

Lecture 05

"Loss function"

\hat{y} → predicted label
 w → weight

y → actual label
 b → bias

- we can call ' w ' and ' b ' as PARAMETERS
- The purpose is to pick the ' w ' and ' b ' in such a way, where logistic regression give $\hat{y} \approx y$
- we want that y or \hat{y} (predicted) should be closer.

Training the Logistic Regression mean:

$$\hat{y} = \text{weight} * \text{input image pixels} + \text{bias}$$

↓
predicted label

- our aim is to pick (w, b) in a way that gives the value of \hat{y} (predicted) should be closer to 1 for the positive class (here mean by cat class in which we are interested)

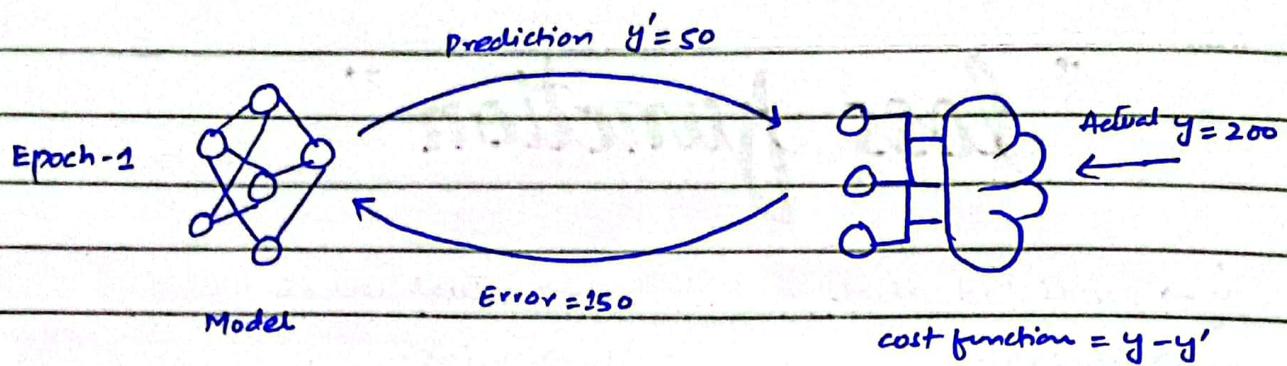
→ Goal

How to find ' w ' and ' b '

- ' w ' and ' b ' can be calculated through loss function.

[Loss function is used to calculate distance between y and \hat{y}]

- Distance in the loss function tell us how far off the predicted output is from the actual output.

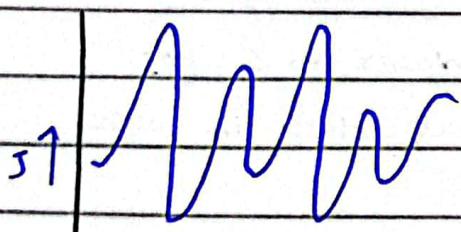


→ Loss function accept y and y' as input and return a number (that represent distance) as an output

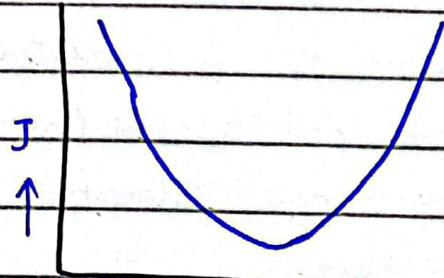
→ If the distance is closer to 0, it means there's a minimal distance b/w y and y' .

→ If that distance value is more than 0, it means there's a huge distance b/w y and y'

Non-convex vs convex functions



- Multiple local minima
- Global minimum harder to find



- Single global minimum
- local minimum is global

→ we can pick any function as a loss function but neural networks prefer convex functions for reliable optimization to find global minima.

→ Loss function use convex functions that have single minimum/maximum point

→ convex function has only one global optimum point

→ In Logistic Regression, we use negative log loss as a loss function

Negative log loss: (cross-Entropy loss OR Binary cross entropy loss)

$$L(y, \hat{y}) = \begin{cases} \rightarrow \hat{y} \approx y \Rightarrow L(y, \hat{y}) \rightarrow 0 \\ \rightarrow \hat{y} \neq y \Rightarrow L(y, \hat{y}) \rightarrow \infty \end{cases}$$

$$L(y, \hat{y}) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})] - \text{eq (1)}$$

case: 01

$y=0, \hat{y}=0$ we can say $y \approx \hat{y}$

$$\begin{aligned} L(y, \hat{y}) &= -[0 \log \hat{y} + (1-0) \log(1-0)] \\ &= -[0 + 1 \log 1] \\ &= -[1 \log 1] \quad \because \log_{10} 1 = 0 \\ &= -[\log 1] \quad \because 10^0 = 1 \\ &= -0 \end{aligned}$$

