

## # Introduction to deep learning

- Each neuron in first layer of network represents a particular parameter
- Every neuron of a intermediate layer is connected to each neuron of prev layer

## # Supervised learning

- For image recognition usually convoluted neural networks are used
- For sequenced data like audio which is sent over time as a 1D time series we use
  - RNN - Recurrent neural network
- For translation of texts - also have time component i.e. are sequenced as characters so we use advanced RNN
- Self driving - custom / hybrid neural networks

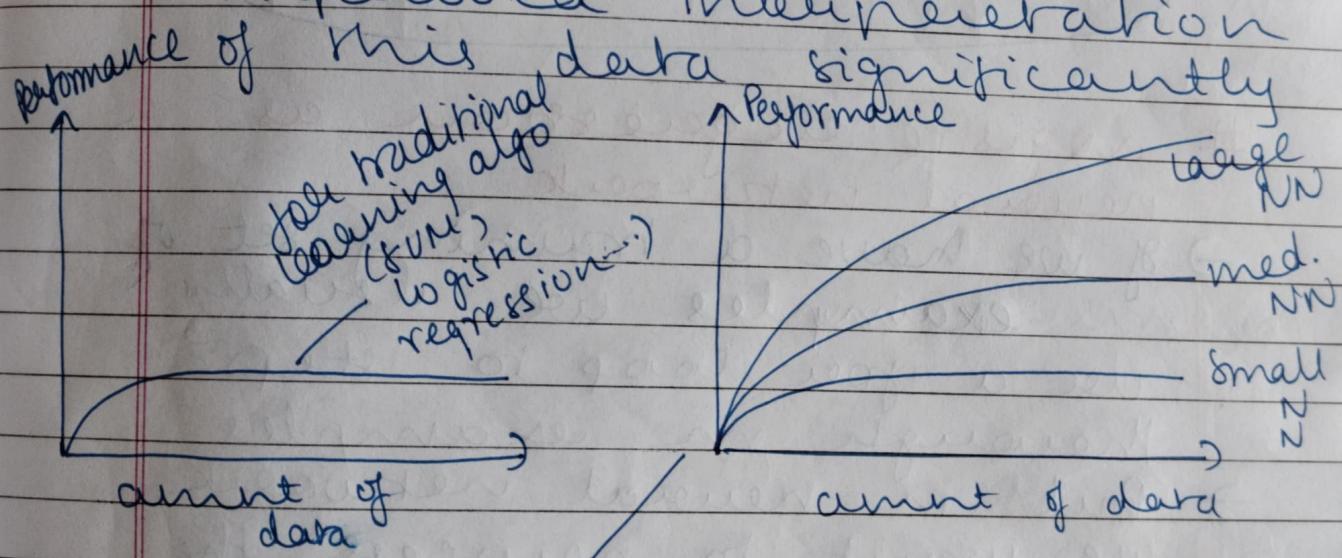
## Structured data -

- Databases of data
  - Each column will give u a parameter array in the computer ex upto price

- the table will include initial inputs, desired output and all the middle parameters
- All data has very well defined meaning

### Unstructured data

- Eg - raw audio, images, text
- unstructured not well defined (eg if u give image then data is pixel values)
- More difficult to interpret unstructured data
- Neural networks have improved interpretation



We need a larger

neural network and more data for better machine learning

- most reliable method to get better performance is to train a bigger network on fewer more data at it

- This only works upto a point as we run out of data or network is so big that it takes too long to train
- All data used must be labelled
- m :- no of training examples
- if u don't have large training set how you train it determines performance
- Algorithmic changes also affect deep learning performance (Eg: switching from sigmoid to ReLU)

Idea → code → experiment

## # Logistic regression as a neural network

- If we have a training set of  $m$  examples we usually use a for loop to step through  $m$  examples
- But in neural network we want to process the training set by training without using an explicit for loop.
- When we organize computation of neural network we have a forward propagation step followed

by a backward reorganization step

### Logistic regression

- form of binary classification
- for example - for a recognizing a cat image
  - 1 → cat or 0 → no cat
- In a computer to store images 3 separate matrices are stored corresponding to the red green and blue channels
- If input is  $64 \times 64$  pixels then we have 3  $64 \times 64$  matrices corresponding to red, green, blue pixel intensity values
- To turn pixel intensity values into a feature vector we need all the pixel values into an input feature vector  $x$
- To understand we define a feature vector  $x$  corresponding,

$$x = \begin{bmatrix} R_1 \\ R_2 \\ \vdots \\ R_{64} \\ G_1 \\ G_2 \\ \vdots \\ G_{64} \\ B_1 \\ B_2 \\ \vdots \\ B_{64} \end{bmatrix}$$

