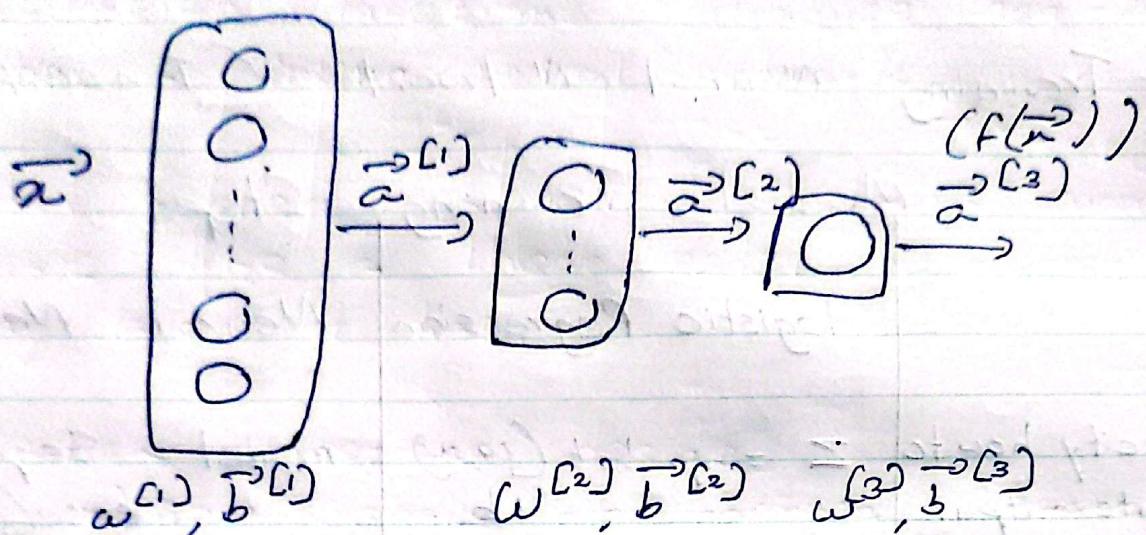


Week 02

01] Training Neural Network (Tensorflow)

Model	Training Steps
logistic Regression	Neural Network
① Specify how to compute output given input x and parameters w, b (define model)	$z = np.dot(w, x) + b$ $f(x) = 1 / (1 + np.exp(-z))$
$f_{\vec{w}, b}(\vec{x}) = ?$	<code>model = Sequential([Dense(...), Dense(...), Dense(...)])</code>
② Specify loss and cost	logistic loss
$L(f_{\vec{w}, b}(\vec{x}), y)$	$loss = -y * np.log(f(x)) - (1-y) * np.log(1-f(x))$
$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^n L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$	<code>model.compile(loss = BinaryCrossentropy())</code>
③ Train on data to minimize $J(\vec{w}, b)$	$w = w - \alpha \nabla J(\vec{w}, b)$ $b = b - \alpha \nabla J(\vec{w}, b)$
	<code>model.fit(x, y, epochs=100)</code>

1) Create the model



```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
```

model = Sequential [

```
Dense(Units= 25, activation = 'Sigmoid'),
Dense(Units = 15, activation = 'Sigmoid'),
Dense(Units = 1, activation = 'Sigmoid'), ]
```

2) Loss and Cost functions

i) binary Classification

$$L(f(\vec{x}), y) = -y \log(f(\vec{x})) - (1-y) \log(1-f(\vec{x}))$$

* Compare prediction vs target

* Logistic loss also known as binary cross entropy

from tensorflow.keras.losses import
BinaryCrossentropy, MeanSquaredError

model.compile(loss = BinaryCrossentropy())

i) Regression (predicting numbers and not categories)

model.compile($\text{loss} = \text{Mean SquaredError}()$)

3) Gradient descent

repeat {

$$w_j^{[L]} = w_j^{[L]} - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b_j^{[L]} = b_j^{[L]} - \alpha \frac{\partial}{\partial b_j} J(\vec{w}, b)$$

}

* Compute derivatives for gradient descent using "back propagation"

model.fit(X, y, epochs = 100)

02] Choosing activation functions

- * Activation functions introduce non-linearity into a ~~linear~~ neural network.
- * Without them, a neural network would behave like a linear model, no matter how many layers it has.

