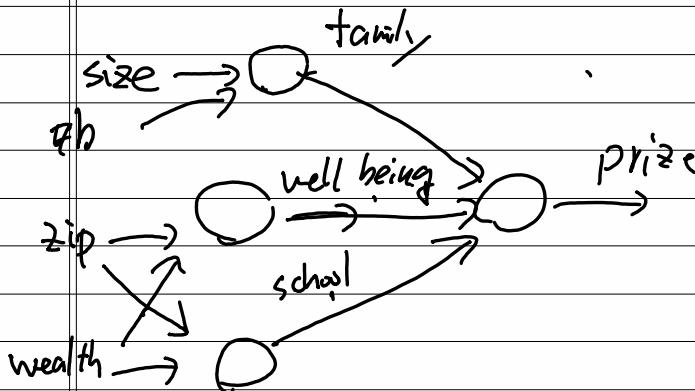


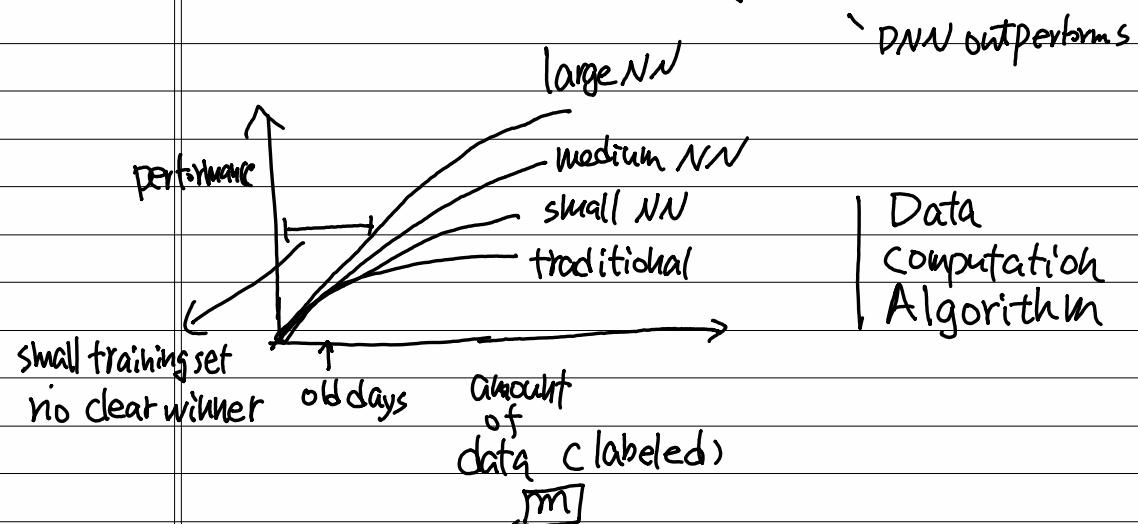
Coursera.



good intuition

but usually full-connected

structured data / unstructured data
table image, video, audio...



W2

m training samples $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})\}$

$$X = \begin{bmatrix} | & | \\ x^{(1)} & \dots & x^{(m)} \\ | & | \end{bmatrix}_{n \times m} \text{ features}$$

m

training sample

$$X \in \mathbb{R}^{n \times m}$$

$$X.\text{shape} = (n, m)$$

$$Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}_{1 \times m}$$

$$Y \in \mathbb{R}^m$$

stack in column
for easier implementation

logistic regression

Given X , want $\hat{y} = P(y=1|x)$

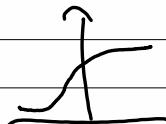
$$X \in \mathbb{R}^{n \times m}$$

parameter: $w \in \mathbb{R}^n$, $b \in \mathbb{R}$

sigmoid function

σ

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



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

$$\Rightarrow \sigma(z) = \frac{1}{1 + 0} = 1$$

$$\begin{aligned} z \downarrow \sigma(z) &= 1 \\ \text{negative} &= \frac{1}{1+e^{\infty}} \\ &= 0 \end{aligned}$$

COST FUNC

cost
func

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

want $\hat{y} \approx y^{(i)}$

loss:

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

loss
func

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

$$y=1 \rightarrow L(\hat{y}, y) = -\log(\hat{y}) \quad \text{want } -\log(\hat{y}) \downarrow$$

