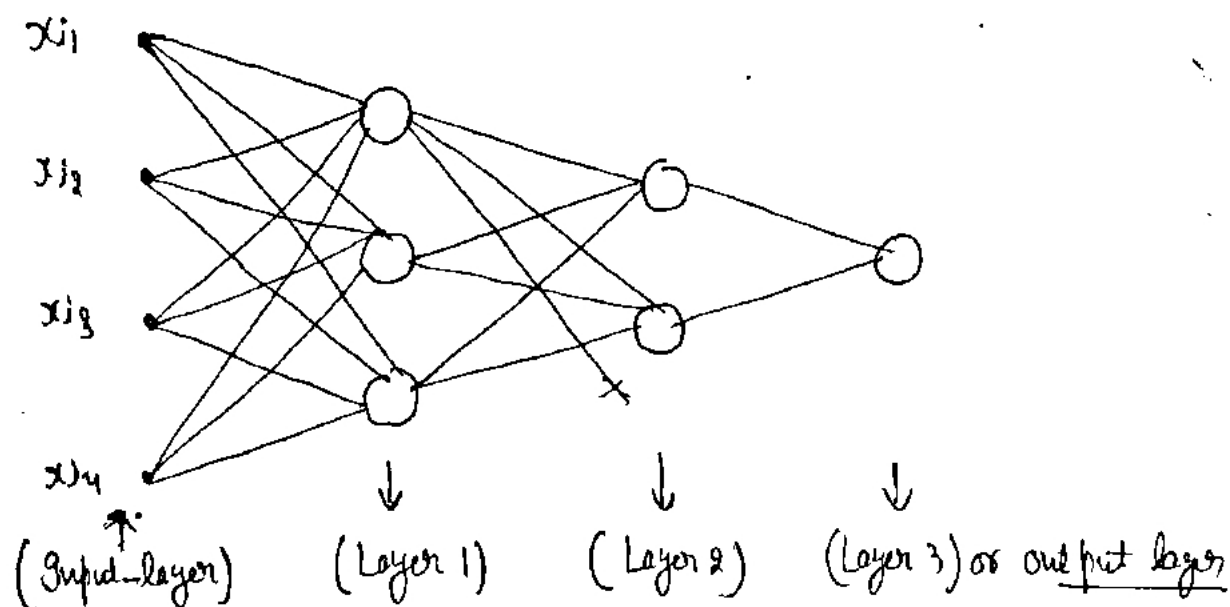→ "Notation" of Neural Network (Multilayer Perceptron) :→

We have seen that Multilayered structure is an extremely powerful idea to build extraordinarily powerful models, while you can avoid overfitting.

lets look at some notations, so that it will become easy for us to follow, what is happening.

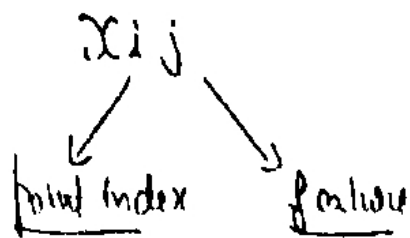let's understand the notation with the help of an example.



(Input layer)      (Layer 1)      (Layer 2)      (Layer 3) or output layer

let's assume our Dataset is $D = \{x_i, y_i\}$, lets take a regression problem because they are easy to understand. let's assume $x_i$ is a four dimensional point.