} no of values  
 $= 64 \times 64 \times 3 = 12288$

$n = n_x = 12288$

- In binary classification our goal is to learn a classifier that can input image represented by feature vector  $x$  and predict  $y$  either as 1 or 0
- A training set with  $m$  examples is represented as  $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$

$$\rightarrow X = \begin{bmatrix} & & & & & \\ & & & & & \\ & & x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \text{more} \\ \text{compact} \\ \text{notation} \\ \text{for training} \\ \text{set} \end{array} \leftarrow m \rightarrow$$

$$X \in \mathbb{R}^{(n_x) \times m}$$

$X$ . shape  $(n_x, m)$

command

for shape of matrix in python

$$\text{output } \rightarrow y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

$$y \in \mathbb{R}^{1 \times m}$$

$$y \text{ shape } = (1 \times m)$$

→ Hence given  $X$  we need an algorithm that would give output  $\hat{y} = P(y=1 \text{ given } x)$

probability

$$x \in \mathbb{R}^{n_x} \quad (0 \leq \hat{y} \leq 1)$$

$x$  dimensional vector

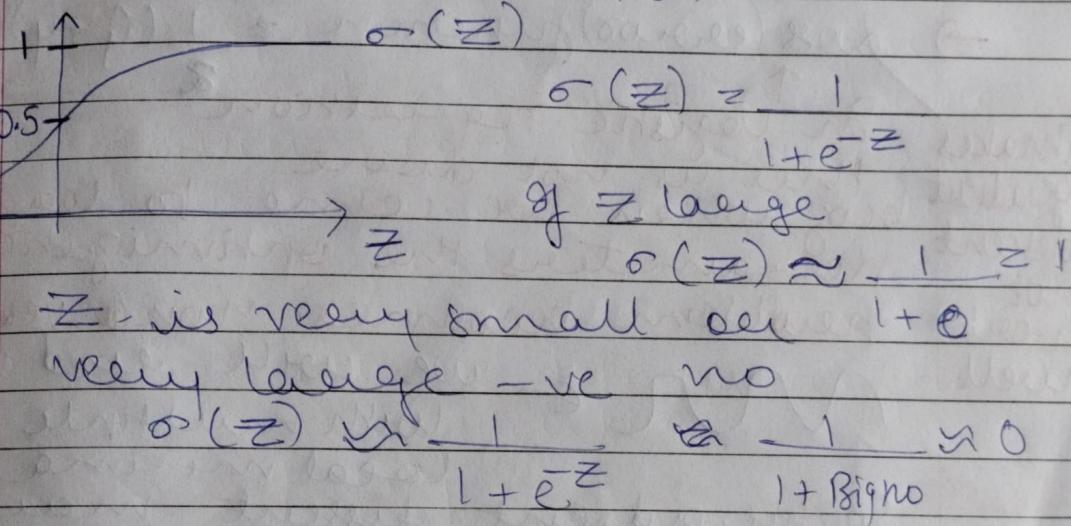
Parameters  $\rightarrow w \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$   
 Output  $\hat{y} = w^T x + b$ , — doesn't work

linear func of input  $x$  what we would use if doing linear regression  
 not a very good algorithm bcoz we want chance that  $y$  is equal to 1

$\hat{y}$  should be b/w 0 and 1

$w^T x + b$  — can be  $> 1$  or  $-ve$  which doesn't make sense

So in logistic regression our output is  $\hat{y} = \sigma(w^T x + b)$



→ When we program neural networks we usually keep  $w$  &  $b$  separate where  $b$  corresponds to an inter-spectrum

→ In some conventions we define extra feature  $x_0 = 1$

$$x \in \mathbb{R}^{n+1} \quad \hat{y} = \sigma(\theta^T x)$$

11

Just in a minute  
of taking  
notebook

$$\Theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \vdots \\ \theta_n \end{bmatrix} \quad \left\{ \begin{array}{l} b \leftarrow \\ w \in \end{array} \right. \quad \begin{array}{l} \text{Easier} \\ \text{to keep} \\ b \text{ & } w \\ \text{as separate} \end{array}$$

