

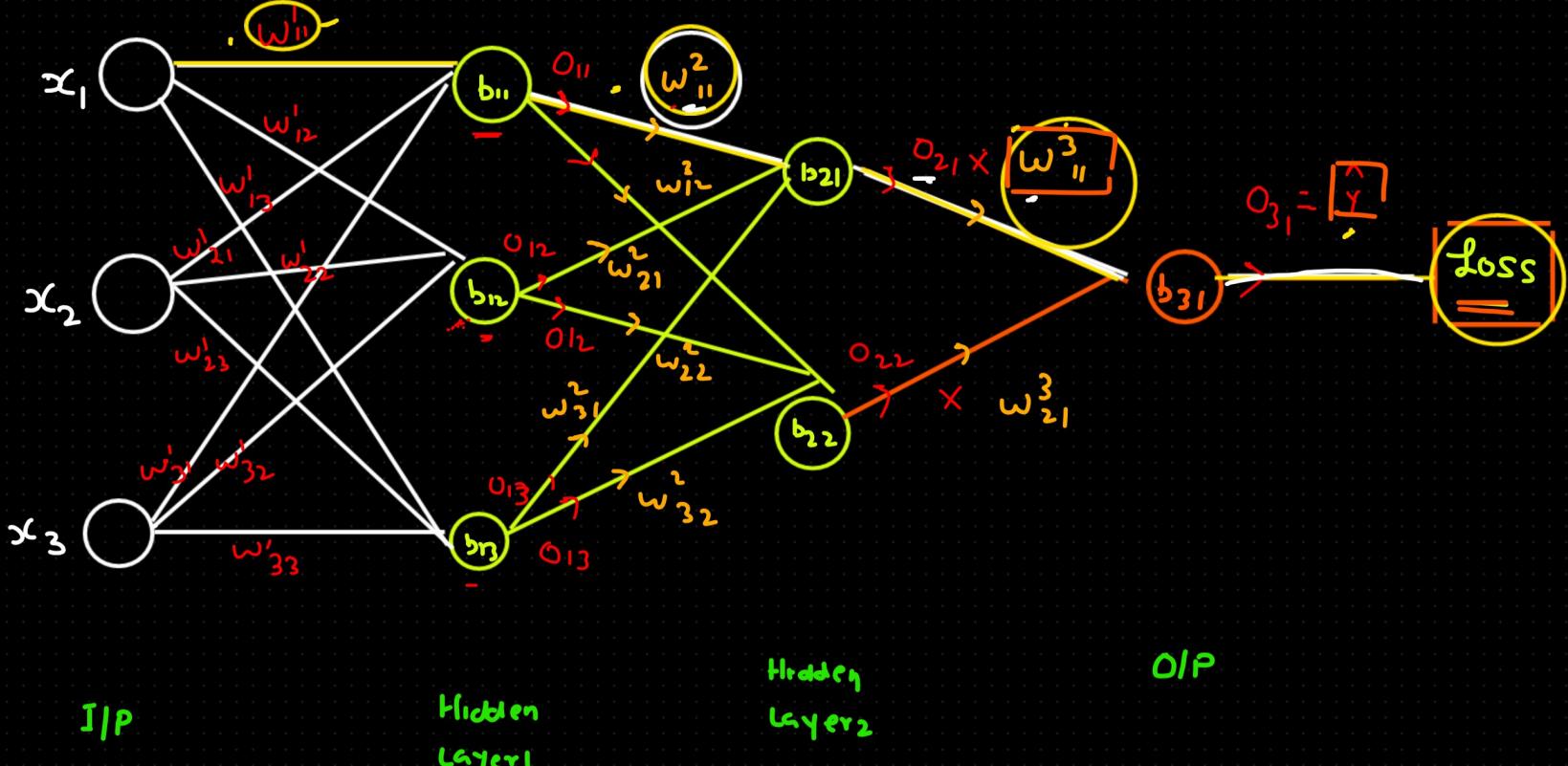
Please note that:

Technical name of neural network is called Logical Architecture of Neural Network. The same is shown below:

--> Purpose of Back Propagation is to update the trainable parameters i.e. weights & bias

1

Starting with back propagation quick revision



$$\text{Trainable Parameter } 3 \times 3 + 3 = 12$$

$$(3 \times 2) + 2 = 8$$

$$(2 \times 1) + 1 = 3 = 23$$

Now let's revise the BP \Rightarrow to update the trainable Parameter (Weight + Biases)

