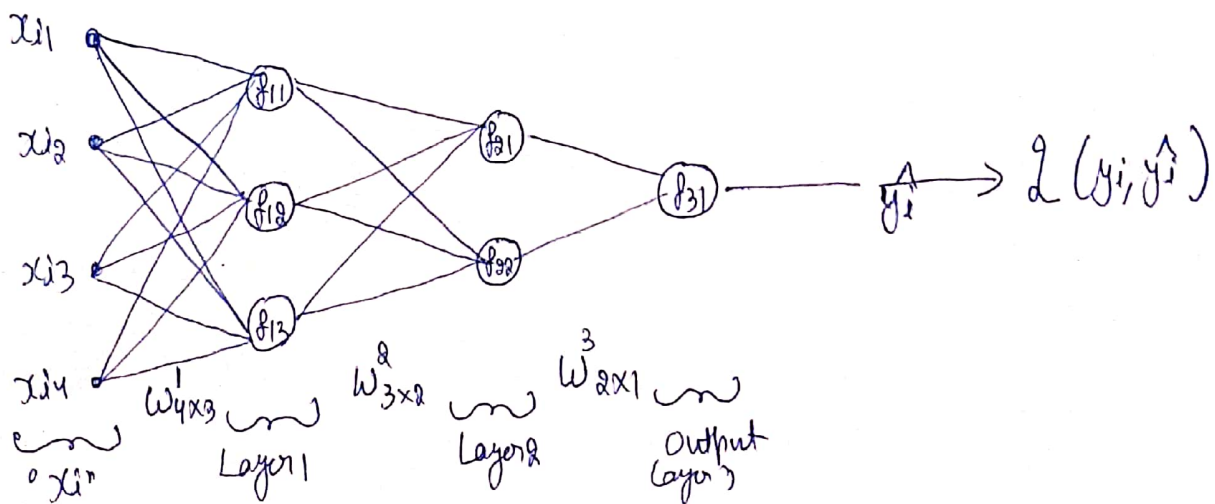The most important thing in "neural networks" is MLP (Multilayer Perceptron)

let's say we have a dataset "D" comprising of $x_i$'s & $y_i$'s

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

$x_i \in \mathbb{R}^4$ } How we want to train an MLP and this
$y_i \in \mathbb{R}$ } problem is a standard Regression problem

In standard regression problem, the simple loss function that we have is "squared loss"

let's draw an MLP with the notations we saw earlier :→



output is $\hat{y_i}$

We need a loss function, just like we saw in a single neuron model. So, $\hat{y_i}$ goes into a loss function represented by 'L'

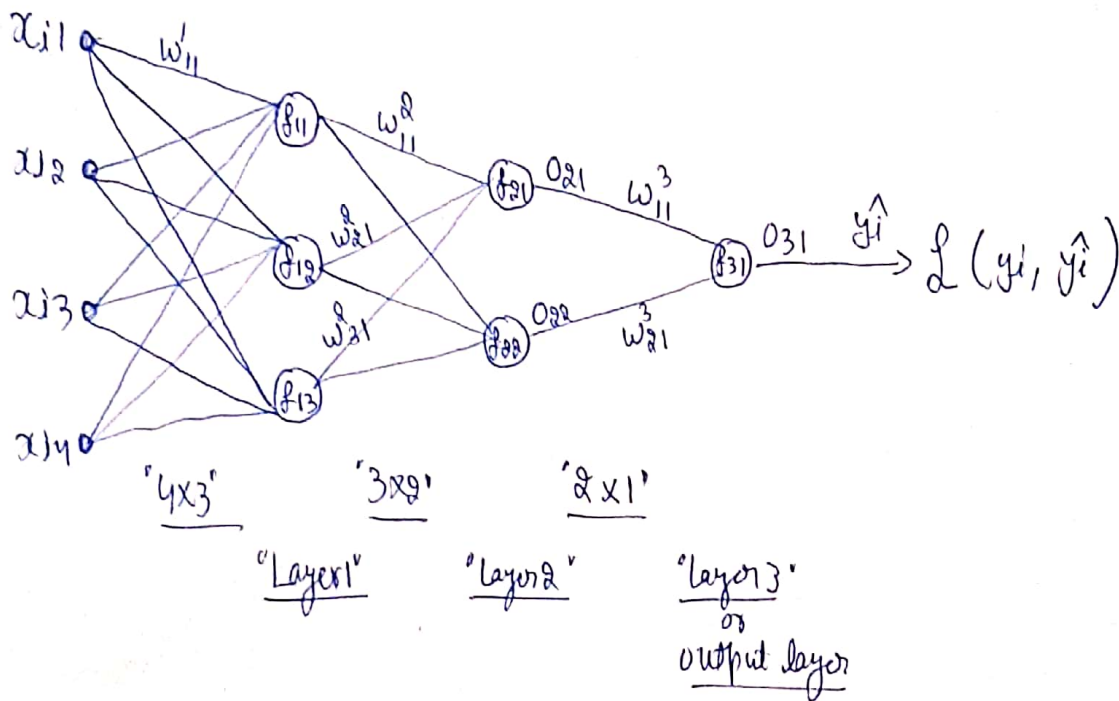$y_i$'s are from the dataset itself & $\hat{y_i}$ is the predicted values