## # Logistic regression cost function

$$\rightarrow \hat{y}^{(i)} = \sigma(w^T x^{(i)} + b) \quad \text{where } \sigma(z^{(i)}) = \frac{1}{1 + e^{-z^{(i)}}}$$

$$z^{(i)} = w^T x^{(i)} + b$$

Given  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

we want  $\hat{y}^{(i)}$  vs  $y^{(i)}$  ideal

$$\rightarrow \text{loss (error) function} = \frac{1}{2} (\hat{y} - y)^2$$

makes  
gradient  
descent  
work  
well

$\uparrow$  in logistic regression  
This is not done  
bcz when we come to learn  
parameters the optimization  
problem becomes nonconvex.

$\curvearrowleft$  we will end up  
with multiple  
local minima  
so gradient descent may  
not find a global optimum

$\text{U} \curvearrowleft$  what we should  
use

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

$$\text{if } y = 1 : \mathcal{L}(\hat{y}, y) = -\log \hat{y}$$

we need

$-\log \hat{y}$  small  
as possible

$\therefore \log \hat{y}$  large as possible

$\therefore \hat{y}$  should be large

$$\text{if } y = 0 : \mathcal{L}(\hat{y}, y) = -\log(1 - \hat{y})$$

should be as small  
as possible

$1 - \hat{y}$  should be large

$\therefore$  we want  $\hat{y}$  as small  
as possible.

$$\text{cost function} = J(w, b)$$

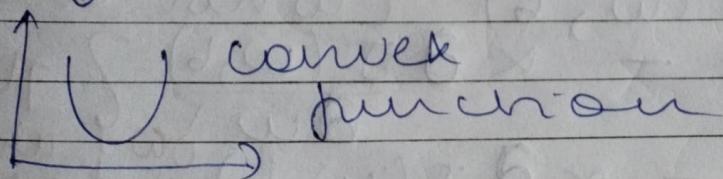
measures how well =  $\frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$   
we do over  
entire training set

$$= -\frac{1}{m} \sum_{i=1}^m \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

$\Rightarrow$  so in training we are going  
to try to find parameters  
 $w$  &  $b$  that minimize overall  
cost func  $J$

## # gradient descent

$\rightarrow$  we need  $w, b$  that minimize  
 $J(w, b)$

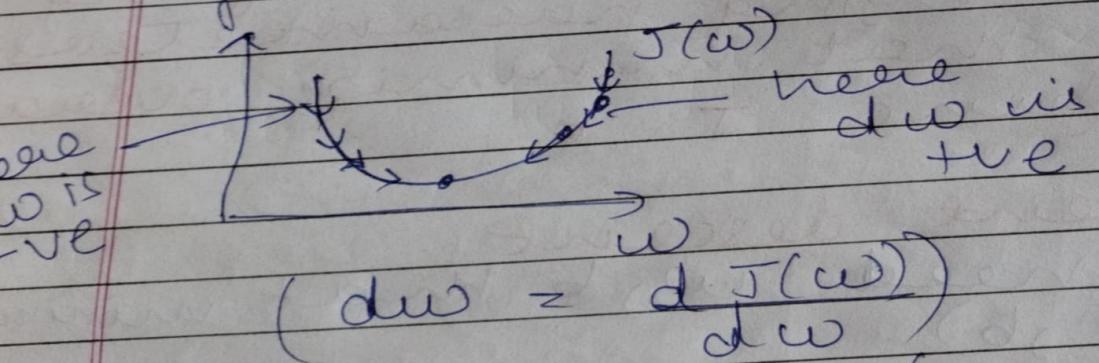


- To find a good value of  $w$ , b we first initialize them to some initial value
  - Gradient descent takes a step in steepest downward direction again and again until we converge to global maximum close to zero
  - We keep updating  $w$  to make it closer to zero
- $w = w - \alpha \frac{d J(w)}{d w}$
- learning rate that controls how big step we take in each iteration

$$w = w - \alpha dw$$

$dw$  represents  $\frac{d J(w)}{d w}$  in code

- if  $dw$  is -ve we take  $\downarrow w$
- if  $dw$  is +ve we  $\uparrow w$



Since  $J(w, b)$  ( $w$  not only parameter)  
 $w = w - \alpha \frac{\partial J(w, b)}{\partial w}$  ? These are actual  $\partial w$  in code