Optimizer = (Gradient Descent) =

$$W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial w}$$

Loss is dependent on the predicted Output \hat{y} that is represented by the Output O_{31} which in turn is dependent on some trainable parameters like: $(W_{311}), (W_{321})$ and b_{31}

$$\frac{\partial L}{\partial w} \Rightarrow \frac{\partial L}{\partial w_{11}^3}$$

$$\begin{array}{c} \text{Loss} \rightarrow \hat{y} \rightarrow \\ \hat{y} = O_{31} \\ \hline \end{array}$$

Here just for easier demonstration only 1 trainable parameter is taken

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_{11}^3} \Rightarrow ?$$

Here for instance we can clearly observe that the chain dependency is formed which mathematically is represented using chain rule differentiation

$$w_{11}^3 \rightarrow O_{21} \rightarrow w_{11}^2$$

(Chain rule of differentiation)

Just visualize that how these dependencies are being calculated

$$\text{Loss} \rightarrow \hat{y}(O_{31}) \rightarrow$$

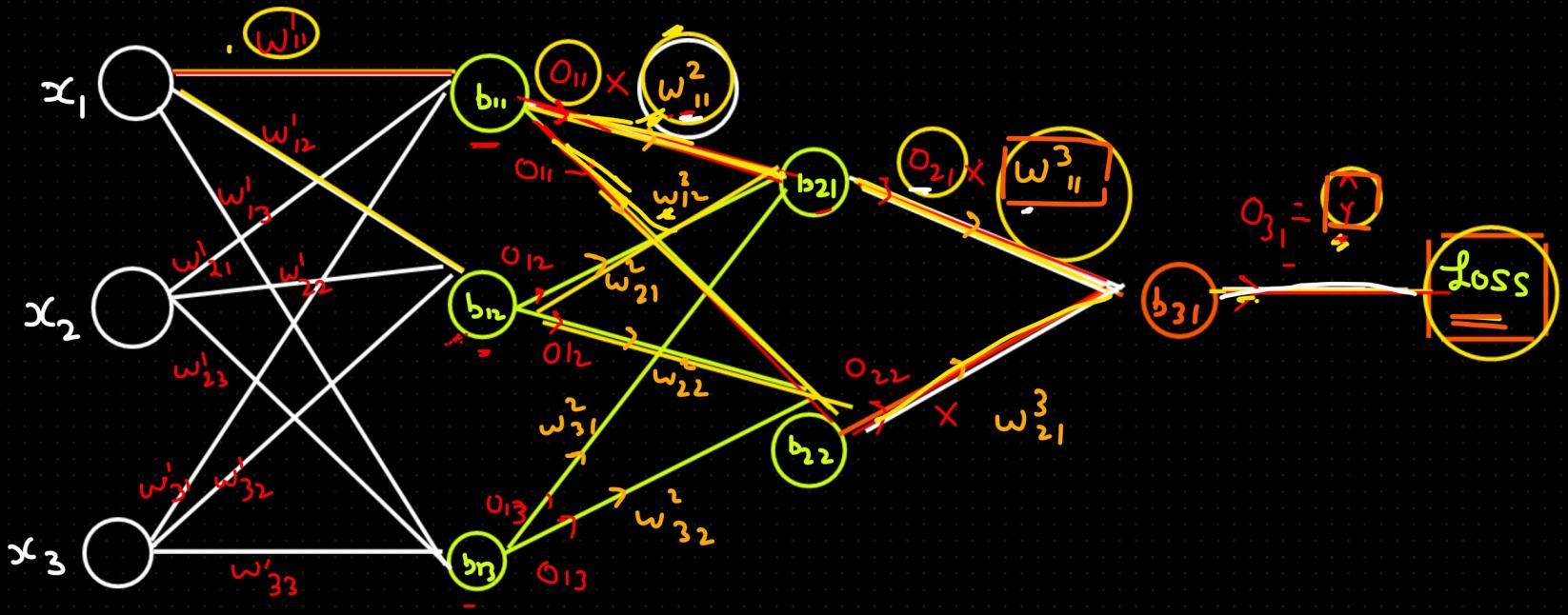
$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial w_{11}^2}$$

$$\text{Loss} \rightarrow \hat{y} \rightarrow w_{11}^3 \rightarrow O_{21} \rightarrow w_{11}^2 \rightarrow O_{11} \rightarrow w_{11}^1$$

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial w_{11}^3} \times \frac{\partial w_{11}^3}{\partial O_{11}} \times \frac{\partial O_{11}}{\partial w_{11}^1}$$

W_{new} and W_{old} are weights. Row (η) represents the learning parameter. Where L represents the Loss which is dependent on weight trainable parameter. That's why partial derivation of L is done wrt W .