P.T.O

Weights in the first layer are represented using a 4×3 matrix

$$\left[ W^1_{4\times3} \right]$$

In layer second, they are represented using $W^2_{3\times2}$

In layer third, they are represented using $W^3_{8\times1}$

The weights & output associated with the MLP as per the notation we saw earlier is as given below :-



$X_{i1}$   $W^1_{11}$   $f_{11}$   $W^2_{11}$   $f_{21}$ $O_{21}$   $W^3_{11}$

$X_{i2}$   $f_{12}$   $W^2_{21}$   $f_{31}$ $O_{31}$ $\hat{y_i} \rightarrow L(y_i, \hat{y_i})$

$X_{i3}$   $W^2_{31}$   $f_{22}$ $O_{22}$   $W^3_{21}$

$X_{i4}$   $f_{13}$

'4×3'   '3×2'   '2×1'

"Layer1"   "Layer2"   "Layer3"
or
output layer

Let's now try to pose this problem as an optimization problem:-

Training this perceptron (MLP) given a dataset $D = \{x_i, y_i\}$ basically means, we need to determine the weights associated with the edges.

So, during training we need to determine the weights & they are represented as

$$\left[ W^1_{4\times3}, \quad W^2_{3\times2} \quad \& \quad W^3_{2\times1} \right]$$

So in total we have to determine these $20$ weights + training basically means, computing these $20$ weights.

Well now go step by step.

Step no.1 :- Define the loss function + the loss function here is the squared loss.

$$\mathcal{L} = \underbrace{\sum_{i=1}^{n} (y_i - \hat{y_i})^2}_{\text{Squared loss}} + \text{regularizer}$$

(Note :- Regularizer is on the weights)

Note :- Any weight can be represented at $\boxed{W_{ij}^{K}}$

$L_2$ regularizer is defined on the weights as :-

$$\sum_{i,j,k} (W_{ij}^{K})^2$$

$L_1$ regularizer is defined on the weights as :-

$$\sum_{i,j,k} |W_{ij}^{K}| \longrightarrow \text{absolute value}$$

just ignore regularizer for now, to make optimization problem simpler.

Let us define $\mathcal{L}i$ which is basically the squared error for a single training point.

$$\mathcal{L}i = (y_i - \hat{y_i})^2 \qquad \text{Now } \mathcal{L} \text{ can be written as :-}$$

$$\boxed{\mathcal{L} = \sum_{i=1}^{n} \mathcal{L}i + \text{regularizer}}$$

$$L = \sum_{i=1}^{n} L_i + \text{regularizer}.$$

Step 2: form an optimization problem.

Now the optimization problem can be written as:

$$\begin{bmatrix} \min_{w^1, w^2, w^3} L \begin{cases} \to \text{Squared loss} \\ \to \text{regularizer} \end{cases} \end{bmatrix}$$

Minimize the squared distance between the "predicted value" and the "actual value", for all the weight matrices.

Here "L" has two parts (squared loss & regularization)

In other words we can write this optimization problem as:

Minimize Loss, for all the possible weights $(w_{ij}^k)$

$$\downarrow$$

$$\begin{bmatrix} \text{for all possible values of} \\ i, j + k \end{bmatrix}$$

$$\boxed{\min_{w_{ij}^k} L}$$

Where 'L' is nothing but the sum of "square loss" & "regularizer"

Step 3: Apply "Gradient-Descent" to solve the optimization problem.

a) Initialization of weights. $w_{ij}^k$ randomly.

b) We need to iteratively keep on changing these weights

$$\left[ (w_{ij}^k)_{new} = (w_{ij}^k)_{old} - \eta \frac{\partial L}{\partial w_{ij}^k} \right] \to \text{update rule}$$

© Perform updates till convergence.