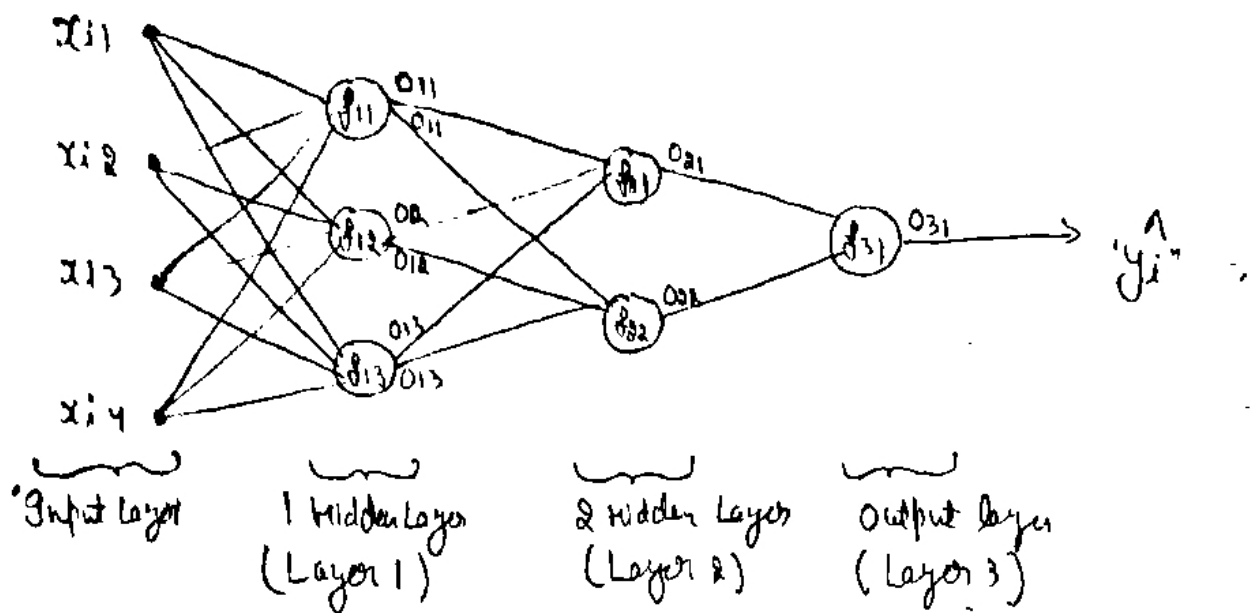$$x_i \in \mathbb{R}^4 \; ; \; \& \; y_i \in \mathbb{R}$$

$x_i$ is a 4-d point and $y_i$ is a 'real value'. So, what we are solving here is a regression problem.

we are representing each of these features as $x_{ij}$ where 'i' is the point index (which point it is) & 'j' is the corresponding feature.

$X_{ij}$

$\downarrow$ input index    $\downarrow$ feature

Let's now connect all the inputs in the input layer with all the neurons in the first Hidden layer or layer 1 & so-on as shown:



In layer "1" we have 3 neurons & In each neuron we have a function.

Each function in layer 1 is represented as '$f_{1j}$'

We have 3 neurons in layer 1, therefore corresponding to each neuron, we have a function

Layer 1 = { $f_{11}$   $f_{12}$   $f_{13}$   }

Layer1 $\swarrow$ function 1 --- $\downarrow$ Layer1 $\searrow$ function 3

Now let's assume we have 2 neurons in "layer 2" & all the outputs from "Layer 1" are connected to each neuron in "Layer 2" as shown above.

Let's now represent the functions in Layer 2

Layer 2 = { $f_{21}$ , $f_{22}$ }

Layer 2 $\swarrow$ fun 1   Layer 2 $\searrow$ fun 2

& let's assume we have 1 neuron in 'layer' 3 or the output layer, and all the outputs from 'layer2' are connected to 'layer3' as we have seen.

The function in layer 3, since we have a single neuron, is represented as '$f_{31}$'.

$$\text{Layer 3} = \{ f_{31} \}$$

'$\hat{y}_i$' is the output that we get, when we input '$x_i$'

The ideal output we want is '$y_i$', but the model will output '$\hat{y}_i$'.    So, if '$y_i$' is very similar to '$\hat{y}_i$', we have done an extremely good job.

let's now give notation for the functions used in each of the neurons

In generalised terms, the function is represented by

$$\boxed{f_{ij}}$$ , $\left\{ \begin{array}{l} \text{'i' represents the layer number} \\ \& \\ \text{'j' represents the function index} \end{array} \right\}$

When we give inputs to a neuron (to activation function) we will get an output, (corresponding output).

