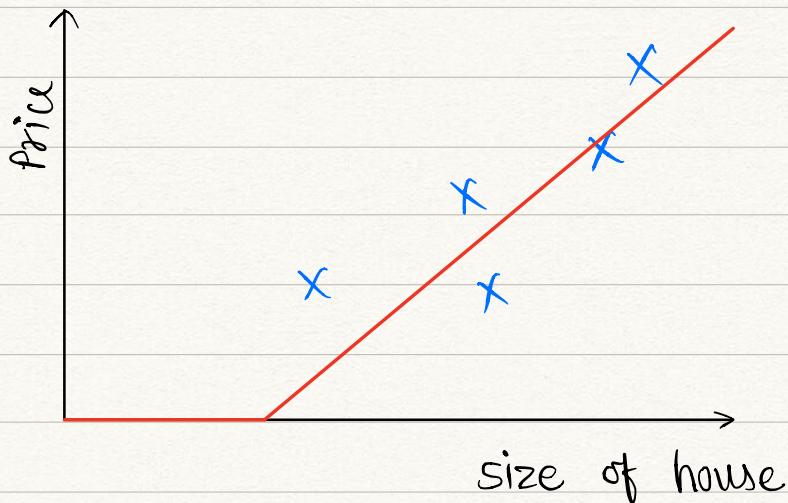


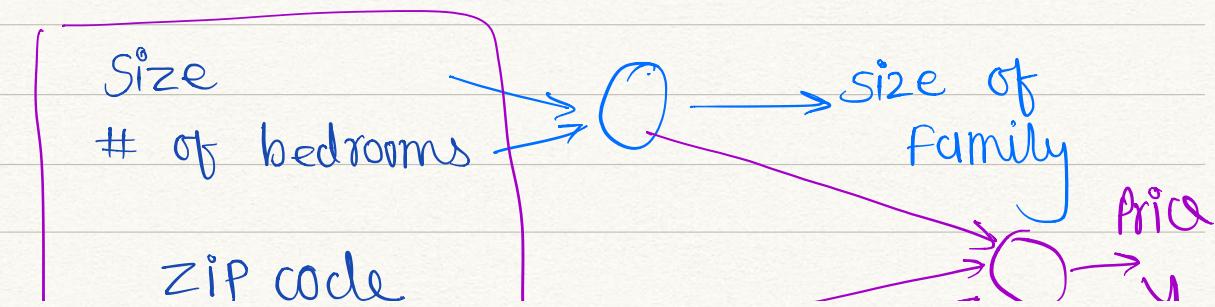
Q. What is Neural Network?

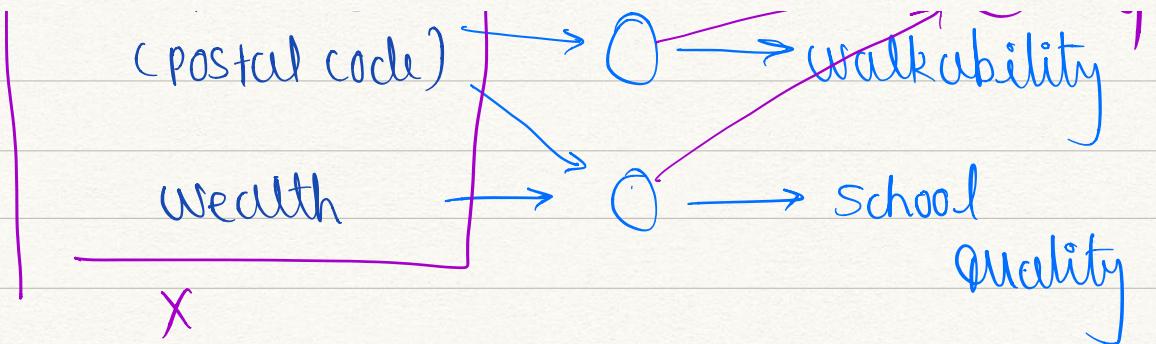


size \xrightarrow{x} \xrightarrow{y} price } Single Neuron network

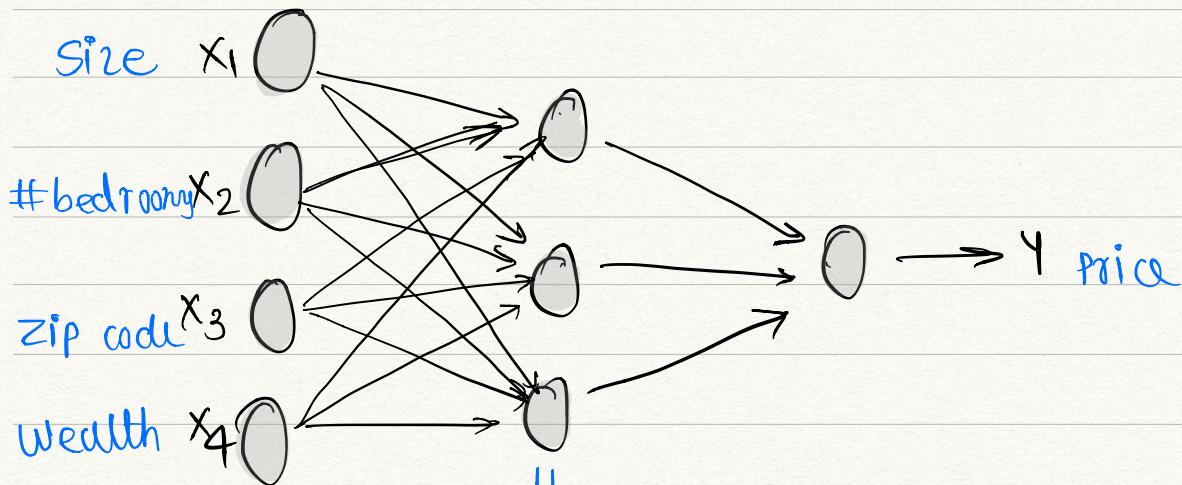
Neuron } goes to zero sometimes &
takes off as st. line.
 \downarrow

ReLU Function
Rectified linear unit





give X & Y , Network will figure out middle Neurons.



- Takes all four features as input
- Densely connected

* Supervised Learning:

| X | Y | Application |
|---------------|--------------------------------|--------------------|
| Home features | Price | Real Estate |
| Ad. User info | Click on Ad _i (col) | Online Advertising |

Image
Audio
English

Object
Text Transcript
Chinese

Photo tagging
Speech recognition
Machine Translation

Autonomous driving
cluster hybrid

Image, Radar info

Position of other car

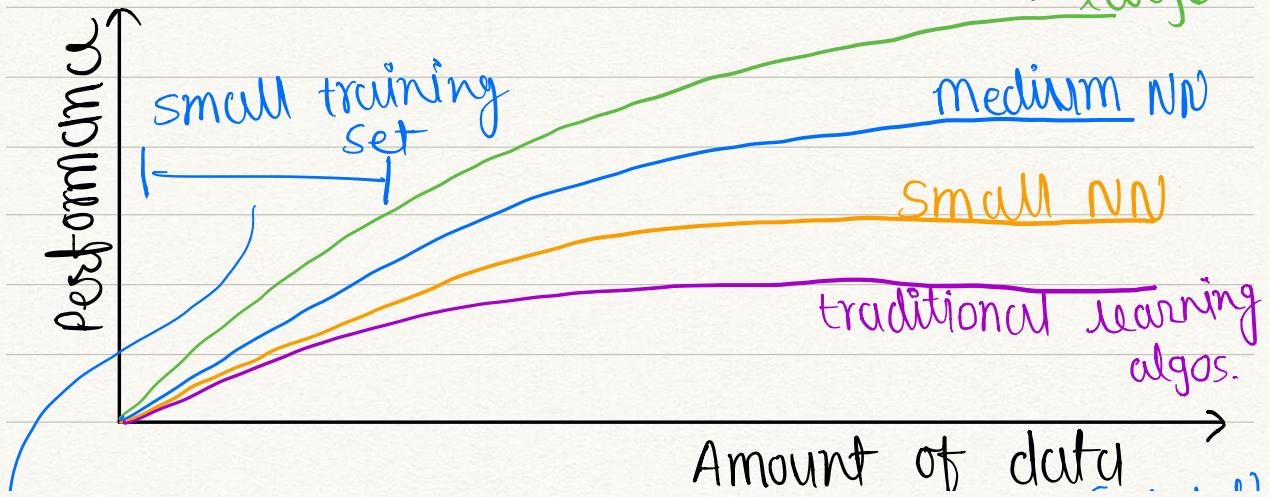
Structured Data
→ organised tabular

Unstructured Data

- Audio
- Image
- Text
↓

Hard for computer

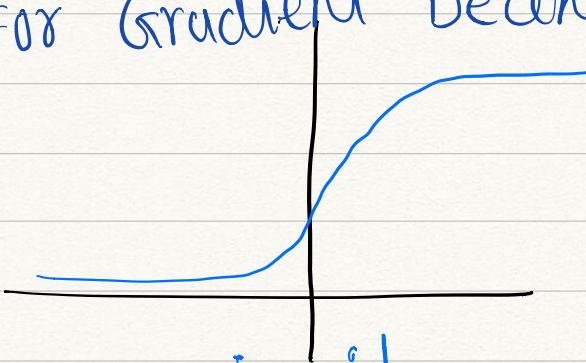
Why Deep learning is taking off?