\hat{y} large

$$y=0 \rightarrow L(\hat{y}, y) = -\log(1-\hat{y}) \quad \text{want } \log(1-\hat{y}) \uparrow$$

\hat{y} small

$$\hat{y} = (0, 1)$$

cost func

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

GRADIENT DESCENT

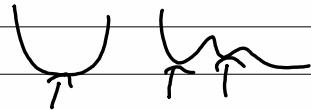
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

convex non-convex

$J(w)$

$$w := w - \alpha \frac{d J(w)}{dw}$$

learning rate



$J(w, b)$

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

$$\frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \frac{d J(w, b)}{db}$$

(partial derivative symbol)

∂ = partial (for multiple var)

d = derivative

DERIVATIVE

$$\text{slope} = \frac{h}{w}$$

$$f(a) = a^2$$

$$f(a) = 3a$$

$$\frac{d f(a)}{da} = 2a$$

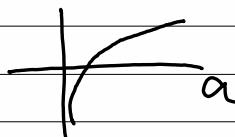
increase a by 0.1

$$y \uparrow 2 \times 0.1 = 0.2$$

$$\frac{d f(a)}{da} = 3$$

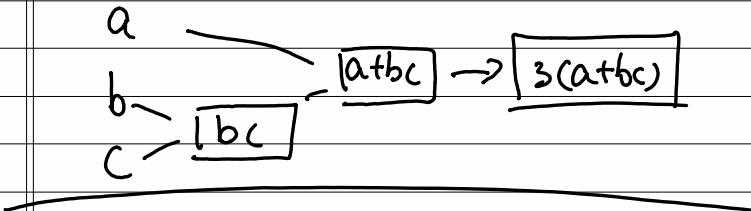
nudge $a \rightarrow$, y increase $3a$

$$f(a) = \log(a) \quad \frac{d}{da} f(a) = \frac{1}{a}$$

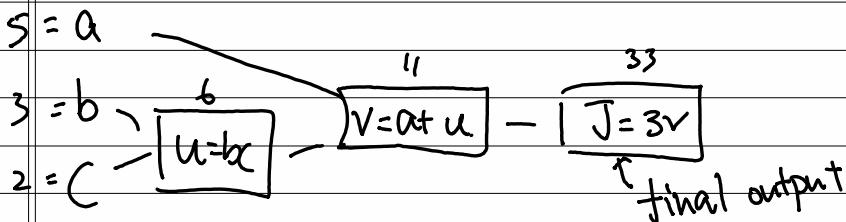


Computation Graph

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



Derivatives with computation graph



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

$$\frac{dJ}{da} = ?$$

$$v = 11 \rightarrow 11.001$$

$$a = 5 \rightarrow 5.001 \rightarrow \frac{dv}{da} = 1$$

$$J = 33 \rightarrow 33.003$$

$$v = 11 \rightarrow 11.001$$

$$\Delta v = 0.001 \times 3 = 0.003$$

$$J = 33 \rightarrow 33.003$$

$$\frac{dJ}{da} = \frac{dJ}{dv} \frac{dv}{da} \quad \text{chain rule}$$

Coding
convention

$$\frac{d \text{finaloutputvar}}{d \text{var}} \rightarrow \underline{d \text{var}}$$

- derivative of final output var.

backpropagation

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

$$S = a$$

$$3 = b$$

$$2 = c$$



$$11$$

$$33$$

$$V = a + u$$

$$dv = 3$$

$$u = bc$$

$$du = 3$$

$$u = 6 \rightarrow 6.001$$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$\frac{dJ}{du} = \frac{dJ}{dv} \cdot \frac{dv}{du} = 3.$$

$$3 \quad 1$$

$$b = 3 \rightarrow 3.001$$

$$\frac{dJ}{db} = \frac{dJ}{dv} \cdot \frac{dv}{du} \cdot \frac{du}{db} = 6$$

$$(\begin{matrix} 1 & 1 \\ 3 & 1 \end{matrix}) \quad 2$$

$$u = 6 \rightarrow 6.002$$

$$v = 11 \rightarrow 11.002$$

$$J = 33 \rightarrow 33.006$$

$$\frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} = 9$$

$$(\begin{matrix} 1 \\ 3 \end{matrix}) \quad 3$$

Logistic regression gradient descent

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

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

single sample

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

$$\begin{aligned} Z &= w_1 x_1 + w_2 x_2 + b \rightarrow a = \sigma(Z) \rightarrow L(a, y) \\ \frac{dZ}{dz} = \frac{dL}{dz} &= \frac{dL(a,y)}{da} \\ \frac{da}{dZ} &= -\frac{y}{a} + \frac{1-y}{1-a} \\ &= \underbrace{\frac{a(1-a)}{a}}_{\text{sigmoid derivative}} \end{aligned}$$