& output notation is same as the function notation.

"Output of $f_{11}$" $\rightarrow$ $\begin{bmatrix} \text{'}O_{11}\text{'} \\ \text{'}O_{12}\text{'} \\ \text{'}O_{13}\text{'} \end{bmatrix}$    "Output of $f_{21}$" $\rightarrow$ $\begin{bmatrix} \text{'}O_{21}\text{'} \\ \text{'}O_{22}\text{'} \end{bmatrix}$

"Output of $f_{12}$" $\rightarrow$

"Output of $f_{13}$" $\rightarrow$    "Output of $f_{22}$" $\rightarrow$

"output of $f_{31}$" $\rightarrow$ $\begin{bmatrix} \text{'}O_{31}\text{'} \end{bmatrix}$

Note :- output of neuron is having same numbering as neuron itself

In generalised terms, output is represented as :→

$O_{ij}$ { output of ith layer & jth neuron in that layer }

When we look at our perceptron, we see that there are weights associated with the edges & we are doing this to differentiate 'one input' from 'another'.

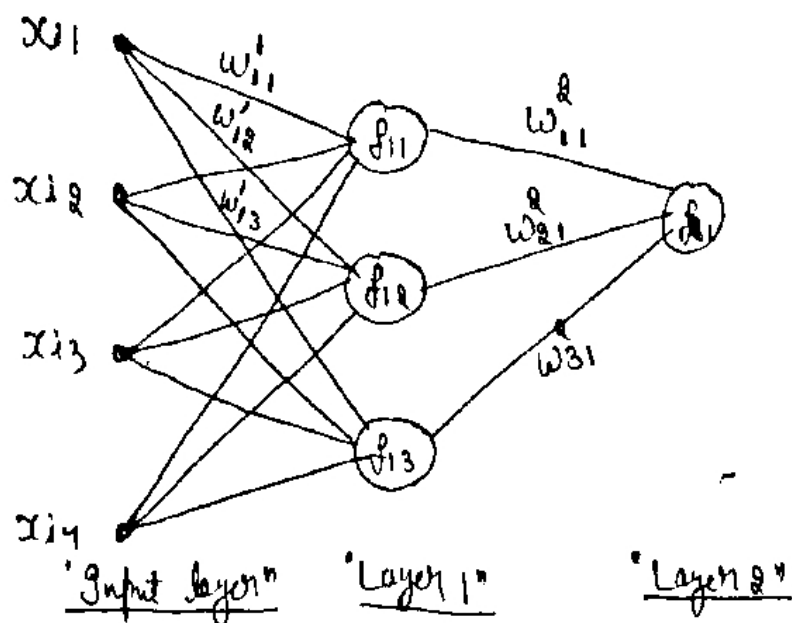'Weights' on the edges are represented in a generalised way as :→

$$W_{ij}^{k}$$

'k' → what is the next layer.
　　　　　↳ next layer is put as a 'superscript'
'i' → is 'from'
'j' → is 'to'