→ Scale drives Deep learning (labelled) progress

Position of NN is not clearly defined.

For Gradient Decent



sigmoid



ReLU

Decay of learning rate as we go forward.

No decay of learning rate

change in Algorithms.

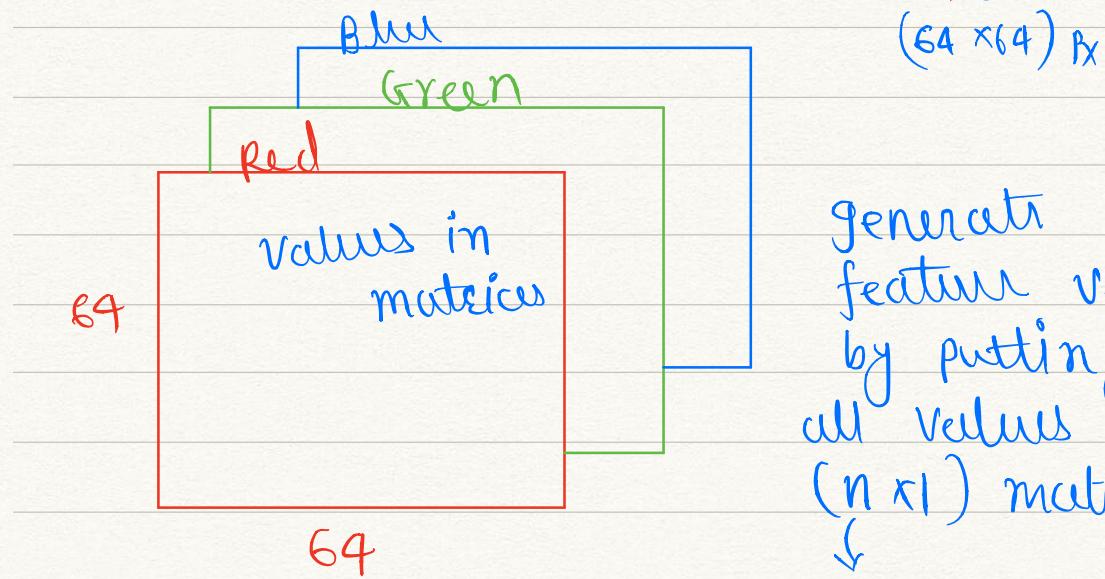
⇒ Advances in Computation.

We can use for loop to go through entire training set but it will take lot of time complexity.

→ In Neural Networks we don't use for loops to go through training sets

Binary Classification :

(e.g.) Cat or No Cat in image



generating
feature vector
by putting
all values in
(n x 1) matrix
↓
64 x 64 x 3

Here : (x, y)

$$x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

m training Examples: $(x_1, y_1), (x_2, y_2), \dots$

$$X = \begin{bmatrix} & & & & \\ & & & & \\ x_1 & x_2 & x_3 & \dots & x_m \\ & & & & \\ & & & & \end{bmatrix} \begin{matrix} \\ \\ \\ n_x \\ \hline m \text{ columns} \end{matrix}$$

$$Y = [y_1 \ y_2 \ \dots \ y_m]_{1 \times m}$$

Logistic Regression:

Given x , $\hat{y} = P(y=1|x)$
 $x \in \mathbb{R}^{n_x}$

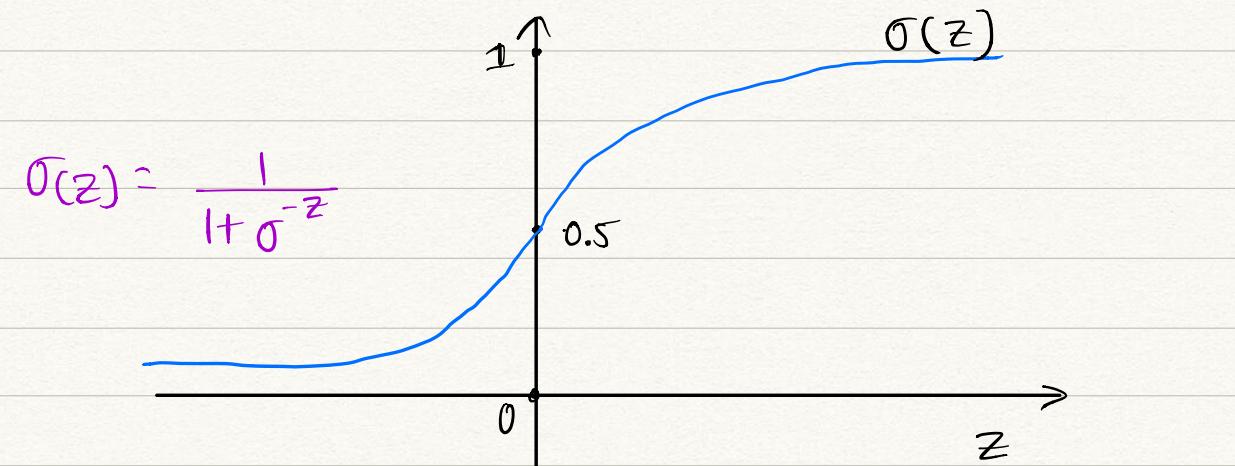
Parameter: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

Output $\hat{y} = \frac{w^T x + b}{\text{not good}}$

can't be binary (0 or 1)

so apply sigmoid function

$$\text{output } \hat{y} = \sigma(w^T x + b) \quad \downarrow z$$



In logistic Regression:

We try to find w & b so that \hat{y} become good estimate of chance of y being equal to 1.

Another Notation:

$$x_0 = 1, \quad x \in \mathbb{R}^{n_{x+1}}$$

$$\Theta = \begin{bmatrix} \hat{y} \\ \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{bmatrix} = \sigma(\Theta^T x) \rightarrow b \quad w$$

Logistic Regression Cost Function:

$$L(\hat{y}, y) = \frac{1}{2} \cancel{C} \cancel{(\hat{y}-y)^2}$$

optimization problem becomes non convex. might have local optimum so gradient decent might not find global optimum.

loss function

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

small as possible

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

if $y=0$

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

large $\Rightarrow \hat{y} \rightarrow \text{large}$

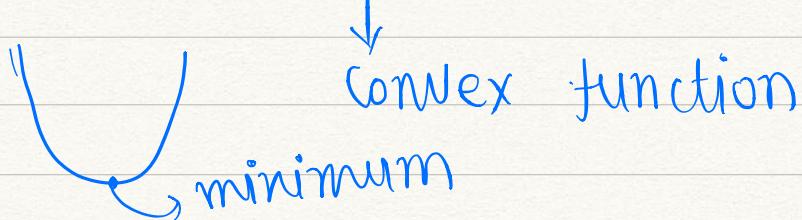
large $\rightarrow \hat{y} \rightarrow \underline{\text{small}}$

Cost Function:

$$\begin{aligned} \text{overall } J(w, b) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\ &= \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})] \end{aligned}$$

Gradient Descent:

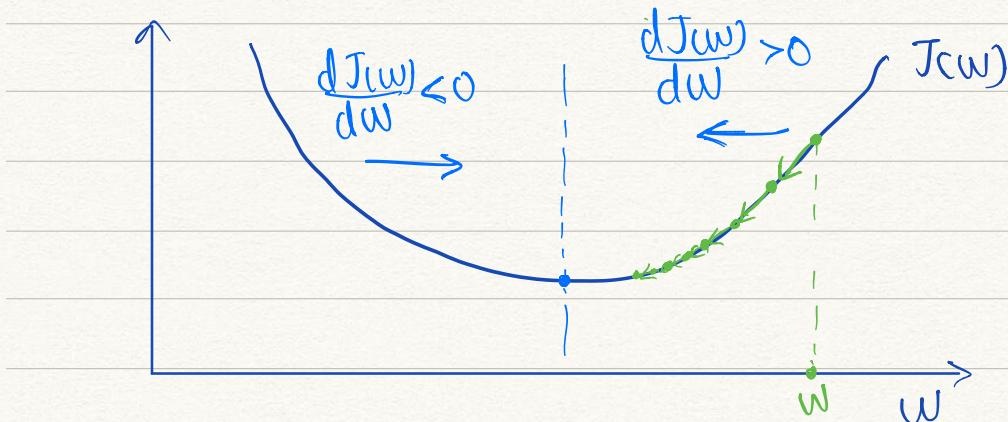
We want to find w, b for which $J(w, b)$ is minimum.