$$\text{Loss} = 0$$

↳ There is no distance b/w y and \hat{y}

case: 02

$$y=1, \hat{y}=0.1$$

$$\begin{aligned} L(y, \hat{y}) &= -[1 \log(0.1) + (1-1) \log(1-0.1)] \\ &= -[1 \log(0.1) + 0] \\ &= -[1 \log 0.1] \end{aligned}$$

$$= -[-1]$$

$$\therefore \log_{10} 0.1 = \log_{10} 10^{-1} = -1$$

$$= 1$$

There's a significant distance b/w y and \hat{y} , that means the values of w and b should be optimized.

Cost function:

→ It is the combined loss of all training examples

For one sample:

$$L(y, \hat{y}) = -[y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

For 'M' samples: Add $L(\overset{(i)}{\hat{y}}, \overset{(i)}{y})$ for all and take average

\curvearrowleft cost function

$$J(w) = \frac{1}{M} \sum_{i=1}^M L(\overset{(i)}{\hat{y}}, \overset{(i)}{y})$$
$$\therefore \hat{y} = \sigma(w^T x + b)$$

$$J(w, b) = -\frac{1}{M} [y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

$$J(w, b) = -\frac{1}{M} \left[y \log \{\sigma(w^T x + b)\} + (1-y) \log (1 - \{\sigma(w^T x + b)\}) \right]$$

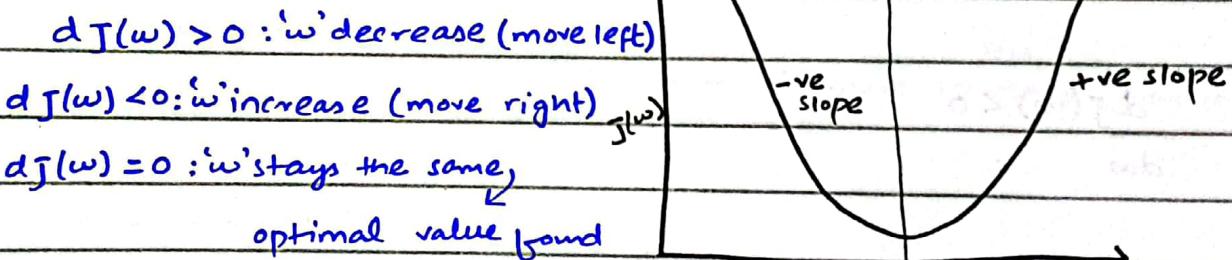
$\overset{(i)}{y} \approx \overset{(i)}{\hat{y}}$ (predicted outputs are approximately equals to actual outputs)

Lecture : 06

Gradient Descent

- Gradient Descent is used to optimize the parameters (w & b)
- Gradient means → slope or derivatives or rate of change
- Descent means → to move ⁱⁿ downward direction

Gradient descent is an optimization algorithm which iteratively moves in the direction of the negative gradient to find the minimum value of loss function.



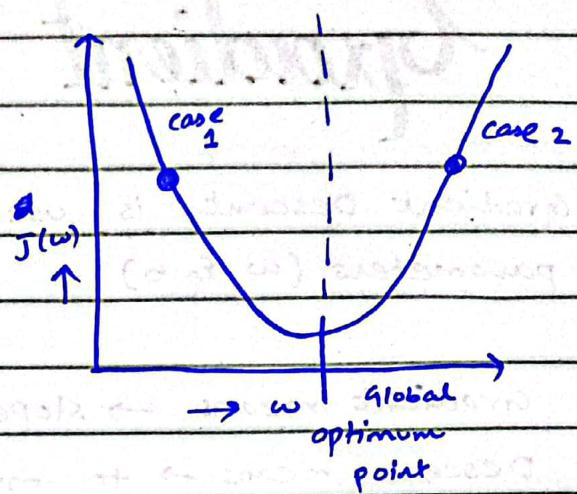
$J(w) \rightarrow$ cost of all training examples

$$w = w - \frac{d}{dw} J(w)$$

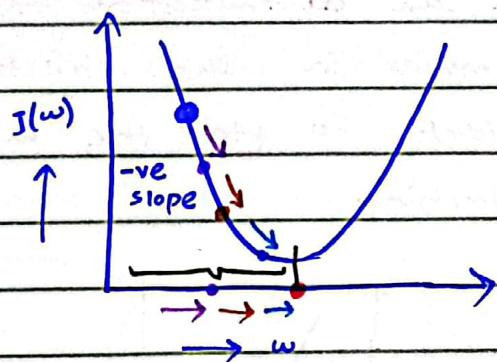
(change in loss per unit change
in weight)
updated as

