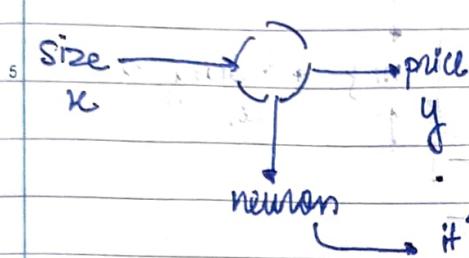


Introduction to neural networks



it just acts
as a func

ReLU

Rectified linear units

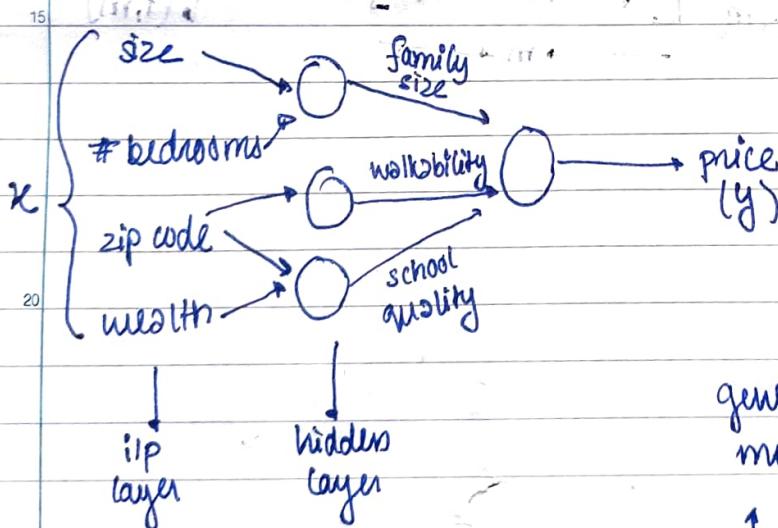
mean taking
a max of 0

hence

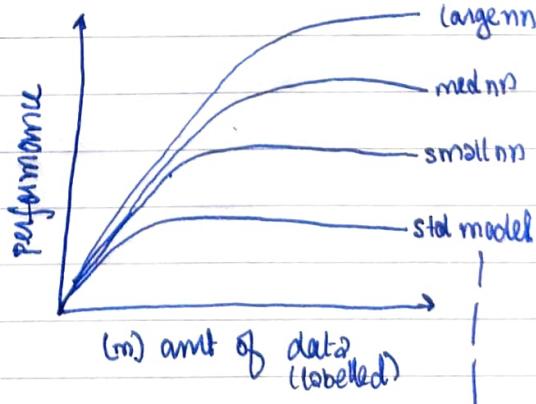
this is what

ReLU func looks like

10 a neural network is built by stacking multiple neurons w/ diff features as i/p e.g.



generally w/ traditional models,



like normal
supervised model
w/o neural network

Note:

Traditionally we write \rightarrow vertical vectors.

$$x = \begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^m \end{bmatrix} \quad m \text{ f. } \rightarrow \text{training ex. no.}$$

$y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^m \end{bmatrix} \quad (m, 1) \quad \rightarrow n_x \leftarrow \text{feature no.}$

but in neural network if we do horizontal stacking computations becomes easier

$$x = \begin{bmatrix} x^1 & x^2 & \cdots & x^m \end{bmatrix} \quad m \rightarrow \text{feat. no.}$$

$y = [y^1 \ y^2 \ \dots \ y^m] \quad (1, m)$

$x_{\text{shape}} = (n_x, m)$

horizontal representation

Date: _____

H_ologistic Regression as Neural Network

Generally,

given,

$x \in R^m$, then $\hat{y} = p(y=1 | x)$

Let say we has 64×64 px

$$\Rightarrow x = \underbrace{64 \times 64}_{\text{dim}} \text{ vektor}$$

$$(x^1, x^2, \dots, x^m)$$

\vec{r}_x dimension
vector

x^1, x^2

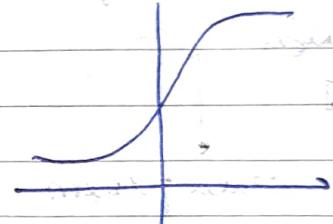
each
is a
feature.