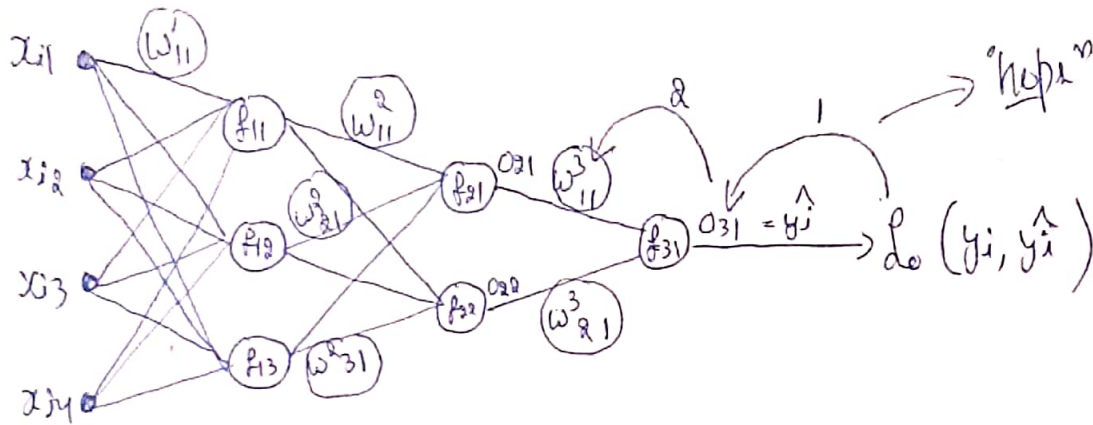Convergence means, you keep changing $w^k_{ij}$ till the time $(w^k_{ij})_{new} + (w^k_{ij})_{old}$ are almost very very close.

Now the next important step in this is to compute its partial derivative.

$$\boxed{\dfrac{\partial L}{\partial w^k_{ij}}}$$

Let's now compute these partial derivatives :



for simplicity let's just compute the derivatives of $L_o$ with respect to the weights mentioned on these edges

$$\frac{\partial L}{\partial w^3_{11}}$$

Here $w^3_{11}$ goes into $f_{31}$ & it impacts $O_{31}$ & $O_{31}$ goes into $L$ & it impacts $L$

$$\therefore \frac{\partial L}{\partial w^3_{11}} = \frac{\partial L}{\partial O_{31}} * \frac{\partial O_{31}}{\partial w^3_{11}}$$

Similarly let's now compute $\dfrac{\partial L}{\partial w^3_{21}}$

$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial w_{21}^3}$$

Here we are just following the path to compute there ["partial derivatives".]

So, we have computed the partial derivatives w·r·t to weights on third layer.

Now let's compute partial derivatives w·r·t to the $w_{3\times2}^2$ matrix (weight matrix)

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial w_{11}}$$

$$\frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial w_{21}^2}$$

$$\frac{\partial L}{\partial w_{31}^2} = \frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial w_{31}^2}$$

These are the weights in $W_2$

Although there are 6 weights in total, we have used only three weights.
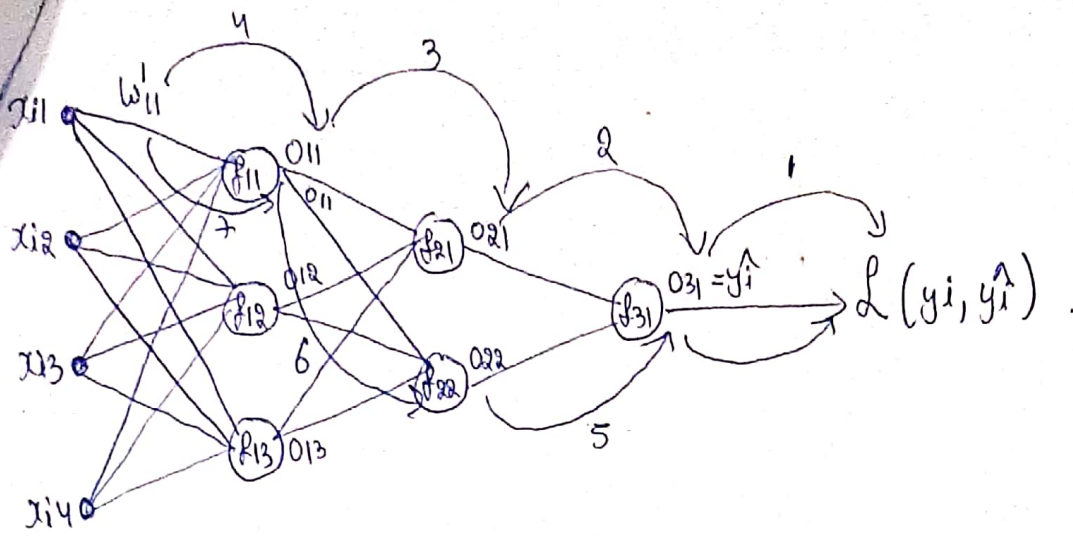Similarly we can compute other weights also.

