}} \boxed{\frac{\partial h}{\partial x}}$$

$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}$$

$$\frac{\partial k}{\partial x} = \frac{\partial k}{\partial h} * \left\{ \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x} \right\}$$

$$\frac{\partial O_{31}}{\partial w'_{11}} = \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w'_{11}} + \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w'_{11}}$$

$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial O_{31}} \cdot \boxed{\frac{\partial O_{31}}{\partial w'_{11}}}$$



$$\partial w^1_{11} = \frac{\partial L}{\partial O_{31}} \cdot \left\{ \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \frac{\partial O_{11}}{\partial w'_{11}} \; \text{(+)} \; \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w'_{11}} \right\}$$

-----------------------------------------------------------------------------------------------------

Memoisation :

saving stuff we're recomputing



$$\frac{\partial L}{\partial w^3_{11}} = \boxed{\frac{\partial L}{\partial O_{31}}} \cdot \frac{\partial O_{31}}{\partial w^3_{11}} \quad \leftarrow \text{chain rule}$$

$$\frac{\partial L}{\partial w^3_{21}} = \boxed{\frac{\partial L}{\partial O_{31}}} \cdot \frac{\partial O_{31}}{\partial w^3_{21}} \quad \leftarrow \text{chain rule}$$

compute once & reuse

$W^2$

$$\frac{\partial L}{\partial W_{11}^3} = \boxed{\frac{\partial L}{\partial O_{31}}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W_{11}^2}$$

$$\frac{\partial L}{\partial W_{21}^2} = \boxed{\frac{\partial L}{\partial O_{31}}} \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W_{21}^2}$$

$$\frac{\partial L}{\partial W_{31}^2} = \boxed{\frac{\partial L}{\partial O_{31}}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W_{31}^2}$$
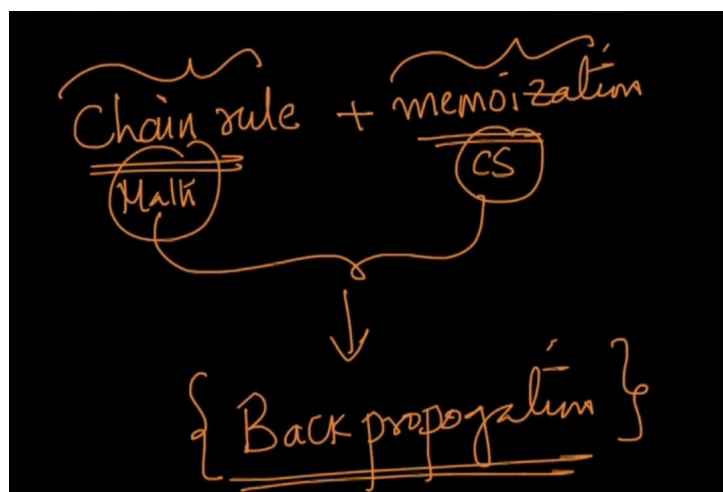
a lot of stuff is used again and again



$W^2$

$$\frac{\partial L}{\partial W_{11}^3} = \boxed{\frac{\partial L}{\partial O_{31}}} \cdot \boxed{\frac{\partial O_{31}}{\partial O_{21}}} \cdot \frac{\partial O_{21}}{\partial W_{11}^2}$$

$$\frac{\partial L}{\partial W_{21}^2} = \boxed{\frac{\partial L}{\partial O_{31}}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W_{21}^2}$$

$$\frac{\partial L}{\partial W_{31}^2} = \boxed{\frac{\partial L}{\partial O_{31}}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial W_{31}^2}$$