param: $w \in \mathbb{R}^n$

$$b \in \mathbb{R}$$

then we know,

$$y = \text{Sigmoid}(w^T x + b)$$



sometime it can also be denoted as,

$$n_0 = 1, \quad g = p(y=1|x)$$

when θ is the weight vector = $\theta^T x$

1960 Jones established & began to use the Williams and

but this makes it clear that $b = 60^\circ$

implementing nn
a tad bit more
difficult hence

$$\Theta = \begin{bmatrix} 0.3 \\ \theta' \\ \vdots \\ \theta^{n_x} \end{bmatrix} \rightarrow b = \Theta^0$$

$w = \theta^1, \dots, \theta^{n_r}$

we use prev notation
keeping w & b separate

Cost function

for creating loss func. we could,

its a convex loss func

$$L(g, y) = \frac{1}{2} (g - y)^2 \rightarrow \text{but this causes a convex optimization}$$

* used ~~target~~ ~~easy~~ in convex landscape

non-convex model is linear func of its params

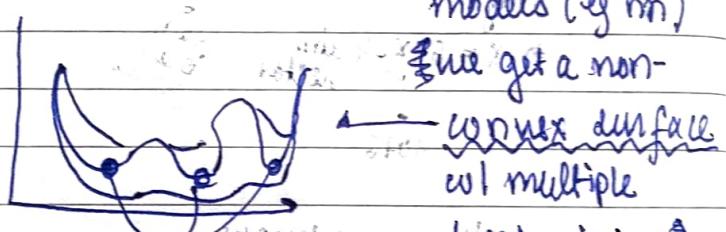
but useful only if

causes problem in non-linear param func

local minima's

→ problem not because of MSE

but bcz of presence of multiple local minima



hence creating a non-convex optimization problem

→ MSE is used in non-convex application even if it might cause problem because its simple & reliable and can give pretty good values even if we tweak the model a bit

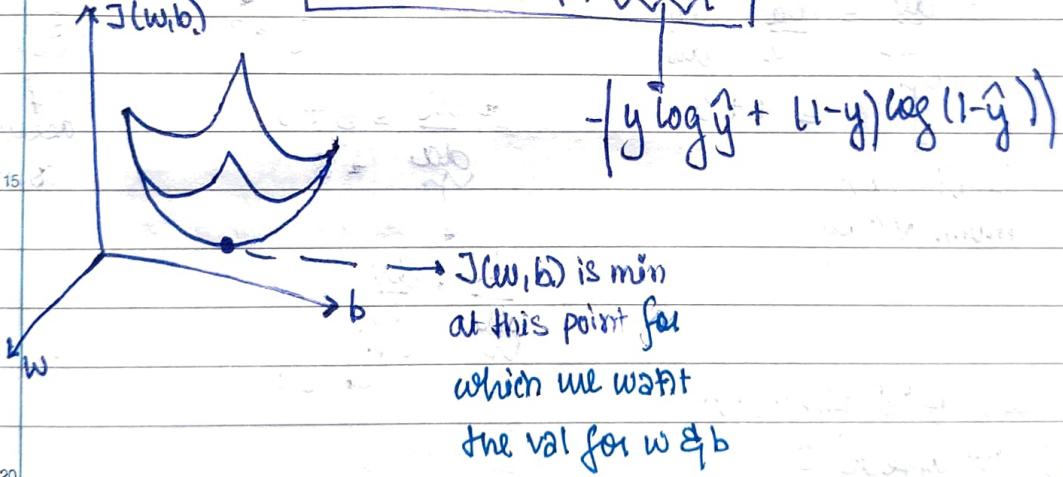
but MSE not used widely,
we use

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

done w/ optimization

the loss func was defined unit single training eg
 cost func measures how well we doing in entire dataset