Now let's compute partial derivatives w·r·t to the first weight matrix.

Let's compute partial derivative $\frac{\partial L}{}$ w·r·t $w_{11}^1$, This is little bit tricky

P.T.O

$x_{i1}$ , $w'_{11}$

4

3

2

1

$O_{11}$

$f_{11}$

$O_{11}$

$x_{i2}$

7

$f_{12}$

$O_{12}$

$f_{21}$ $O_{21}$

6

$f_{22}$ $O_{22}$

$f_{31}$ $O_{31} = \hat{y_i}$

$\rightarrow \mathcal{L}(y_i, \hat{y_i})$ .

$x_{i3}$

$f_{13}$ $O_{13}$

5

$x_{i4}$

"$x_i$"
Input

we are Computing $\dfrac{\partial L}{\partial' w_{11}}$ .

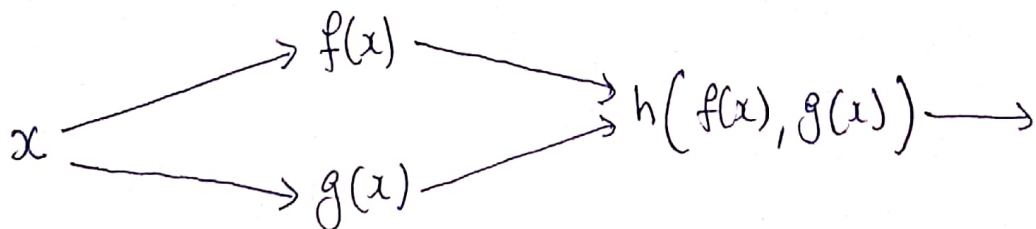Here we can see $w'_{11}$ goes into '$f_{11}$' & impacts $O_{11}$ but $O_{11}$ instead of going to a single neuron, goes as an input to two neurons. $(f_{21})$ & $(f_{22})$ & impacts both "$O_{21}$" & "$O_{22}$"

So, there are basically two paths from $w'_{11}$ to $\mathcal{L}$.

∴ we have to choose another rule of chain rule process.

This can be Computed using a simple rule in Calculus.

let's say we have a Variable '$x$' which goes & impacts two functions $f(x)$ & $g(x)$ & now let's assume we have a function '$h$' which is impacted by both $f(x)$ & $g(x)$ as shown below :-



$x$ $\rightarrow f(x) \rightarrow$
$\rightarrow g(x) \rightarrow$
$\rightarrow h(f(x), g(x)) \longrightarrow$

'$x$' is an input to $f(x)$ which generates some output & this output goes into function '$h$', similarly '$x$' is an input to $g(x)$ which generates some output & this output again goes into function '$h$'.

So, here we have two paths from ['x'] to ['h'].

$$\left[ \frac{\partial h}{\partial x} = \boxed{\frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x}} + \boxed{\frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}} \right) \right]$$ "second path"

↓
"One path"

This is becoz "h" is a function of "f" + "g"

But here in our case, the path from $W'_{11}$ is divergeing at ['O11'] & later these two paths are converging at ['O31'].

Here we can see step 1 is common in both these paths.

These paths can be analysed as is



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

That is the structure that we actually have in the MLP.