→ It will start at some initial point and then go down on the steepest descent in direction of $\nabla J(w, b)$.

→ It will converge to global minimum in case of convex function.

Gradient Descent for 1 D:



Repeat { $w: w - \alpha \frac{dJ(w)}{dw}$ }

$\alpha \rightarrow$ warning rate

$\frac{dJ(w)}{dw} \Rightarrow$ in code " dw "

Coding term: $w \leftarrow w - \alpha dw$

Actual in 2D:

$$w := w - \alpha \frac{d J(w, b)}{dw} \rightarrow \text{partial derivative}$$

$$b := b - \alpha \frac{d J_{CW,b})}{db}$$

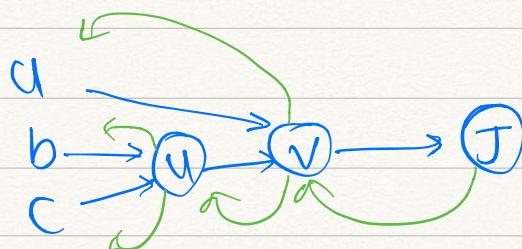
Computation Graph:

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

$$y = bc$$

$$J = 3V$$

$$V = g + u$$



$$u = bc$$

$$V = ct + bc$$

$$J=3V$$

$$\frac{dJ}{dV} = 3$$

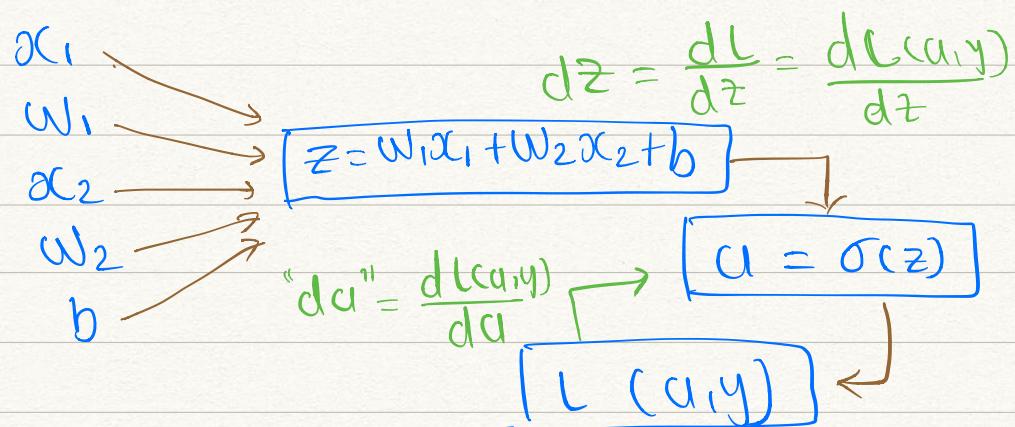
$$\frac{dJ}{du} = 3$$

Logistic Regression Gradient Descent:

$$z = \mathbf{w}^T \mathbf{x} + b$$

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

computational graph:



$$da = \left(\frac{-y}{a}\right) + \left(\frac{1-y}{1-a}\right)$$

$$dz = \frac{dL(a,y)}{dz} = \frac{dL}{da} \cdot \frac{da}{dz}$$

$$\frac{dL}{dw_1} = dw_1 = x_1 dz \quad dw_2 = x_2 dz$$

$$db = dz$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha b$$

Logistic Regression on m examples:

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

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

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

$$dw_i^{(i)} \rightarrow (x^{(i)}, y^{(i)})$$

Algorithm :

$$J=0, dw_1=0, dw_2=0, db=0$$

$$\left. \begin{array}{l} \text{for } i=1 \text{ to } m \\ \quad z^{(i)} = w^T x^{(i)} + b \\ \quad a^{(i)} = \sigma(z^{(i)}) \\ \quad J = -[y^{(i)} \log(a^{(i)}) + (1-y^{(i)}) \log(1-a^{(i)})] \end{array} \right\}$$

for loop

for loop

$$\left\{ \begin{array}{l} dz^{(i)} = a^{(i)} - y^{(i)} \\ dw_1 + = x_1^{(i)} dz^{(i)} \\ dw_2 + = x_2^{(i)} dz^{(i)} \\ db + = dz^{(i)} \end{array} \right.$$

$\uparrow n=2$

$J | = m$

$dw_1 | = m$

$dw_2 | = m$

$db | = m$

↓
accumulator

$$dw_1 = \frac{\partial J}{\partial w_1}$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

Vectorization: get rid of for loop.

Logistic Regression Cost Function (Adv.)

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

$y=0$ $P(y|x) = 1-\hat{y}$

$$P(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

function with
above boundary

conditions.

$$\log p(y|x) = y \underbrace{\log(\hat{y}) + (1-y)\log(1-\hat{y})}_{L(y, \hat{y})}$$

Cost Function:

$$P(\text{labeled in training set}) = \prod_{i=1}^n p(y^{(i)}|x^{(i)})$$

We have to maximize it.

Put log & maximize it.

$$\log P(\cdot) = \sum_{i=1}^n \underbrace{\log p(y^{(i)}|x^{(i)})}_{-L(\hat{y}^{(i)}, y^{(i)})}$$

Maximum likelihood Estimate:

choose parameter that
maximize $\log P(\cdot)$

$$= -\sum_{i=1}^n L(\hat{y}^{(i)}, y^{(i)})$$

Cost: $J(w, b) = \frac{1}{m} \sum_{i=1}^n L(\hat{y}^{(i)}, y^{(i)})$

\downarrow so get rid of negative sign.

$\frac{1}{m} \rightarrow$ for better scale

By minimizing cost function, we are carrying out max. likely hood estimation over logistic Reg. model under the assumption that our training examples are IID (Identical, Randomly distri.)

Vectorization is useful in getting rid of for loop to run algorithm faster.

$$z = w^T x + b \quad w = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix}_{R^{n_x}} \quad x = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix}_{R^{n_x}}$$

Non Vectorized:

```
z=0  
for i in range(n-x)  
    z+= w(i) * x(i)  
z+=b
```

Vectorized:

$$z = np. \underbrace{\text{dot}(w, x)}_{w^T x} + b \quad (\text{numpy})$$

Avoid for loops whenever possible

$$U = AV$$
$$U_i = \sum_j A_{ij} V_j$$

$$U = np. \text{dot}(A, V)$$

$$U = np. \text{zeros}(cn, 1)$$

for i --

for $j = \dots$
 $v_i = A_{ij} \cdot v_j$

Vectors & Matrix valued Function:

e.g. Need to apply exponential operation, on every element of matrix.

$$v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \quad u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

$u = \text{np.zeros}((n, 1))$
for i in range(n):
 $u_{(i)} = \text{math.exp}(v_{(i)})$

Using Numpy:

```
import numpy as np.  
u = np.exp(v)
```

e.g. np.log, maximum
abs

Logistic Regression Derivatives:

$$d\mathbf{w} = \text{np. zeros } ((n-x, 1))$$

$$d\mathbf{w}^+ = \mathbf{x}^{(i)} dz^{(i)}$$

$$d\mathbf{w}^- = m$$

Vectorizing Logistic Regression :

$$\begin{aligned} z^{(1)} &= \mathbf{w}^T \mathbf{x}^{(1)} + b & z^{(2)} &= \mathbf{w}^T \mathbf{x}^{(2)} + b & z^{(3)} &= \mathbf{w}^T \mathbf{x}^{(3)} + b \\ u^{(1)} &= \sigma(z^{(1)}) & u^{(2)} &= \sigma(z^{(2)}) & u^{(3)} &= \sigma(z^{(3)}) \end{aligned}$$

$$\mathbf{z} = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = \mathbf{w}^T \mathbf{x} + [b \ b \ b \ \dots \ b] \quad 1 \times m$$

$$\mathbf{z} = \text{np. dot}(\mathbf{w}^T \mathbf{x}) + b$$

Broadcasting  real number 1×1

$$\mathbf{A} = [u^{(1)} \ u^{(2)} \ \dots \ u^{(m)}] = \sigma(\mathbf{z})$$

Vectorising logistic Reg.'s Gradient Output:

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = \dots$$

$$\begin{aligned} dz &= [dz^{(1)} \quad dz^{(2)} \quad \dots \quad dz^{(m)}] \\ A &= [a^{(1)} \quad a^{(2)} \quad \dots \quad a^{(m)}] \\ y &= [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}] \end{aligned}$$

$$dz = A - y$$

$$\begin{aligned} dW &= 0 \\ dW^+ &= x^{(1)} dz^{(1)} \\ dW^+ &= x^{(2)} dz^{(2)} \\ &\vdots \\ dW^+ &= m \end{aligned}$$

$$\begin{aligned} db &= 0 \\ db^+ &= dz^{(1)} \\ db^+ &= dz^{(2)} \\ &\vdots \\ db^+ &= m \end{aligned}$$

$$dW = \frac{1}{m} \times dz^T$$

$$db = \frac{1}{m} \sum_{i=1}^m dz_{(i)}$$

$$= \frac{1}{m} \begin{bmatrix} x_1 & x_2 & \dots & x_m \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix} = \frac{1}{m} \text{ np.sum}(dz)$$

Summing all up:

for i in range(10000):

$$\begin{aligned}z &= w^T x + b \\&= \text{np. dot}(w^T x) + b\end{aligned}$$

$$A = \sigma(z)$$

$$\begin{aligned}dz &= A - Y \\dw &= \frac{1}{m} \times dz^T\end{aligned}$$

$$db = \frac{1}{m} \text{ np. sum}(dz)$$

$$w := w - \alpha dw$$

$$b := b - \alpha db$$

Makes your code run faster.

e.g.

| | Apple | Beef | Eggs | Potatoes | |
|-------|-------|------|------|----------|--|
| carb | 56 | 0 | 4.4 | 68 | |
| prot. | 1.2 | 104 | 52 | 8 | |
| Fat | 1.8 | 135 | 99 | 0.9 | |

Calculate % of calories from C, P & F
for each of 4.

$$\Rightarrow \text{Cal} = A \cdot \text{sum}(A, \text{axis}=0)$$

$$\underline{\text{Percentage}} = 100 * A / (\text{cal} \cdot \text{reshape}(1, 4))$$

no need to use it. Redundant

e.g. $\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \sim \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$

$$\begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix} + \begin{pmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{pmatrix}$$

↓ Python will convert it to

$$\begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$$

||

$$\begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

General Principle:

$$\begin{array}{ccc} (m,n) & + & (1,n) \longrightarrow (m,n) \\ \text{matrix} & \times & / \\ & & (m,1) \longrightarrow (m,n) \end{array}$$

Operate element wise.

→ Sum values Vertically : axis = 0
Operation Horizontally : axis = 1

Tips:

Never use rank 1 array.

`a = np.random.randn(5)`

`a.shape = (5,)`

rank 1 array

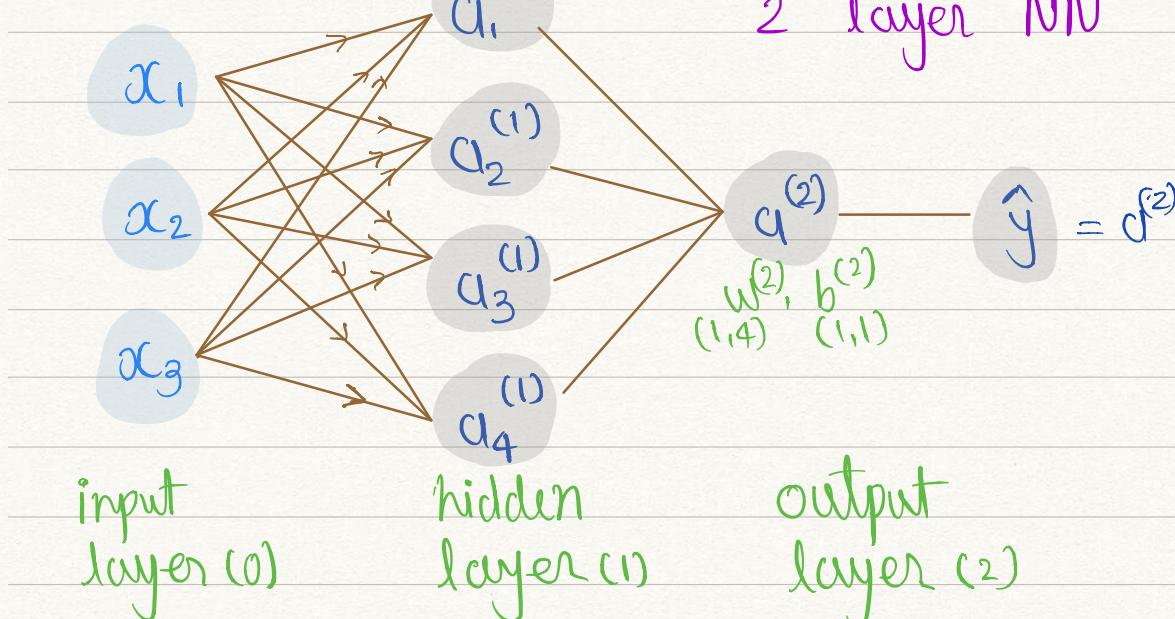
use properly defined array:

`a = np.random.randn(5,1)`

`a = np.random.randn(1,5)`

Neural Network Representation:

$$u^{(0)} = x$$

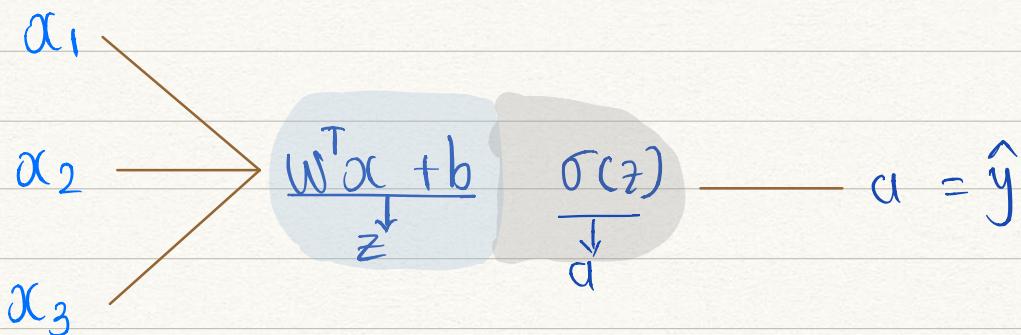


$$u^{(1)} = \begin{bmatrix} c_1^{(1)} \\ c_2^{(1)} \\ \vdots \\ c_4^{(1)} \end{bmatrix}$$

$\hat{y} = u^{(2)}$ because in logistic reg. we only get single o/p from N.N.

