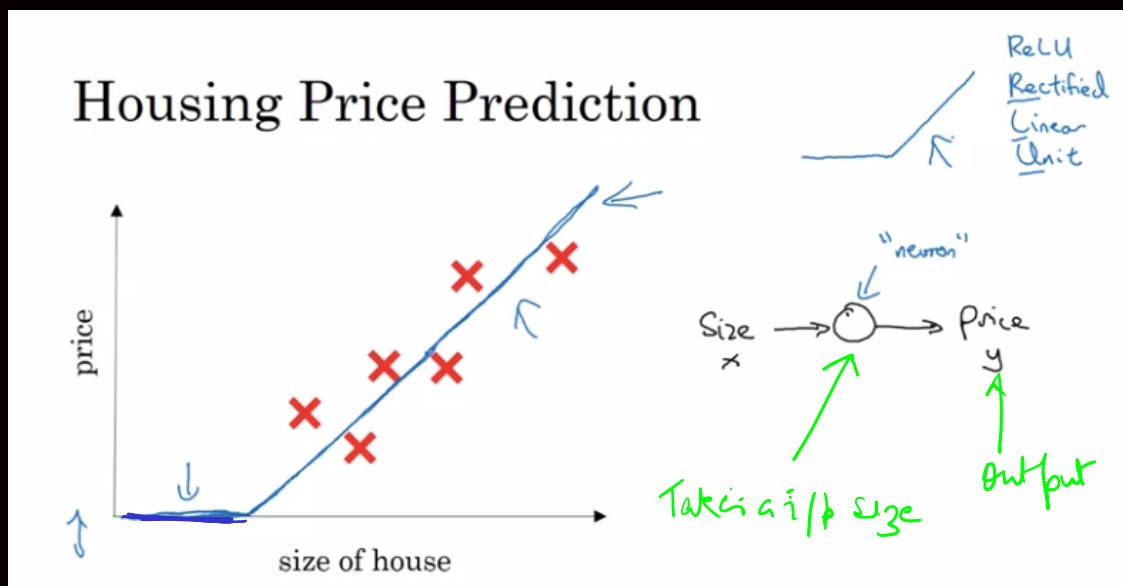
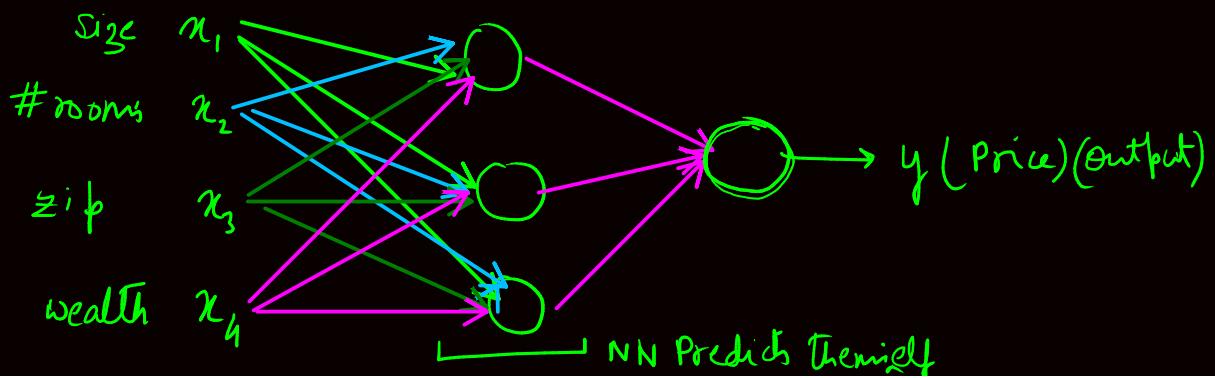
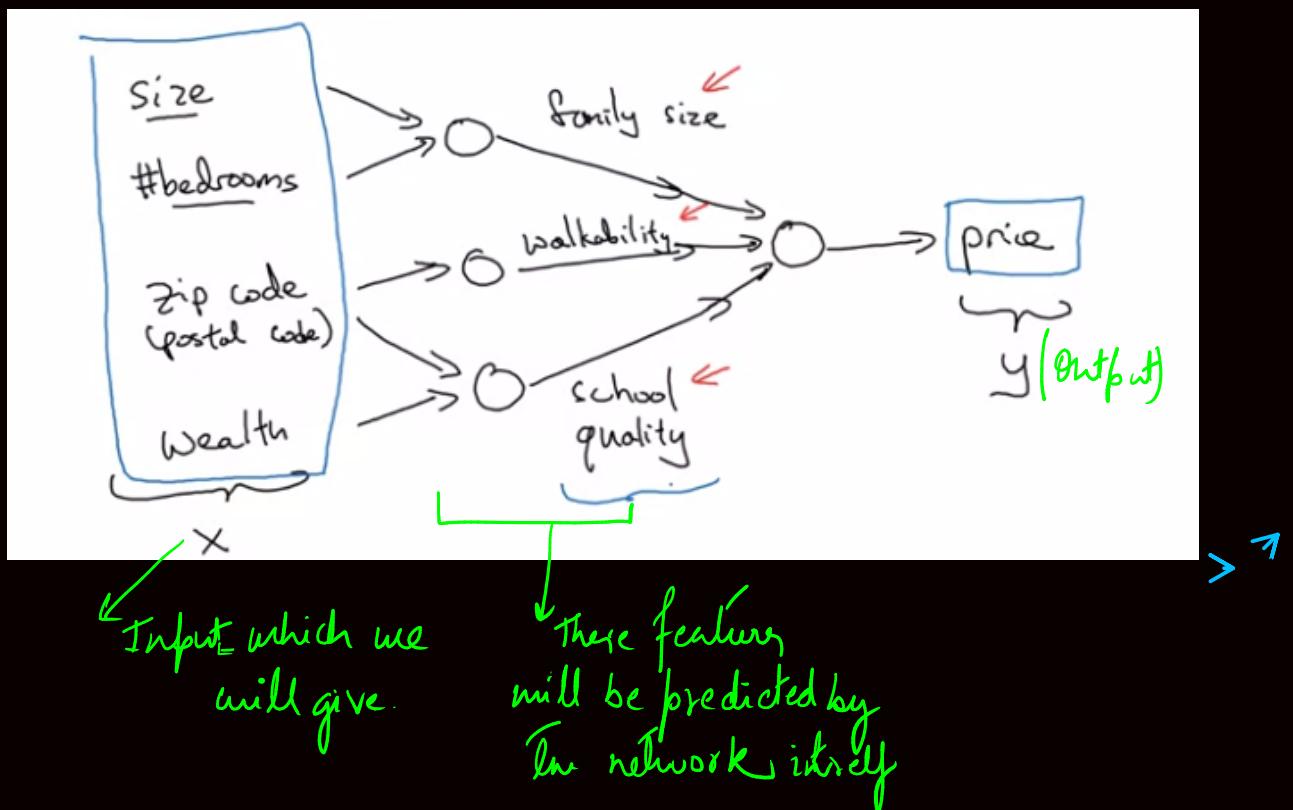


Week 1

Neural Networks & Deep Learning



of smallest Possible Neural Network

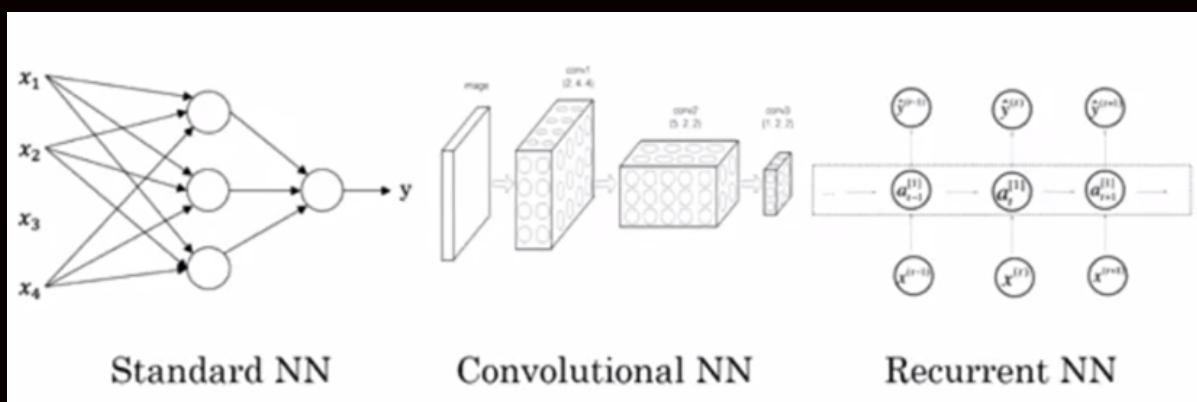


Supervised Learning

Applications

Input(x)	Output (y)	Application
Home features	Price	Real Estate
Ad, user info	Click on ad? (0/1)	Online Advertising
Image	Object (1,...,1000)	Photo tagging
Audio	Text transcript	Speech recognition
English	Chinese	Machine translation
Image, Radar info	Position of other cars	Autonomous driving