Now we have to figure out how to compute $\frac{\partial h}{\partial x}$ → for that we have two paths

$$\frac{\partial R}{\partial x} = \frac{\partial R}{\partial h} * \left[ \frac{\partial h}{\partial f} * \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} * \frac{\partial g}{\partial x} \right]$$

$$\left[ \frac{\partial h}{\partial x} = \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x} \right]$$
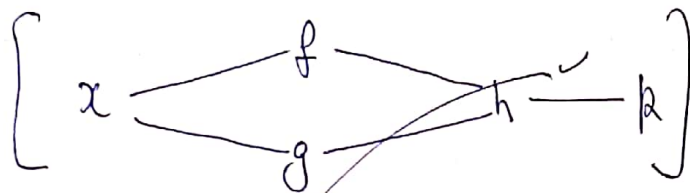
Now let's use that similar concept to find $\frac{\partial L}{\partial w'_{11}}$

$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial O_{31}} * \boxed{\frac{\partial O_{31}}{\partial w_{11}}}$$ ⟶ Now we need to compute this

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

$$\therefore \frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial O_{31}} * \left\{ \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w'_{11}} + \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w'_{11}} \right\}$$

similarly we can compute <u>other weights also belonging to first wei-</u>
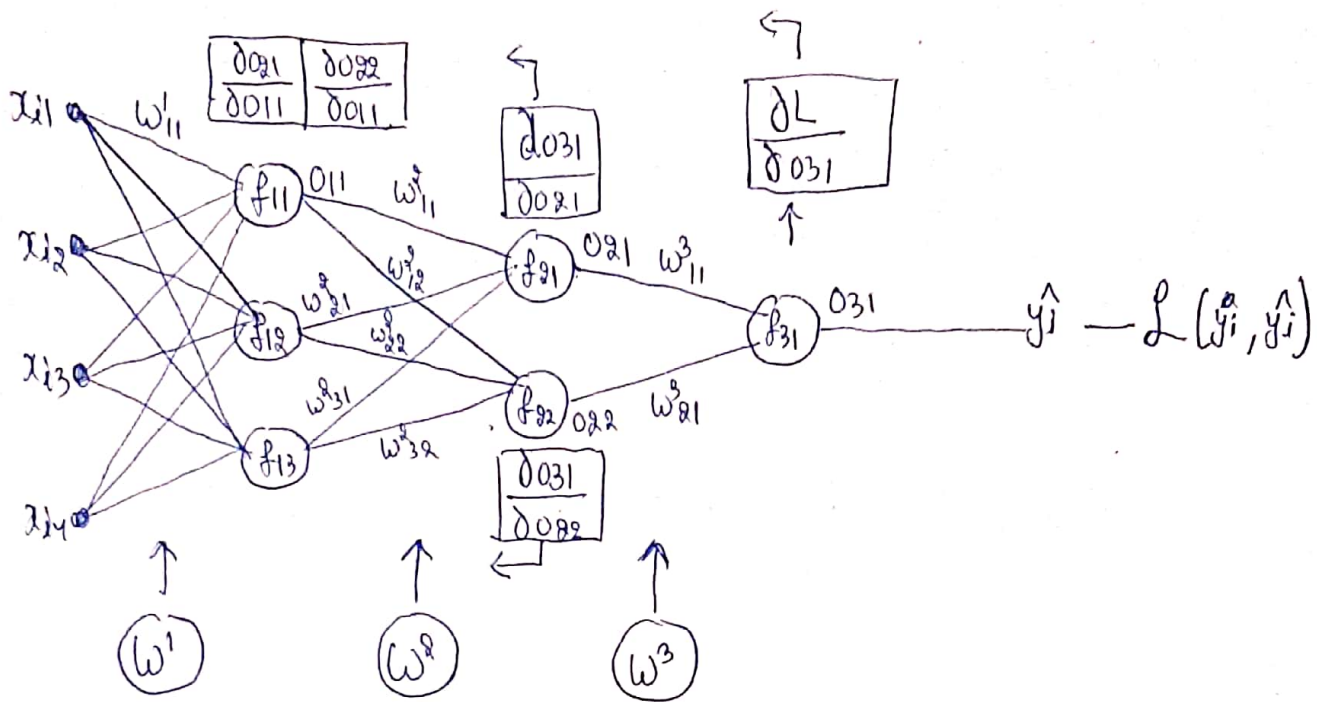<u>ght matrix "</u>]

Now once we have computed all their derivatives we can quickly update the weights associated with the edges: as,

$$\boxed{(w^{k}_{ij})_{new} = (w^{k}_{ij})_{old} - \eta \frac{\partial L}{\partial w^{k}_{ij} \, old}}$$

↳ This is how we train an MLP.

