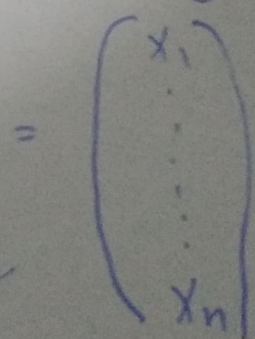
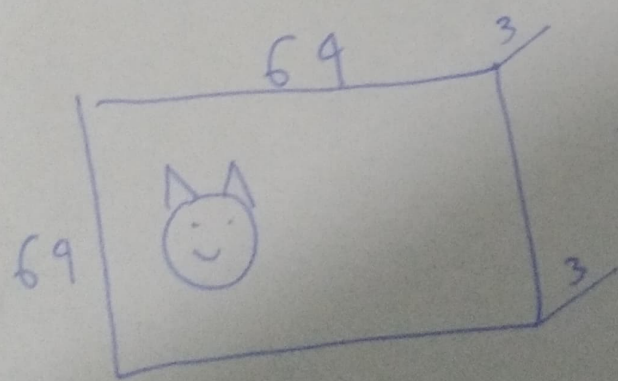


Deep Learning:-

- Computational power
- Data
- Algorithms.

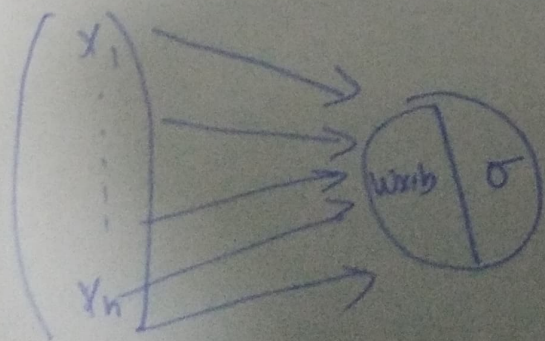
⊙ How can logistic regression be seen as neural network?

Goal:- Find cats in images ($1 \rightarrow$ presence of cats)



$0 \rightarrow$ absence of cats.

\sim Flattened the image.



$$\hat{y} = \sigma(\theta^T x)$$

$$= \sigma(w x + b)$$

$(1, 12288)$ $(12288, 1)$

1) Initialize w, b

2) Find the optimal w, b

Use $\hat{y} = \sigma(w x + b)$ to predict

* Define a loss function:

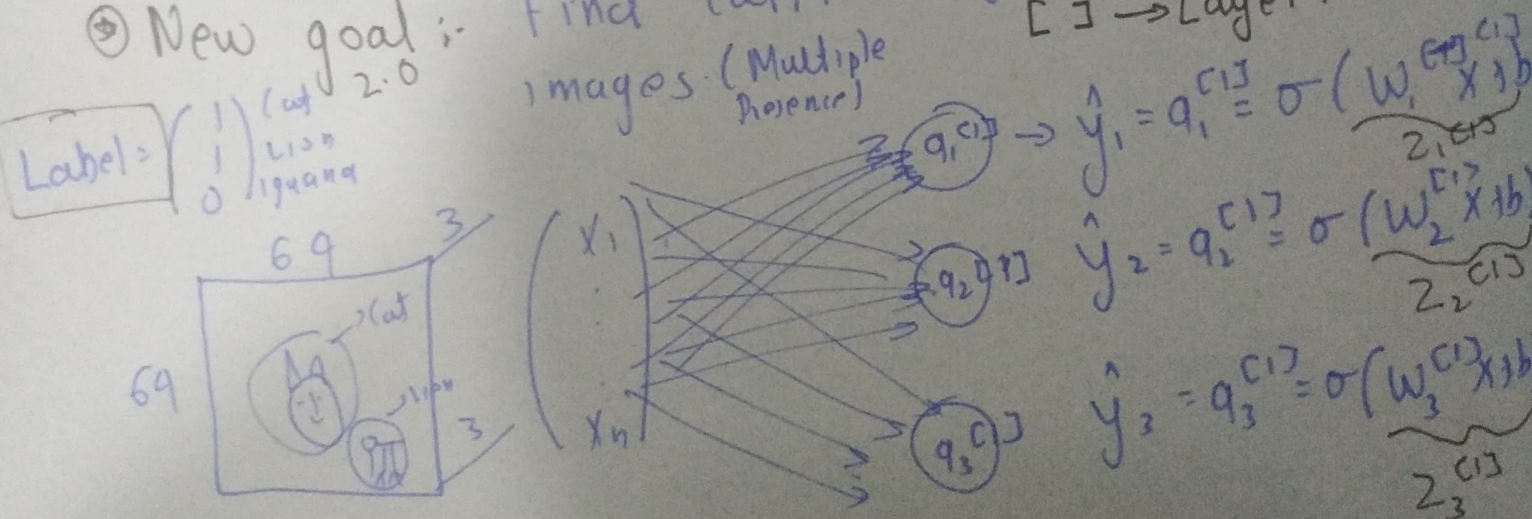
$$\min \mathcal{L} = -[y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

$$\left[\begin{aligned} w &= w - \alpha \frac{\partial \mathcal{L}}{\partial w} \\ b &= b - \alpha \frac{\partial \mathcal{L}}{\partial b} \end{aligned} \right] \rightarrow \text{Gradient descent.}$$