} Standard NN
} CNN
} RNN (Sequence data)
} Custom Hybrid Network

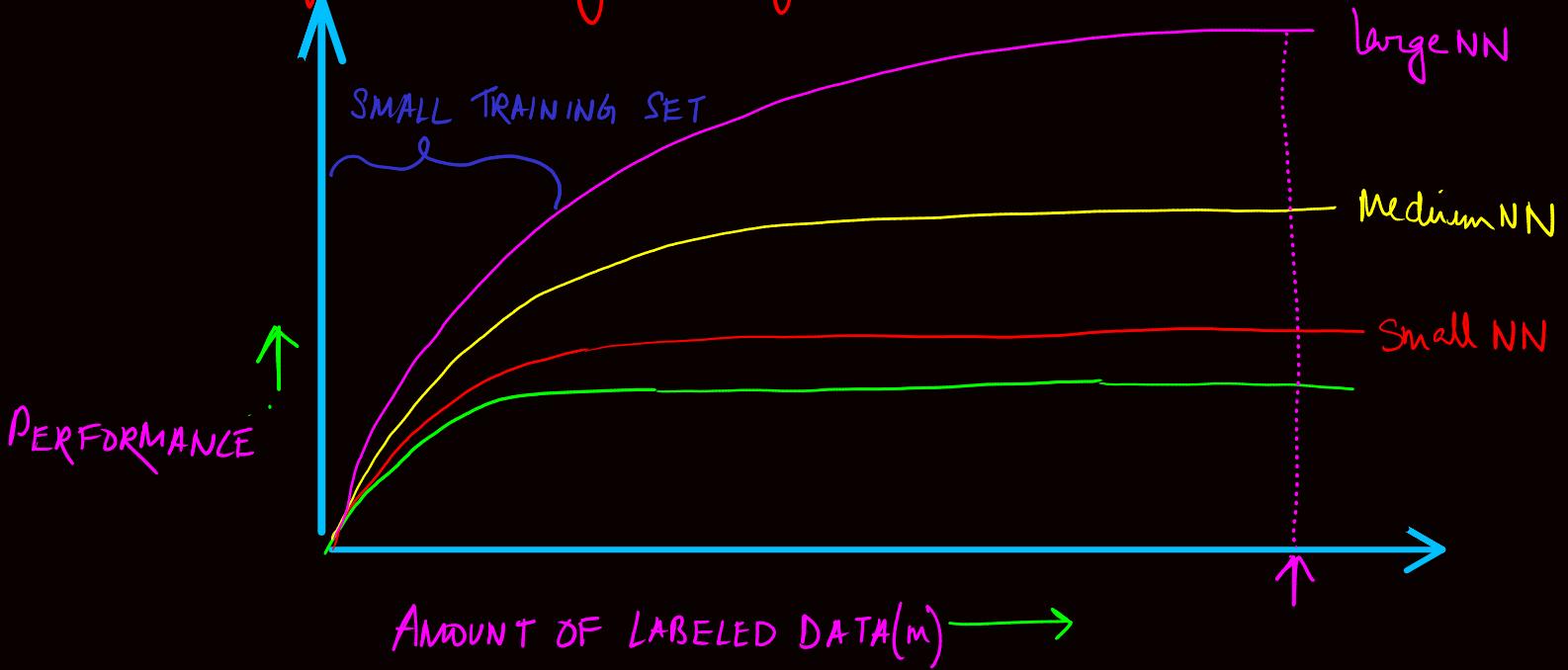


Other Classification

Structured Data				Unstructured Data	
Size	#bedrooms	...	Price (1000\$)		
2104	3		400	Audio	Image
1600	3		330		
2400	3		369		
:	:		:		
3000	4		540		
User Age	Ad Id	...	Click	Four scores and seven years ago...	
41	93242		1	Text	
80	93287		0		
18	87312		1		
:	:		:		
27	71244		1		

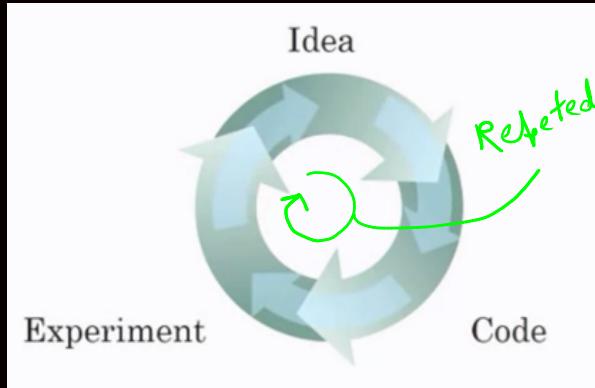
NN focuses more on
Unstructured Data

Why deep learning is taking off now?



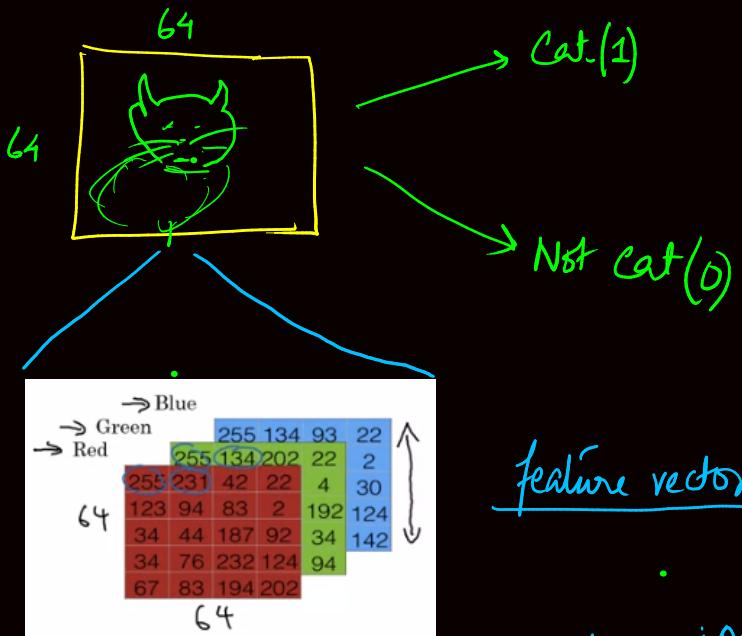
Driving forces:

- Data.
(↑ in data)
- Computation
(Hardware)
- Algorithms.
E.g.



Week 2

BINARY CLASSIFICATION



feature vector: $x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix}$

dimension: $64 \times 64 \times 3 = 12288$

$n_x = 12288$

Training Example = (x, y) ; $x \in \mathbb{R}^{n_x}$, $y \in \{0, 1\}$

In training example = $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$X = \left[\begin{array}{cccc} & \downarrow \text{1 feature vector (1st image)} & & \\ n^{(1)} & n^{(2)} & \dots & n^{(m)} \\ | & | & & | \\ \hline m & & & \end{array} \right] \quad ; \quad X \in \mathbb{R}^{n_x \times m}$

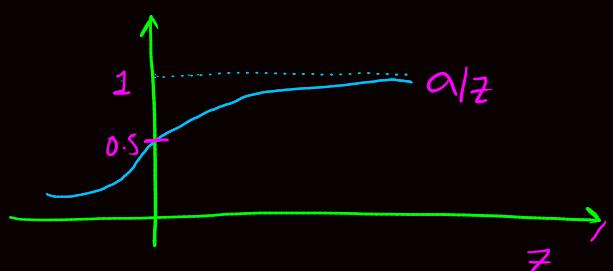
$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]_{1 \times m}$$

LOGISTIC REGRESSION

Given x , want $\hat{y} = P(y=1|x)$; $0 < \hat{y} < 1$

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

$$\text{Output } \hat{y} = \sigma(\underbrace{w^T x + b}_z)$$



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

when $z \rightarrow \infty$, $\sigma(z) = 1$
 $z \rightarrow -\infty$, $\sigma(z) = 0$

COST FUNCTION

Prediction on the training example i^{th} :

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

Loss funⁿ (error funⁿ):

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

If $y=1$ $L(\hat{y}, y) = -\log \hat{y} \rightarrow \hat{y}$ must be large

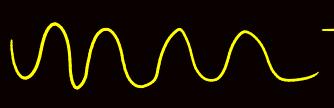
$y=0$; $L(\hat{y}, y) = -\log(1-\hat{y}) \rightarrow \hat{y}$ be be as small as possible.

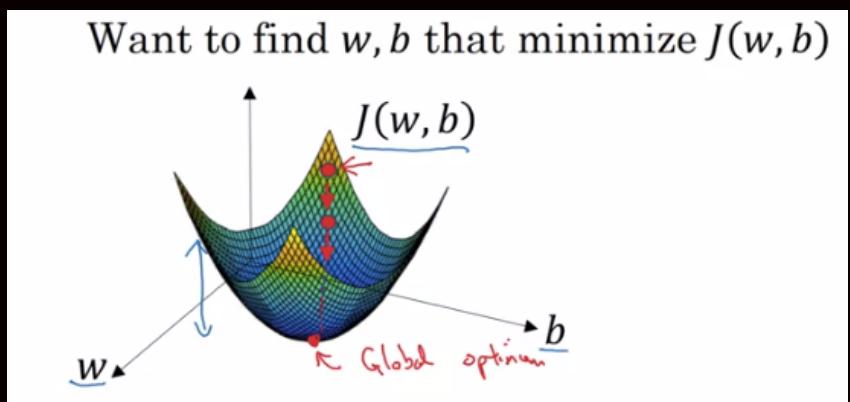
Loss funⁿ \rightarrow Single Training Example

Cost funⁿ: (for the entire training set)

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

* Convex funⁿ  → one global minima

Non convex funⁿ  → multiple global minima.