$$\left[w = w - \frac{d}{dw} J(w) \right] \text{Gradient descent Algorithm}$$

→ choose random 'w' value



Case: 01



$$\frac{dJ(w)}{dw} < 0$$

$$\rightarrow w = w - \frac{dJ(w)}{dw}$$

$$w = w - (-\text{ve number})$$

$$w = w + \text{num}$$

$$\rightarrow w = w - (-\text{ve number})$$

$$w = w + \text{num}$$

$$w = w + (0.1)$$

$$\rightarrow w = w - 0$$

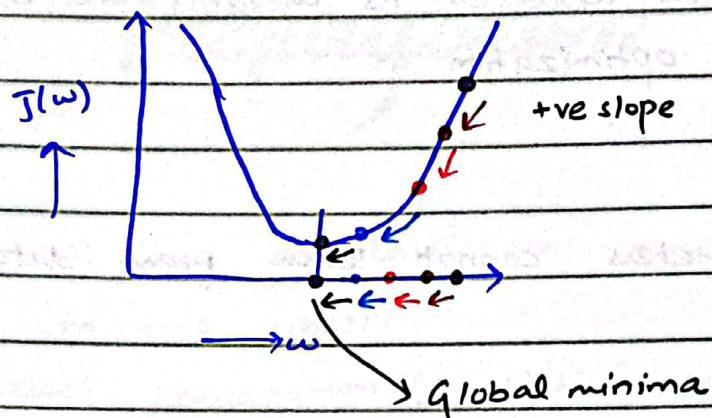
w = w (reach at global minimum point)

(weight is converged now)

(loss is minimized now)

$$J(\omega) \approx 0 \quad (\text{minimum loss})$$

case: 02



$$dJ(\omega) > 0$$

$$\rightarrow \omega = \omega - \frac{d}{d\omega} J(\omega)$$

$$\omega = \omega - (+\text{ve num})$$

$$\rightarrow \omega = \omega - (+\text{ve num})$$

$$\omega = \omega - (\text{very small number} \approx 0)$$

$$\omega = \omega - 0 \rightarrow$$

- * ω is optimized

- * cost is minimized

- * $y \approx \hat{y}$

- * Neural Network is trained now

$\omega = \omega - \alpha \frac{d}{d\omega} J(\omega)$
--

Learning rate (α) :-

The learning rate is a hyperparameter that controls the step size at which a model updates its weights and bias during optimization.

→ hyperparameters cannot learn from data.

Lecture : 07 (part: 1+2)

forward and backward

Propagation

→ Computations and calculations in neural network can be in two phases

- 1) Forward propagation (predictions, loss function)
- 2) Backward propagation (optimize parameters)

→ For understanding forward and backward propagation we use computational graph.

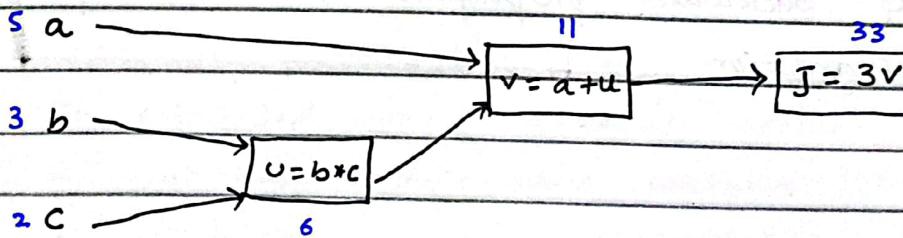
computational graph:

→ Explain different computations in the form of one graph node.

Let:

$$J(a, b, c) = 3(a + b * c)$$

↙
input parameters

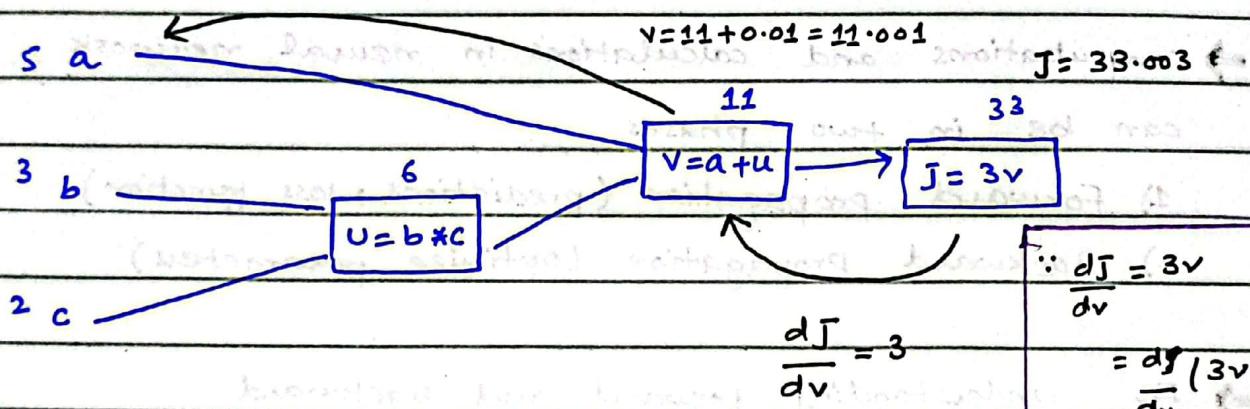


$$\begin{aligned}
 \frac{dV}{da} &= \frac{d}{da}(a+u) \\
 &= \frac{d}{da}(a) + \frac{d}{da}(u) \\
 &= 1 + 0 \\
 &= 1
 \end{aligned}$$

$\frac{dJ}{da}$ (unit change in a, b, c effect J)

→ (changing the unit of the independent variable and observing its effect on the dependent variable) Backward propagation

$$s+0.01 = s.01 \quad \frac{dv}{da} = 1 \quad v = s.001 + 6 = 11.001$$



$$\begin{aligned}
 \rightarrow \frac{dJ}{da} &= \frac{dv}{da} \cdot \frac{dJ}{dv} \\
 &= 1 \times 3 \\
 &= 3
 \end{aligned}$$

$$\begin{aligned}
 \because \frac{dJ}{dv} &= \frac{dJ}{dv}(3v) \\
 &= 3 \frac{dJ}{dv}(v) \\
 &= 3(1) \\
 &= 3
 \end{aligned}$$

Forward propagation: Inference, loss find

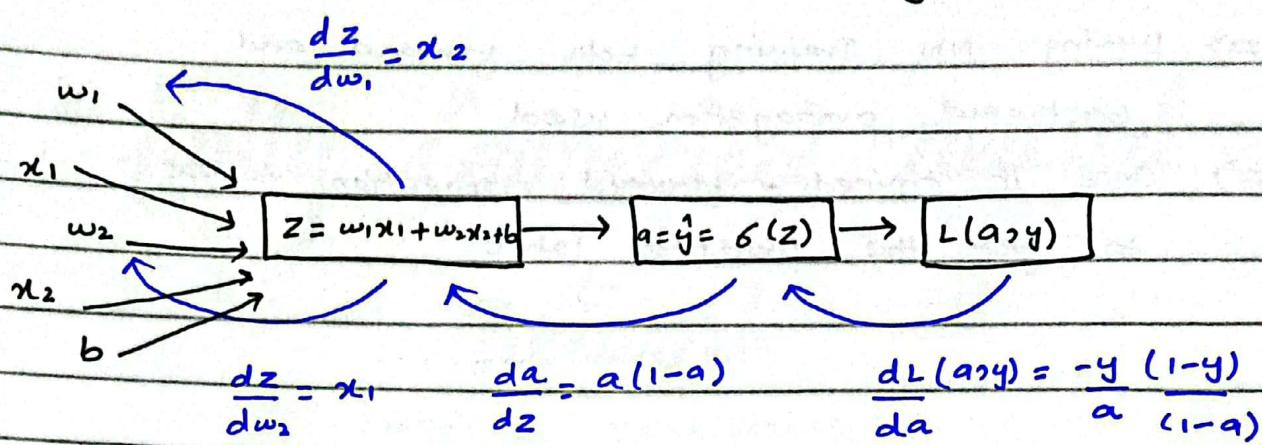
Backward propagation: Adjust parameters to minimize the loss

Forward & Backward propagation in terms of
Logistic Regression

$$Z = \omega^T x + b$$

$$\hat{y} = \sigma(Z)$$

$$\hat{y} \approx y$$



- Backward propagation provides relationship between parameters and loss
- If it's positive means 'w' increase 'loss' increase
- for negative means 'w' increase 'loss' decrease

$$\begin{aligned}\frac{dL}{dw_1} &= \frac{dZ}{dw_1} \times \frac{da}{dZ} \times \frac{dL}{da} \\ &= x_2 \cdot a(1-a) \cdot -\frac{y(1-y)}{a(1-a)} \\ &= x_2(a-y)\end{aligned}$$