In the same way we will construct optimizer or gradient descent formula for bias trainable parameter. In the above formula just replace W with b (bias)

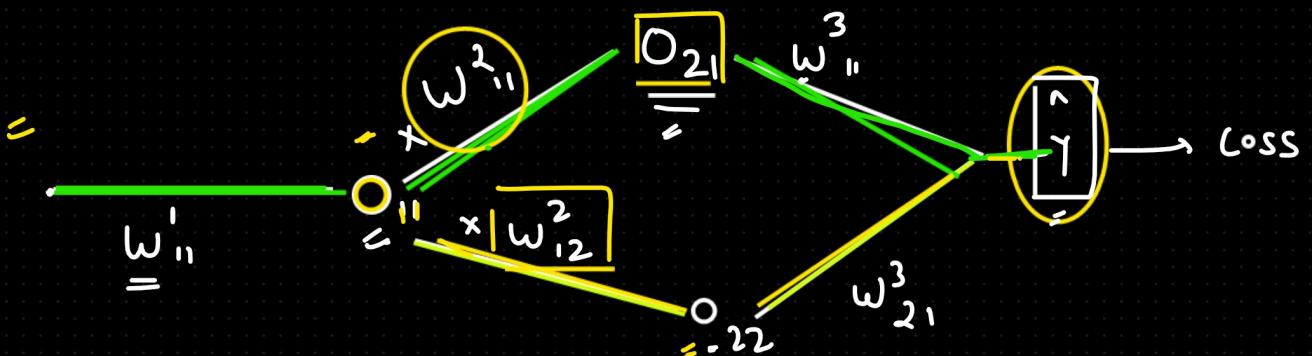


I/P

Hidden
Layer1

Hidden
Layer2

O/P



$\left[\frac{\partial L}{\partial w_{ii}} \right] =$ will take both Path
and for final out will add on
both Path

Above for easier demonstration purpose we have only considered a Single path. In real time all possible path along with all the trainable parameters (weights and bias) will be taken into consideration.

$$= \frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial \vec{q}} \left[\frac{\partial o_{31}}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{11}} + \frac{\partial o_{31}}{\partial o_{22}} \times \frac{\partial o_{22}}{\partial o_{11}} \times \right. \\ \left. \frac{\partial o_{11}}{\partial w_{11}} \right]$$

{
 Loss fn
 activation fn
 optimizer

Loss fn \Rightarrow Regression \Rightarrow MSE, MAE, huberloss (MSE and MAE already covered in the last lecture)
Classification \Rightarrow Binary, Categorical, Sparse cross entropy

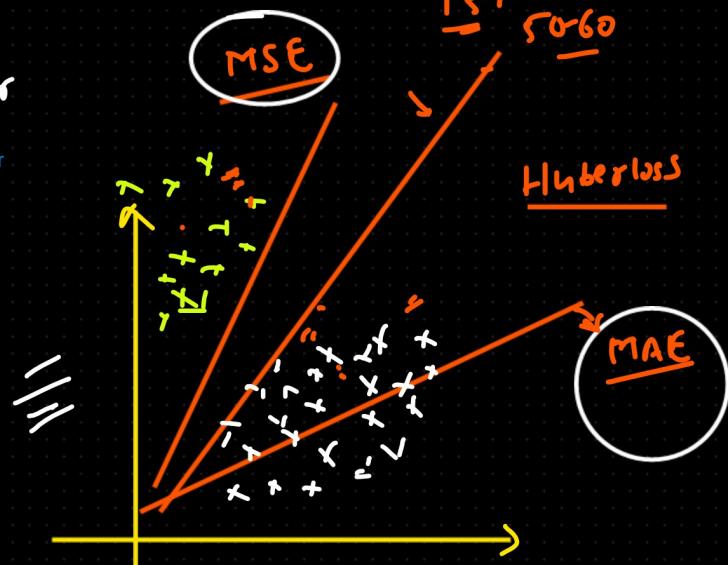
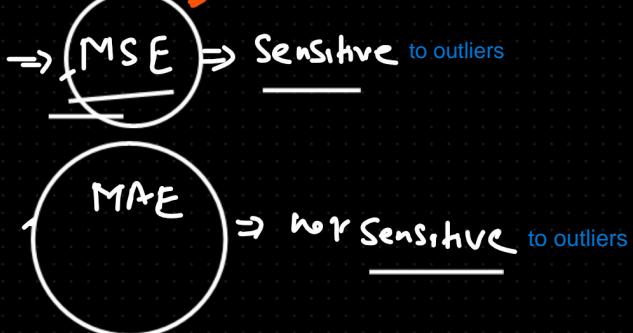
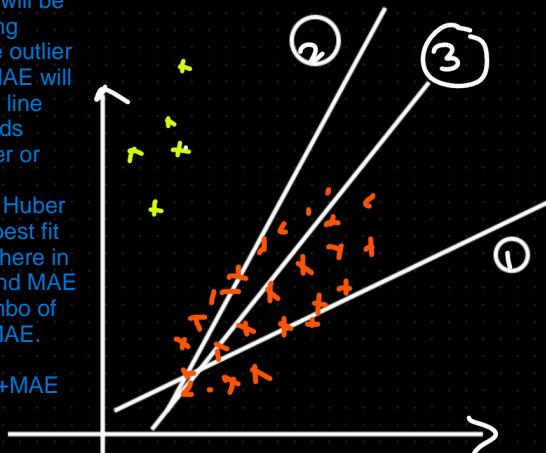
MSE+MAE