$$\frac{dL}{dw_i}$$

$$\begin{aligned} &= \frac{dL}{dz} \cdot \frac{dz}{dw_i} \\ &= x_i \cdot \frac{da}{dz} \end{aligned}$$

In training example

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

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

To calculate $\rightarrow d_{w_1}^{(i)}, d_{w_2}^{(i)}, db^{(i)}$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)})$$

$$\frac{\partial}{\partial w_2} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_2} L(a^{(i)}, y^{(i)})$$

$$d_{w_1}^{(i)} = (x_1^{(i)}, y^{(i)})$$

▷ Vectorization

CODE

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

$$J(w,b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)}) \quad \frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial b}$$

for each sample for $i = 1 \text{ to } m$

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

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

$$dw_1 += x_1 \cdot (a - y)$$

$$o^{(i)} = \sigma(z)$$

$$dw_2 += x_2 \cdot (a - y)$$

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

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

$$dw_1 += x_1 dz^{(i)}$$

$$dw_2 += x_2 dz^{(i)}$$

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

$$J /= m$$

$$dw_1 /= m$$

$$dw_2 /= m$$

$$db /= m$$

Vectorization / vectorized logistic regression

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

$$\begin{aligned} X &\quad X \in \mathbb{R}^{n \times m} \\ \begin{bmatrix} \hat{z}^{(1)} & \dots & \hat{z}^{(m)} \end{bmatrix} &= \begin{matrix} T \\ w^T X + [b \dots b] \end{matrix} \\ 1 \times n &\quad n \times n \quad 1 \times m \end{aligned}$$

$$= \begin{bmatrix} w^T x^{(1)} + b & w^T x^{(2)} + b & w^T x^{(3)} + b & \dots \end{bmatrix} \quad]$$

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

$\underbrace{\phantom{np.\text{dot}(w^T, x) + b}}$

\hat{z} scalar

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

$\underbrace{\phantom{a^{(i)} - y^{(i)}}}$

CODE

one iter of gradient descent

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

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

$$d\hat{z} = A - Y$$

$$dw = \frac{1}{m} X d\hat{z}^T$$

$$db = \frac{1}{m} np \text{sum}(d\hat{z})$$

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

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

Broadcasting

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \dots$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}_{(1, n)} \rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}_{(m, 1)} = \dots$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m, n)} + \begin{bmatrix} 100 \\ 200 \\ 300 \end{bmatrix}_{(m, 3)} \rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m, n)} + m \cdot n = \dots$$

$$\begin{array}{rcl} (m, n) & + & (l, n) \sim (m, l) \\ \times & & \text{broadcast} \\ \div & & (m, 1) \sim (m, n) \\ & & \text{broadcasting} \end{array}$$

DEBUG

Avoid "rank-1" array.

randn(5) - DON'T USE

randn(5, 1) - column vector

randn(1, 5) - row vector

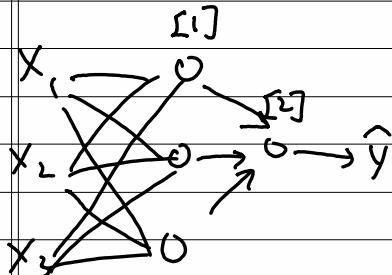
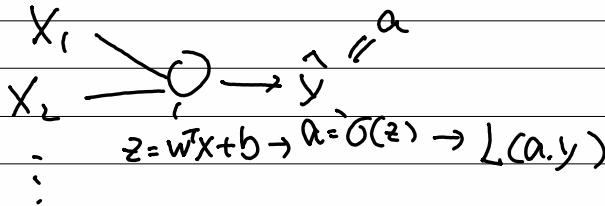
* → assert(a.shape == (5, 1))

↳ (5, 1).reshape((5, 1))