$b = b - \alpha \frac{\partial J(w, b)}{\partial b}$   $\partial b$  in code

## # Derivatives

$$\text{if } \frac{df(a)}{da} = x$$

\* if  $a$  is increased by  $\Delta a$   
then  $f(a)$  increases by  $x\Delta a$

## # Computation graph

- Computations of a neural network are organized in terms of a forward pass or forward propagation step followed by a in which we compute o/p of neural network followed by a backward pass or back propagation step which we use to compute gradients or derivatives
- computation graph explains why its ordered this way
- suppose we have a func

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

$$bc = v$$

$$v = a + u \quad J = 3u$$

$$a \xrightarrow{dJ/da}$$

$$b \rightarrow [v = bc] \quad c \rightarrow$$

$$[v = a + u] \rightarrow [J = 3v]$$

$$dJ \xrightarrow{dv}$$

- Computation is done left to right per function  $J$
- for derivatives computation is from right to left

## # Computing Derivatives

$$\rightarrow \frac{dJ}{dv} = \frac{\Delta J}{\Delta v} = 3$$

$$\rightarrow \frac{dJ}{da} = \frac{\Delta J}{\Delta a} = 3 = \frac{dJ}{dv} \times \frac{dv}{da}$$

$$\rightarrow \frac{dv}{da} = \frac{\Delta v}{\Delta a} = 1$$

chain rule

$\rightarrow$  When writing codes to implement backpropagation we will find some final output variable that we really care about (e.g.: -  $J$ )  $\Rightarrow$  it's really the last node of the computational graph.

$\rightarrow$  We will be calculating  $\frac{d}{d \text{ var}}$  final o/p var

In code we represent  $\frac{d}{d \text{ var}}$   
by ("dvar")  $\uparrow$  (var = an intermediate variable)

$$\rightarrow \frac{dJ}{db} = \frac{dJ}{dv} \cdot \frac{dv}{db} = 3 \times c = 3c$$

## # Logistic regression & gradient descent

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$L(a, y) = -(y \log(a) + (1-y) \log(1-a))$$

$$x_1 \rightarrow \\ w_1 \rightarrow$$

$$x_2 \rightarrow z = w_1 x_1 + w_2 x_2 + b \rightarrow \\ w_2 \rightarrow \frac{dz}{da} = \frac{dL}{dz} \cdot \frac{dL(a,y)}{da} \leftarrow \\ b \rightarrow \text{decode } \frac{dz}{da} \quad dz = a - y$$

$$\hat{y} = a = \sigma(z) \rightarrow \alpha(a, y) \quad \begin{matrix} \frac{dL}{da} \\ \frac{da}{dz} \end{matrix} \\ \text{decode } \frac{da}{dz} = \frac{dL(a, y)}{da} \leftarrow \\ \frac{da}{dz} = \frac{-y + 1-y}{a \quad 1-a}$$

$$dw_1 = x_1 dz \quad dw_2 = x_2 dz \\ db = dz$$

$$\text{decode } w_1 = w_1 - \alpha dw_1 \\ w_2 = w_2 - \alpha dw_2 \\ b = b - \alpha db$$

# gradient descent for m examples

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

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \\ \sigma(w^T x^{(i)} + b) \\ \frac{\partial J(w, b)}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} \alpha(a^{(i)}, y^{(i)}) \\ \underbrace{\qquad\qquad\qquad}_{d w_1^{(i)}}$$

Let's initialize  $J = 0$ ,  $d w_1 = 0$ ,  $d w_2 = 0$ ,  $db = 0$

$\rightarrow$  For  $i = 1$  to  $m$  (for loop)

$$z^{(i)} = w^T x^{(i)} + b \\ a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log(a^{(i)}) + (1-y^{(i)}) \log(1-a^{(i)})]$$

$$d\tilde{z}^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 + = x_1^{(i)} d\tilde{z}^{(i)}$$

$$dw_2 + = x_2^{(i)} d\tilde{z}^{(i)}$$

$$db + = dz^{(i)}$$

$$J / = m, dw_1 / = m; dw_2 / = m \text{ layer}$$

$$db / = m$$

$n = 2$   
only  
2 features  
ie 2 neurons in 1st layer

we have  
reduced  
by

in  
one of  
m  
examples

$$w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

→ Vectorization will help  
us do this without for  
loops