$w^{(1)} \rightarrow (4 \times 3) \rightarrow 4 \text{ output unit}$
 3 input unit

Computing NN's Output:



$a^{(l)}$ \rightarrow layer number

a_i \rightarrow node in layer

$$\begin{aligned} z_1^{(1)} &= W_1^{(1)T} x + b_1^{(1)} \\ a_1^{(1)} &= \sigma(z_1^{(1)}) \end{aligned} \quad \left. \begin{array}{l} \text{layer 1, node 1} \\ \text{computation} \end{array} \right.$$

$$\begin{aligned} z_2^{(1)} &= W_2^{(1)T} x + b_2^{(1)} \\ a_2^{(1)} &= \sigma(z_2^{(1)}) \end{aligned} \quad \left. \begin{array}{l} \text{layer 1, node 2} \\ \text{computation} \end{array} \right.$$

Vectorize the whole process:

$$\left[\begin{array}{c} -W_1^{(1)T} - \\ -W_2^{(1)T} - \\ -W_3^{(1)T} - \\ -W_4^{(1)T} - \end{array} \right] \left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right] + \left[\begin{array}{c} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{array} \right] = \left[\begin{array}{c} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{array} \right] = \underline{z^{(1)}}$$

$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ \vdots \\ a_4^{(1)} \end{bmatrix} = \sigma(z^{(1)})$$

given input x :

$$z^{(1)}_{(4 \times 1)} = w^{(1)}_{(4 \times 3)} x_{(3 \times 1)} + b^{(1)}_{(4,1)}$$

$$a^{(1)}_{(4,1)} = \sigma(z^{(1)})_{(4,1)}$$

Algorithm:

for $i=1$ to m :

$$\begin{aligned} z^{(1)(i)} &= w^{(1)} \cdot x^{(i)} + b^{(1)} \\ a^{(1)(i)} &= \sigma(z^{(1)(i)}) \\ z^{(2)(i)} &= w^{(2)} a^{(1)(i)} + b^{(2)} \\ a^{(2)(i)} &= \sigma(z^{(2)(i)}) \end{aligned}$$

$\left. \begin{array}{l} \text{4 lines} \\ \text{of for} \\ \text{loop} \end{array} \right\}$

$$x = \begin{bmatrix} | x^{(1)} | & | x^{(2)} | & | x^{(3)} | & \dots & | x^{(m)} | \end{bmatrix}^T$$

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

training examples.

$$A^{(1)} = \underbrace{\left[\begin{array}{c|c|c|c} 1 & & & \\ \hline u^{(1)(1)} & u^{(1)(2)} & \cdots & u^{(1)(m)} \end{array} \right]}_{\text{hidden units}}$$

Explanation for vectorized implementation:

$$z^{(1)(1)} = w^{(1)} x^{(1)} + b^{(1)}$$

$$z^{(1)(2)} = w^{(1)} x^{(2)} + b^{(1)}$$

$$w^{(1)} x^{(1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$w^{(1)} x^{(2)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$w^{(1)} \left[\begin{array}{c|c|c} 1 & x^{(1)} & x^{(2)} \\ \hline x^{(1)} & \vdots & \vdots \\ \hline \end{array} \right] = \left[\begin{array}{c|c|c} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{array} \right] = \left[\begin{array}{c|c|c} z^{(1)(1)} & z^{(1)(2)} & \dots \\ \hline +b^{(1)} & +b^{(1)} & \dots \\ \hline z^{(1)} & z^{(1)} & \dots \end{array} \right]$$

↓
Set of all training examples

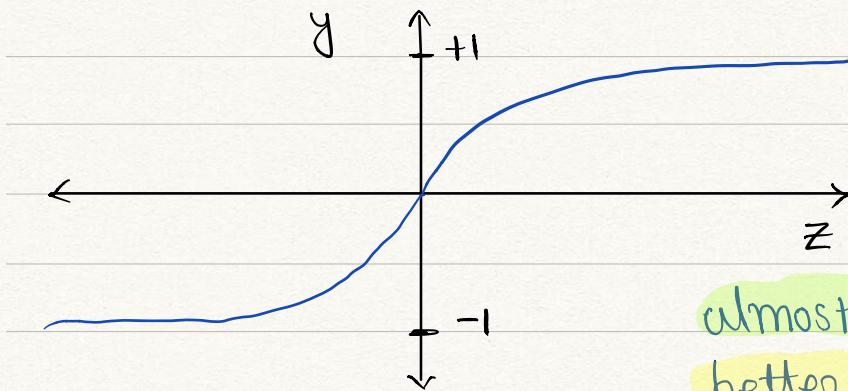
$$\left. \begin{array}{l} z^{(1)} = w^{(1)} X + b^{(1)} \\ A^{(1)} = \sigma(z^{(1)}) \\ z^{(2)} = w^{(2)} A^{(1)} + b^{(2)} \\ A^{(2)} = \sigma(z^{(2)}) \end{array} \right\} \begin{array}{l} \text{vectorization} \\ \text{of 4 lines} \\ \text{in for loop} \end{array}$$

Activation Functions:

Sigmoid Function:

$$\sigma = \frac{1}{1 + e^{-z}}$$

tanh Function:



$$y = \tanh(z)$$

$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

almost always works
better than sigmoid

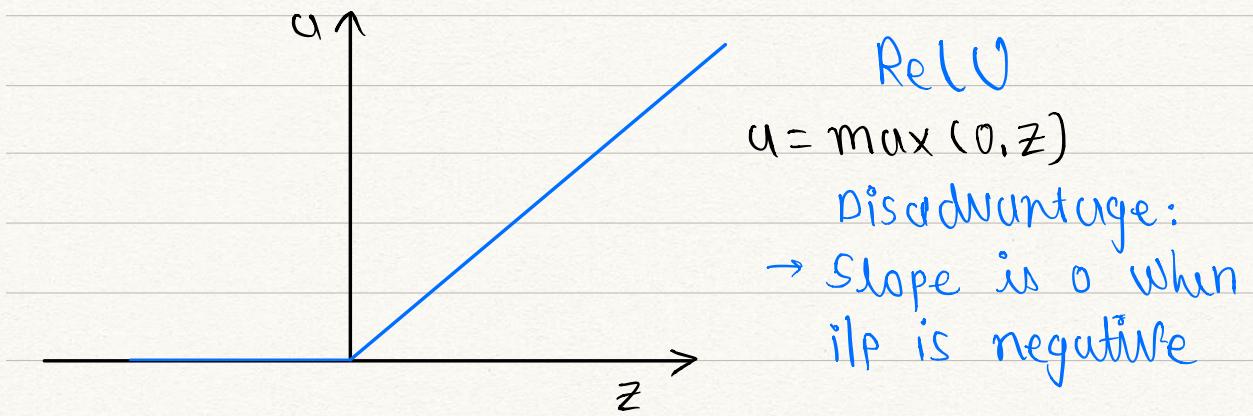
- have close to zero mean
- centering the data
- easier for further calculation

→ when you have binary classification output layer should be sigmoid.
keep hidden layers as tanh.

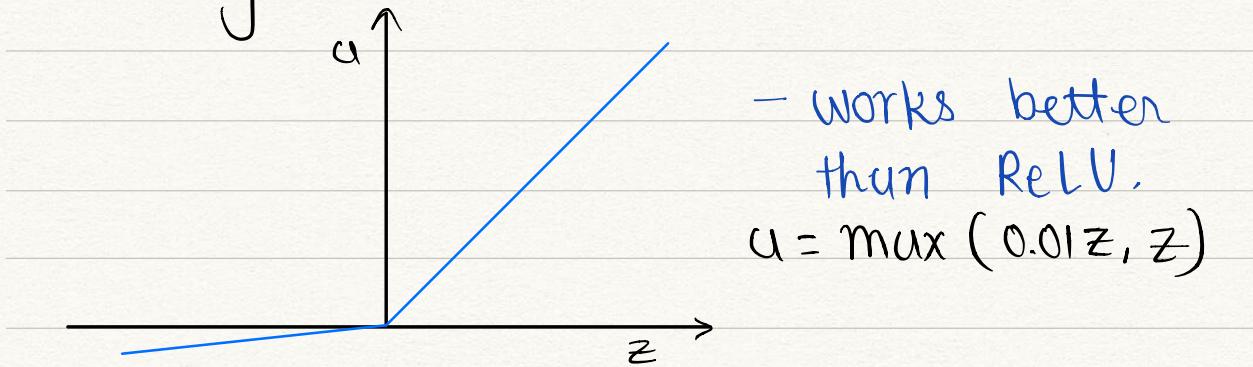
→ You can have different activation function for different layers.

→ When you have very large or very small i/p, slope of function is close to zero, so it can slow down gradient descent.

Rectified Linear Unit :



Leaky ReLU:



Why do we need non linear Activation Functions?

linear activation function:

$$g(z) = z$$

↳ identity function.

$$\begin{aligned} u^{(1)} &= z^{(1)} = w^{(1)}x + b^{(1)} \\ u^{(2)} &= z^{(2)} = w^{(2)}u^{(1)} + b^{(2)} \\ u^{(2)} &= w^{(2)}(w^{(1)}x + b^{(1)}) + b^{(2)} \\ &= \frac{(w^{(2)}w^{(1)})}{w^{(1)}}x + \frac{(w^{(2)}b^{(1)} + b^{(2)})}{w^{(1)}} \end{aligned}$$
$$u^{(2)} = w'x + b'$$

No matter how many layer NN has
its just doing linear activation.
No benefit of doing this.

You can use linear activation on output layer, not recommended at hidden layer.