* Neuron = Linear + Activation

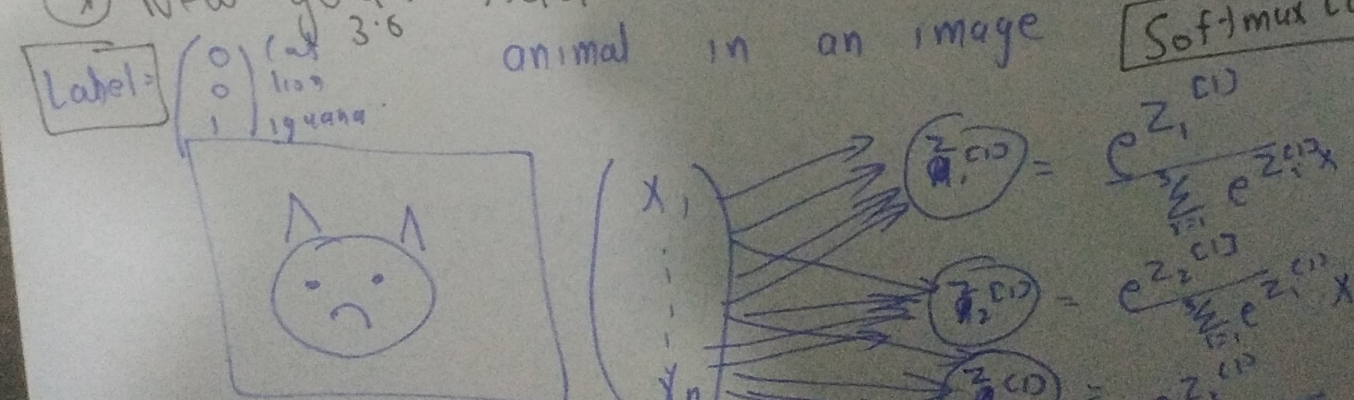
* Model = Architecture + Parameters

* New goal: Find cat/lion/iguana in images. (Multiple presence) $[] \rightarrow \text{Layer.}$



* New goal: Add constraint unique animal in an image

Softmax layer.



Redesign loss function for Goal 2.0

$$\mathcal{L} = - \sum_{k=1}^3 \left[y_k \log y_k + (1 - y_k) \log (1 - y_k) \right]$$

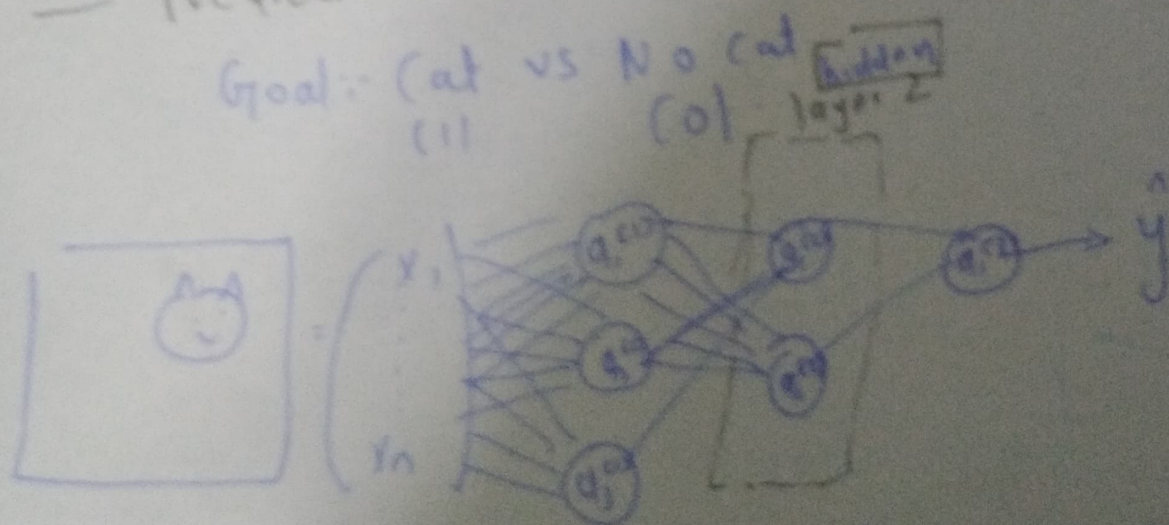
Redesign loss function for Goal 3.0

Class-Entropy loss

$$\mathcal{L} = - \sum_{k=1}^3 y_k \log \hat{y}_k$$

II Neural Networks:-

Goal:- Cat vs No cat



→ Layer is a cluster of neuron not connected to each other.

→ Hidden:- Output and Input are hidden from this layer.