$x_{i1}$  $x_{i2}$  $x_{i3}$  $x_{i4}$

$w_{11}'$  $w_{12}'$  $w_{13}'$  $f_{11}$  $f_{12}$  $f_{13}$  $w_{11}^{2}$  $w_{21}^{2}$  $f_{1}$  $w_{31}$
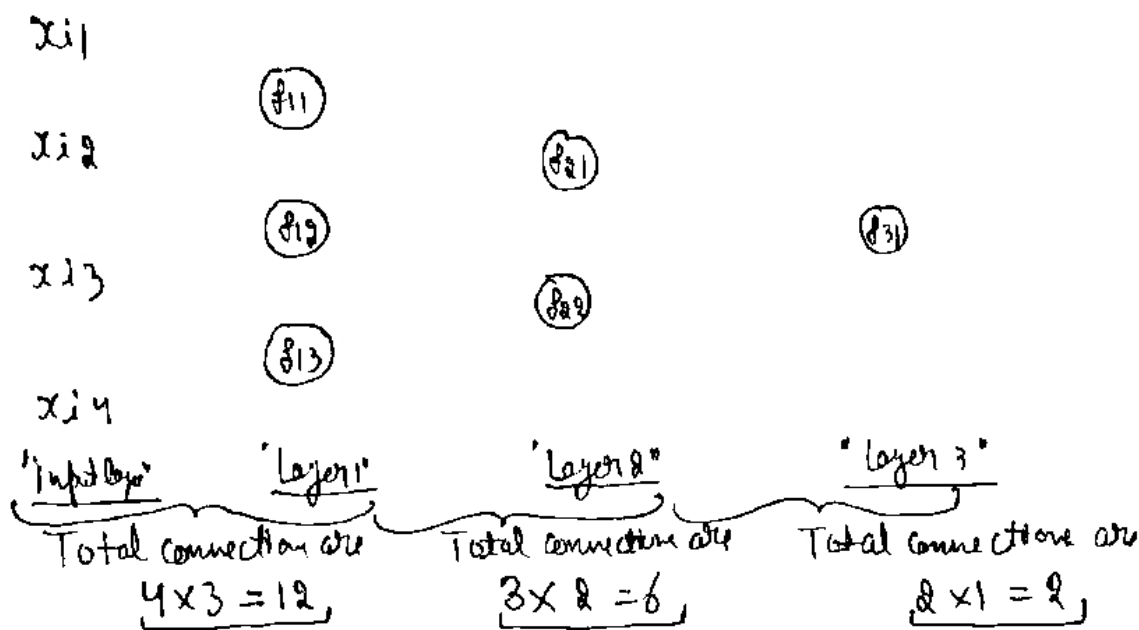
& so on

"Input layer"　"Layer 1"　"Layer 2"

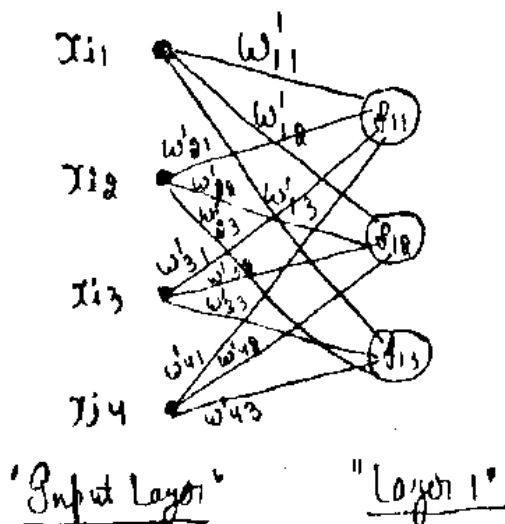This is how we number each of the ['edge-values'], ['neuron function'] & the outputs from ['each of the neurons']

Since in the 'network' we have connected all the 'possible pairs' hence this network is called a ['fully connected neural network'].

Note :- Connections go from 'one layer' to the 'next layer' & not within the layer.

$x_{i1}$

$(f_{11})$

$x_{i2}$

$(f_{21})$

$(f_{12})$

$(f_{31})$

$x_{i3}$

$(f_{22})$

$(f_{13})$

$x_{i4}$

'Input layer'    'Layer 1'         'Layer 2'          'Layer 3'

Total connection are    Total connection are    Total connection are

$4 \times 3 = 12$,       $3 \times 2 = 6$,        $2 \times 1 = 2$,

All the possible connections we are making, hence it is called a 'fully connected neural network'.

lets consider the 'first layer' & the 'input layer' only



$x_{i1}$   $W_{11}^1$

$W_{12}^1$   $(f_{11})$

$W_{21}^1$
$W_{22}^1$   $W_{13}^1$

$x_{i2}$

$W_{31}^1$   $(f_{12})$

$W_{32}^1$
$W_{33}^1$

$x_{i3}$

$W_{41}^1$   $W_{42}^1$   $(f_{13})$

$W_{43}$

$x_{i4}$

'Input Layer'         'Layer 1'