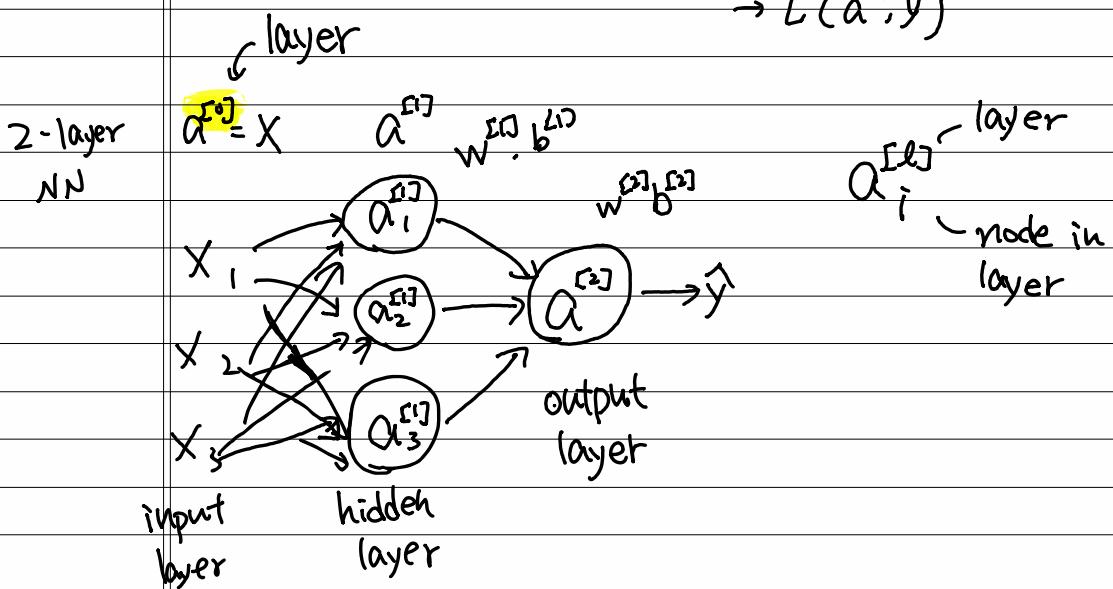
NN

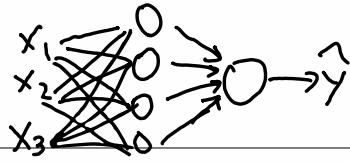
WEEK 3

Logistic regressionish



$$\begin{aligned} z^{[1]} &= w^{[1]} x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \\ &\rightarrow a^{[2]} = \sigma(z^{[2]}) \\ &\rightarrow L(a^{[2]}, y) \end{aligned}$$





Computing
NN output

$$z_1^{(1)} = w_1^{(1)T} x + b_1^{(1)}, \quad a_1^{(1)} = \sigma(z_1^{(1)})$$

hen roh

vectorize

$$z_4^{(1)} = w_4^{(1)T} x + b_4^{(1)}, \quad a_4^{(1)} = \sigma(z_4^{(1)})$$

$$\begin{bmatrix} -w_1^{(1)} \\ -w_2^{(1)} \\ \vdots \\ -w_4^{(1)} \end{bmatrix}_{4 \times 3} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} b_1^{(1)} \\ \vdots \\ b_4^{(1)} \end{bmatrix} = \begin{bmatrix} w_1^{(1)T} \cdot x + b_1^{(1)} \\ \vdots \\ w_4^{(1)T} \cdot x + b_4^{(1)} \end{bmatrix}$$

4×3

$$W^{(1)}$$

$$b^{(1)}$$

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

CODE

$$z^{(1)} = W^{(1)T} X + b^{(1)}$$

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

m samples

$$x \rightarrow a^{(1)} = y$$

$$x^{(1)} \rightarrow a^{(1)(1)} = y^{(1)}$$

$$\vdots$$

$$x^{(m)} \rightarrow a^{(1)(m)} = y^{(m)}$$

$$z^{(1)} = W^{(1)T} a^{(1)} + b^{(1)}$$

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

for $i = 1$ to m

$$z^{(1)(i)} = W^{(1)T} x^{(i)} + b^{(1)}$$

\vdots

vectorize
in training

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

$n \times m$.