Huber loss

$$\left\{ \begin{array}{l} \frac{1}{2} (y - \hat{y})^2 \leftarrow \text{thru Parl' in presence of outliers} \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 \leftarrow \text{if D.P. is outlier} \end{array} \right.$$

See below snippet to understand the formula better

In case of MSE best fit line will be more tending towards the outlier whereas, MAE will give best fit line more towards more denser or non outlier datapoints. Huber loss gives best fit line somewhere in b/w MSE and MAE. Kind of combo of MSE and MAE. Huber Loss=MSE+MAE



defines the loss function piecewise

$$L_\delta(a) = \begin{cases} \frac{1}{2} a^2 & \text{for } |a| < \delta, \\ \delta \cdot (|a| - \frac{1}{2} \delta), & \text{otherwise.} \end{cases}$$

This function is quadratic for small values of a , and linear for large values, with equal values and slopes of the different sections at the two points where $|a| = \delta$. The variable a often refers to the residuals, that is to the difference between the observed and predicted values $a = y - f(x)$, so the former can be expanded to [2]

$$L_\delta(y, f(x)) = \begin{cases} \frac{1}{2} (y - f(x))^2 & \text{for } |y - f(x)| \leq \delta, \\ \delta \cdot (|y - f(x)| - \frac{1}{2} \delta), & \text{otherwise.} \end{cases}$$

The Huber loss is the convolution of the absolute value function with the rectangular function, scaled and translated. Thus it "smoothes out" the former's corner at the origin.

alpha δ represents the loss that is difference b/w actual value and predicted value

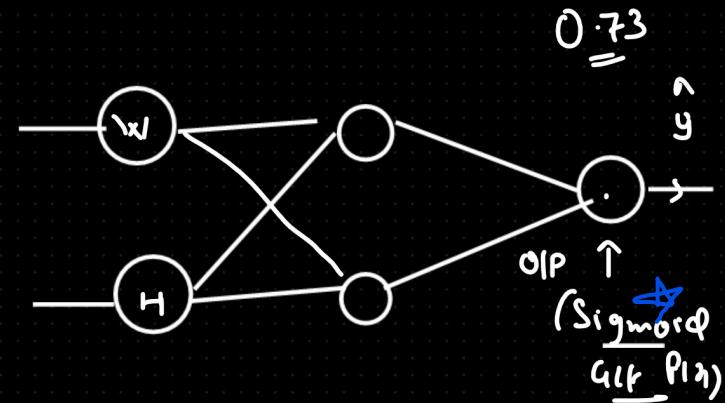
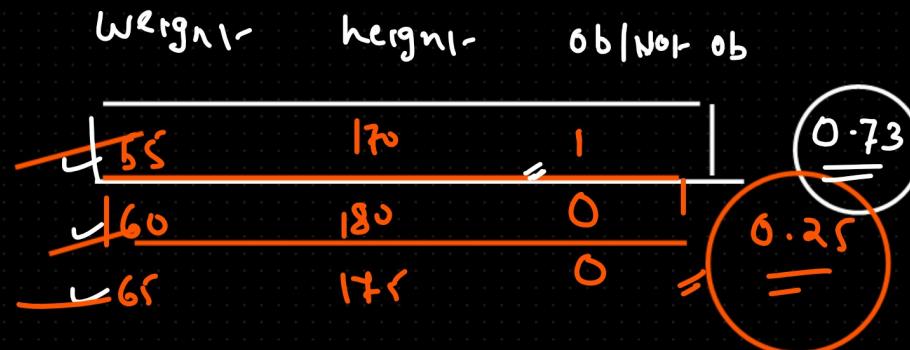
Dell δ Hyper-parameter is decided based on outlier