Here we have in total 12 weights b/w ['layer 1'] & ['input layer']

beez we have 'four' inputs in the input layer & we have '3' neurons in 'layer 1'

$4 \times 3 = 12$

PTO

The weights ( 12 weights) between 'Input layer' and 'Layer-1' are represented in matrix form as :-

$$W^1 = \begin{bmatrix} W^1_{11} & W^1_{12} & W^1_{13} & \boxed{W^1_{14}} \\ W^1_{21} & W^1_{22} & W^1_{23} \\ W^1_{31} & W^1_{32} & W^1_{33} \\ W^1_{41} & W^1_{42} & W^1_{43} \end{bmatrix}_{4 \times 3}$$

This matrix of weights can be represented using $4 \times 3$ matrix.

We can represent all of these '12 weights' using a matrix of size '$4 \times 3$'.

Similarly weights on the edges b/w 'Layer1' & 'Layer2' can be represented by a matrix $W^2$ of size '$3 \times 2$'

$$W^2 = \begin{bmatrix} W^2_{11} & W^2_{12} \\ W^2_{21} & W^2_{22} \\ W^2_{31} & W^2_{32} \end{bmatrix}_{3 \times 2} \qquad W^3 = \begin{bmatrix} W^3_{11} \\ W^3_{21} \end{bmatrix}_{2 \times 1}$$

So, all these connections can be represented using matrices.

P.T.O

→ "Training a single Neuron Model" :->

We first understand, how to train a single neuron model, before we go to training an MLP (Multi layer Perceptron).

'Training' means to find the ["best weights"] best edge weights in a neural network model. why traing data (Q Train).

Single neuron model is 'Perceptron'. & 'logistic Regression'

"Perceptron & Logistic Regression" → ["Single Neuron Models"] for classification

Similarly if we think of "Linear-Regression". It can also be represented as a [single neuron model] for 'regression'
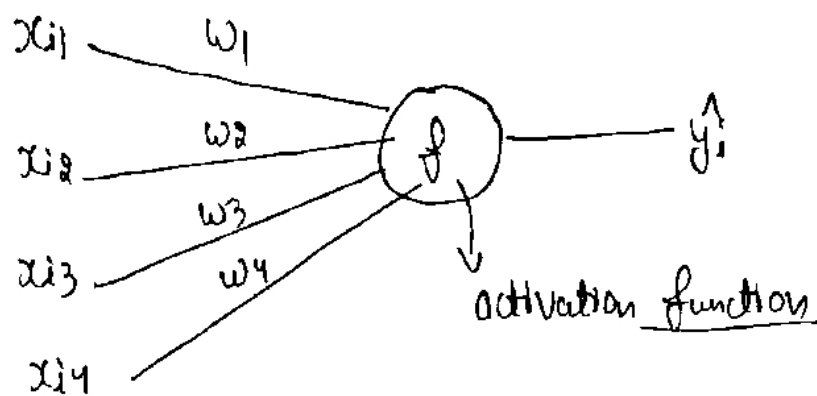
"Linear-Regression" ⟶ ["Single neuron model for Regression"]

Now we will use regression as an example instead of classification becoz, the mathematics & calculus become easy. So whatever we learn for 'linear-regression' can easily be extended to 'logistic-regression', 'perceptron' etc.

We have already seen how a "linear-regression" model can be represented as an "optimization problem".

So, we revisit all these concepts of linear-regression as they are very important in understanding the training of a neural network.

Let's now understand 'Linear Regression problem" from the perspective of a 'neuron'. (Single neuron model)



activation function

In a linear regression, the problem it is

$$\hat{y_i} = \sum_{j=1}^{d} w_j x_{ij} \quad \bigg| \quad D_{Tra} = \{x_i, y_i\} \quad x_i \in \mathbb{R}^d$$
$$y_i \in \mathbb{R}$$

In linear regression, we have.

$$\hat{y_i} = \bar{w}^T x_i$$

(This is the relationship b/w 'x_i' & $\hat{y}$)

So, given 'x_i' & 'y_i' our job is to find the optimal 'w'.