Activation Functions decide

- * how signals flow forward (forward propagation)
- * How gradients flow backward (back propagation / during learning)

Common Activation Functions

1] ReLU (Rectified Linear Unit)

$$f(z) = \max(0, z)$$

Derivative :

$$f'(z) \begin{cases} 1 & z > 0 \ (0, \infty) \\ 0 & z \leq 0 \end{cases}$$

- * Very fast to compute
- * Avoids vanishing gradients (for positive inputs)
- * Outputs only positive values
- * Neurons can die if outputs become 0 for all inputs (then gradient = 0)

2] Sigmoid

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

Derivative:

$$f'(z) = f(z) \times (1 - f(z))$$

- *) Output Range $(0, 1)$
- *) Smooth, S shaped curve
- *) Output can be interpreted as probability
- *) Causes vanishing gradient when used in hidden layers
- *) Gradients get tiny when z is far from 0

3] Linear

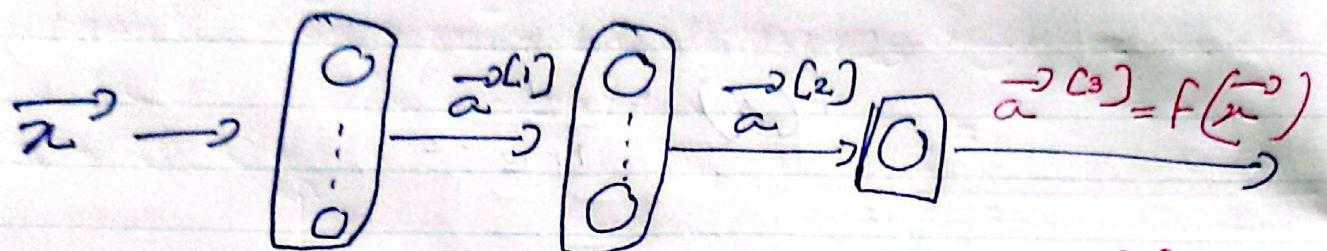
$$f(z) = z$$

Derivative:

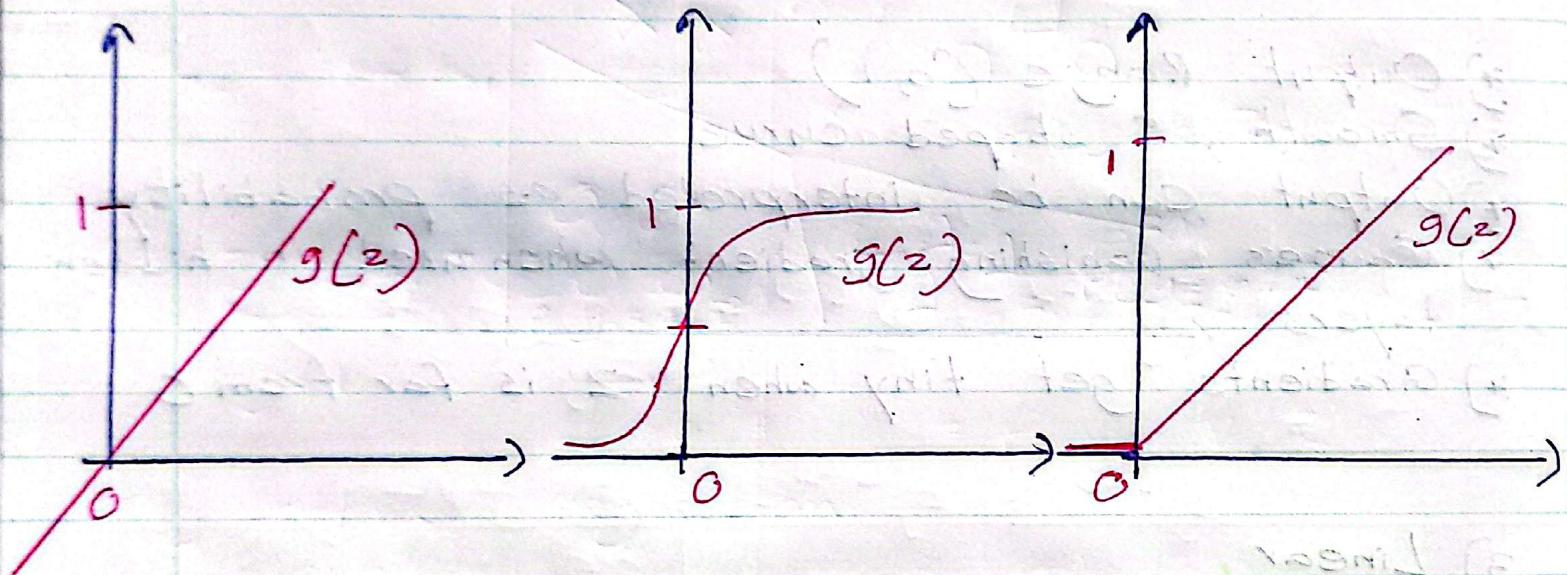
$$f'(z) = 1$$

- *) No non-linearity (No activation function)
- *) Output range $(-\infty, \infty)$

3) Output layer



$$f(\vec{z}) = a_1 c_1 + a_2 c_2 + a_3 c_3 = g(z_1, z_2, z_3)$$



Regression
Linear activation
 $y = +/-$

Binary Classification
Sigmoid activation
 $y = 0/1$

Regression
ReLU
 $y = 0 \text{ or } +$

4) Hidden layer

* Most Common choice is ReLU because it avoids vanishing gradient, fast and effective

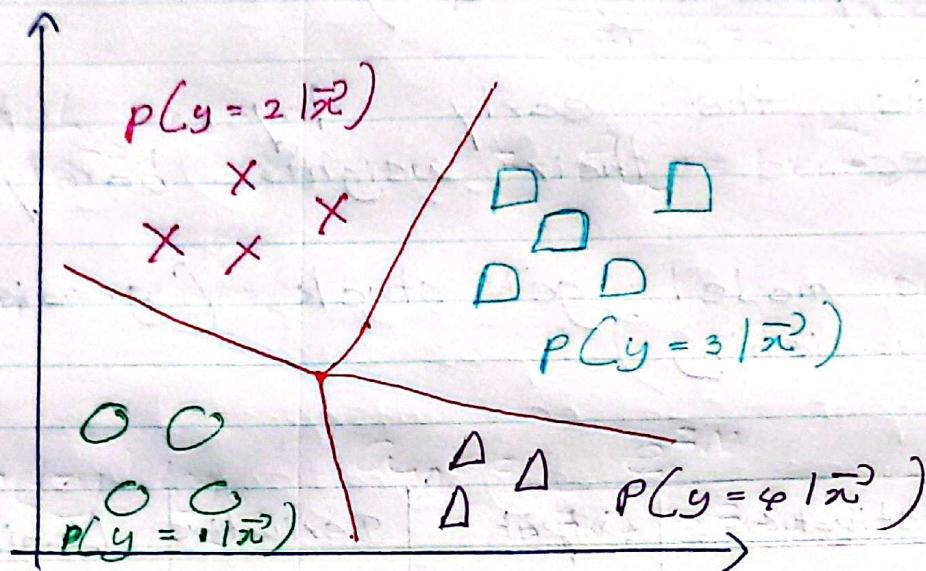
Vanishing gradient

- * As gradients are propagated backward through many layers, they can become smaller and smaller (approaching zero)
- * Then the early layers stop learning because their weights barely change
- * The model get stuck (gradient descent vanishing)

Activation	Output	Derivative	Vanishing Gradient descent
Sigmoid	(0, 1)	(0, 0.25)	High
ReLU	(0, ∞)	(0, 1)	Low
Tanh	(-1, 1)	(0, 1)	medium

03] Multi Class

* If target y can take on more than two possible values it a multi class classification



04] Softmax regression

② Logistic regression (2 possible output values)

$$o \alpha_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1|\vec{x})$$

$$\times \alpha_2 = 1 - \alpha_1 = P(y=0|\vec{x})$$

④ Softmax regression (4 possible outputs)
 $y = 1, 2, 3, 4$

$Z_1 = \vec{w}_1 \cdot \vec{x} + b_1, \quad a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$
○ X □ D