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m h(\hat{y}_i, y_i)$$



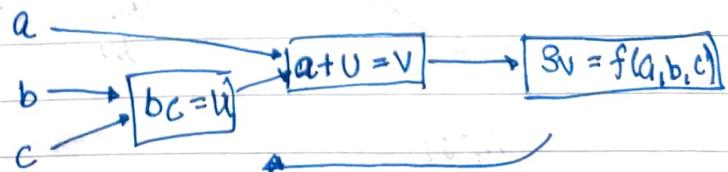
Computation graphs

let's say we have,

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

$$\begin{aligned} bc &= v \\ a+v &= a+bc = v \end{aligned}$$

$$3v = f(a, b, c)$$



if we go back we get a derivative

lets do some back propagation,

$$J = 3v$$

$$\Rightarrow \frac{dJ}{dv} = 3 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{1 step back}$$

now if we increase 'a' how does that affect J,

$$a = 5 \rightarrow 5.001$$

$$v = 11 \rightarrow 11.001$$

$\left. \begin{array}{l} \\ \end{array} \right\}$ from here we see

$$\text{initial } J \Rightarrow J = 3v \rightarrow 33 \rightarrow 33.003 \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{now } a \uparrow \Rightarrow J \uparrow \times 3$$

$$\Rightarrow \frac{dJ}{da} = 3$$

$$\frac{dJ}{da} = \frac{da}{da} + \frac{du}{da}$$

$$J = 3v = 3a + 3v$$

$\left. \begin{array}{l} \\ \end{array} \right\}$
if $a \uparrow$ how
much $v \uparrow$ by

$$\Rightarrow \frac{dJ}{da} = 3 + 0$$

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

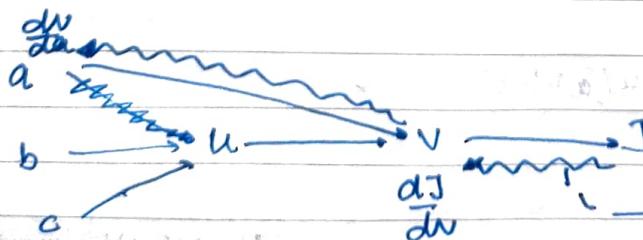
assumed
cost.

used if 'a' affects 'v'

or 'v' affects J

$$\rightarrow a \rightarrow v \rightarrow J$$

so till now we have done



→ back propagation

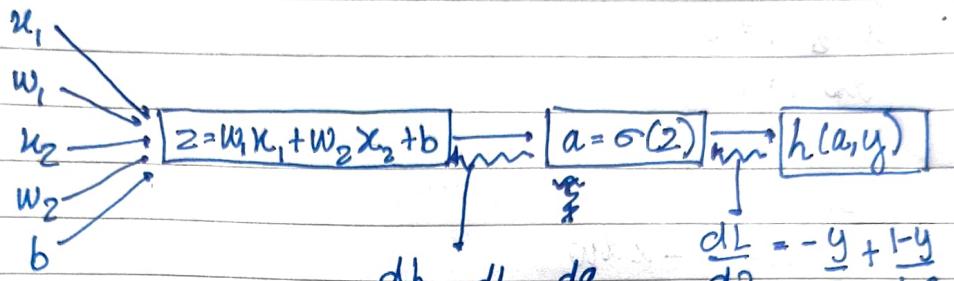
likewise,

$$\frac{dv}{du} = 3 \quad \text{if } \frac{dJ}{dv} = 3 \quad \Rightarrow \frac{dJ}{du} = 3$$

then,

$$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db} \quad \Rightarrow \frac{dJ}{db} = 3 \cdot 3c \quad \text{if } \frac{dJ}{dc} = 3b$$

Now using derivatives to back propagate on logistic reg.