―――――

P.T.O

# "Backpropagation"



The diagram shows a neural network with inputs $x_{i1}, x_{i2}, x_{i3}, x_{i4}$ feeding through hidden layers with nodes $f_{11}, f_{12}, f_{13}$ (with weight $w^1_{11}$), then $f_{21}, f_{22}$ (weights $w^2_{11}, w^2_{12}, w^2_{21}, w^2_{22}, w^2_{31}, w^2_{32}$), then $f_{31}$ (weights $w^3_{11}, w^3_{21}$), producing output $O_{31} \longrightarrow \hat{y_i} \longrightarrow \mathcal{L}(\hat{y_i}, y_i)$.

Boxes shown:
$$\boxed{\dfrac{\partial O_{21}}{\partial O_{11}} \quad \dfrac{\partial O_{22}}{\partial O_{11}}}$$
$$\boxed{\dfrac{\partial O_{31}}{\partial O_{21}}}$$
$$\boxed{\dfrac{\partial L}{\partial O_{31}}}$$
$$\boxed{\dfrac{\partial O_{31}}{\partial O_{22}}}$$

Intermediate outputs: $O_{11}, O_{21}, O_{22}, O_{31}$. Weights gradients $W^1, W^2, W^3$.

**Note :->** Backpropagation is one of the most important algorithm in whole of deep learning.

let's now see how it works :->

Imagine we have a dataset $D$ of points $\{x_i, y_i\}$

Now we will input every datapoint "$x_i$" thro' the network just like any optimization problem, let's go step by step.

① Initialize $W^k_{ij}$ (weights) all the weights randomly.

  Note :-> There are various ways of initializing the weights.

② for each "$x_i$" in $D$.

  [ Send each "$x_i$" through the network & here each "$x_i$" is having four features as it is a four dimensional data ]

(a) pass $x_i$ forward through the network

⤷ This is called forward propagation.
(becoz we are taking the data & sending it forward in the network)

So, we input one data point & we get the loss. w.g $(y_i + \hat{y_i})$

Now using this loss we can compute all the derivatives to update the weights.

(b) Compute loss, $\mathcal{L}(\hat{y_i}, y_i)$

(c) Compute all the derivates using chain rule

(d) Update weights from end of the network to the start. ⟶ This is called "back propagation"

So what we did, we send our input & passes it forward through the network & computed the loss, this is called forward propagation. Now we will do a backward propagation to update the weights from end of the network to the start.

The weights can be updated as's

$$(w^3_{11})_{new} = (w^3_{11})_{old} - \eta \left[ \frac{\partial L}{\partial w^3_{11}} \right]_{(w^3_{11} old)}$$

↓
This can be computed using the chain rule.

$$\left[ \frac{\partial L}{\partial w^3_{11}} = \frac{\partial L}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial w^3_{11}} \right]$$

Now once we have these derivatives, we can simply update the weights.

When we went from input to output, that is called.
["forward propagation".]

& when we are coming back, we are updating the weights until loss function & this is called ["Backward propagation".]

Note :- In forward propagation, we are sending the input & trying to compute the output which is $\hat{y}$ & In backward propagation, we are using the error or the loss $L(y_i, \hat{y}_i)$ to update the weights, so that next time, we send the input, these weights will be better tuned to give "lesser loss" or "lesser error".

③ Repeat step no. 2 till convergence.

Note :- At the end of step 2, all the data points have been seen once by the network.

There is a possibility that at the end of step no 2, still we have not converged, so, we have to repeat step no. 2 till convergence.

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

There is an important term called "epoch". So one "epoch" of training means that we have input all the points in our dataset once through the "neural network".