Loss fn for the classification

Please don't get confuse b/w loss and activation funⁿ. Sigmoid is the activation function that may be used in the case of Binary classification. Whereas, loss function used for binary classification is called Binary cross entropy

1 Binary Cross entropy \Rightarrow 2 classes log loss \Rightarrow logistic regression

Loss function for Binary class classification



* Single value $\underline{\text{Loss}} = \boxed{-y \log(\hat{y}) - (1-y) \log(1-\hat{y})}$

* for the whole data $\underline{\text{Loss fn}} = \frac{1}{n} \sum_{i=1}^n [-y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)]$

Assignment
Perform the derivation (done)
of log loss

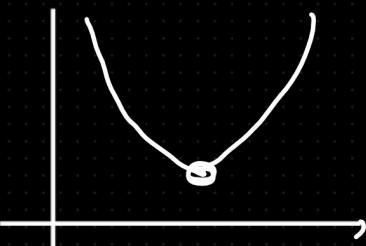
$$\begin{aligned} & -1 \log(0.73) - (1-1) \log(1-0.73) \\ & -1 \times \log(0.73) - 0 \times \log(0.27) \\ & \boxed{-\log(0.73)} = \boxed{-(0.13667714)} \end{aligned}$$

$$-\underline{0} \times \log(0.25) - \underline{(1-0)} \times \log(\underline{1-0.25}) = -\underline{1} \times \log(0.75) = -\log(0.75) = \underline{\underline{0.12}}$$

Advantage

we can use G.O. for this loss

since it is differentiable because of convex curve



Disadvantage

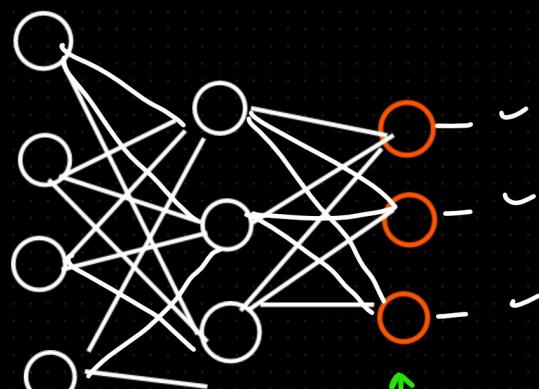
Categorical cross entropy \Rightarrow Multiclass classification \Rightarrow More than 2 class

Iris dataset
example is taken
Sepal length,
sepal width, petal
length and petal
width

S_L	S_W	P_L	P_W	Class
				Setosa
				Virginica
				Versicolor

One hot encoded value for
the target class

Since in multiclass Target will be classified into more than 2 classes so as a standard step we will require to apply One Hot Encoding on top of it



Soft Max \Rightarrow

$$\hat{y} = \text{act}(\text{Summation}) \xrightarrow{\text{Softmax}} \frac{e^{z_i}}{\sum e^{z_i}}$$

$$\Rightarrow \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\frac{e^{z_2}}{\dots} \quad \frac{e^{z_3}}{\dots}$$

Categorical
cross entropy

K represents the number of classes in which target will be classified

$$\text{Loss} = - \sum_{i=1}^K y_i \log(\hat{y}_i)$$

$$\text{Loss} = - \sum_{i=1}^3 y_i \log(\hat{y}_i)$$

Please note that here only loss function is being discussed. For calculating the Cost just take the average of the individual losses

$$- y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - y_3 \log(\hat{y}_3)$$

Since we have 3 classes that is why we are having 3 nodes in the Output layer



Applying one hot encoding to the Target column

	Verginica	Setsosa	Versicolor
1	1	0	0
2	0	1	0
3	0	0	1

this is assumption

$$\begin{bmatrix} 0.3 & 0.5 & 0.2 \end{bmatrix} = \hat{y}$$

$$\begin{bmatrix} 1 & 0 & 0 \end{bmatrix} = y$$

wrt first records we are getting $\hat{y} = [0.3 0.5 0.2]$ whereas actual $y = [1 0 0]$. Applying loss function formula on \hat{y}

$$\left[-y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - y_3 \log(\hat{y}_3) \right] = -1 \times \log(0.3) - 0 \times \log(0.5) - 0 \times \log(0.2)$$