Now let's go back to a single neuron framework about

Here input to "f" is $\sum_{j=1}^{y} w_j x_{ij}$

becoz $x_{i1}$ gets multiplied with $w_1$

$x_{i2}$ " " with $w_2$

$x_{i3}$ " " with $w_3$

& $x_{i4}$ " " with $w_4$

& then we sum up all of them. as $\sum_{j=1}^{y} w_j x_{ij}$

∵ the input that we get is '$w^T x_i$'

In a 'linear regression' the output is also '$w^T x_i$' which means that function ['$f$'] is called an [identity function] in case of a linear regression

An identity function is, if we input "$z$" to a function & same "$z$" is returned as an output then it is an 'identity function'

$$[ f(z) = z ]$$

So, this is in the case of 'Linear regression'. becoz in linear regression the input that we give to the activation function is $w^T x_i$ & the output is also $w^T x_i$, so if both input & output of a function is same, then the function is known as an 'identity function'.

Note :- 
- In case of a logistic regression, '$f$' is a sigmoid function
- In case of a simple perceptron, '$f$' is a thresholding function.
- In case of a linear regression, '$f$' is an identity function

Here we always represent the function as '$f$' because whatever we derive, it should be generalized.

Second, thing, In linear regression, when we represented it as an optimization problem.

So, optimization problem of "Linear Regression" can be written as is

—find the $w's$ that minimizes the loss function $\phi$ in linear regression we have used square loss

$$\underset{w's}{\min} \sum_{i=1}^{n} (y_i - \hat{y_i})^2 + \text{regularizer}$$

$n \rightarrow$ no of points in training data

where $\hat{y_i} = w^T x_i$
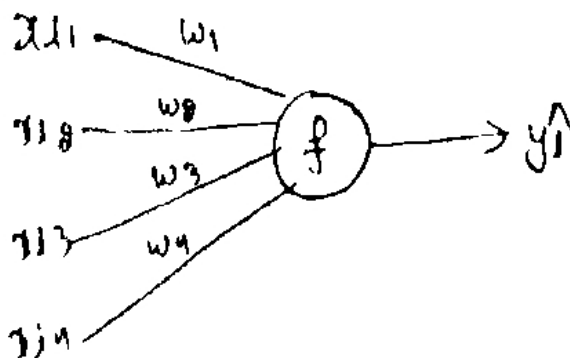
This can can be written as :

$$\left\{ \underset{w's}{\min} \sum_{i=1}^{n} (y_i - w^T x_i)^2 + \text{regularizer} \right\}$$

$\downarrow$ squared error

It could be any $L_1$ or $L_2$

This is the optimization problem for "Linear Regression"

Even in a single neuron model, this is the optimization problem that we solve.

Let's try to connect this optimization problem with a single neuron model :→



$D = \{x_i, y_i\}_{i=1}^{n}$

$\downarrow$ training data

where $x_i \in \mathbb{R}^y$

& $y_i \in \mathbb{R}$

This means we are solving a Regression problem

Let's go step by step

Step 1 :- Define a 'loss-function'

$$\mathcal{L} = \sum_{i=1}^{n} (y_i - \hat{y_i})^2 + reg.$$

When we input 'xi' through this network, we are going to get a 'ŷi' which is different from 'yi', we want to penalize it. The farther away "ŷi" is from 'yi', we want to penalize it. In our loss function.

Let us also define a term called 'Lᵢ'

$$\mathcal{L}_i = (y_i - \hat{y_i})^2$$

Note :- squared loss of one point, we are calling it as 'Lᵢ'.

& the sum of all the squared losses, is called 'L'

$$\boxed{\mathcal{L} = \sum_{i=1}^{n} \mathcal{L}_i}$$

'L' is loss across all the training data & 'Lᵢ' is the loss across each training point

Also $\boxed{\hat{y_i} = \bar{w}^T x_i}$

We have our ŷi, Now on top of it, we apply the loss function, we are computing using [yi] & [ŷi]. we get ŷi from the model & 'yi' is from the actual data & all we care about is $\mathcal{L}(y_i, \hat{y_i})$.

we have our $x_i$'s & our goal is to find the optimal $w$'s (weights). This is what it namely means.