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

$$A^{[1]} = G\{Z^{[1]}\}$$

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = G\{Z^{[2]}\}$$

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

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

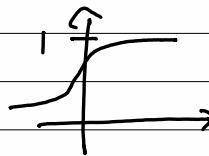
\nwarrow different
Node
in
layers

\nearrow different
sample

$$Z^{1} = W^{[1]} X^{(1)}$$

ACTIVATION

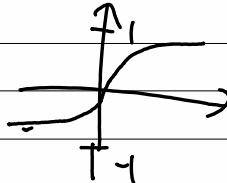
$$G = \frac{1}{1+e^{-z}}$$



prob only for output
 $[0, 1]$

$$g = \tanh(z)$$

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



almost always better
for hidden layers

on:

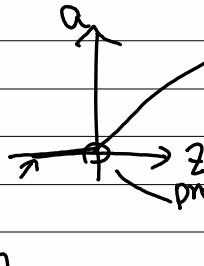
△ large $z \rightarrow$ small gradient \rightarrow slow

DEFAULT

ReLU

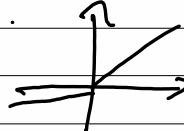
Rectified Linear Unit

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



pro: learn faster!
- bcⁿ not affected by
small slope
△ usually quickly hⁿ des
problem in theory only w/ $z > 0$

Leaky ReLU



$$\alpha = \max(0, 0.01z, z)$$

Derivative of Activation Function

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

prime

$$\begin{aligned} g'(z) &= \frac{1}{1+e^{-z}} \left(-\frac{1}{1+e^{-z}} \right) \\ g \text{ prime of } z &= g(z) \cdot (1-g(z)) \end{aligned}$$

$$\begin{aligned} \frac{d}{dz} e^{-z} \\ = -z \cdot e^{-z} \end{aligned}$$

* Power rule applies when numerator is
the power of derivative

$$\frac{d}{dx} x^n = n \cdot x^{n-1}$$

$$\begin{aligned} \frac{d}{dx} f(x)^n &= \frac{d}{df(x)} f(x)^n \cdot \frac{df(x)}{dx} \\ &= n \cdot f(x)^{n-1} \cdot f'(x) \end{aligned}$$

$$\begin{aligned} &= -\frac{d}{dz} (1+e^{-z}) \times (1+e^{-z})^2 \\ &= -e^{-z} \cdot (1+e^{-z})^{-2} \\ &= \frac{-e^{-z}}{(1+e^{-z})^2} \end{aligned}$$

$$\begin{aligned} \frac{d}{dx} f(x)^{-1} &= -f(x)^{-2} \cdot f'(x) \\ &= \frac{-f'(x)}{f(x)^2} \end{aligned}$$

$$\begin{aligned} \frac{d}{dx} \frac{\theta(x)}{f(x)} \\ = \frac{d}{dx} \theta(x) \cdot f(x)^{-1} \end{aligned}$$

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

$$\frac{d}{dz} \tanh(z) = 1 - (\tanh(z))^2$$

ReLU

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

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

implement NN

Forward Prop

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

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

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

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

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

Back Prop

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

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

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

$$dz^{[1]} = w^{[1]T} dz^{[2]} * g'(z^{[1]})$$

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

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

Initialization

\times all $w = 0$

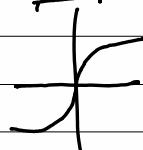
- all hidden unit will be the same

$$w^{[1]} = \text{np.random.randn}(C_{2,2}) \xrightarrow{\text{under}} 0, 0.1 \xrightarrow{\text{in this range}}$$

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

:

$$b^{[2]} = 0$$



W3 Assignment

planar

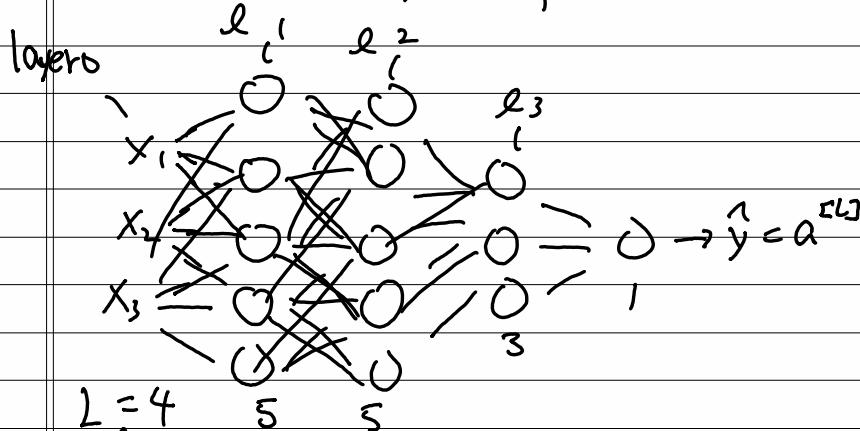
$$X \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad 2 \times m$$