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

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

$$= a - y$$

then

$$\frac{dh}{dw_i} = \frac{\partial h}{\partial z} \cdot \frac{\partial z}{\partial w_i}$$

$$= x_i \frac{dh}{dz}$$

then updates for gradient descent becomes,

$$w_i = w_i - \alpha \frac{dh}{dw_i}$$

Gradient for m examples

$$J=0, dw_i^T = 0, db=0$$

for i in m:

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

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

$$J = -(y^{(i)} \log \hat{y}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_i^T = x_i^T dz^{(i)}$$

$$dw_i^T = x_n^i dz^{(i)}$$

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

$$\begin{aligned} J &= m \\ \frac{\partial J}{\partial w_1} &= \frac{\partial w_1}{\partial w_1} J = m \\ \frac{\partial J}{\partial w_2} &= \frac{\partial w_2}{\partial w_2} J = m \\ \frac{\partial J}{\partial b} &= \frac{\partial b}{\partial b} J = m \end{aligned}$$

then,

$$w = w - \alpha \frac{\partial J}{\partial w}$$

$$b = b - \alpha \frac{\partial J}{\partial b}$$

if there are multiple feature then we would have to iterate over them but this causes our algorithms to run slowly
hence,

we use vectorization

before for.

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

we did,

for i in range () :

$$z += w[i] * x[i]$$

~~z += b~~ ~~for each iteration we do this~~

vectorized format do,

$$z = np.dot(w, x) + b$$

Note:

both

GPU } \rightarrow have SIMD instru.
CPU }

single ins multi data \rightarrow basically parallelism

more used by vectorization

\Rightarrow lesser time taken

Hist function

lets say,

$$\text{if } y=1 \therefore P(y|x) = \hat{y}$$

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

} only 2 $(y=0, y=1)$::
binary classification

these 2 eqn
can be written as

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

$$\text{at } y=1 \quad \hat{y}^1 \cdot (1-\hat{y})^{(1-1)}$$

$$= \hat{y}$$

$$\text{so at } y=0 \quad \hat{y}^0 \cdot (1-\hat{y})^{(1-0)}$$

$$= 1-\hat{y}$$

same
as before

Now,

taking log

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

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

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

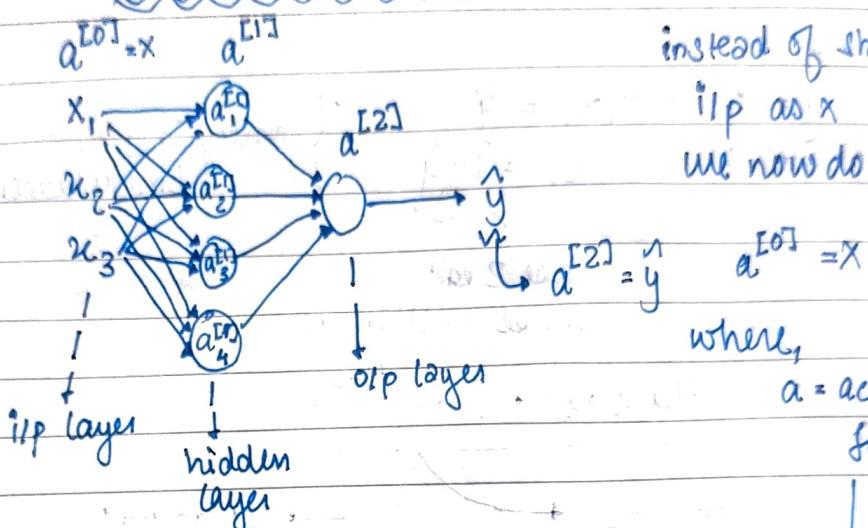
$$\Rightarrow \downarrow L(\hat{y}, y) \rightarrow \uparrow \log(P(y|x))$$