 $= P(y=1 | \vec{x}) 0.30$

$Z_2 = \vec{w}_2 \cdot \vec{x} + b_2, \quad a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$

 $= P(y=2 | \vec{x}) 0.20$

$Z_3 = \vec{w}_3 \cdot \vec{x} + b_3, \quad a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$

 $= P(y=3 | \vec{x}) 0.15$

$Z_4 = \vec{w}_4 \cdot \vec{x} + b_4, \quad a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$

 $= P(y=4 | \vec{x}) 0.35$

④

Softmax regression (N possible outputs)

$$z_j = \vec{w}_j \cdot \vec{x} + b_j \quad (j = 1, \dots, N)$$

$a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} = P(y=j | \vec{x})$

$(a_1 + a_2 + \dots + a_N = 1)$

Cost

Logistic Regression

$$z = \vec{\omega} \cdot \vec{x} + b$$

$$a_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1 | \vec{x})$$

$$a_2 = 1 - a_1 = P(y=0 | \vec{x})$$

$$\text{loss} = \underbrace{-y \log a_1}_{\text{if } y=1} - \underbrace{(1-y) \log}_{(1-a_1)} \underbrace{(1-a_1)}_{\text{if } y=0}$$

$$J(\vec{\omega}, b) = \text{average loss}$$

Softmax Regression

$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}}$$

$$= P(y=1 | \vec{x})$$

⋮

$$a_N = \frac{e^{z_N}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}}$$

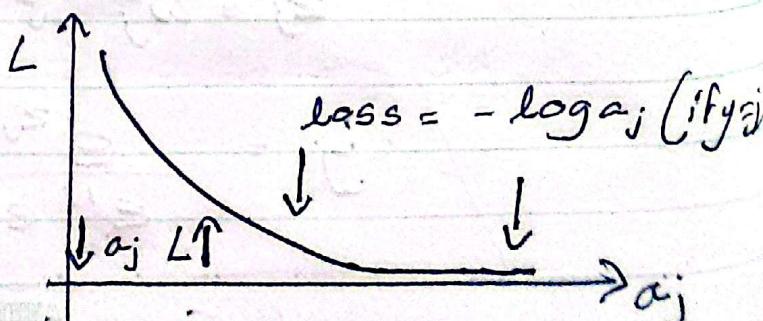
$$= P(y=N | \vec{x})$$

Crossentropy loss

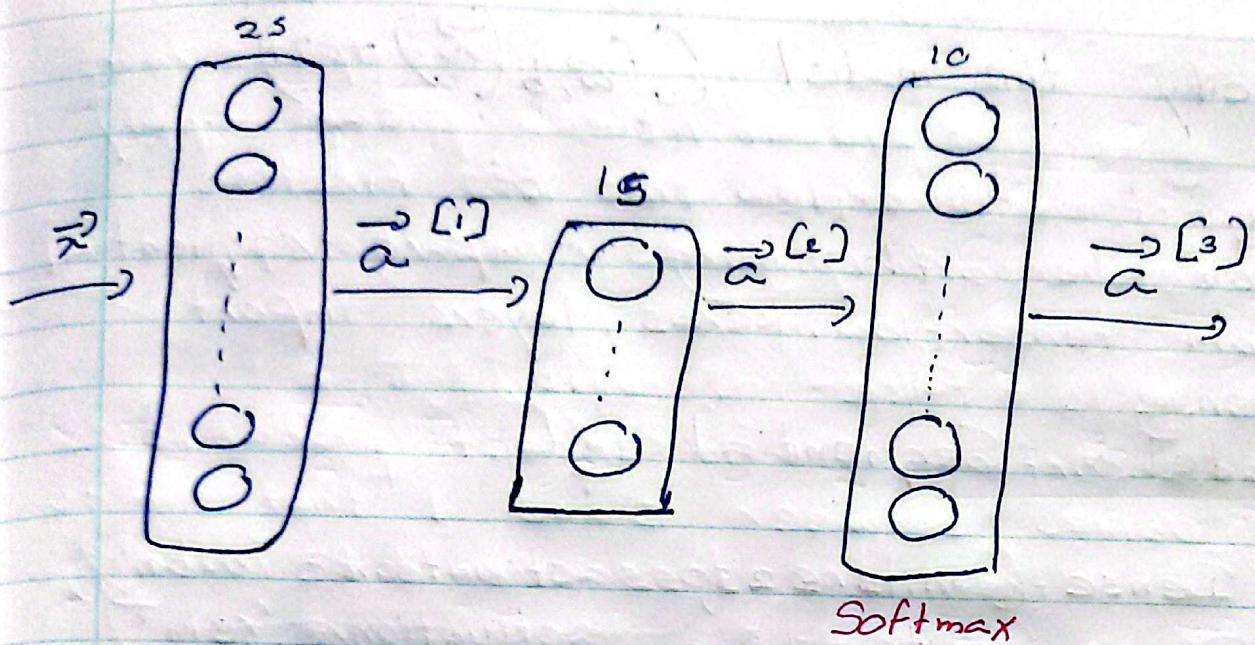
$$\text{loss}(a_1, \dots, a_N, y)$$

$$= \begin{cases} -\log a_1 & (\text{if } y=1) \\ -\log a_2 & (\text{if } y=2) \end{cases}$$

$$-\log a_N & (\text{if } y=N)$$



05) Neural Network with Softmax Output



$$z_i^{[3]} = \vec{w}_i \cdot \vec{a}^{[2]} + b_i^{[3]}$$

$$a_i^{[3]} = \frac{e^{z_i^{[3]}}}{e^{z_1^{[3]}} + \dots + e^{z_k^{[3]}}} = P(y=1 | \vec{a}^{[3]})$$

$$z_{10}^{[3]} = \vec{w}_{10} \cdot \vec{a}^{[2]} + b_{10}^{[3]}$$

$$a_{10}^{[3]} = \frac{e^{z_{10}^{[3]}}}{e^{z_1^{[3]}} + \dots + e^{z_{10}^{[3]}}} = P(y=10 | \vec{a}^{[3]})$$