for slightly more memory, we get a huge speedup



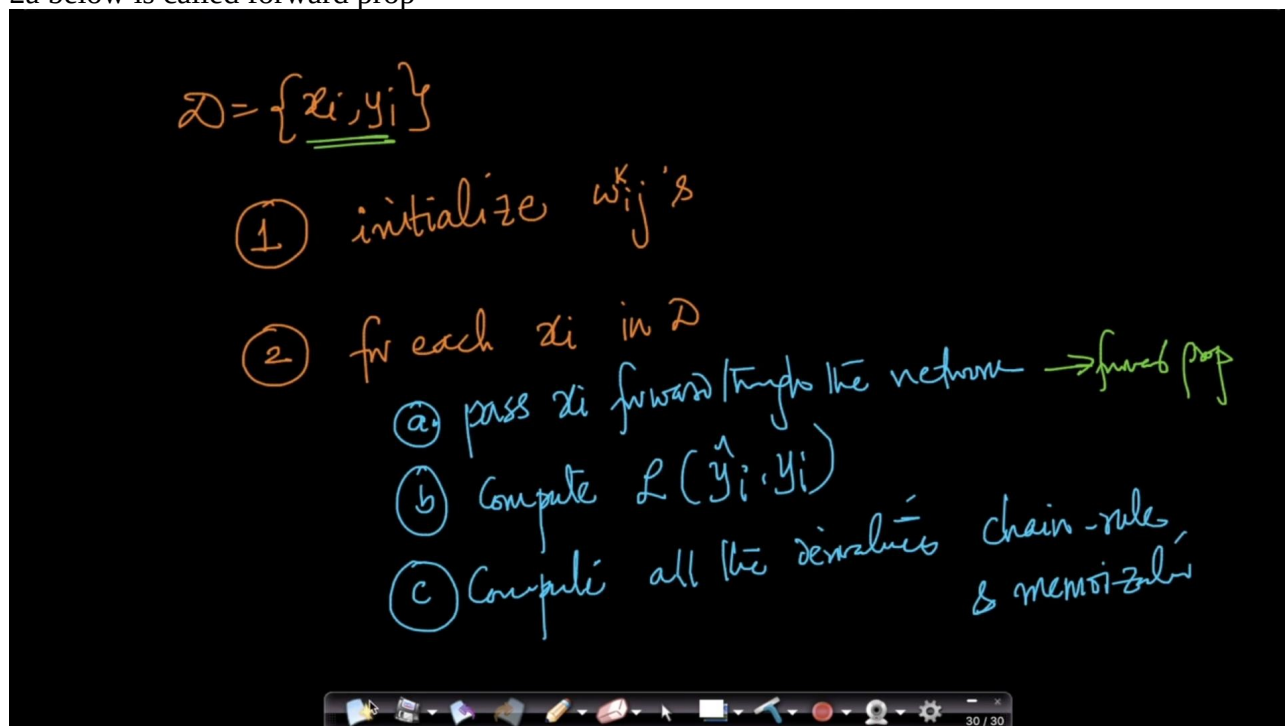Chain rule + memoization

Math        CS

↓

{ Back propogation }

backprop explained :
2a below is called forward prop

$$\mathcal{D} = \{x_i, y_i\}$$

① initialize $w_{ij}^k$'s

② for each $x_i$ in $\mathcal{D}$
   ⓐ pass $x_i$ forward/through the network → forward prop
   ⓑ Compute $\mathcal{L}(\hat{y}_i, y_i)$
   ⓒ Compute all the derivatives  chain-rule,
                                   & memoization

$$\{ ⓓ \text{ Update weights from end of the network to the start} \}$$

backward prop

③ repeat step 2 till convergence } → Lr. reg & optimization

$$\left(w_{ij}^k\right)_{new} \approx \left(w_{ij}^k\right)_{old}$$

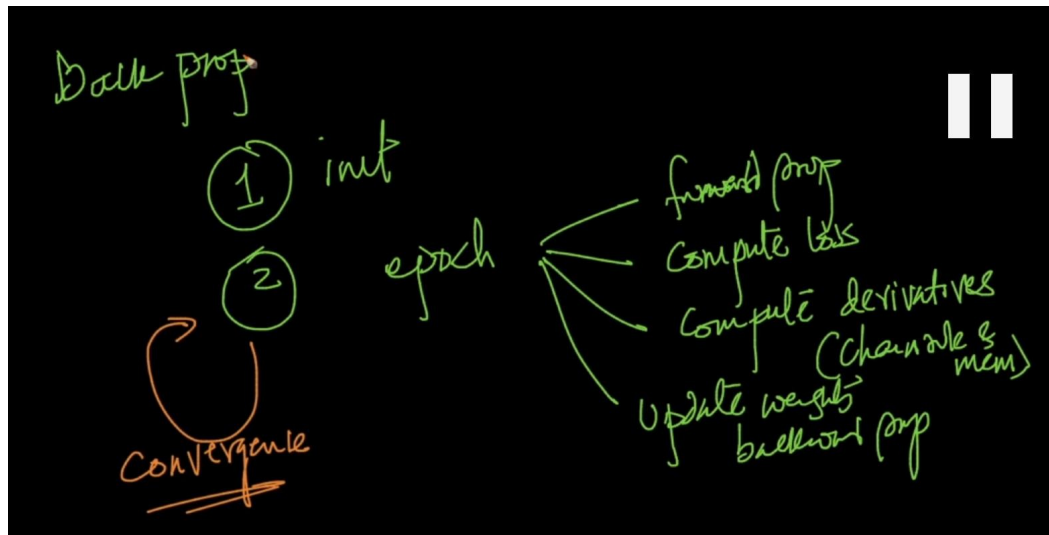epoch means you've passed all the points in the dataset once through neural network

