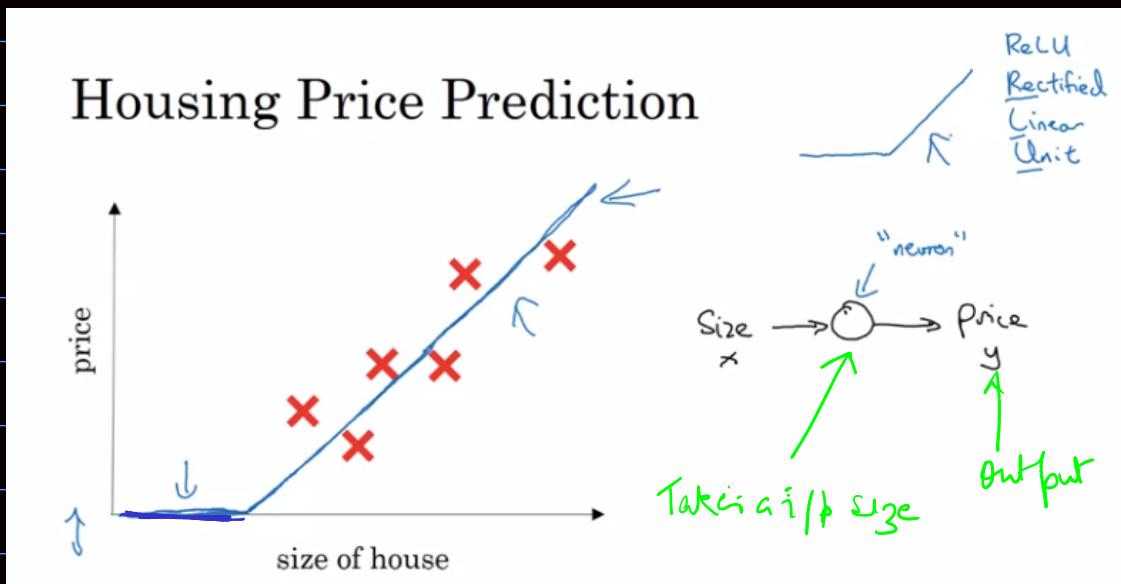
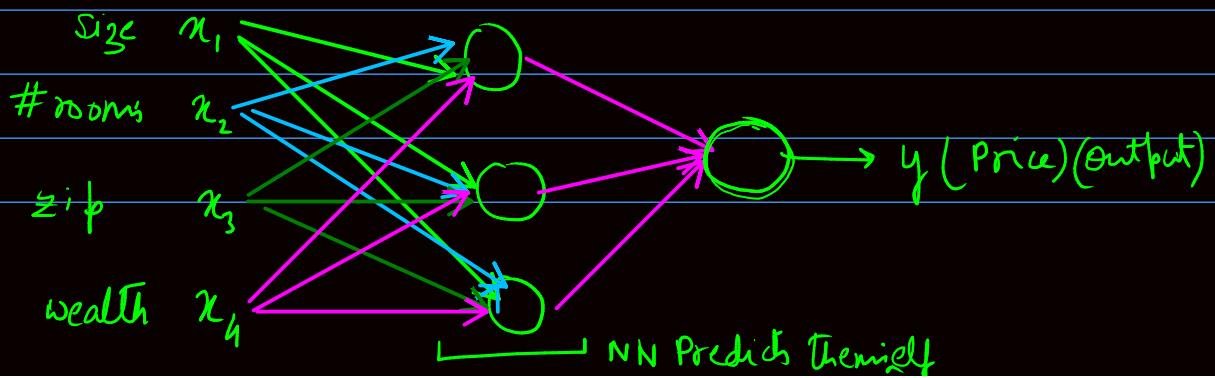
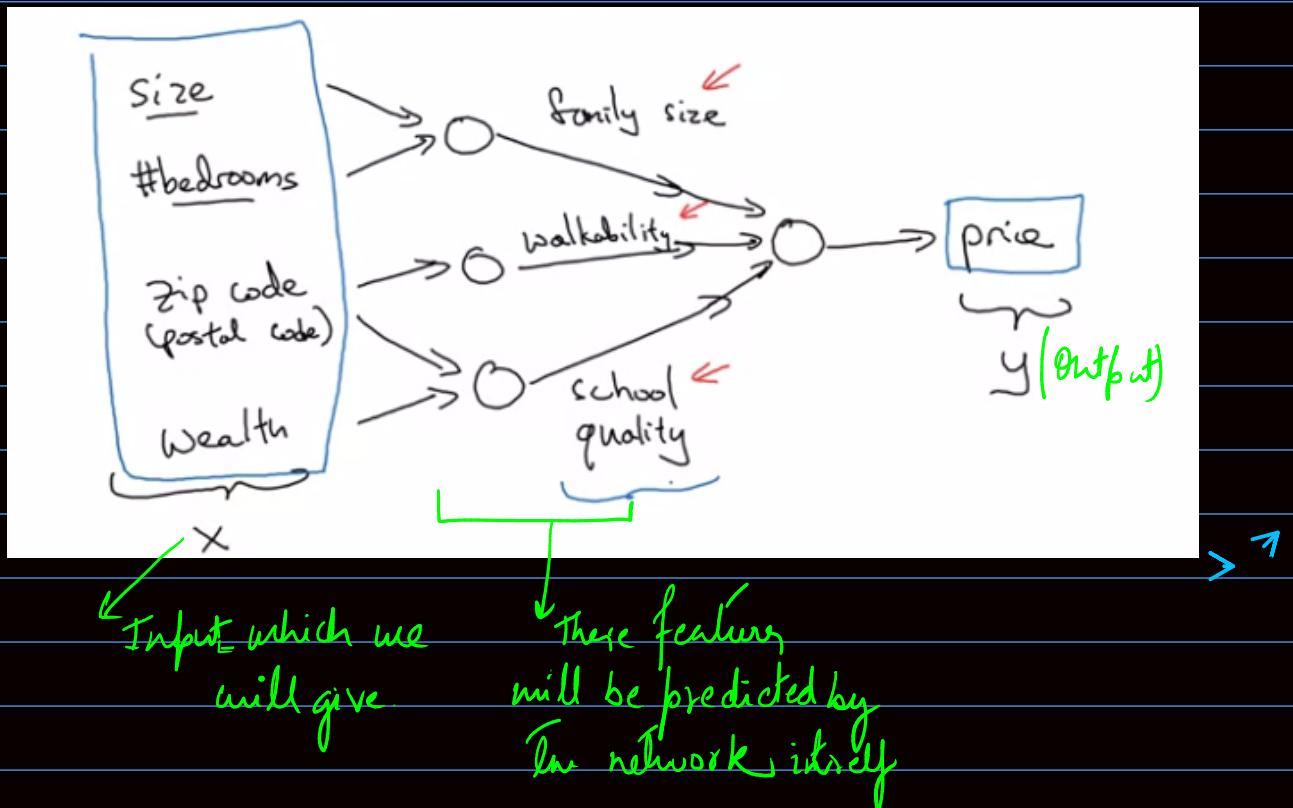