Derivatives of Activation Functions:

When we implement back prop. on NN
We need to compute slope or the
derivative of activation function.

Sigmoid Activation Function:

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

$$\begin{aligned} g'(z) &= \frac{d}{dz} g(z) = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right) \\ &= g(z) \cdot (1 - g(z)) \\ &= a(1-a) \end{aligned}$$

tanh activation:

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

$$\begin{aligned} g'(z) &= \frac{d}{dz} (g(z)) = 1 - (\tanh(z))^2 \\ &= 1 - (g(z))^2 \\ &= \underline{\underline{1 - a^2}} \end{aligned}$$

ReLU & Leaky ReLU:

$$g(z) = \max(0, z) \rightarrow \text{ReLU}$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

~~undefined " $z=0$ "~~

↳ due to discontinuity.

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & z < 0 \\ 1 & z \geq 0 \end{cases}$$

~~undefined $z=0$~~

Gradient Decent for NNs:

Parameters: $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}$
 $(n^{(1)}, n^{(0)})$ $(n^{(1)}, 1)$ $(n^{(2)}, n^{(1)})$ $(n^{(2)}, 1)$
 $n_x = n^{(0)}, n^{(1)}, n^{(2)} = 1$

Cost Function: $J(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)})$

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

\hat{y} $\uparrow a^{(2)}$

Repeat : { compute predictions:

$$\partial w^{(1)} = (\hat{y}^{(i)}, i=1, 2, \dots, m) \quad \partial b^{(1)} - \partial J$$

$$\text{update: } \begin{aligned} \hat{\mathbf{w}}^{(1)} &= \mathbf{w}^{(1)} - \frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(1)}} \\ \hat{\mathbf{b}}^{(1)} &= \mathbf{b}^{(1)} - \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(1)}} \end{aligned}$$

Formulas for Computing derivatives:

Forward Propagation:

$$\begin{aligned} \mathbf{z}^{(1)} &= \mathbf{w}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \\ \mathbf{A}^{(1)} &= g^{(1)}(\mathbf{z}^{(1)}) \\ \mathbf{z}^{(2)} &= \mathbf{w}^{(2)} \mathbf{A}^{(1)} + \mathbf{b}^{(2)} \\ \mathbf{A}^{(2)} &= g^{(2)}(\mathbf{z}^{(2)}) = \sigma(\mathbf{z}^{(2)}) \end{aligned}$$

Back Propagation:

$$dz^{(2)} = A^{(2)} - y$$

$$y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$d\mathbf{w}^{(2)} = \frac{1}{m} dz^{(2)} \mathbf{A}^{(1)T}$$

$$db^{(2)} = \frac{1}{m} \text{np.sum}(dz^{(2)}, \text{axis}=1, \text{keepdims=T})$$

$$(n^{(2)}, 1)$$

$$dz^{(1)} = \underbrace{\mathbf{w}^{(2)T} dz^{(2)}}_{(n^{(1)}, m)} * \underbrace{g^{(1)'}(\mathbf{z}^{(1)})}_{\text{elementwise product}} (n^{(1)}, m)$$

$$db^{(1)} = \frac{1}{m} \cdot \text{np.sum}(dz^{(1)}, \text{axis}=1, \text{keepdims=T})$$

$$d\mathbf{w}^{(1)} = \frac{1}{m} d\mathbf{z}^{(1)} \mathbf{x}^T$$

Backpropagation Intuition:

Computing Gradient:

$$\begin{aligned} X & \xrightarrow{\mathbf{W}} z = \mathbf{w}^T X + b \\ \mathbf{w}, b & \quad dz \\ dw, db & \quad \frac{dL(a,y)}{da} \\ &= \frac{d}{da} (-y \log a - (1-y) \log (1-a)) \end{aligned}$$

$$da = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\begin{aligned} dz &= da \cdot g'(z) \xrightarrow{g(z) \rightarrow \text{sigmoid}} \frac{1}{\sigma(z)} \\ a &= \sigma(z) \\ \frac{\partial L}{\partial z} &= \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \xrightarrow{\frac{d}{dz}(\sigma(z))} \end{aligned}$$

$$dw = dz \cdot x$$

$$db = dz$$

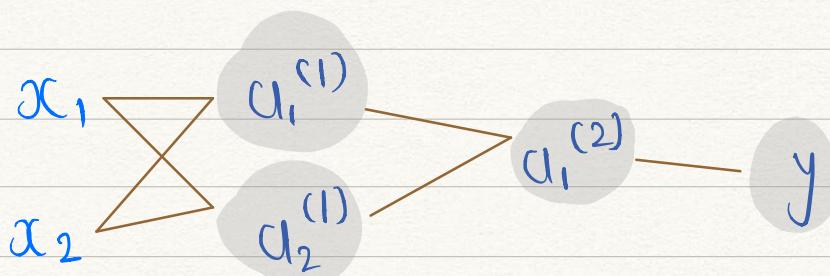
Diagram illustrating a two-layer neural network. The forward pass consists of two layers:
 $z^{(1)} = w^{(1)T}x + b^{(1)}$ and $a^{(1)} = \sigma(z^{(1)})$. The backward pass (backpropagation) consists of two layers:
 $z^{(2)} = w^{(2)T}a^{(1)} + b^{(2)}$, $a^{(2)} = \sigma(z^{(2)})$, and $da^{(2)}$.
 The error function is $L(a^{(2)}, y)$.
 The diagram shows the flow of gradients from the output layer back through the hidden layer, involving terms $dz^{(2)}$, $w^{(2)}$, $b^{(2)}$, and $da^{(2)}$.
 The total gradient for the weight update is $dw^{(2)} = dz^{(2)} \cdot a^{(1)T}$.
 Below the diagram, the element-wise product rule for the gradient calculation is given as:

$$d z^{(1)} = \frac{(n^{(1)}, n^{(2)})}{(n^{(2)}, 1)} \left(\frac{d z^{(1)}}{n^{(2), 1}} \right) * \underbrace{g^{(1)'}(z^{(1)})}_{\text{element wise product}}$$

$$= \frac{(n^{(2)})}{(n^{(2)}, n^{(1)})} \left(\frac{d z^{(1)}}{n^{(2), 1}} \right) * g^{(1)'}(z^{(1)})$$

Random Initialization:

What happens if you initialize weights to zero?



$$w^{(1)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{(1)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$u_1^{(1)} = u_2^{(1)}$$

$$dz^{(1)} = dz^{(2)}$$

$$dw = \begin{pmatrix} u & v \\ u & v \end{pmatrix}$$

$$W^{(1)} = W^{(1)} - 2dw$$

$$W^{(1)} = \begin{bmatrix} - & - & - & - \\ - & - & - & - \end{bmatrix}$$

identical rows.

If you initialize with zero, hidden layer will be symmetric and we will keep on computing same function, regardless of no. of layers.

So we do...

$$W^{(1)} = \text{np.random.randn}(2, 2) * 0.01$$

$$b^{(1)} = \text{np.zeros}((2, 1))$$

multiply with small no.

we do not have symmetry breaking problem

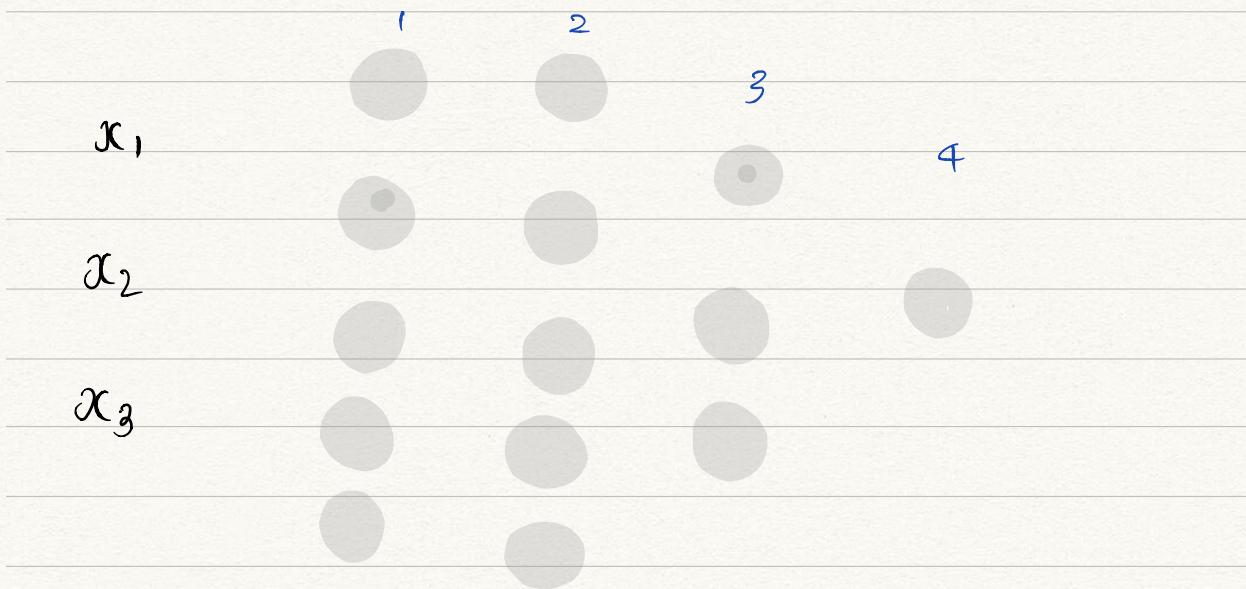
$$W^{(2)} = \text{random}$$

$$b^{(2)} = 0$$

→ If w values are too large , z will be very big or very small. in which case we will not get good gradient decent because it will have dw close to zero.
(end up at flat part of activation function)