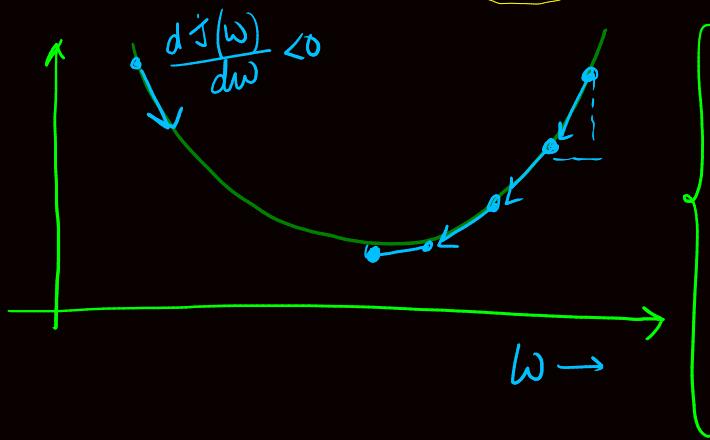


Gradient Descent

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

learning rate α

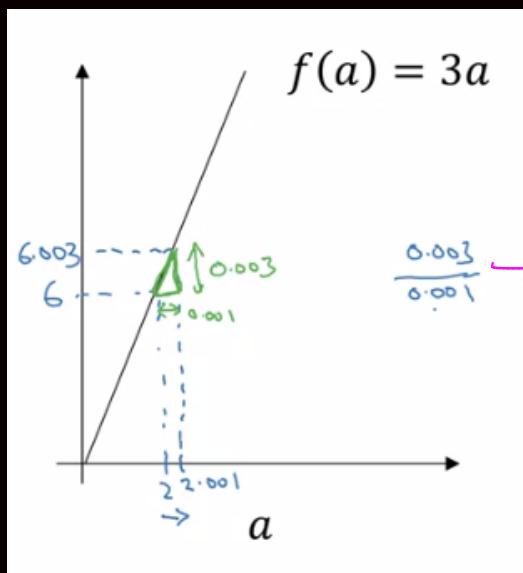
$$\Rightarrow \omega = \omega - \alpha dw$$



Updates

$$\begin{aligned} \omega &= \omega - \alpha \frac{dJ(\omega, b)}{dw} \\ b &= b - \alpha \frac{dJ(\omega, b)}{db} \end{aligned}$$

Derivatives



Slope

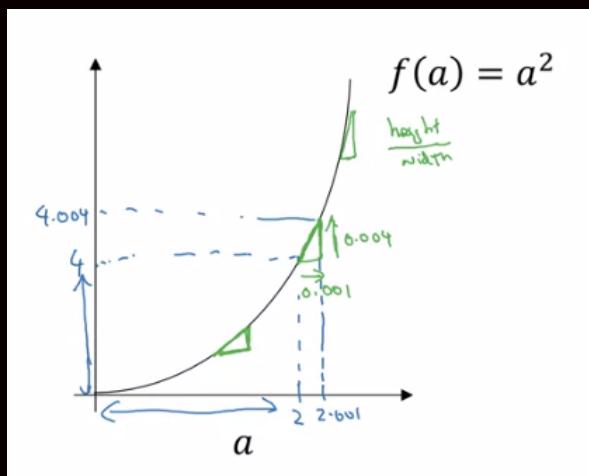
$$\frac{df(a)}{da} = 3 = \frac{d}{da} f(a)$$

i.e if we tweak the variable a 1 time, $f(a)$ is tweaked 3 times that value

$$a = 5, f(a) = 15$$

$$h = 0.001, f(a+h) = 15.003$$

When slope is diff throughout the function



$$\frac{df(a)}{da} = 2a$$

$$\frac{df(a)}{da} = 4 \text{ when } a = 2$$

$$\frac{df(a)}{da} = 10 \text{ when } a = 5$$

Similar fun's are

$$\textcircled{1} \quad f(a) = a^3$$

$$\frac{df(a)}{da} = 3a^2$$

$$\textcircled{2} \quad \log a = f(a)$$

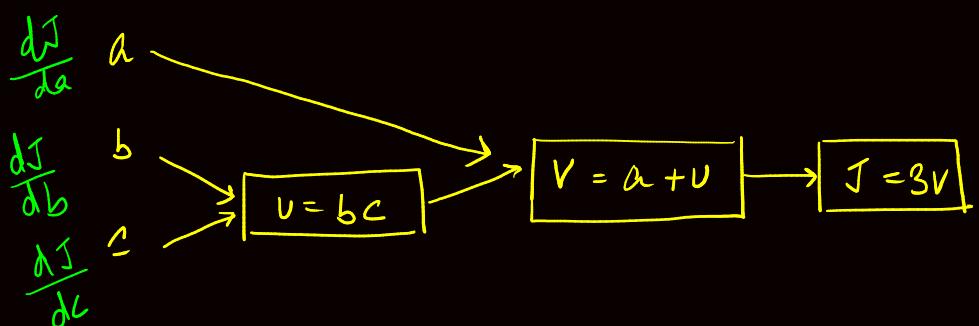
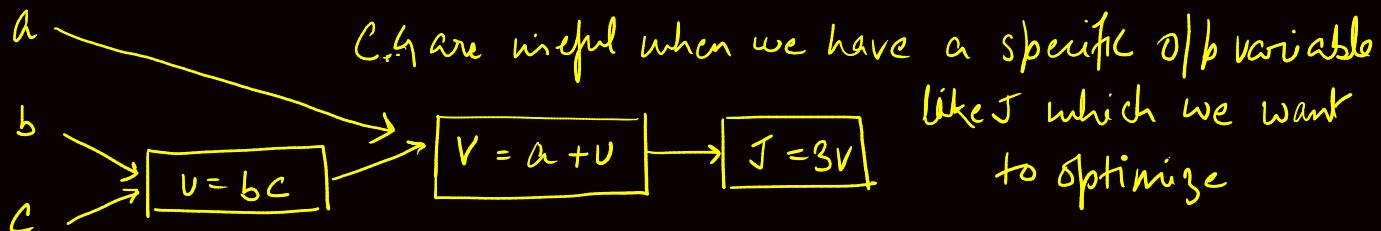
$$\frac{df(a)}{da} = \frac{1}{a}$$

* Derivative of a fun' = Slope of a fun'

COMPUTATION GRAPH

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

$$U = bc, V = a + U, J = 3V$$



$$\begin{aligned} \frac{dJ}{da} &= J = 3V \\ \frac{dJ}{db} &= V = 11 \rightarrow 11.001 \{ \text{ny} \} \\ \frac{dJ}{dc} &= J = 33 \rightarrow 33.003 \{ 3x \} \end{aligned}$$

$$\therefore \frac{dJ}{dv} = \frac{d(3v)}{dv} = 3$$

Chain Rule

$$\begin{aligned} \frac{dJ}{da} &= a = 5 \rightarrow 5.001 \\ &= v = 11 \rightarrow 11.001 \\ &= J = 33 \rightarrow 33.001 \end{aligned}$$

a effekt \rightarrow v effekt \rightarrow J

$$\boxed{\frac{dJ}{da} = \frac{dJ}{dv} \times \frac{dv}{da}}$$

$$= 3 \times 1$$

Final O/b Variable to Optimize = J

$$\frac{d(\text{Final O/b Variable})}{d(van)} = dJ/dvan \leftarrow "dvan"$$

↑

Smart Notation for code

$$\star \frac{dJ}{dvan}$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial b}$$

$$b = 3 \rightarrow 3.001$$

$$v = b \cdot c \rightarrow 6 \rightarrow 6.002$$

$$= \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} \cdot \frac{\partial u}{\partial b}$$

$$= 3 \times 1 \times 2$$

$$= 6$$

$$\frac{\partial J}{\partial c} = 6 \quad (\text{Same})$$

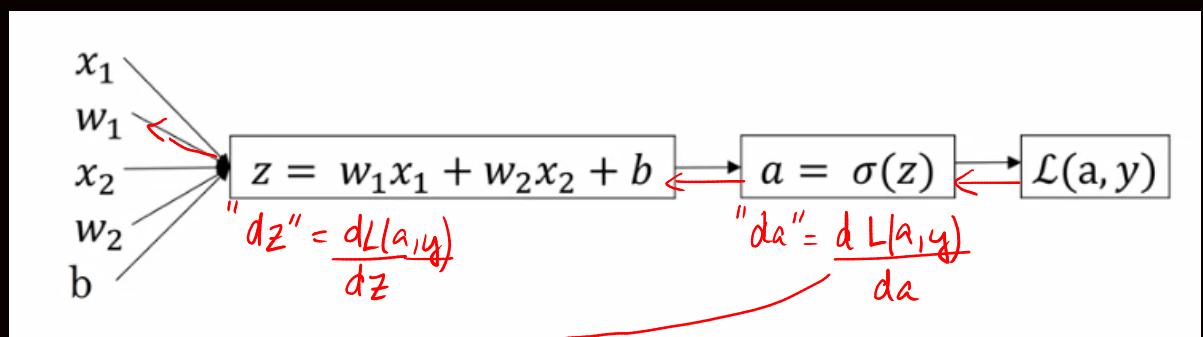
* The most good way to find derivatives is to do a right to left computation. by chain rule

LOGISTIC REGRESSION GRADIENT DESCENT

① For 1 training example

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

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



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

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

$$= (a - y)$$

$$* \frac{\partial L}{\partial w_1} = \overbrace{\frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z}}^{\text{d}z} \times \frac{\partial z}{\partial w_1}$$

$(a - y) \times n_1$
 $\underbrace{d_z}_{d_w} \times n_1$

$$\therefore \frac{\partial L}{\partial w_1} = "d w_1" = n_1 d z \quad , \frac{\partial L}{\partial w_2} = n_2 d w_2 \quad , \frac{d L}{d b} = d z$$

Updates learning rate

$$w_1 = w_1 - \alpha d w_1 ; \quad w_2 = w_2 - \alpha d w_2 ; \quad b = b - \alpha d b$$

② For 'm' training example

Recap $z = w^T x + b$ $L(a, y) = -(y \log a + (1-y) \log(1-a))$

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

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

Steps

① Initialization

$$J=0, d w_1=0, d w_2=0, d b=0$$

For $i=1$ to m

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

$$a^{(i)} = a(z^{(i)})$$

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

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

$$d w_1 += x_1^{(i)} d z^{(i)}$$

This increases in only 2 features w_1, w_2

$$d\omega_2 = \omega_2^{(i)} dz^{(i)}$$

$$db = dz^{(L)}$$

$$\begin{aligned} J/m \\ d\omega_1/m \\ d\omega_2/m \\ db/m \end{aligned}$$

This technique is very tedious

∴ For such calculation we adopt what we know as vectorization,

VECTORIZATION

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

In python

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

```
a = np.random.rand(1000000)
b = np.random.rand(1000000)

tic = time.time()
c = np.dot(a,b)
toc = time.time()

print(c)
print("Vectorized version:" + str(1000*(toc-tic)) + "ms")

c = 0
tic = time.time()
for i in range(1000000):
    c += a[i]*b[i]
toc = time.time()

print(c)
print("For loop:" + str(1000*(toc-tic)) + "ms")
```

250286.989866
Vectorized version:1.5027523040771484ms
250286.989866
For loop:474.29513931274414ms

for loops are much slower.