$a_i^{[l]}$ where,
 i = layer no.
 i = node in layer

Date : _____

in counting
layer
we don't
count
i/p layers

Neural Network



instead of showing
i/p as x
we now do,

where,
 a = activation
func

here,

$a^{[1]}$ is 4 dim vector
 \therefore it has 4 diff nodes in
its layer

refer to val the
current layer is
passing to subsequent
layer

The hidden layer has params attached to it,

$$\begin{matrix} w^{[1]} & \& b^{[1]} \\ \uparrow & & \downarrow \\ (4, 3) \text{ mat} & & (4, 1) \text{ vect} \end{matrix}$$

3 \because 3 i/p

4 \because 4 nodes in the layer

The o/p layer also has params,

$$\begin{matrix} w^{[2]} & \& b^{[2]} \\ \uparrow & & \downarrow \\ (1, 4) & & (1, 1) \end{matrix}$$

4 \because hidden
layer has
4 nodes

1 \because 1 o/p node

for o/p
node, the
hidden
node
acts as i/p

Let's take a look at a single node at the hidden layer
div into two parts

$$z_1^{[1]} = w_1^T x + b_1^{[1]} \quad \text{---> Node} \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

doing this for all layers w1 for loop
is very inefficient
instead we would

$$\begin{bmatrix} w_1^{[1]} \\ w_2^{[1]} \\ w_3^{[1]} \\ w_4^{[1]} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$w^{[1]}$
 $(4,3)$

$$\therefore z^{[1]} = w^{[1]} x + b^{[1]}$$

$(4,1) \quad (4,3) \quad (3,1) \quad (4,1)$

in case of multiple training examples

$a_2^{[2]}[i] = g^{[2]}$
layer 2
example no.

number
of training eg were
stacked

for $i=1$ to m :

$$z^{[1]}[i] = w^{[1]} x^{[i]} + b^{[1]}$$

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

$$z^{[2]}[i] = \dots$$

$$x = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \rightarrow (n_x, m)$$

Why stacking training examples vertically is good ??

let see,

$$w^{[1]} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad \& \quad x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, \quad x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

then according to matrix mul

from

$$w^{[1]} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad w^{[1]} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

& if we do

$$w^{[1]} x = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 2 & 2 \\ 1 & 1 & 1 \end{bmatrix}$$

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

$$\text{now if we add } b^{[1]} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

so above then python using broadcasting adds $b^{[1]}$ to each col

Hence,

vectorizing in column help makes calc easier & faster

$$\text{i.e. } z^{[1]} = w^{[1]} x + b^{[1]} \rightarrow A^{[0]}$$

helps vectorize

for $i=1$ to m :

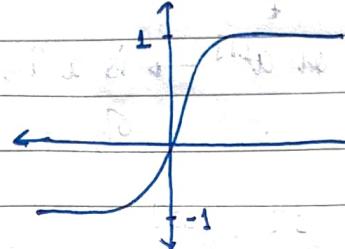
$$z^{[1](i)} = w^{[1]} x^{(i)} + b^{[1]}$$

Activation Functions

Sigmoid func might not be best.

tanh → mainly use for hidden layer

* Sigmoid → o/p layer
to get $y \in \{0, 1\}$
rather than b/w $\{-1, 1\}$
in binary classification



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

better: data centered
→ mean about 0 making learning for next layer easy

downside of both → 2 is \gg or \ll

then gradient / derivative of slope ≈ 0
→ slows down gradient descent

Instead what's very popular is

ReLU

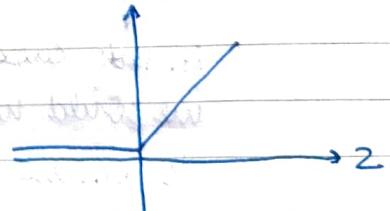
$$a = \max(0, z)$$

disadv ← derivative / slope = 0

when $z = \text{ne}$

slope = 1

when $z = \text{pos}$



hence

leaky ReLU

$$a = \max(0.01z)$$



Why do neural network need Non-linear Activation function?