the number of times we pass them through the network is called epoch

$$\mathcal{D} \rightarrow \{x_i, y_i\} \rightarrow \text{epoch}$$
$$5 \text{ times} \rightarrow 5 \text{ epochs}$$

passing = computing loss and updating weights

we run for multiple epochs irl

how do we pick each of the point $x_i$ in D : we should pick points uniformly at random
pick a point and do step 2



repeat step 2 till convergence

"back prop is a multi epoch training methodology where we leverage chain rule and memoisation to update weights"

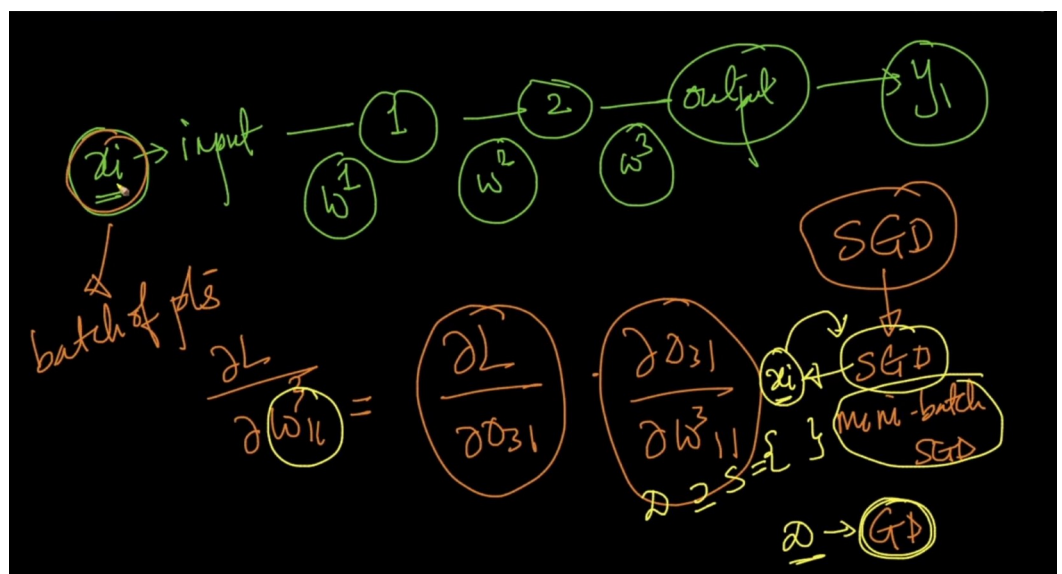backprop only works if activation functions are differentiable

if it's easy and fast differentiable, you can speedup the training of the neural network using backprop algo

-------------------------------------------------------------------------------------------------------

in SGD you take one point at a time to compute derivates to update weights
in mini batch sgd you take a set of points
in gradient descent you take all of the points

----------------------------------------------------------------------------------------------------

below is sort of like sgd because we're sending one point each



the derivates above are computed using just one point
this is one approach

----------------------------------------------------------------------------------------------------

instead of sending 1 point, we send a set of points
the size S is called mini batch size

we can also do gradient descent based approach where we send all datapoints

the big problem is keeping all datapoints in RAM and computing derivates using the whole data can be extremely time consuming

people usually take one point at a time or mini batch based SGD(most popular)

-------------------------------------------------------------------------------------------------------

example :

the possible ram sizes are shown below
we'll have to run the loop 100 times to complete 1 epoch
this is more efficient than running the loop 10k times with one point each