GPU [] → SIMD → Single Instruction multiple data
CPU

AVOID USING EXPLICIT FOR LOOPS whenever possible

Examples

Suppose

$$\textcircled{1} \quad u = Av$$

$u = np.\text{dot}(A, v) \rightarrow$ Vectorized.

$$U_i = \sum_j A_{ij} U_j$$

$U = np.zeros((n, 1))$

for $i \dots$:

$$U[i] += A[i][j] * v[j]$$

} Not Vectorized

\textcircled{2}

Say you need to apply the exponential operation on every element of a matrix/vector.

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

```
→ u = np.zeros((n, 1))
for i in range(n):
    u[i]=math.exp(v[i])
```

import numpy as np
 $u = np.exp(u)$
 Vectorized // faster.

\textcircled{3}

Original For Loop Technique.

$$z = 0, dw_1 = 0, dw_2 = 0, db = 0$$

Alternate

$$dw = np.zeros(n_x, 1)$$

For $i = 1$ to m

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

$$a^{(i)} = a(z^{(i)})$$

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

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

$$dw_1 += X_1^{(i)} dz^{(i)}$$

$$dw += n^{(i)} dz^{(i)}$$

$$d\omega_2 = n_2^{(i)} dz^{(i)}$$

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

J/m

$d\omega_1/m$
 $d\omega_2/m$
 db/m

$d\omega/m$

VECTORIZING LOGISTIC REGRESSION

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

$$X = \left[\begin{array}{c|c|c|c|c|c} & & & & & \\ \hline X^{(1)} & n^{(2)} & \dots & - & X^{(m)} & \\ \hline & & & & & \\ \hline \end{array} \right]_{(n \times m)} \quad w = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = \underbrace{w^T X}_{\downarrow} + [b \ b \ \dots \ b]_{1 \times m}$$

$$[z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = \underbrace{[w^T x^{(1)} + b]}_{\text{Row vector}} \ \underbrace{[w^T x^{(2)} + b]}_{\text{Row vector}} \ \dots \ \underbrace{[w^T x^{(m)} + b]}_{\text{Row vector}}_{1 \times m}$$

$$Z = \text{np-dot}(w^T, X) + b.$$

→ Real w broadcasted to $(1 \times m)$ matrix
to add to the previous term

$$\text{Silly, } A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(Z)$$

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

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

$$A = \begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(n)} \end{bmatrix} \quad Y = \begin{bmatrix} y^{(1)} & \dots & y^{(m)} \end{bmatrix}$$

$$\delta Z = A - Y = \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & \dots & a^{(n)} - y^{(n)} \end{bmatrix}$$

* $\boxed{\delta b = \frac{1}{m} \text{np.sum}(\delta Z)}$

* $\boxed{\delta w = \frac{1}{m} X \times \delta Z^T}$

$$= \frac{1}{m} \begin{bmatrix} | & | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & & \\ | & | & | & \dots & | \end{bmatrix}_{n \times m}^T \begin{bmatrix} \delta z^{(1)} \\ \delta z^{(2)} \\ \vdots \\ \delta z^{(m)} \end{bmatrix}_{m \times 1}$$

$$= \frac{1}{m} [x^{(1)} \delta z^{(1)} + \dots + x^{(m)} \delta z^{(m)}]_{n \times 1}$$

\rightarrow {epochs}

* Code for iter in range(1000):

$$\left. \begin{array}{l} Z = \text{np.dot}(w.T, X) + b \\ A = \sigma(Z) \\ \delta Z = A - Y \\ \delta w = \frac{1}{m} X \times \delta Z^T \end{array} \right\}$$

Vectorized. imp of Gradient Descent

$$\delta b = \frac{1}{m} \text{np.sum}(\delta Z)$$

$$w = w - \alpha \delta w$$

$$b = b - \alpha \delta b$$

}

BROADCASTING IN PYTHON

Calories from Carbs, Proteins, Fats in 100g of different foods:

	Apples	Beef	Eggs	Potatoes
Carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	135.0	99.0	0.9

→ 21.94%

g. Find the % of calories for carb, Protein & Fat for each of the food items individually

```
In [1]: import numpy as np
In [2]: A = np.array([[56.0, 0.0, 4.4, 68.0],
                  [1.2, 104.0, 52.0, 8.0],
                  [1.8, 135.0, 99.0, 0.9]])
print(A)
[[ 56.    0.    4.4   68. ]
 [  1.2  104.   52.    8. ]
 [  1.8  135.   99.    0.9]]
In [3]: cal = A.sum(axis=0)
print(cal)
[ 59.  239.  155.4  76.9]
In [4]: percentage = 100 * A/cal.reshape(1,4)
print(percentage)
[[94.91525424  0.          2.83140283  88.42652796]
 [ 2.03389831 43.51464435 33.46203346 10.40312094]
 [ 3.05084746 56.48535565 63.70656371  1.17035111]]
```

* ①
$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \xrightarrow{\text{Python Broadcasting}} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \cdot \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix}$$

②
$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} \cdot = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 265 & 306 \end{bmatrix}$$

$1 \times n \rightarrow m \times n$

etc..

$$\begin{array}{c}
 * (m, n) \\
 \text{matrix} \\
 + \quad (1, n) \rightarrow (m, n) \\
 * \quad (n, 1) \rightarrow (m, n) \\
 /
 \end{array}$$

* Tip: Don't use data structure.

```

* a = np.random.randn(5)
  a.shape = (5,)           } Don't use a=a.reshape((5,1))
  "rank 1 array"          } if such Rank 1 arrays come
                           } ↑
                           a = a.reshape((5,1))    ↗
a = np.random.randn(5, 1) → a.shape = (5,1)      column vector ✓
a = np.random.randn(1, 5) → a.shape = (1,5)      row vector. ✓
assert(a.shape == (5,1)) ←
  
```

* Logistic Regression cost fun

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

$$\begin{array}{ll}
 * \text{If } y = 1: & p(y|x) = \hat{y} \\
 \text{If } y = 0: & p(y|x) = 1 - \hat{y}
 \end{array}$$

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

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

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

$$\begin{array}{l}
 \text{maximize} \uparrow \log(p(y|x)) = \underset{\text{to minimize the loss}}{\underset{\uparrow}{L(\hat{y}, y)}}
 \end{array}$$

Cost function examples

$$\log P(\text{label } i \text{ in training set}) = \underbrace{\log \prod_{i=1}^m P(y^{(i)} | x^{(i)})}_{= -L(\hat{y}^{(i)}, y^{(i)})}$$

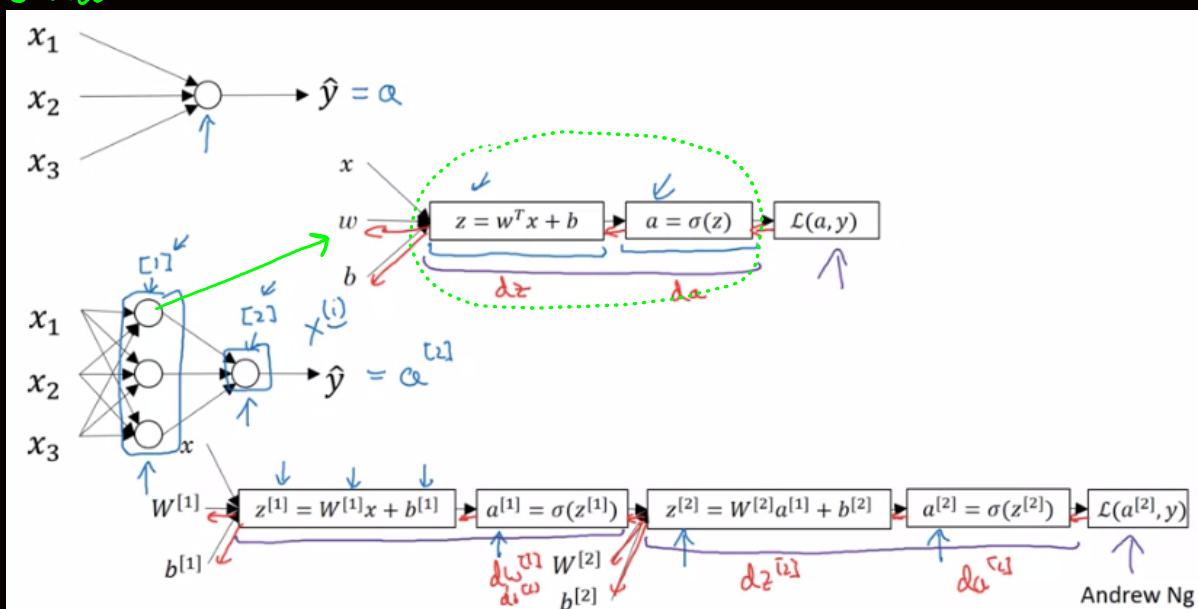
By maximum likelihood estimation

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

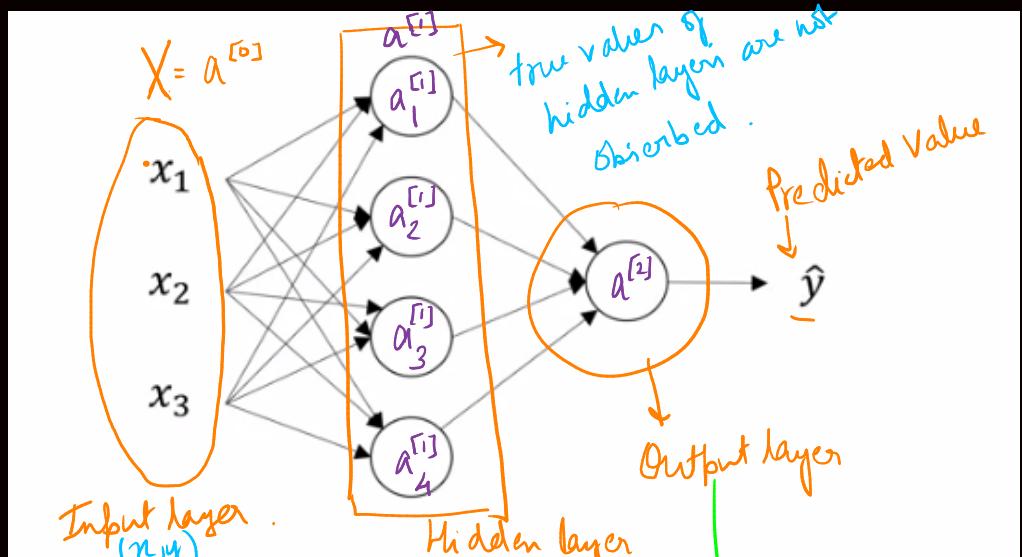
minimize //

WEEK 3

Overview



2-layer Network [ifp excluded]