logistic regression:

$$a_1^{[3]} = g(z_1^{[3]}) \quad a_2^{[3]} = g(z_2^{[3]})$$

Soft max regression

$$\vec{a}^{[3]} = (a_1^{[3]}, \dots, a_{10}^{[3]}) = g(z_1^{[3]}, \dots, z_{10}^{[3]})$$

Default Version (Not recommended)

1) Specify the model ($F_{\vec{w}, b}(\vec{x}) = ?$)

```
import tensorflow as tf  
from tensorflow.keras import sequential  
from tensorflow.keras.layers import Dense  
model = Sequential [  
    Dense (units = 25, activation = 'relu'),  
    Dense (units = 15, activation = 'relu'),  
    Dense (units = 10, activation = 'softmax')  
]
```

2) Specify loss and Cost ($L(F_{\vec{w}, b}(\vec{x}), y)$)

```
from tensorflow.keras.losses import  
SparseCategoricalCrossentropy
```

```
model.compile (loss = SparseCategorical  
Crossentropy ())
```

3) Train on data to minimize $J(\vec{w}, b)$

```
model.fit (X, Y, epochs = 100)
```

06] Improved Implementation of Softmax

```
model = Sequential ([  
    Dense (units = 25, activation = 'relu'),  
    Dense (units = 15, activation = 'relu'),  
    Dense (units = 10, activation = 'linear')])
```

```
from tensorflow.keras.losses import  
SparseCategoricalCrossEntropy
```

```
model.compile (..., loss = SparseCategorical  
Cross entropy  
(from_logits = True))
```

```
model.fit (X, Y, epochs = 100)
```

```
logits = model (X)  
f_X = tf.nn.softmax (logits)
```

- *] In neural networks, before applying an activation function the output of the last linear layer is called the logit.
- *] If you apply an activation you will get probabilities (a).
- *] If you don't apply you just have the raw Scores (z)
 - $z = \text{logits}$
 - $a = \text{activated outputs (probabilities)}$

*) from_logits = True :

This tells us my model output raw logits, not probabilities, please apply Softmax inside the loss function before computing cross entropy