$$Y \in \mathbb{R}^m \quad 1 \times m \text{ (label)}$$

w4

DNN - Deep NN

shallow vs deep



$$L = 4 \quad 5 \quad 5$$

$$n^{[l]} = \# \text{unit in layer } l$$

$$h^{[1]}=5, h^{[2]}=5, h^{[3]}=3, n^{[4]}=1$$

$$h^{[1]} = n_x = 3$$

$$a^{[l]} = \text{activation in layer } l$$

$$a^{[l]} = g^{[l]}(z^{[l]}) \quad w^{[l]} = \text{weights for } z^{[l]}$$

z

Forward Prop

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

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

$$z^{[l]}$$

↓

vectorize

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

$$\{ z^{[l-1]} \dots z^{[0]} \}$$

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

$$z^{[l]}$$

↓

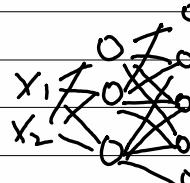
$$\hat{y} = g(z^{[L]}) = A^{[L]}$$

DEBUG DNN

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

(3, 1)

(2, 1)



so, w should be (3, 2)

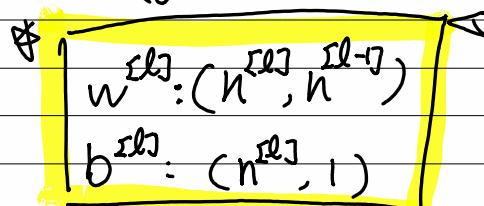
$$\Rightarrow (h^{[1]}, 1) = (h^{[0]}, h^{[1]})(h^{[0]}, 1)$$

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

(5, 1)

(3, 1)

so, (5, 3)



$$dw^{[l]} : (n^{[l]}, n^{[l-1]})$$

$$db^{[l]} : (n^{[l]}, 1)$$

VECTORIZED

$$z^{[l]} = w^{[l]} \cdot x + b^{[l]}$$

($n^{[l-1]} \times 1$) ($n^{[l]} \times n^{[l-1]}$) ($n^{[l]} \times 1$) ($n^{[l]} \times 1$)

$$\sum^{[l]} = w^{[l]} \cdot x + b^{[l]}$$

($n^{[l]} \times m$) ($n^{[l]} \times n^{[l-1]}$) ($n^{[l]} \times m$) ($n^{[l]} \times 1$)

Python broadcasting

\downarrow
($n^{[l]} \times m$)