image

www.nature.com/scientificreports/ | (2022) 12:1030 | Article number: 1030

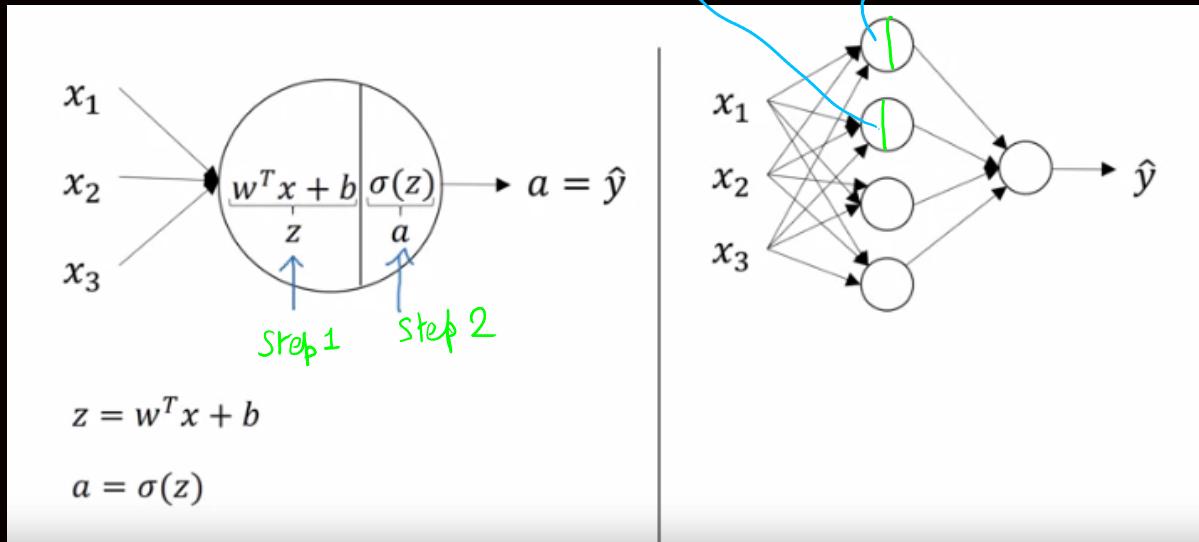
ω \downarrow $b^{[1]}$, $b^{[2]}$

$$\omega^{[2]}, b^{[2]}$$

$$z_1 = \omega_2 x + b_2$$

$$q_2 = \vartheta(z_2)$$

$$= \omega_1^{[1]T} x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]})$$



$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]}) \\ z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]}) \end{aligned}$$

$$Z = \underbrace{\begin{bmatrix} -\omega_1^{[1]T} \\ -\omega_2^{[1]T} \\ -\omega_3^{[1]T} \\ -\omega_4^{[1]T} \end{bmatrix}}_{W^{[1]}} \underbrace{\begin{bmatrix} n_1 \\ n_2 \\ n_3 \end{bmatrix}}_{3 \times 1} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}}_{b^{[1]}}$$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]})$$

Given input x :

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

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

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

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

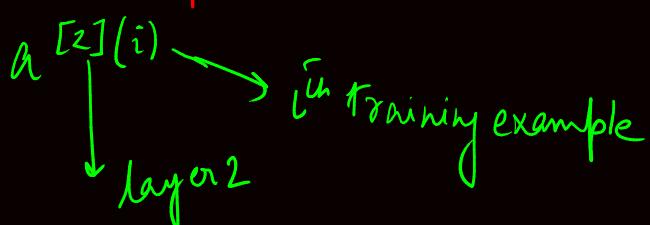
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

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

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

$$(1,1) \quad (1,1)$$

For m training example.



Vectorized Solution

training example

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

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

$$A^{[1]} = \sigma(Z^{[1]})$$

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

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

where, no. of training examples

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

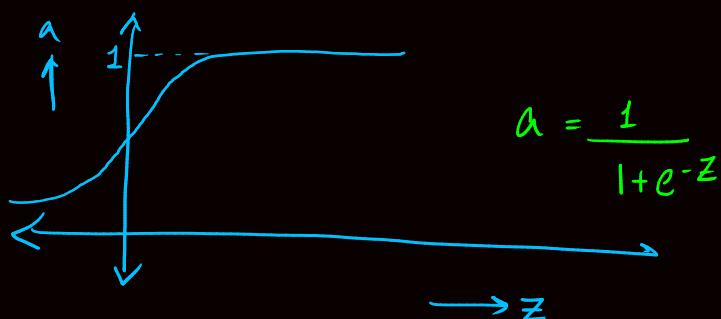
$A^{[1]} = \begin{bmatrix} | & | & | & | \\ a_1^{1} & a^{[1](2)} & \dots & a^{[1](m)} \\ a_2^{1} \\ \vdots \end{bmatrix}$

← No. of training examples → ↑ No. of nodes in that particular layer.

ACTIVATION FUNCTIONS

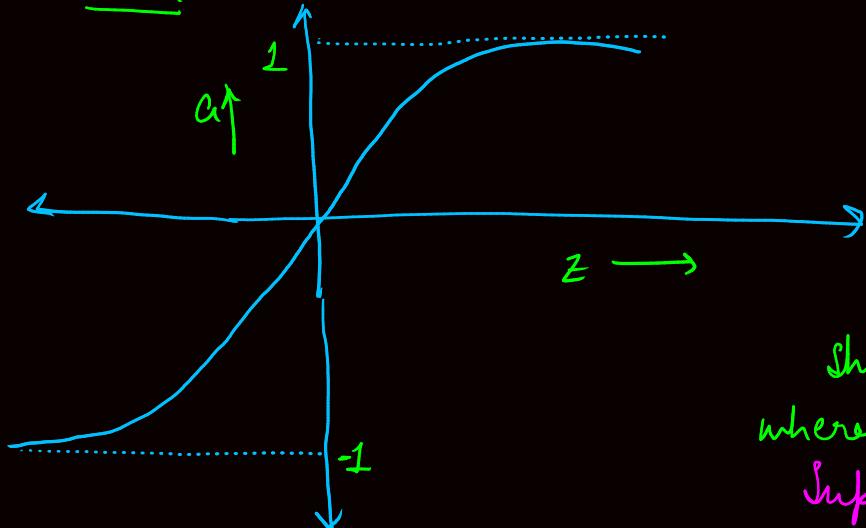
Sigmoid

The tanh activation usually works better than sigmoid activation function for hidden units because the mean of its output is closer to zero, and so it centers the data better for the next layer.



Can be used for binary classification task

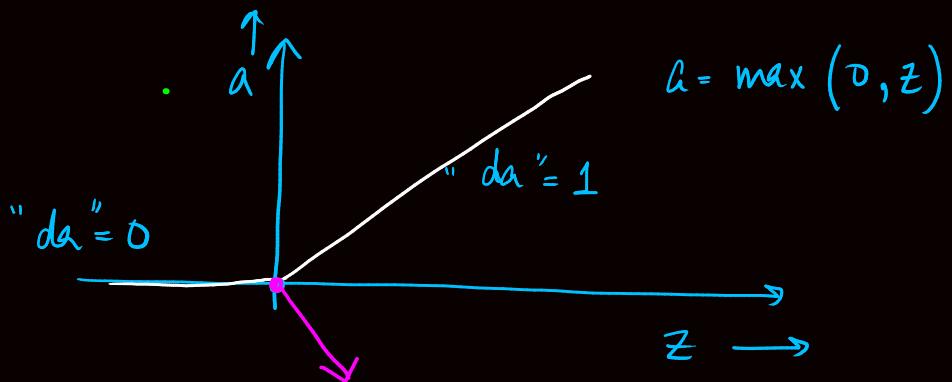
Tanh



"Zero mean"
shifted version of sigmoid
where $-1 < a < 1$
Superior tanh sigmoid

* Drawback of sigmoid & tanh is that if the activation function value is very large or very small, the slope ≈ 0 , thus learning becomes very slow.

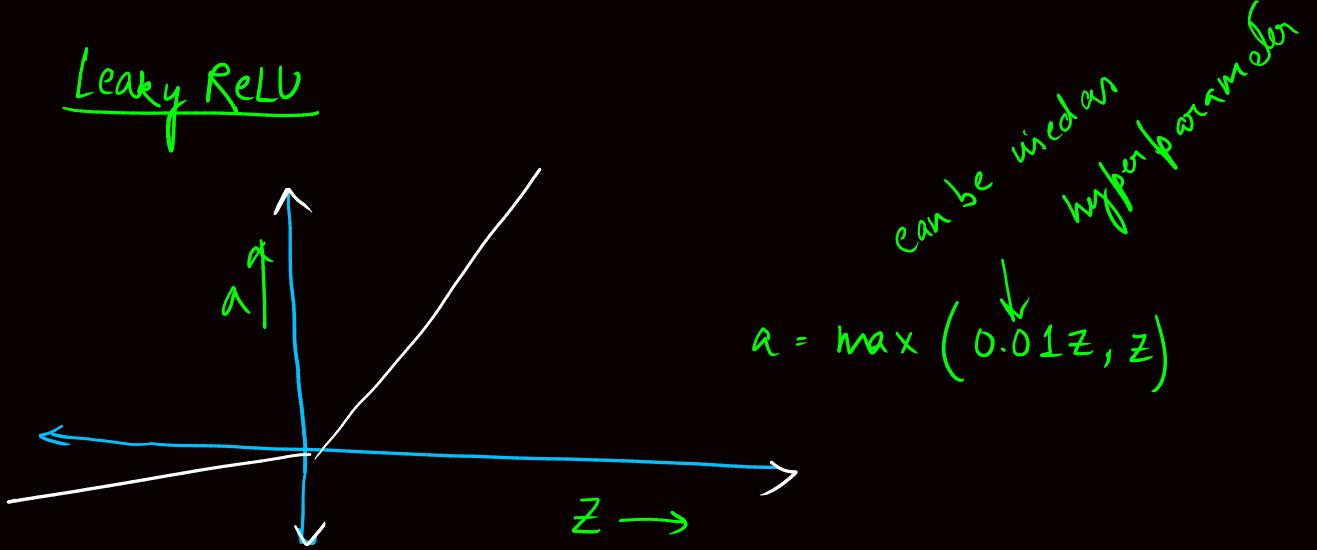
ReLU (Rectified Linear Unit)



Technically this point is undefined but in computer the value of this point accounts for 0.00000.....1

- * If our model is for binary classification, let the opp layer activation funⁿ by sigmoid & all the rest layers, ReLU must be the default choice.