Deep L-layer Neural Network

Total layers = hidden layers + output layers



$L = 4$ (No. of layers)

$n^{(l)}$ = No. of units in layer l

$$n^{(1)} = 5 \quad n^{(2)} = 5 \quad n^{(3)} = 3 \quad n^{(L)} = 1$$
$$n^{(0)} = n_x = 3$$

$a^{(l)} = g^{(l)}(z^{(l)}) \rightarrow$ activation in layer l
 $w^{(l)} = \text{weights for } z^{(l)}$

Forward Propagation in Deep Net.

General equation:

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$

$$a^{(l)} = g^{(l)}(z^{(l)})$$

Vectorized:

$$\begin{aligned} z^{(1)} &= w^{(1)} A^{(0)} + b^{(1)} \\ A^{(1)} &= g^{(1)}(z^{(1)}) \end{aligned}$$

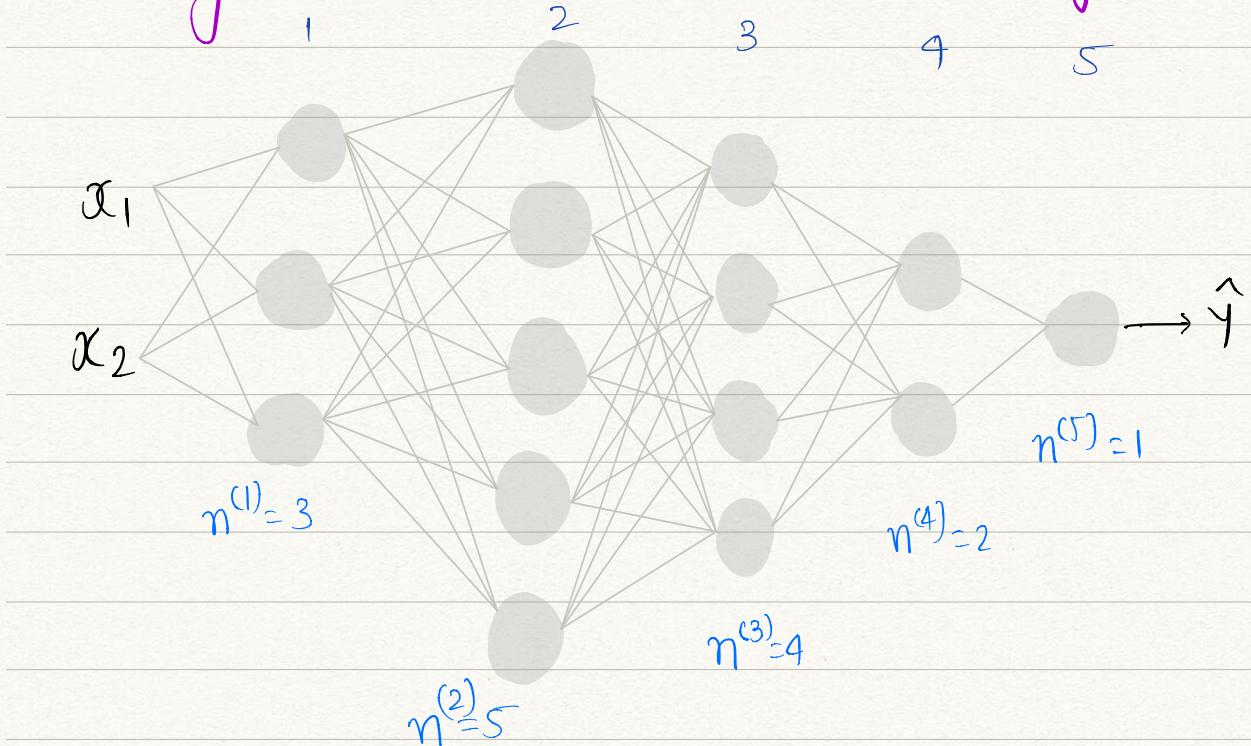
for loop

$$\begin{aligned} z^{(2)} &= w^{(2)} A^{(1)} + b^{(2)} \\ A^{(2)} &= g^{(2)}(z^{(2)}) \end{aligned}$$

$0 \rightarrow l$

$$Q = g^{(4)}(z^{(4)}) = A^{(4)}$$

Getting Your Matrix dimension right:



$$\vec{z}^{(1)} = W^{(1)} \cdot x + b^{(1)}$$

| | | | |
|----------------|----------------------|----------------|----------------|
| $(3, 1)$ | $(3, 2)$ | $(2, 1)$ | $(3, 1)$ |
| $(n^{(1)}, 1)$ | $(n^{(1)}, n^{(0)})$ | $(n^{(0)}, 1)$ | $(n^{(1)}, 1)$ |

$$W^{(l)} = (n^{(l)}, n^{(l-1)})$$

$$b^{(l)} = (n^{(l)}, 1)$$

$$W^{(1)} = (n^{(1)}, n^{(0)}) = (3, 2)$$

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

$$W^{(2)} = (5, 3)$$

$$b^{(2)} = (5, 1)$$

$$W^{(3)} = (4, 5)$$

$$b^{(3)} = (4, 1)$$

$$W^{(4)} = (2, 4)$$

$$b^{(4)} = (2, 1)$$

$$W^{(5)} = (1, 2)$$

$$b^{(5)} = (1, 1)$$

$$dW^{(l)} = (n^{(l)}, n^{(l-1)}) \cdot db^{(l)} = (n^{(l)}, 1)$$

Vectorized Implementation:

$$z^{(l)} = W^{(l)} \cdot x + b^{(l)}$$

$(n^{(l)}, m)$ \downarrow $(n^{(l)}, n^{(0)})$ \downarrow $(n^{(0)}, m)$ \nearrow $(n^{(0)}, m)$

$$z^{(l)}, A^{(l)} = (n^{(l)}, m) = dz^{(l)}, dA^{(l)}$$

Why deep Representation?

- As we proceed deep into layers, they get more complex.
- Initial hidden layers compute simple functions & as they move deeper they combine previous results & get more complex.

e.g. Audio → low levels $\xrightarrow{\text{compose}}$ Phonems
audio waveform }

Sentences $\xleftarrow{\text{compose}}$ Words $\xleftarrow{\text{compose}}$

Circuit Theory & Deep learning:

Informally:

There are functions you can compute with a "small" L-layer deep neural net that shallower nets require exponentially more hidden units to compute.