⇒ when weight is optimized and loss is minimize then

$$\hat{y} \approx y$$

our Neural Network is Trained

⇒ once the relationship b/w parameters and loss is calculated then gradient descent algorithm is use to update the weights.

$$w_1 = w_1 - \alpha \left(\frac{dL}{dw_1} \right)$$

=> During NN Training both forward and backward propagation used

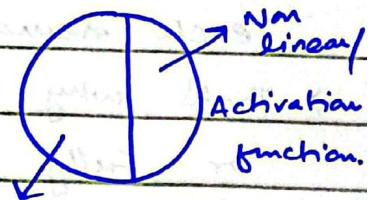
=> Once it trained, forward propagation is used to get the predicted labels

Lecture: 08

fully connected neural network

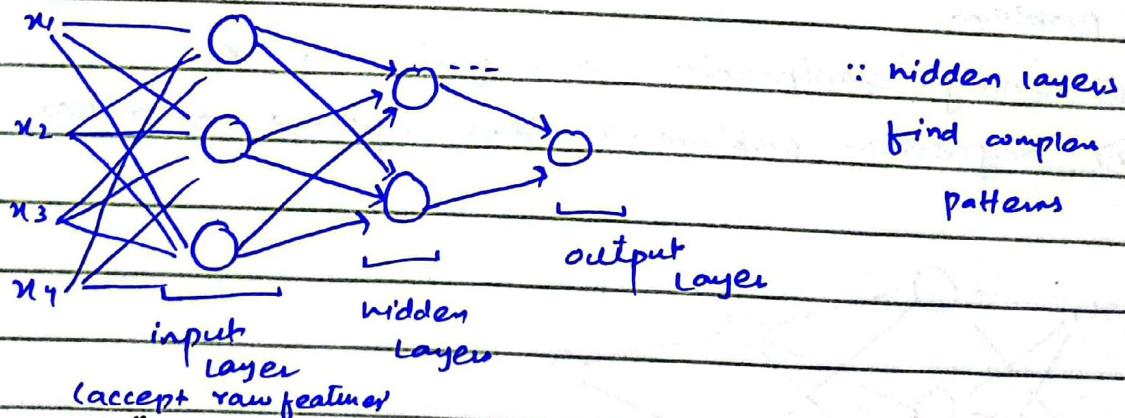
(1980's by Yann LeCun)

- Fully connected neural network is also known as standard neural network or dense neural network.
- Each neuron in the layer is connected to every neuron in subsequent layer
- Each neuron has two parts
 - ① Linear part
 - ② Non Linear / activation function
- activation function depends on the type of problem, it could be (RELU, sigmoid, softmax etc)



1) Layers

It contains input-layer, hidden layers, output layers



2) Activation function

It depends on the problem

- For binary classification sigmoid is used as an activation function in output layer.

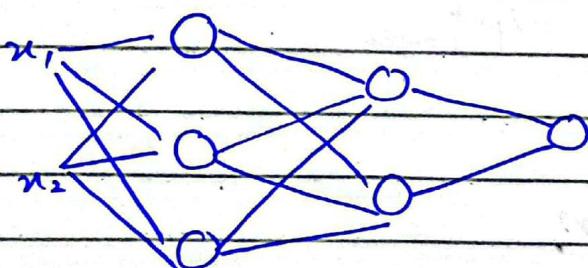
- for regression problem RELU can be use as an activation function in output layer.
- Neuron is also known as activation function / non-linearity
- Mostly, RELU is use as an activation function in input and hidden layer.

3) Structure

- Each neuron in a layer is connected to each neuron in subsequent layer.
- That's why it is known as Dense neural network or Fully connected neural network or standard neural network.

4) Forward Propagation

- In forward propagation predicted labels and loss is calculated
- loss function is chosen on the basis of problem
 - ① Binary classification \rightarrow Negative log loss
 - ② Regression problem \rightarrow Mean squared Error loss



	(rows)	(columns)
	Hidden units	X input features
$z^{[1]} = w^{[1]} x + b \dots a^{[1]} = g(z^{[1]})$		<i>W.E.R</i>
$z^{[2]} = w^{[2]} a^{[1]} + b \dots a^{[2]} = g(z^{[2]})$		activation function
$z^{[3]} = w^{[3]} a^{[2]} + b \Rightarrow y^i = a^{[3]} = g(z^{[3]})$		
Then:		
y is compared with \hat{y}		
On Backward propagation, we optimize the weights to reduce the loss.		
when $y \approx \hat{y}$ then		

Applications :-

- Hand-written digits identification (cv problem)
- House price prediction
- For structured data, standard neural network is preferred.
- speech recognition
- NLP (sentiment classification)
- Fraud detection in finance