Leaky ReLU



Why non-linear activation funⁿ?

- * Combining 2 or more linear funⁿ will in turn give an linear funⁿ Thus no learning

$$\begin{aligned} z^{[1]} &= \underline{\underline{z^{[1]}}} = \underline{\underline{W^{[1]}x + b^{[1]}}} \\ a^{[1]} &= \underline{\underline{g^{[1]}(z^{[1]})}} = \underline{\underline{z^{[1]}}} = \underline{\underline{W^{[1]}a^{[1]} + b^{[1]}}} \\ a^{[2]} &= \underline{\underline{a^{[2]}}} = \underline{\underline{W^{[2]}\left(\underline{\underline{W^{[1]}x + b^{[1]}}} + b^{[2]}\right)}} + b^{[2]} \\ &= \underbrace{(W^{[2]}W^{[1]})}_{W^T} x + \underbrace{(W^{[2]}b^{[1]} + b^{[2]})}_{b^T} \\ &= \underline{\underline{W^T x + b^T}} \end{aligned}$$

Andrew Ng

One place where we use linear activation function in regression problem (ML)

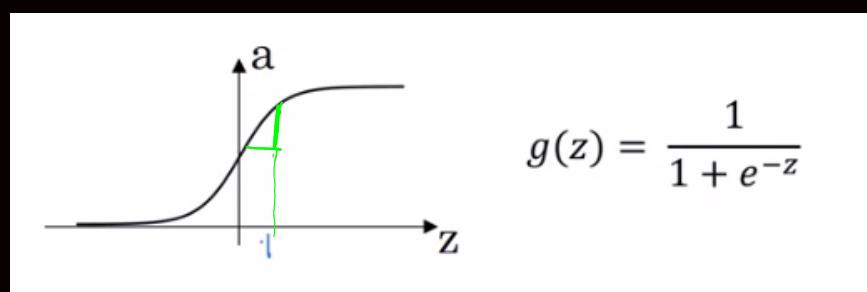
Eg Predict house prices

$$y \in \mathbb{R} \{ \$0, \dots, \$100000 \}$$

Use linear here
But all other hidden layers must use ReLU

Derivative of Activation function

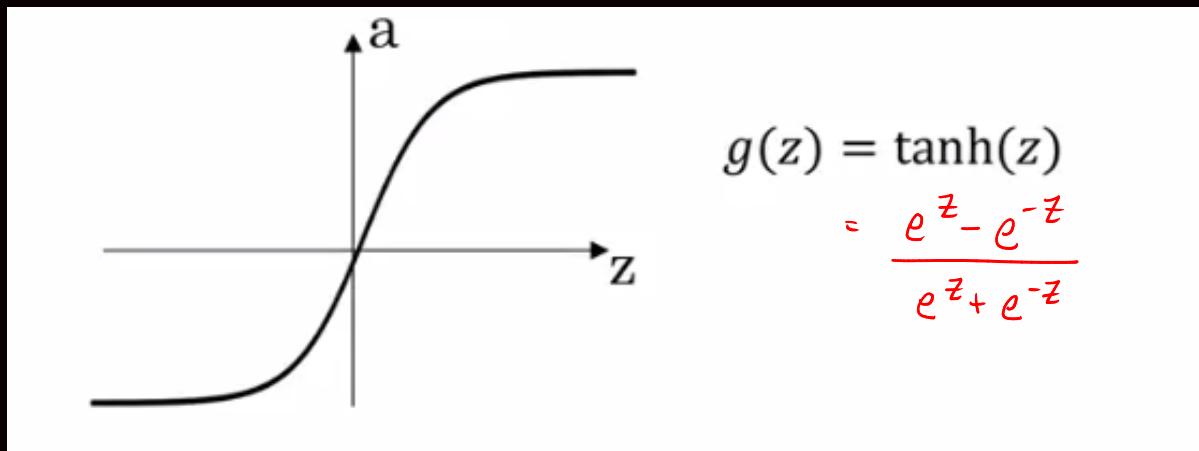
① Sigmoid



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

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

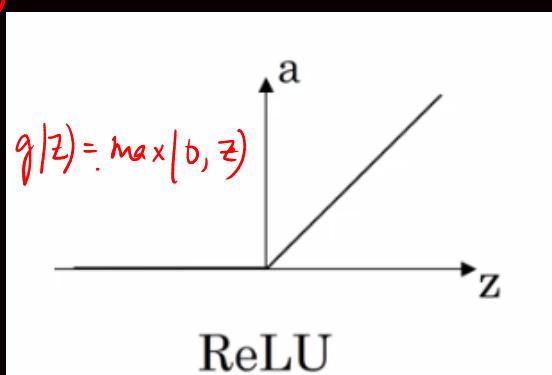
② Tanh



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

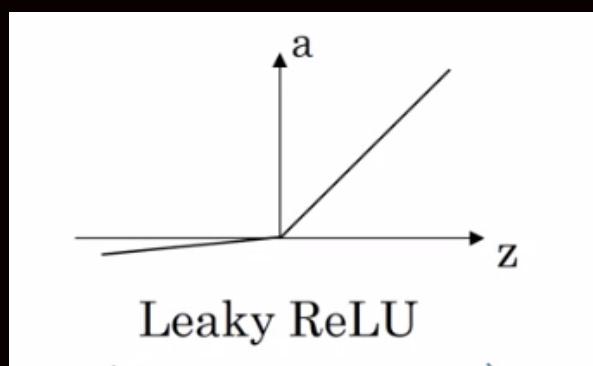
$$a = g(z) = g'(z) = 1 - a^2$$

③ ReLU



$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \\ 0.0000...1 & \text{if } z = 0 \end{cases}$$

④ Leaky ReLU



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

GRADIENT DESCENT FOR NEURAL NETWORK

Parameters = $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$ $n_x = n^{[0]}, n^{[1]}, n^{[2]}$
 $(n^{[1]} \times n^{[0]}) (n^{[1]} \times 1) \quad (1 \times n^{[1]}) \quad [n^{[2]}, 1]$

$$\text{Cost fun} = J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}, y)$$

$\rightarrow a^{[2]}$

Gradient Descent

Till parameters converge

Repeat {

- ① Compute predictions = $(\hat{y}^{(i)}, i = 1, 2, \dots, m)$
- ② $d w^{[1]} = \frac{d J}{d w^{[1]}}, d w^{[2]}, d b^{[1]}, d b^{[2]}$
- ③ $w^{[1]} = w^{[1]} - \alpha d w^{[1]} ; b^{[1]} = b^{[1]} - \alpha d b^{[1]}$
 $w^{[2]} = w^{[2]} - \alpha d w^{[2]} ; b^{[2]} = b^{[2]} - \alpha d b^{[2]}$

} FORMULUS

Forward Propagation

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

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

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

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

Backward Propagation

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

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

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

Automatical changes
any rank arrays to
a vector

$$dz^{[1]} = \underbrace{W^{[2]T} dz^{[2]} *}_{(n^{[1]}, m)} \underbrace{g^{[1]'}(z^{[1]})}_{(n^{[1]}, m)}$$

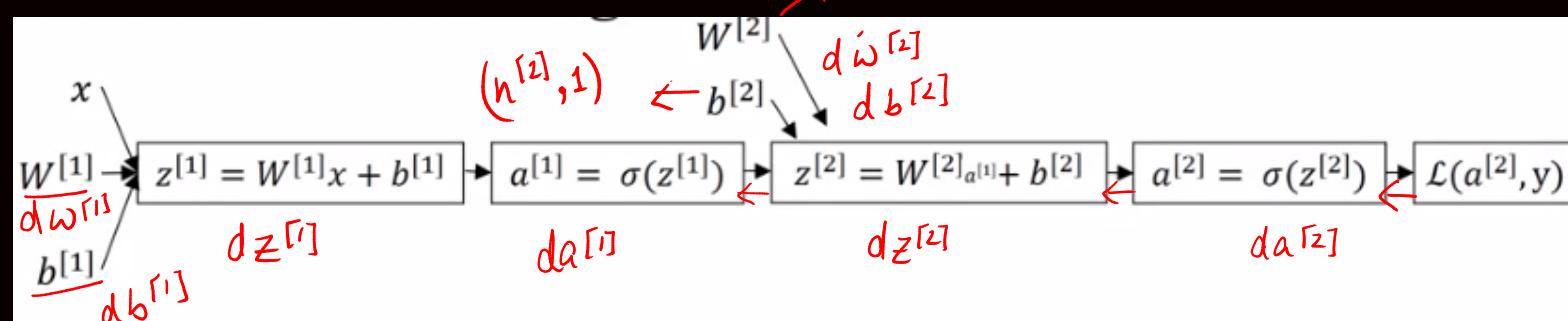
element wise product

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

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

Intuition behind the
Back Prop formulae

ω & $d\omega \rightarrow$ same dimension



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

$$d\omega^{[2]} = dz^{[2]} a^{[1]T} \approx \begin{cases} \text{Logistic Regm.} \\ d\omega = dz x \end{cases}$$

$$db^{[1]} = dz^{[1]}$$

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

$$d\omega^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

Summary of gradient descent

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

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

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

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

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

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

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

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{(n^{[2]}, m)} * \underbrace{g^{[1]'}(Z^{[1]})}_{(n^{[1]}, m)}$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