If in forward prop we did.

$$a^{[1]} = z^{[1]} \text{ instead of } a^{[1]} = g^{[1]}(z^{[1]})$$

g or $a^{[1]}$ → is a linear function of i/p

↓

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = z^{[2]}$$

where

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

doing so makes the model no more complex & expressive than a normal logistic func

linear activation func,

can be used in regression problem where

eg. housing prices

here too it's only used in the o/p layer
& hidden layers use non-linear funcs

in this case too

we could use ReLu

to ensure $\hat{y} \geq 0$

Back propagation

lets say for sigmoid func, $g(z) = \frac{1}{1+e^{-z}}$

we find derivative,

$$\frac{d}{dz} g(z) = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right)$$

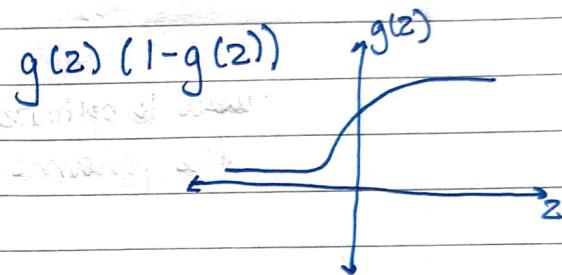
i.e

$$g'(z) = g(z)(1-g(z))$$

lets do a sanity check

say $z \gg$

$$\Rightarrow g(z) \approx 1$$



$$\Rightarrow \frac{d}{dz}(g(z)) = 1(1-1) = 0 \Rightarrow \text{correct}$$

so,

$$g(z) = \tanh z$$

$$\Rightarrow g'(z) = 1 - (\tanh z)^2$$

for ReLU,

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

for leaky ReLU

$$\Rightarrow g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undefined} & \text{if } z = 0 \end{cases} \quad \xrightarrow{\text{---}} 0.01 \text{ if } z < 0$$

in a nn both forward & back propagation happen during training

→ uses iip to make prediction

→ forward prop

→ using loss calc to compute error

diff b/wt w & b

to get gradient

back propagation

used to optimize
the params

uses chain rule

to calc gradients of the
non-linear activation func

Gradient Descent

Params : $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

$$n_x = n^{[0]}, n^{[1]}, n^{[2]} = 1$$

iip feature

hidden units

o/p units

nodes

then,

$$w^{[1]} = (n^{[1]}, n^{[0]})$$

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

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

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

$$\text{if cost func} = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i)$$

for neural networks,

it's better to initialize params randomly rather than = 0

\therefore gradient descent \Rightarrow repeat q

compute $g^{(i)}$, $i=1,2\dots,m$

$$dw^{[1]} = \frac{dJ}{dw^{[1]}}$$

$$db^{[1]} = \frac{dJ}{db^{[1]}}$$

$$w^{[1]} = w^{[1]} - \alpha dw^{[1]}$$

$$b^{[2]} = b^{[1]} - \alpha db^{[1]}$$

$$w^{[2]} =$$

$$b^{[2]} =$$

till convergence

Forward

Back

$$z^{[1]} = w^{[1]}X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(\sum z^{[1]})$$

$$z^{[2]} = w^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(z^{[2]})$$

$$dz^{[2]} = A^{[2]} - y$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

stacked horizontally

$$db^{[2]} = \frac{1}{m} np.sum(dz^{[1]}, axis=0, keepdims=True)$$

axis = 0, 1,
keepdims = True)
to get (n, 1)
not (n,)

$$dz^{[1]} = W^{[2]T} dz^{[2]} \circ g^{[1]}'(z^{[1]})$$

done using np.multiply or * ← -

element-wise prod

$$db^{[1]} = \frac{1}{m} np.sum(dz^{[1]}), \quad dw^{[1]} = \frac{1}{m} dz^{[1]} X^T$$

\rightarrow for any var x , $x \in dx$ have same dims

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

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

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

→ only 1 training eg

in case of 3 layered nn,

$$\left\{ \begin{array}{l} dW^{[2]} = dz^{[2]} \\ \quad \downarrow \\ \quad dh \\ \quad \overline{dz} \end{array} \right.$$

is in case of a

is in case of a

single training of

$$(dW^{[1]} = dZ^{[1]})^T$$

it's basically
a IOTF

but we do it w/ multiple eq

⇒ need to do vectorization

→ While initializing weights we use a random val instead of 0

e.g.