Building Blocks of deep Neural Nets:

Forward & Backward functions:
layer l : $w^{(l)}, b^{(l)}$

Forward :

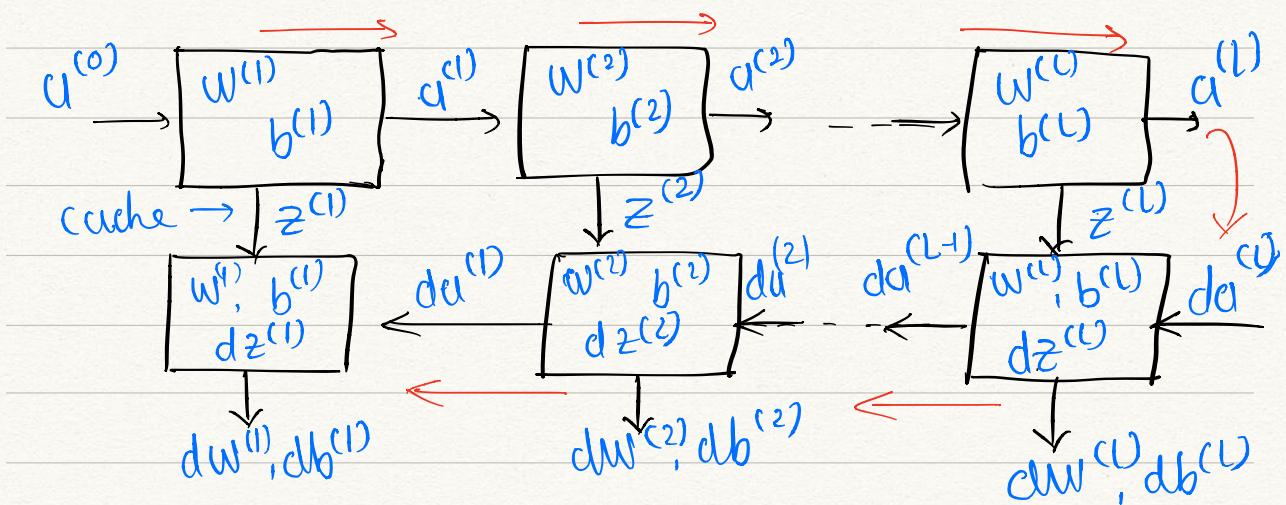
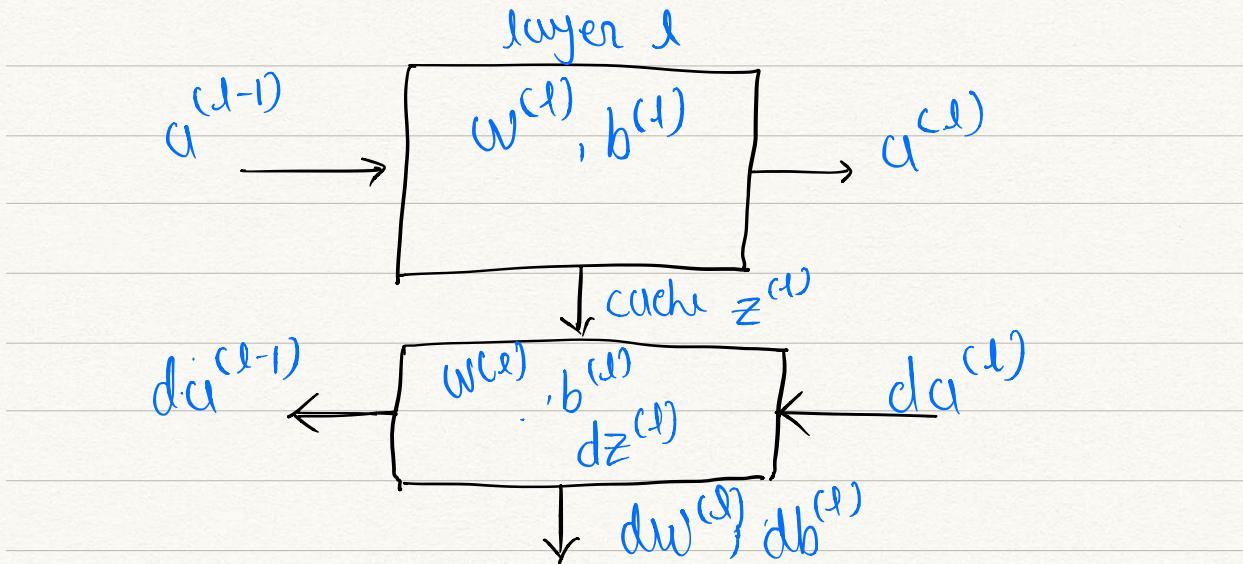
input: $a^{(l-1)}$ o/p: $a^{(l)}$

$$\begin{aligned} z^{(l)} &= w^{(l)} a^{(l-1)} + b^{(l)} \\ a^{(l)} &= g^{(l)}(z^{(l)}) \end{aligned}$$

cache: $z^{(l)}$

Backward:

inp: $da^{(l)}$ o/p: $da^{(l-1)}, dw^{(l)}, db^{(l)}$
cache: $z^{(l)}$



Backward Propagation for layer l :

Input : $da^{(l)}$
 Output : $da^{(l-1)}, dW^{(l)}, db^{(l)}$

$$dz^{(l)} = da^{(l)} * g^{(l)'}(z^{(l)})$$

$$\begin{aligned} dW^{(l)} &= dz^{(l)} \cdot a^{(l-1)} \\ db^{(l)} &= dz^{(l)} \end{aligned}$$

$$da^{(l-1)} = W^{(l)}^T \cdot dz^{(l)}$$

$$dz^{(l)} = W^{(l+1)T} \cdot dz^{(l+1)} * g^{(l)\prime}(z^{(l)})$$

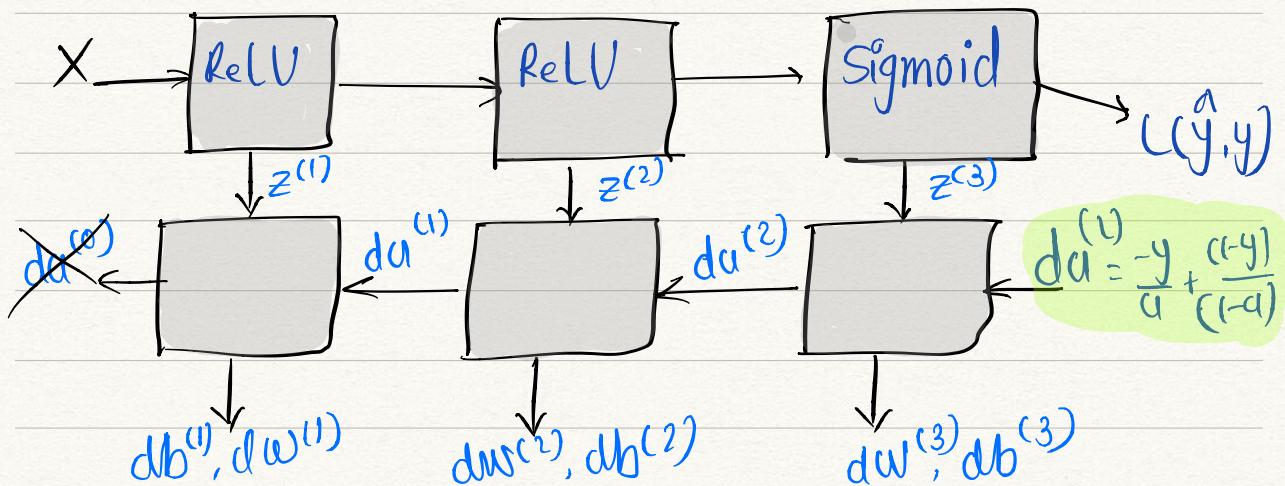
Vectorized:

$$dz^{(l)} = dA^{(l)} * g^{(l)\prime}(z^{(l)})$$

$$dW^{(l)} = \frac{1}{m} dz^{(l)} \cdot A^{(l-1)T}$$

$$db^{(l)} = \frac{1}{m} \text{np.sum}(dz^{(l)}, \text{axis}=1, \text{keepdims}=T)$$

$$dA^{(l-1)} = W^{(l)T} \cdot dz^{(l)}$$



Parameters Vs. Hyperparameters:

Parameters: $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots$

Hyperparameters: Learning Rate (α)

of iterations of gradient descent to carry out.

of hidden layer

of hidden units

Activation function choice

Other: Momentum, mini batch size,
Regularization

Applied deep learning is a very empirical process.

idea \rightarrow code \rightarrow experiment



What does deep learning have to do with brain:

→ Not a lot...



Very simple explanation.

Quiz

Q:1 ~~X~~ cache intermediate value forward to backward

Q:2 all but w^l, b^l, a^l

Q:3 Deeper layer compute

more complex

Q:4

False

Q:5

in (layer_dims)

w: (i, i-1)

Q:6

layer : 4

hidden : 3

Q:7

False ~~X~~ True

Q:8

True

Q:9

$w_1 = (4 \times 4) \quad w^2 = (3, 4) \quad w^3 = (1, 3)$
 $b^1 = (4, 1) \quad b^2 = (3, 1) \quad b^3 = (1, 1)$

Q:10

(n^1, n^{1-1})