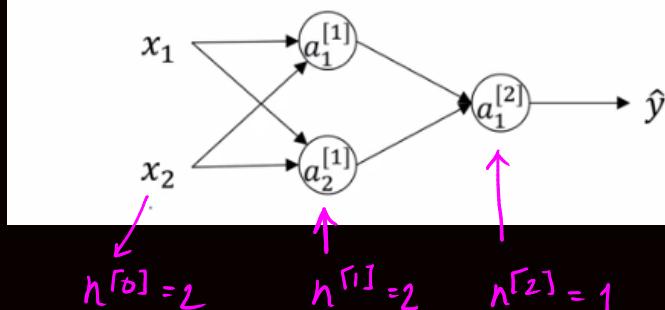
$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$$

Vectorized Implementation

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

RANDOM INITIALIZATION

What happens if you initialize weights to zero?



$$\text{let, } \omega^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\text{Then } a_1^{[1]} = a_2^{[1]} \quad \text{and } dz_1^{[1]} = dz_2^{[1]}$$

$$dw = \begin{bmatrix} v & v \\ u & v \end{bmatrix} \quad w^{[1]} = w^{[1]} - \alpha dw = \begin{bmatrix} \quad \\ \quad \end{bmatrix}$$

$$dw = \partial w$$

No matter how many times we train, all the hidden units will be computing the same function (all the hidden units will be symmetric)

Solution

Initialize with random values:

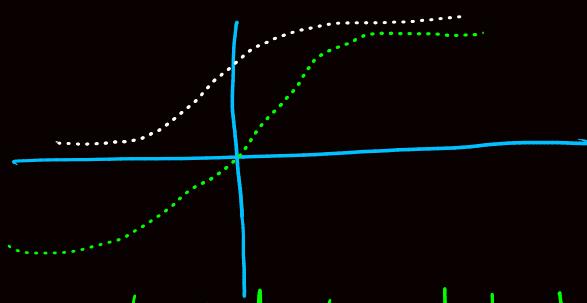
$$w^{[1]} = np.random.randn(2, 2) * 0.01 \quad \text{to initialize } w \text{ by very small values}$$

$$b^{[1]} = np.zeros(1) \quad \{ \text{No symmetry issue} \}$$

$$w^{[2]} = \dots \quad \text{same as } w^{[1]}$$

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

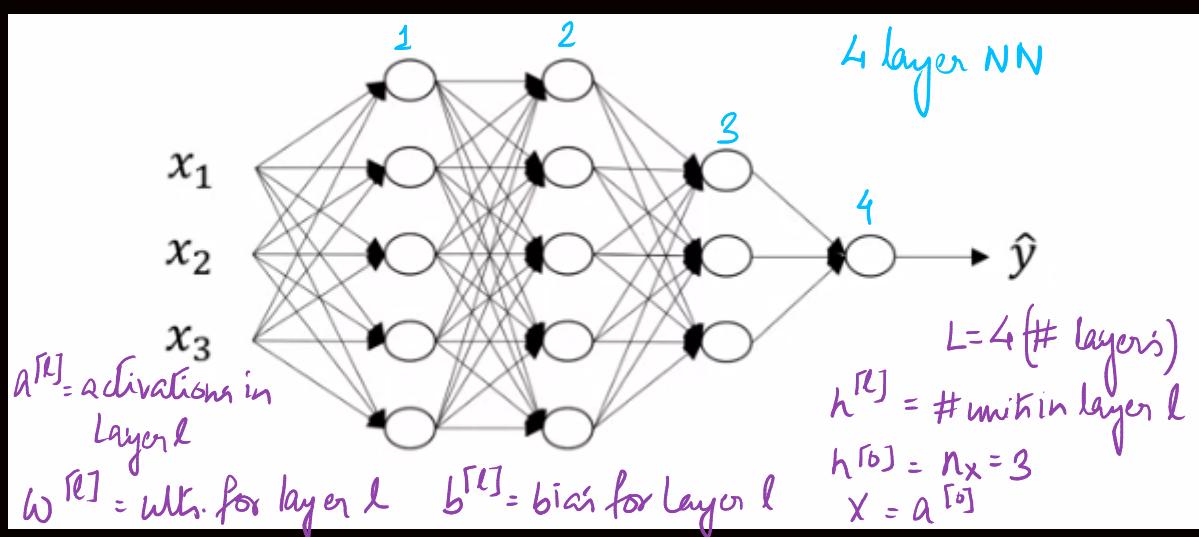
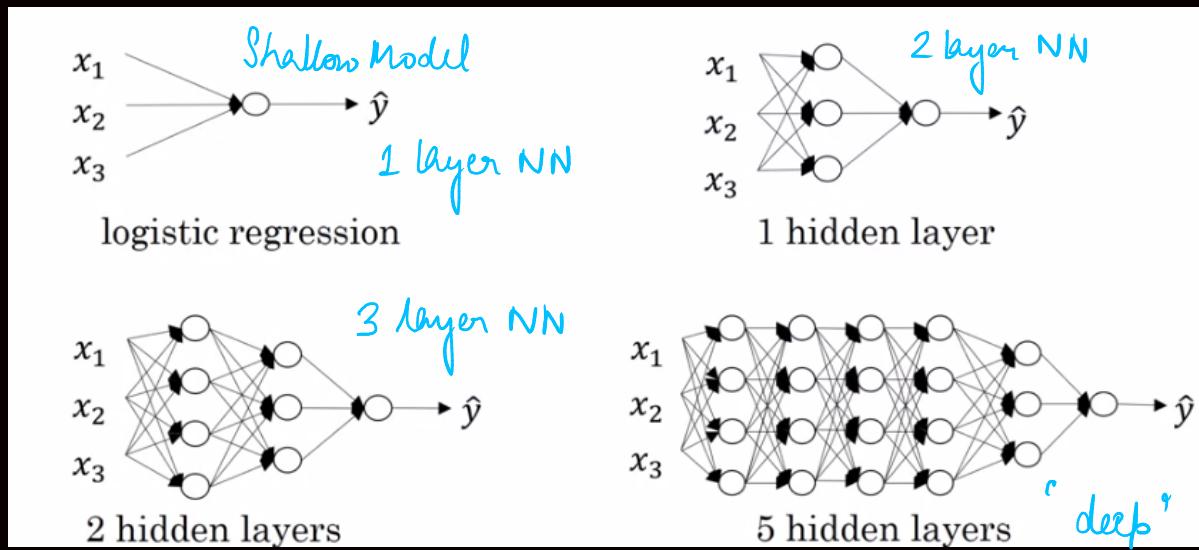
* If w are large $g(x)$ will be big & we will end up in a flat slope if we use tanh or sigmoid as $g(x)$ function



* 0.01 can be used as a hyperparameter

Week 4

Deep L-Layer neural network



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

Forward Propagation

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

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

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

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

$$\vdots \quad \vdots$$

$$a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$$

General = $z^{[\ell]} = \omega^{[\ell]} a^{[\ell-1]} + b^{[\ell]}$

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

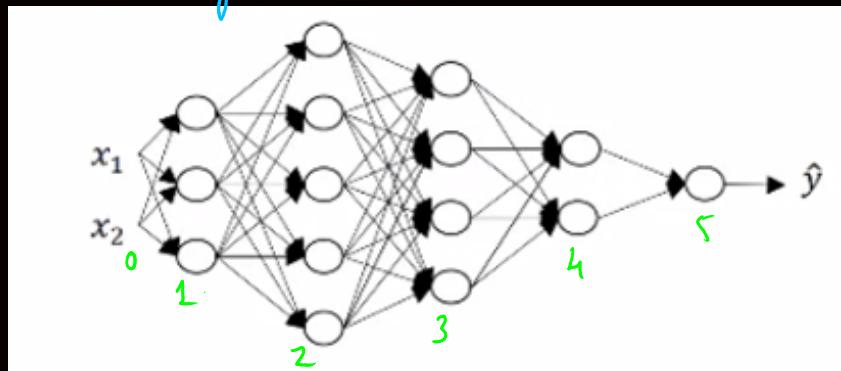
Vectorized

Capital Z ← $Z^{[0]} = W^{[0]} A^{[-1]} + B^{[0]}$
 as defined in
 previous slide

$$A^{[\ell]} = g^{[\ell]}(Z^{[\ell]})$$

$$\hat{y} = g(Z^{[4]}) = A^{[4]}$$

Getting the dimensions right



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

$\begin{matrix} 3 \times 1 \\ \text{---} \\ (3 \times 2) \end{matrix}$	$\begin{matrix} (2 \times 1) \\ \text{---} \\ \vdots \end{matrix}$	$\begin{matrix} (3 \times 1) \\ \text{---} \\ \vdots \end{matrix}$	$\left\{ \begin{array}{l} \omega^{[\ell]} : (n^{[\ell]}, n^{[\ell-1]}) \\ b^{[\ell]} : (n^{[\ell]}, 1) \end{array} \right.$
--	--	--	---