$$z^{[l]}, A^{[l]} : (n^{[l]}, m)$$

$$l=0, A^{[0]} = x = (n^{[0]}, m)$$

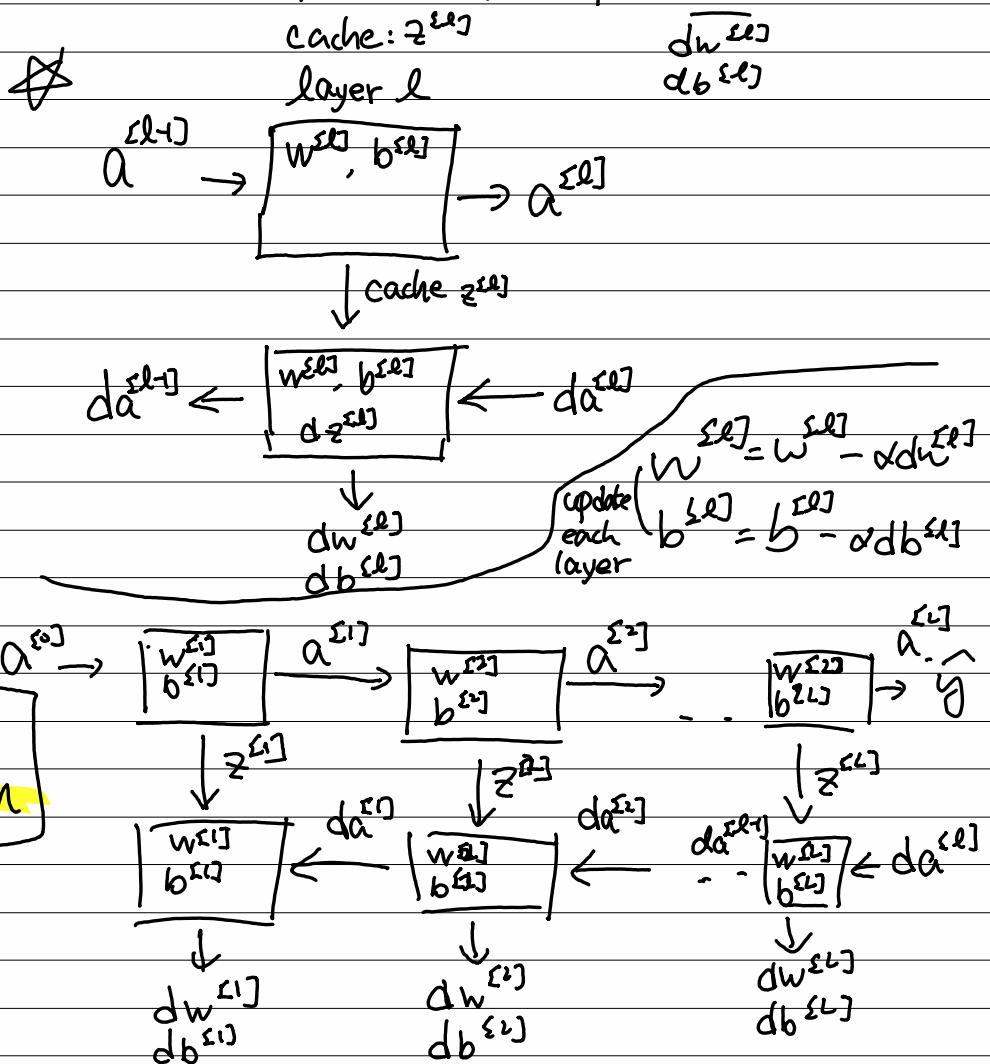
$$dz^{[l]} - dA^{[l]} : (n^{[l]}, m)$$

Forward: input $a^{[l-1]}$, output $a^{[l]}$

$$z^{[l]} = W^{[l]} \cdot a^{[l-1]} + b^{[l]}, \text{ cache } z^{[l]}$$

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

Backward: input $da^{[l]}$, output $da^{[l-1]}$



IMPLEMENTATION

BACKPROP

input: $d\alpha^{[L]}$

output: $d\alpha^{[L-1]}, dW^{[L]}, db^{[L]}$

$$d\tilde{z}^{[L]} = d\alpha^{[L]} * g^{[L]}(\tilde{z}^{[L]})$$

$$dW^{[L]} = d\tilde{z}^{[L]} \cdot \alpha^{[L-1]T}$$

$$db = d\tilde{z}^{[L]}$$

$$d\alpha^{[L-1]} = W^{[L]T} \cdot d\tilde{z}^{[L]}$$

$$d\tilde{z}^{[L-1]} = W^{[L-1]T} d\tilde{z}^{[L]} * g^{[L-1]}(\tilde{z}^{[L]})$$

$$\Delta d\alpha^{[L]} = -\frac{\gamma}{\alpha} + \frac{(1-\gamma)}{1-\alpha}$$

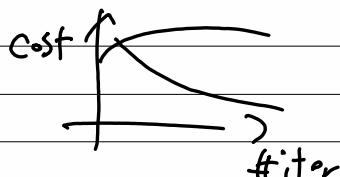
$$dA^{[L]} = \left[-\frac{\gamma}{\alpha} + \frac{(1-\gamma)}{1-\alpha}, \dots, \frac{\gamma^{(m)}}{\alpha^{(m)}} + \frac{(1-\gamma^{(m)})}{1-\alpha^{(m)}} \right]$$

HYPERPARAMETER

- { learning rate
- iteration
- hidden layer L
- hidden unit, $n^{[1]}, n^{[2]}$
- activation func

parameters
 $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots$

look at cost function graph to decide



hyperparameter
trial-and-error.