Step no 2 :⇒ Write the optimization problem.

$$\min_{w_i} \sum_{i=1}^{n} \left(y_i - \underbrace{w^T x_i}_{\hat{y_i}}\right)^2 + reg.$$

So, from neural network perspective

$$\left[\hat{y_i} = f(w^T x_i)\right]$$

In case of linear regression the '$f$' is $\left[\begin{array}{c}\text{identity}\\\text{function}\end{array}\right]$

$$\boxed{\min_{w_i} \sum_{i=1}^{n} \left(y_i - f(w^T x_i)\right)^2 + reg.}$$

$\left[\begin{array}{l}\text{Now if the '}f\text{' is 'I' we get linear regression}\\\text{Is 'Threshold' we get perceptron}\\\text{& if it is sigmoid, we get}\\\qquad\qquad\qquad\text{logistic regression}\end{array}\right]$

This is a generalized optimization problem, we want to solve.

In vector notation, it can be represented as, find a vector '$w^*$' which minimizes the thing.

$$\left[\, \omega^* \; = \; \arg\min_{\omega} \; \sum_{i=1}^{n} \left(y_i - f(\omega^T x_i)\right)^2 + reg. \,\right] \tag{13}$$

Stepno 3 :→ Solve the optimization problem.

(for that we use gradient descent)

ⓐ Initialization of weights. wi's ⟶ random initialization

ⓑ Compute derivative of $L$ wrt 'ω

$$\nabla_{\omega} L \qquad \text{(This is a vector derivative)}$$

$$\nabla_{\omega} L = \begin{bmatrix} \dfrac{\partial L}{\partial \omega_1} \\ \dfrac{\partial L}{\partial \omega_2} \\ \vdots \\ \dfrac{\partial L}{\partial \omega_4} \end{bmatrix} \qquad \because \{\, \omega \in \mathbb{R}^4 \; \text{since} \; x_i \in \mathbb{R}^4 \,\}$$

This is vector representation of 'partial derivative'.

ⓒ Once we compute this, we will iteratively say that.

$$\{\, \omega_{new} = \omega_{old} - \eta * [\nabla_{\omega} L]_{\omega_{old}} \,\} \qquad \rightarrow \text{Vector notation}$$

↳ 'Learning rate'

↓

Now let us write scalar notation of this

$$(\omega_1)_{new} = (\omega_1)_{old} - \eta \left[\dfrac{\partial L}{\partial \omega_i}\right]_{(\omega_i)_{old}}$$

we are computing derivative at a point & derivative is nothing but a _slope_.

we will keep repeatly new & old, we will be having a for loop as

for iterations = 1 to K

$$(w_i)_{new} = (w_i)_{old} - \eta \left[ \frac{\partial L}{\partial w_i} \right]_{(w_i)_{old}}$$

we will run it, till we _converge_.

The major difference between Gradient descent & Stochastic Gradient descent is :→

In both of them, we will compute the partial derivative w.r.t 'w'

In 'GD', $\nabla_w L$ is computed for all the points in our dataset

$$\nabla_w L \longrightarrow x_i's \ \& \ y_i's$$

whereas in case of 'SGD', we will approximate the gradient by taking one point $\{x_i, y_i\}$ or a small set of points (batch of points)
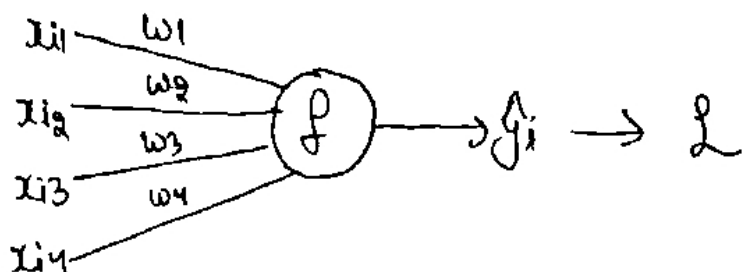
$$\nabla_w L \approx \text{one point } \{x_i, y_i\}$$
or small batch of points

Now let us see how to compute the _derivative_.