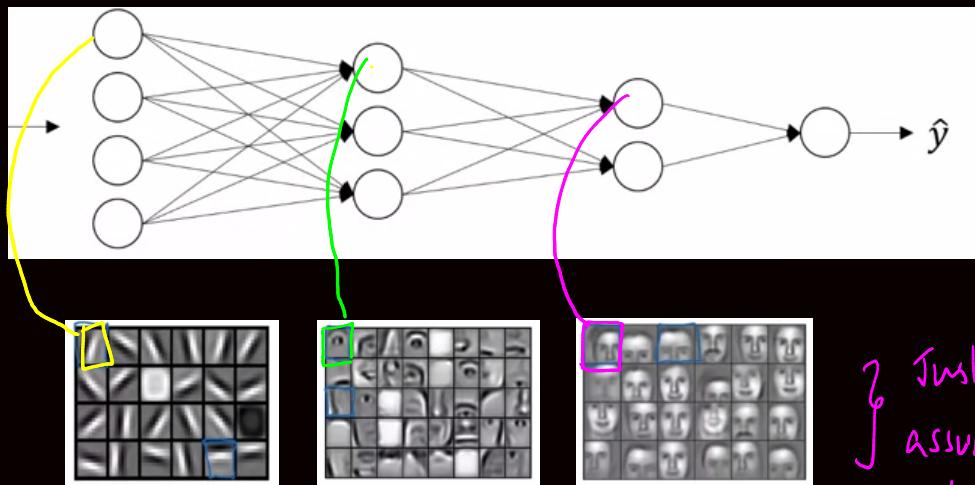
General

$d\omega^{[\ell]} : \text{same as } \omega^{[\ell]}$

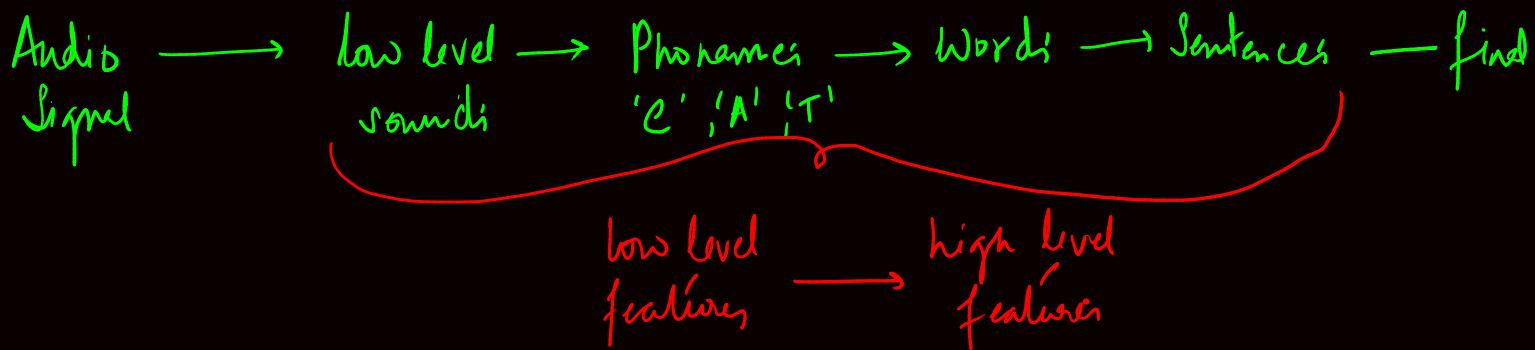
$db^{[\ell]} : \text{same as } b^{[\ell]}$

$Z^{[\ell]}, a^{[\ell]} : (n^{[\ell]}, 1)$

q Why deep representation?



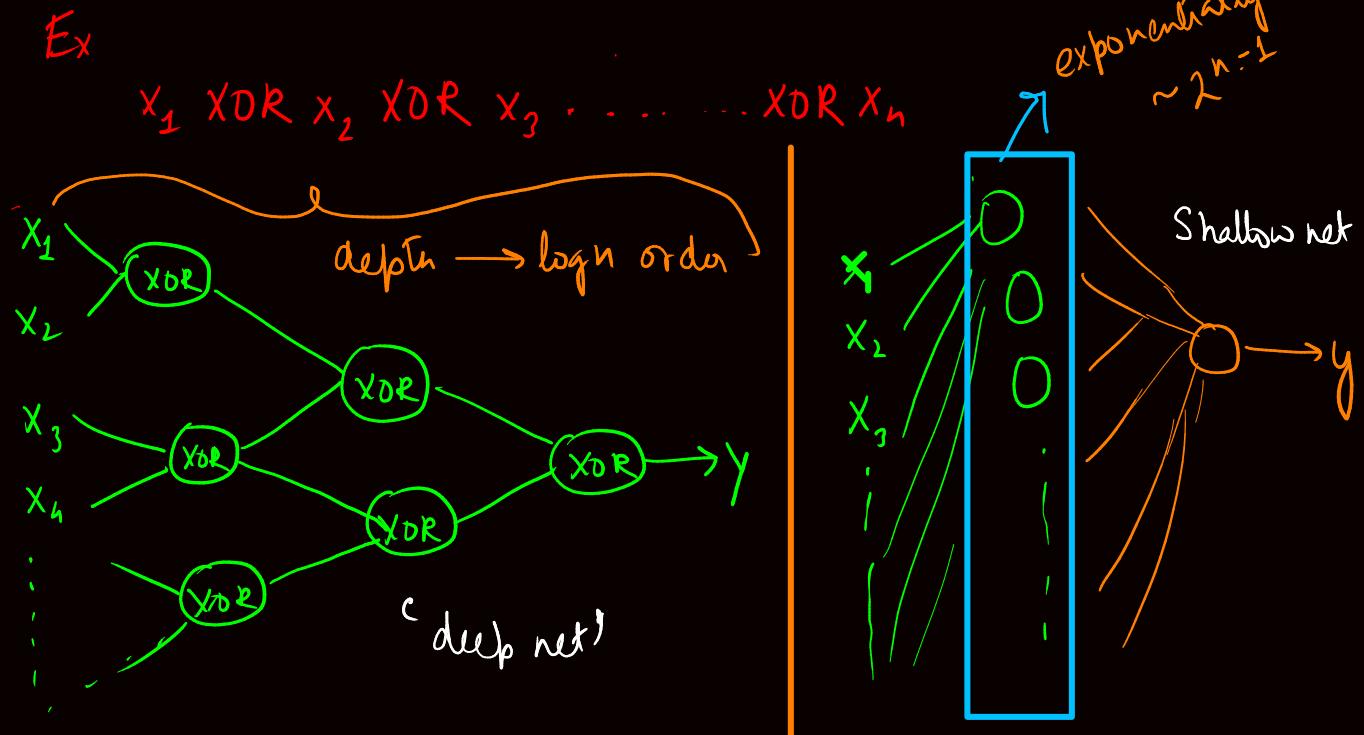
} Just an arbitrary assumption about what a layer does



Circuit theory and deep learning

Informally: There are functions you can compute with a “small” L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

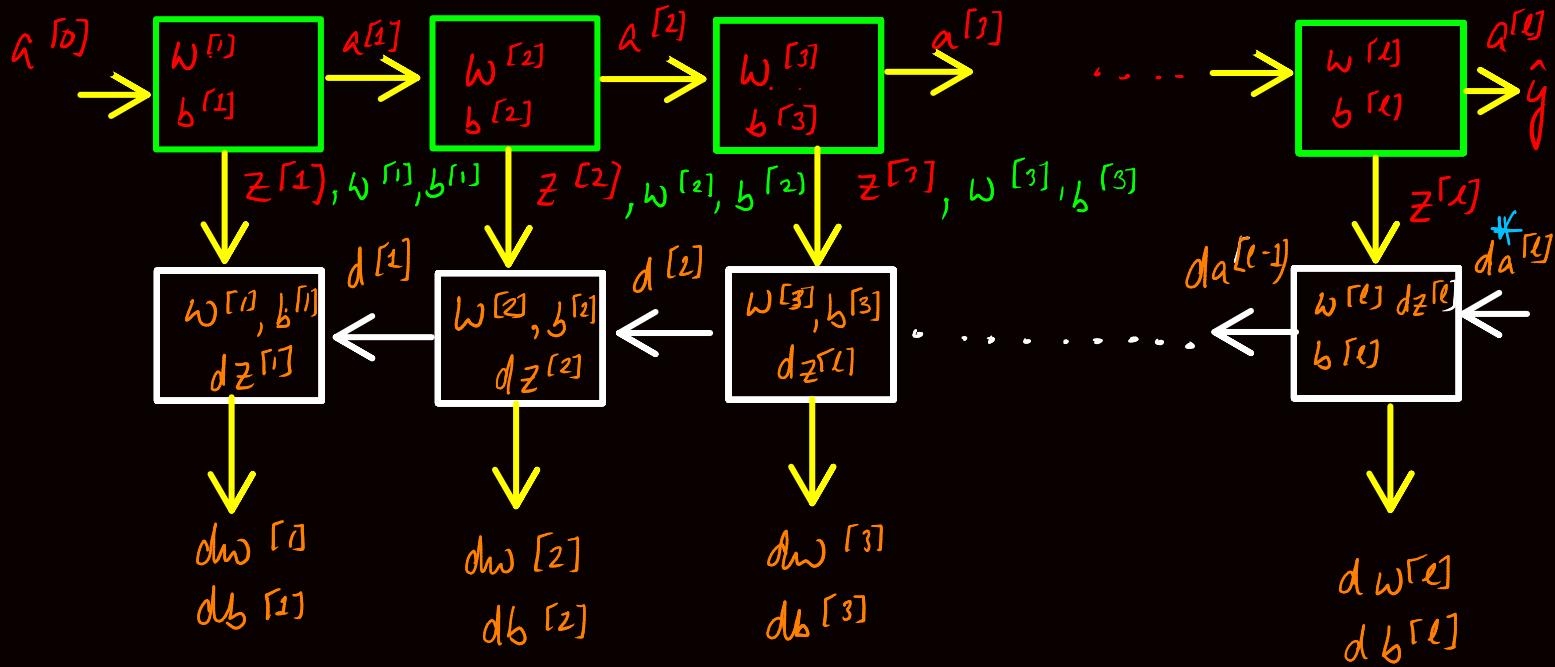
Ex



Building Blocks of Deep Neural Network

* for logistic regression

$$da^{[l]} = \frac{-y}{\alpha} + \frac{(1-y)}{1-\alpha}$$



Formulas:

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

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

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

$$dw^{[l]} = dz^{[l]} * a^{[l-1]T}$$

$$db^{[l]} = dz^{[l]}$$

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

Parameters & Hyperparameters

PARAMETERS: $w^{[l]}, b^{[l]}$

HYPERPARAMETERS: Learning rate, α
 # iterations
 # hidden layer L
 # hidden units, $n^{[1]}, n^{[2]}, \dots$
 Choice of activation funⁿ. etc.

Try with all types of hyperparameters. D.L is an empirical process.