These getting canceled out since respective label encoded values is 0

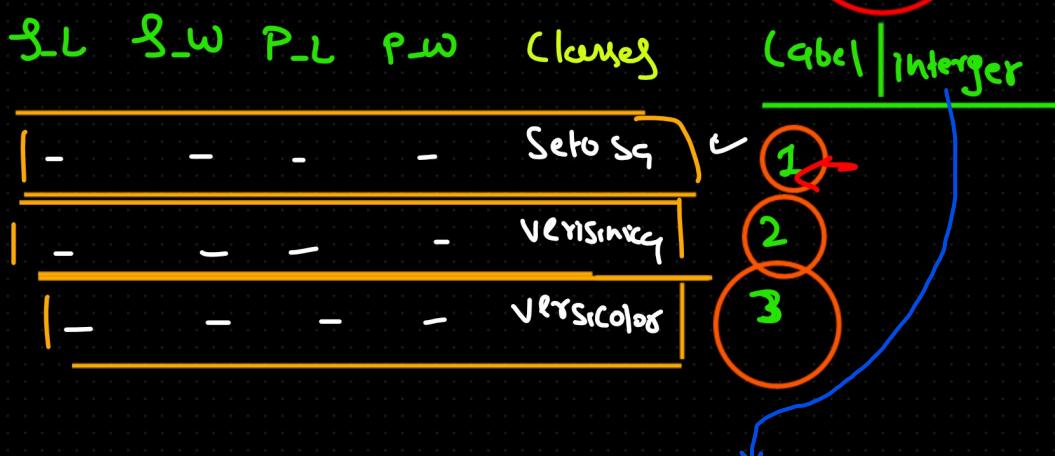
$$= -\log(0.3) = \underline{\underline{0.52}}$$

$$\Rightarrow -0 \times \log(0.3) - 1 \times \log(0.5) - 0 \times \log(0.2)$$

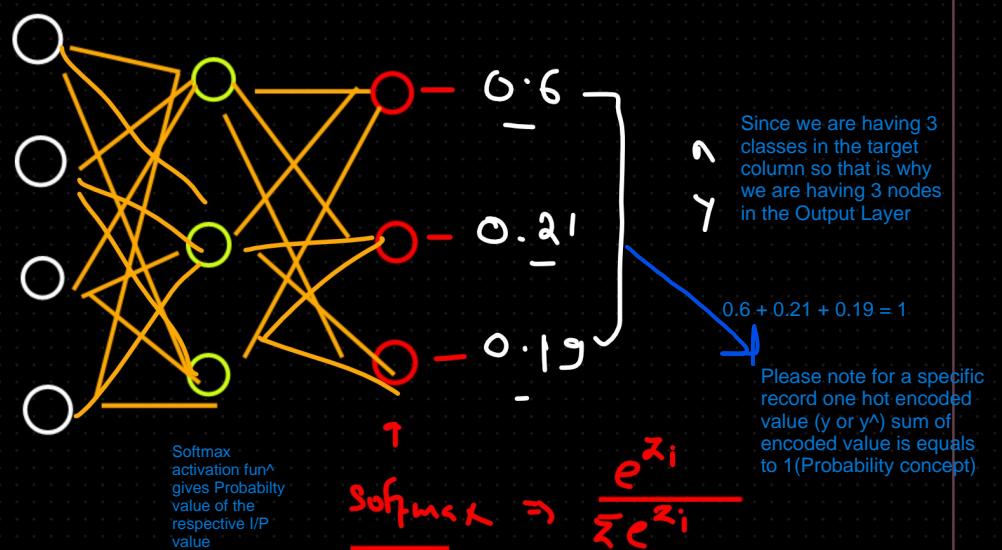
$$= -1 \log(0.5) \Rightarrow -\log(0.5) = \underline{\underline{0.30}}$$

Sparse categorical cross entropy \rightarrow Fast

In above we can clearly observe that in the case of multi class classification there is a lot of sparsity (0s and 1s) which leads to overfitting. To counter that we have Sparse categorical cross entropy



In Categorical cross entropy target is categorically label encoded based on one hot label encoding whereas in Sparse Categorical cross entropy target is Integral encoded



1st row

$$\left[\begin{array}{ccc} 1 & 0 & 0 \\ 0.6 & 0.21 & 0.19 \end{array} \right] \rightarrow \frac{\text{Cat cross entropy}}{?}$$

When classes is less then we use Categorical cross entropy target whereas, when classes are more then we use Sparse Categorical cross entropy

$$\left[\begin{array}{ccc} 1 & 2 & 3 \\ 0.6 & 0.21 & 0.19 \end{array} \right]$$

Actual Vs Predicted

$$-\cancel{y_1 \log(\hat{y}_1)} - \cancel{y_2 \log(\hat{y}_2)} - \cancel{y_3 \log(\hat{y}_3)}$$

$$\Rightarrow -1 \times \log(0.6) \Rightarrow -\log(0.6)$$

For example here we are calculating the loss for the first record so we are directly neglecting the influence of other unnecessary encoded labels. Which earlier in case of Multiclass we neglecting by multiplying with Zero.