Now let us see how to compute the derivative.

$$\S \nabla_w L \S$$

let's draw the perceptron again



$$
\begin{bmatrix}
\text{we are getting } \hat{y_i} \text{ as an output of "} f \text{" \&} \\
\text{on top of that we are applying Loss functi-} \\
\text{on.}
\end{bmatrix}
$$

& In case of regression, this "L" is a 'squared loss'

Now we want to compute $\nabla_w L = \left[ \dfrac{\partial L}{\partial w_1}, \dfrac{\partial L}{\partial w_2}, \dfrac{\partial L}{\partial w_3}, \dfrac{\partial L}{\partial w_4} \right]^T$

This is supposed to be a column vector, so we are taking transpose of it to convert it into a row vector.

Now let us see how to compute each of them

becoz, 'w1', 'w2', 'w3' & 'w4' are independent of each other.

let's compute $\dfrac{\partial L}{\partial w_1}$, if we look at the path from 'L' to $w_1$, we have a function 'f' in b/w & 'a' $w_1$, $\hat{y_i}$ the constant & it is computed from $x_{i1}$

So, here we want to compute derivative of 'L' w.r.t 'w₁'

Let's use the chain rule of differentiation

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f} * \frac{\partial f}{\partial w_1} \longrightarrow \text{chain rule of differentiation}$$

$$\therefore \left\{ \frac{\partial L}{\partial x} = \frac{\partial L}{\partial f_1} * \frac{\partial f_1}{\partial f_2} * \frac{\partial f_2}{\partial x}. \right\} \longrightarrow \text{chain rule of diff.}$$

We are taping the path from 'w₁' to 'L' & whatever variable or function exists in the path, we have to tape the derivative of that w.r.t to 'w₁'

Note ⇒ $x_i's$ & $\hat{y}'s$ are constants

$$L = \sum_{j=1}^{n} (y_1 - \hat{y_i})^2 + reg \qquad \cdot \left( \text{let's ignore regu for a while.} \right)$$

$$+ \hat{y_i} = f(w^T x_i)$$

$$L = \sum_{j=1}^{n} (y_i - f(w^T x_i))^2 \qquad \text{let's write } f(w^T x_i)$$
as f for
simplicity.

$$\frac{\partial L}{\partial f} = \sum_{j=1}^{n} (y_i - f)^2$$

$$\frac{\partial L}{\partial f} = \sum_{j=1}^{n} -2(y_i - f(w^T x_i))$$

Now we need to find $\frac{\partial f}{\partial w_1}$

we prove $f(w^T x_i) = w^T x_i$

Now $\dfrac{\partial f}{\partial w_1} = x_i$   simple

$L$ is basically summation across all the points

let's take $L_i$ for

$$L_i = (y_i - \hat{y}_i)^2 \quad \text{for a single point}$$

$$\frac{\partial L_i}{\partial w_1} = \frac{\partial L_i}{\partial f} * \frac{\partial f}{\partial w_1}$$

$$\Rightarrow -2(y_i - \hat{y}_i) * x_i$$

$$\therefore \frac{\partial L_i}{\partial w_1} = -2x_i (y_i - \hat{y}_i)$$

Now let's find $\dfrac{\partial L}{\partial w_1}$

$$\boxed{\frac{\partial L}{\partial w_1} = \sum_{i=1}^{n} -2x_i (y_i - \hat{y}_i)}$$

Similarly we can compute $\dfrac{\partial L}{\partial w_2}$ ---- $\dfrac{\partial L}{\partial w_4}$

when we have all of these computed, we can easily

compute $\nabla_w L$

once we have $\nabla_w L$, computing

$w_{new}$ is straight forward

as:

$$w_{new} = w_{old} - \eta \left[ \nabla_w L \right]_{w_{old}}$$

↓ learning rate