*) from_logits = False (Default) :

Tensorflow would expect probabilities (values between 0 and 1 that sum to 1), then the loss would be wrong

Numerical round error

*) You set the output layer 'activation = linear' because you already applying softmax handled by loss function, internally, if you set activation = softmax it will double softmax and gives probabilities not raw outputs

*) After training we want to see actual probabilities or class predictions, since we use linear as the output layer it gives logits not probabilities,

So we \rightarrow $\text{softmax}(\text{logits})$ to take the K probabilities

07] Multi Label Classification

- * In multi label classification, each input can belong to more than one class simultaneously

One input → multiple correct classes

Example:

movie genre → Action, comedy both

- * Multi class classification means each input belongs to exactly one class out of many possible classes.

One input → One correct class

- * In multi label classification, output layer consists of one neuron per label
- * Use sigmoid function for each neuron (independent binary output)
- * Use binary cross entropy as loss
- * Each neuron outputs probability between 0 and 1, independently.

08] Adam Algorithm for optimization

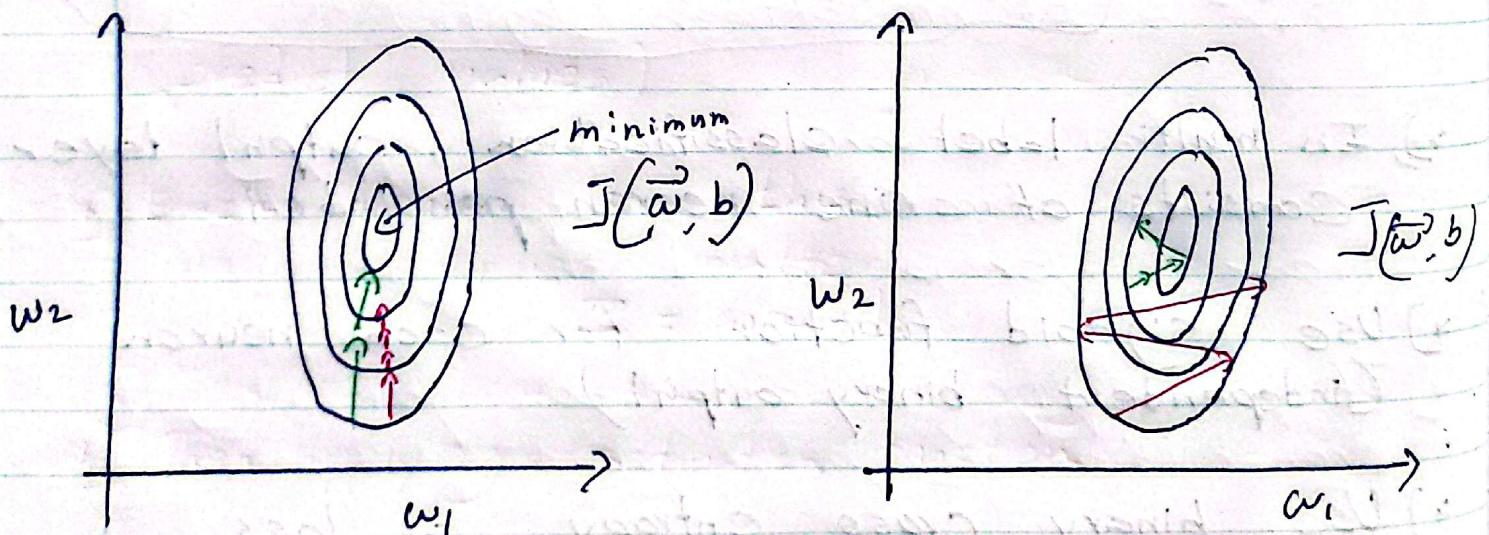
* Adam : Adaptive Moment Estimation
(not just one α)

$$\omega_i = \omega_i - \alpha_i \frac{\partial J(\vec{\omega}, b)}{\partial \omega_i}$$

$$\omega_{i,0} = \omega_{i,0} - \alpha_{i,0} \frac{\partial J(\vec{\omega}, b)}{\partial \omega_{i,0}}$$