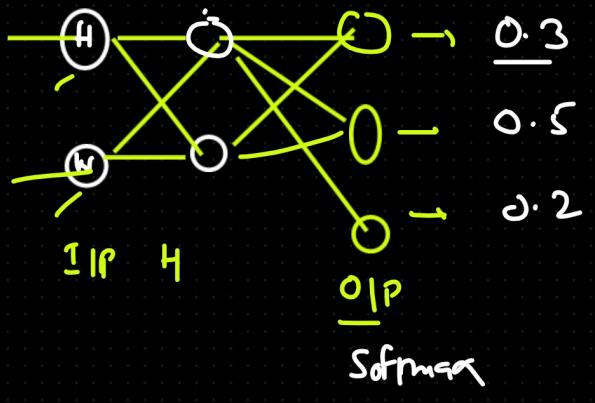
Please note that here,
Labeling does not mean that we will multiply with 2 in case calculation is for 2nd label. We will still multiply by 1 representing actual value of respective label encoded value.
That means that in Sparse categorical cross entropy all the integral encoded target values are internally mapped to 1 when calculating the loss(this consideration is only for the actual value)

$$\left[\begin{array}{ccc} 1 & 2 & 3 \\ 0.6 & 0.21 & 0.19 \end{array} \right]$$

Both are same just that in CCE(Categorical Cross Entropy) we are multiplying the not needed labels with 0 which may seems computational expensive whereas, in SCCE (Sparse Categorical Cross Entropy) we are directly neglecting the effect of not needed labels

In SCCE not needed labels are directly(computational inexpensive) getting rejected whereas in CCE not needed labels are rejected by multiplying with 0(computational expensive)

	Height	Weight	Body size	
→	170	66	M	1
→	175	60	S	2
→	180	75	L	3
→	160	70	M	4
→	155	50	L	5



One-hot encoding

CCA

CCA -

SCCA -

$$= - \sum_{i=1}^3 y_i \log(\hat{y}_i)$$

$$= - y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - y_3 \log(\hat{y}_3)$$

$$= [1 \ 0 \ 0]$$

$$[0.3 \ 0.5 \ 0.2]$$

$$= -1 \times \log(0.3) - \frac{0 \times \log(0.5)}{2} - \frac{0 \times \log(0.2)}{2}$$