$w^{[1]} = np.random.randint(1, 2) \times 0.01$ in case of bias
it doesn't matter

2x2 matrix

naining a
eg

$\times 10^{-1}$ & not
large no.

\therefore on calc activation vs

$2^{[1]}$ turns to be very big
 \Rightarrow we end at

$z^{[1]}$ turns to be very big
 \Rightarrow we end at
flat parts in func
causing gradient
descent to be
very slow

Deep Neural Network

let $l = \text{no. of layers}$

$n^{[l]}$ = no. of units in layer l

$\Rightarrow n^{[1]} = 5 = 5 \text{ nodes in layer 1}$

& $n^{[0]} = n_x = \text{no. of ip node}$

main ip

$a^{[l]}$ = activation func in layer $l = g^{[l]}(z^{[l]})$

$w^{[l]}$ = weights for $z^{[l]}$

Forward Propagation in dnn

same as previously,

general eqn,

$$z^{[l]} = w^{[l]} A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(z^{[l]})$$

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = w^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(z^{[2]})$$

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

:

Getting the dimensions right

for first layers in forward prop we have → normal implementation

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

part of first layer

(nodes in layer, 1) → (no. of neurons, no. of training eg / nodes in previous layer)

(no. of neuron, ith ip to this) in this layer

⇒ general formula,

$A^{[l]} = Z^{[l]} = (n^{[l]}, 1)$	$A^{[l]} = Z^{[l]} = (n^{[l]}, m)$
$W^{[l]} = (n^{[l]}, n^{[l-1]})$	$W^{[l]} = (n^{[l]}, m)$
$X^{[l]} = (n^{[l]}, 1) = b^{[l]} = Z^{[l]}$	$X^{[l]} = (n^{[l]}, m)$

basically $A^{[0]}, l=0$

Vectorized implementation matrix

∴ examples are stacked vertically the dims for Z, X change,

$Z^{[l]} = (n^{[l]}, m)$	$Z^{[l]} = (n^{[l]}, 1)$
$X^{[l]} = (n^{[l]}, m)$	$X^{[l]} = (n^{[l]}, 1)$
$W^{[l]} = (n^{[l]}, n^{[l-1]})$	$W^{[l]} = (n^{[l]}, m)$
$b^{[l]} = (n^{[l]}, 1)$	$b^{[l]} = 1$

training eg

$(n^{[l]}, m) \leftarrow$ to but during '+' it gets duplicated / broadcasted

in back prop,

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

$$\frac{dW^{[l]}}{m} = \frac{1}{m} dz^{[l]} \cdot a^{[l-1]T}$$

element-wise prod

$$db^{[l]} = \frac{1}{m} \sum_{i=1}^m dz^{[l]}_i$$

dot prod

$$da^{[l-1]} = W^{[l]T} dz^{[l]}$$

→ in vectorized form,
 $\frac{1}{m} \text{np.sum}(dz^{[l]}, \text{axis}=1, \text{dims})$

10 substituting $da^{[l-1]}$

$$\Rightarrow dz^{[l]} = W^{[l+1]T} dz^{[l+1]} * g^{[l+1]'}(z^{[l+1]})$$

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

→ is not vectorized

15 $da^{[l]}$

Parameters & Hyper-parameters

20 params : $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots$

hyperparams : α , #iteration, #L, #hidden unit, choice of activation func, no. of

25 params that control W, b

applying dL is very empirical i.e trial & error-ish