$$b = b - \alpha_b \frac{\partial J(\vec{\omega}, b)}{\partial b}$$

* Faster and efficient than gradient descent



If ω_j (or b) keeps moving in same direction, increase α_j

If ω_j (or b) keeps oscillating reduce α_j

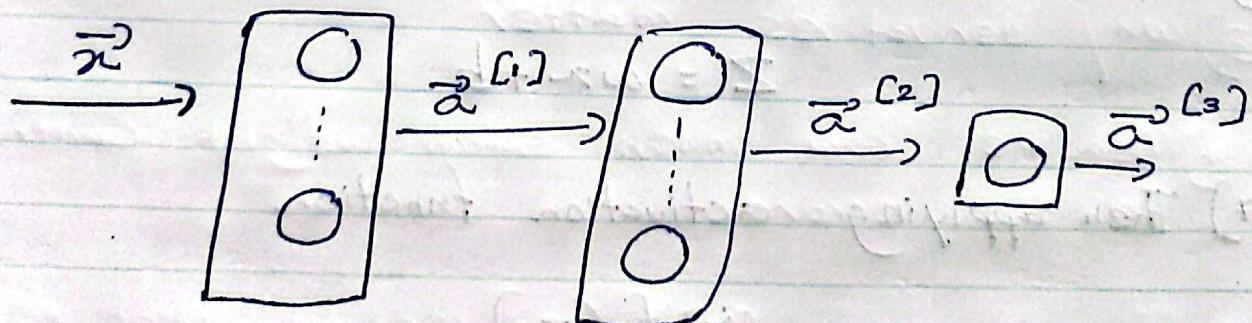
model = Sequential([
 tf.keras.layers.Dense(units=25,
 activation='sigmoid'),
 tf.keras.layers.Dense(units=15, activation='sigmoid'),
 tf.keras.layers.Dense(units=10, activation='linear')
])

Model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=1e-3),
 loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True))

model.fit(X, Y, epochs=100)

Q9] Layer Types

Dense layer



- Each neuron output is a function of all the activation outputs of the previous layer

$$\vec{\alpha}_i^{[2]} = g(\vec{w}_i^{[2]} \cdot \vec{\alpha}^{[1]} + b^{[2]})$$

Convolutional Layer

- * Each neuron only looks at part of the previous layer's outputs.

Why :

- Faster Computation
- Need less training data
(less prone to overfitting)

- 10) How neural network works :

- 1) Forward propagation

- * This happens when data (inputs) flow forward through the network : From input \rightarrow hidden layers \rightarrow output

(predict values)

- * Each neuron does this :

$$Z = w\alpha + b$$

- * Then applying activation function

$$\alpha = f(Z)$$

2) Backpropagation

* This is how the network learns :- by calculating how wrong the prediction was (error) and update weights to reduce it.

Based on calculus (Chain rule)

Step 1 : Compute loss

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

$y \Rightarrow$ true label

$\hat{y} \Rightarrow$ predicted output (from forward propagation)

Step 2 : Compute Derivative of Loss
(Error Signal)

$$\frac{\partial L}{\partial w}$$
 (This tells us how much changing w changes the loss)

Step 3 : Apply chain rule

* we propagate the error backwards layer by layer using chain rule

$$\frac{\delta L}{\delta w} = \frac{\delta L}{\delta a} \cdot \frac{\delta a}{\delta z} \cdot \frac{\delta z}{\delta w}$$

↓ ↓ → δ : derivative
how loss $f'(z)$ of weighted
Changes derivative sum
with activation of activation

Step 4: Update weights

After Computing gradients, we update weights using optimization (Gradient descent, Adam) to reduce the loss.

$$\omega = \omega - \alpha \frac{\partial L}{\partial \omega}$$

ex :-

$$\hat{y} = \sigma(z) := \frac{1}{1+e^{-z}} \quad z = \omega x + b$$

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

$$\frac{\partial L}{\partial \omega} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial \omega}$$

$$\frac{\partial L}{\partial \omega} = (\hat{y} - y) \cdot \alpha$$

3) why derivatives matter

* Measure how much each weight contributed to the error

* Update only in the direction that reduces error

* Learn patterns through repetition over many examples (epochs)