$$= -\log(0.3) \Rightarrow 0.52$$

SCCA
labelling

$$= [1 \ 0.3 \ 0.5 \ 0.2]$$

SCCA vs CCA

$$-1 \times \log(0.5)$$

$$= -\log(0.5)$$

$$= 0.69$$

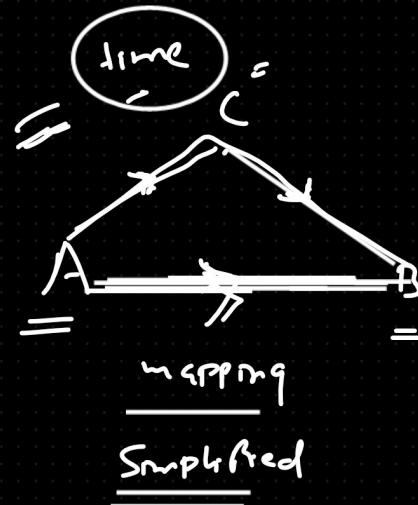
When the classes is low use CCA

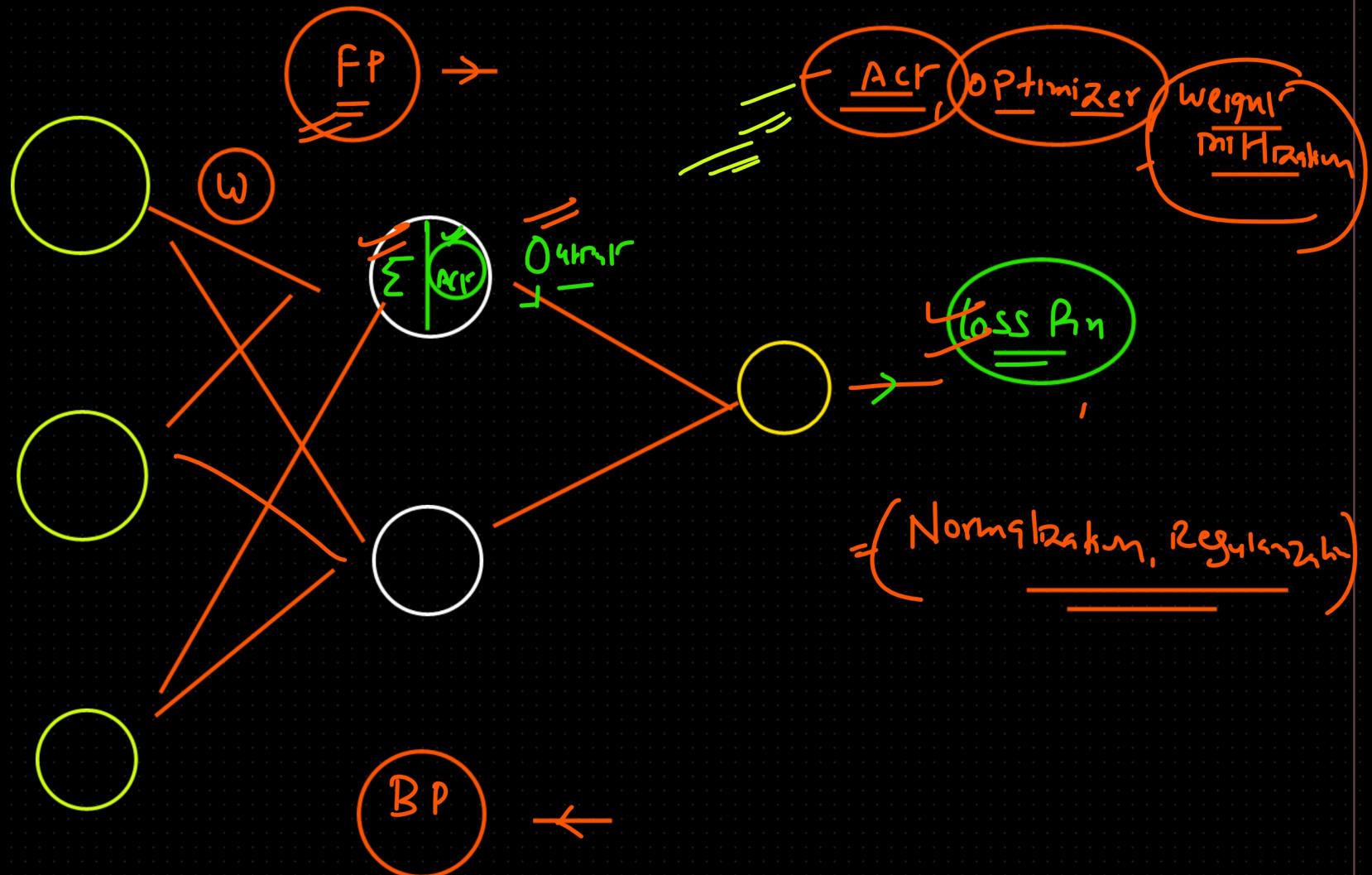
When the classes is high use SCCA

CCA reduce the performance

SCCA increase the Performance

Low data + low class \Rightarrow CCA is fine
large data + high class \Rightarrow SCCA is fine





Linear activation fun

Act Fun

Activation function tells when to trigger the summation part of a neuron. It acts like gate that opens or closes when based on conditions. Please refer to different types as discussed below for better understanding.

1 Binary step Activation

$$y = f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

Activation function Logic

- x represents the summation part of a neuron
- $f(x) = y$ represents the output of a neuron produced after applying activation function $f(x)$

2 Linear activation fun

$$\begin{aligned} y &= f(x) = x \\ f(1) &= 1 \\ f(2) &= 2 \\ f(100) &= 100 \end{aligned}$$

Non Linear activation fun

This will be discussed in depth in the next lecture

1 Sigmoid fun = Logistic fun

$$y = f(x) = \frac{1}{1 + e^{-x}}$$

2 \tanh $\Rightarrow y = f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

3 Relu \Rightarrow $y = f(x) = \max(0, x)$

4 Prelu or Parametric Relu $\Rightarrow f(x) = \max(\alpha x, x)$

5 Softmax $\Rightarrow \text{softmax}(z_i) = \frac{\exp(z_i) \text{ or } e^{z_i}}{\sum \exp(z_i) \text{ or } e^{z_i}}$

#

Regression \Rightarrow Linear & / or Sigmoid fn

Binary classification \Rightarrow Sigmoid, \tanh , (Relu)

Multiclass \Rightarrow Softmax

Why, Where, how = About activation fun^ to be discussed in next lec

Optimizer