$$D \to \{x_i ; y_i\} \longrightarrow \text{Network}$$

$\underbrace{\qquad\qquad\qquad\qquad}$
This is called epoch

[ Passing all the data points through the network.
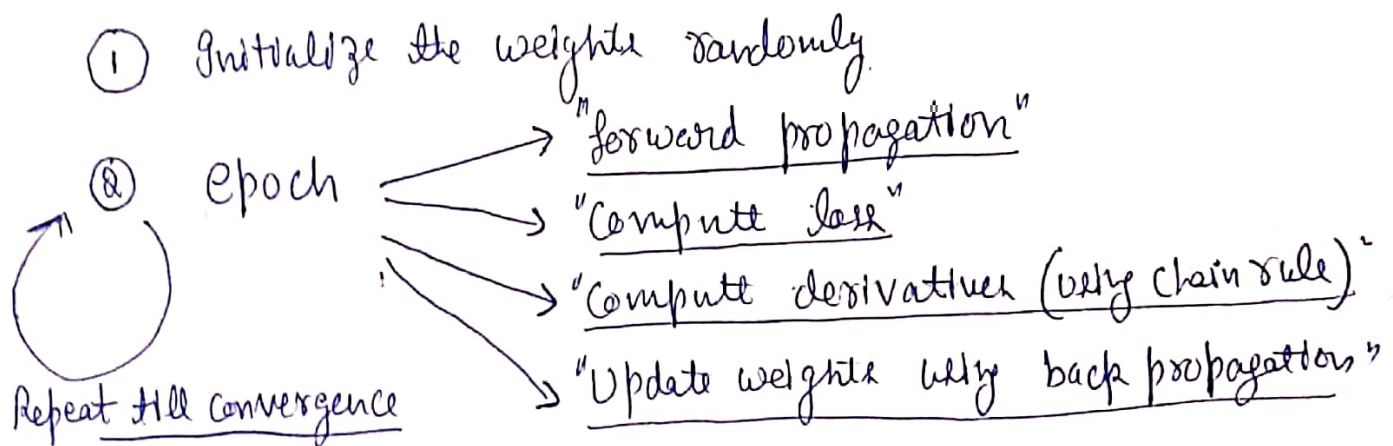
If we pass the datapoints five times through the network, (13)
it is called 5 _epochs_.

Note :> Completing each of step 2 is basically completing
one _epoch_.

In real world, for Neural network training, we run the
model for multiple epoch., which means we pass our dataset
through the network _multiple times_.

Ideally each of the point "$x_i$" in our dataset is picked.
randomly in a uniform way. ( we pick our _points_ _uniformly_ )
at _random_.

Back propagation intuitively is as explained below :>

① Initialize the weights randomly.

② epoch < "forward propagation"
          "Compute loss"
          "Compute derivatives (using chain rule)"
          "Update weights using back propagation"

Repeat till convergence

Note :> Back propagation is a multi-epoch training methodology
where we leverage chain rule to update weights.

vvvimp
Note :> "Back propagation" is only applicable if the activation
functions are differentiable.

P.T.O

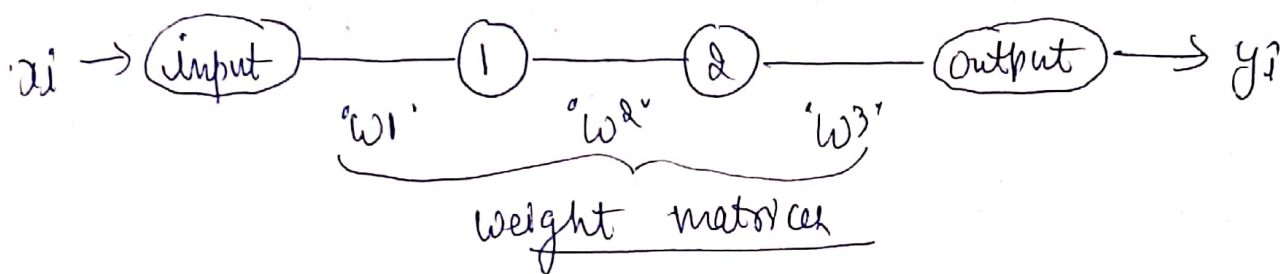<u>Note :></u> if the activation-function is easily differentiable than we can speed up the training of the "Neural Network" why backpropagation.

Second important thing is the batch-size that can increase the training rate.

If we look at the MLP that we just have, it comprises of two hidden layers apart from input & output layers

$xi \rightarrow$ (input) ———— (1) ———— (2) ———— (output) $\rightarrow yi$

'w1'    'w2'    'w3'

weight matrices

In the steps we saw earlier, we are sending one point through the network at a time forward, computing the loss, computing the derivatives and updating the weights. That takes more time.

So, instead of sending one single point through the network why can't we send a batch of points through the network

<u>Note :></u> In 'SGD' we take one point at a time to compute the derivatives to update the weights.

whereas in 'mini-batch' 'SGD' instead of taking one point at a time, we are taking a set of points to compute the derivatives

————— ठ —————