⑧

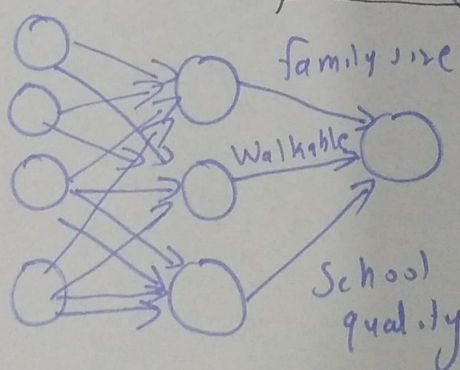
First Layer are going

⑨ Layers tend to learn different feature representations (complex when combined) which help in the final classification task

↳ Deeper layer learn more complex features.

⑩ House price prediction:-

#1
bedroom
Size
Zip code
Wealth



hidden feature which had to be hand encoded in machine learning but in deep learning can be learned if give enough data.

↳ In practice we use fully connected layer which means we connect every input of a layer to every output of ~~first~~ ^{forward} layer.

↳ End-to-End learning:-

↳ Input → [Layers] → Ground Truth

Some Black Box of learned features.

(129)

Forward Propagation equations:-

Matrix = Shape

$$b^{[1]} = (3, 1)$$

$$b^{[2]} = (2, 1)$$

$$b^{[3]} = (1, 1)$$

$$z^{[1]} = W^{[1]} \cdot X + b^{[1]} \rightarrow$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]} \cdot a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

$$z^{[3]} = W^{[3]} \cdot a^{[2]} + b^{[3]} \quad a^{[3]} = (1, 1)$$

$$y = a^{[3]} = \sigma(z^{[3]})$$

We put everything into matrices.

In general

$$W^{[l]} = (a^{[l-1]} \text{ shape}, a^{[l-1]} \text{ shape})$$

$$z^{[l]} = (\text{H of neurons in layer}, m)$$

$$a^{[l]} = (\text{H of neuron in layer}, m)$$

$$b^{[l]} = (a^{[l-1]} \text{ shape}, 1)$$

where m is size of batch.

$$W^{[1]} = (3, n)$$

$$X = (n, 1)$$

$$z^{[1]} = (3, 1)$$

$$a^{[1]} = (3, 1)$$

$$W^{[2]} = (2, 3)$$

$$z^{[2]} = (2, 1)$$

$$a^{[2]} = (2, 1)$$

$$W^{[3]} = (1, 2)$$

* What happens for an input batch of m examples?

Id of example

$$X = \begin{pmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{pmatrix}$$

$[] \rightarrow$ layer

$() \rightarrow$ id of example.

$(3, n)$

(n, m)

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$(3, 1)$

$(3, m)$

$$= \begin{bmatrix} | & | & & | \\ z^{[1]}(1) & z^{[1]}(2) & \dots & z^{[1]}(m) \\ | & | & & | \end{bmatrix}$$

↳ * Even though for b to add to add to $W^{[1]} X$ which is of shape $(3, m)$, b has to be also $(3, m)$ but we were using $(3, 1)$. The reason is that we will use broadcasting. It is a fact that parameter remain same but operation is parallelized so b would automatically through broadcasting would become $(3, m)$ where every row value has been repeated m times.

(131)

Optimizing $W^{[1]}, W^{[2]}, W^{[3]}, b^{[1]}, b^{[2]}, b^{[3]}$
(Continuing same example)

Define loss/cost function

Single example

m examples

$$J(y, \hat{y}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}^{(i)}$$

Batch gradient descent.

$$\text{with } \mathcal{L}^{(i)} = -[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})]$$

Backward Propagation

We have to compute this

$$\forall l = 1 \dots 3$$

$$W^{[l]} = W^{[l]} - \alpha \frac{\partial J}{\partial W^{[l]}}$$

$$b^{[l]} = b^{[l]} - \alpha \frac{\partial J}{\partial b^{[l]}}$$

We will start with

later layers now: and then use chain rule

$$y = a^{[3]}$$

$$\begin{aligned} \textcircled{1} \quad \frac{\partial J}{\partial W^{[3]}} &= \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{[3]}} \cdot \frac{\partial z^{[3]}}{\partial W^{[3]}} \\ \textcircled{2} \quad \frac{\partial J}{\partial W^{[2]}} &= \frac{\partial J}{\partial z^{[3]}} \cdot \frac{\partial z^{[3]}}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial W^{[2]}} \end{aligned}$$

Will go into detail in next lecture

Early layers with chain rule

③

$$\frac{\partial J}{\partial w^{[1]}}$$

=

$$\boxed{\frac{\partial J}{\partial z^{[2]}}}$$

132

from

$$\frac{\partial J}{\partial w^{[2]}}$$

$$\frac{\partial \pi^{[2]}}{\partial a^{[1]}}$$

$$\frac{\partial a^{[1]}}{\partial z^{[1]}}$$

$$\frac{\partial z^{[1]}}{\partial w^{[1]}}$$