Week 1

Neural Networks & Deep Learning



smallest Possible Neural Network

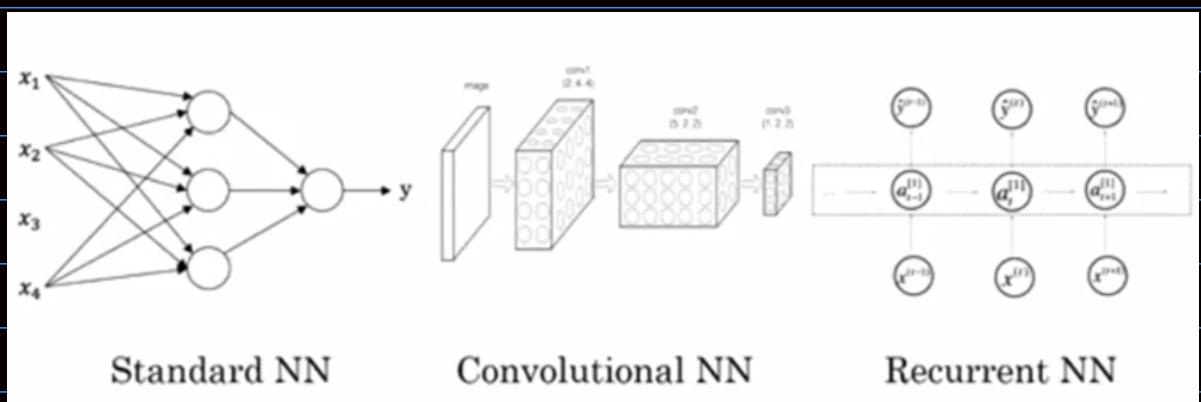


Supervised learning

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

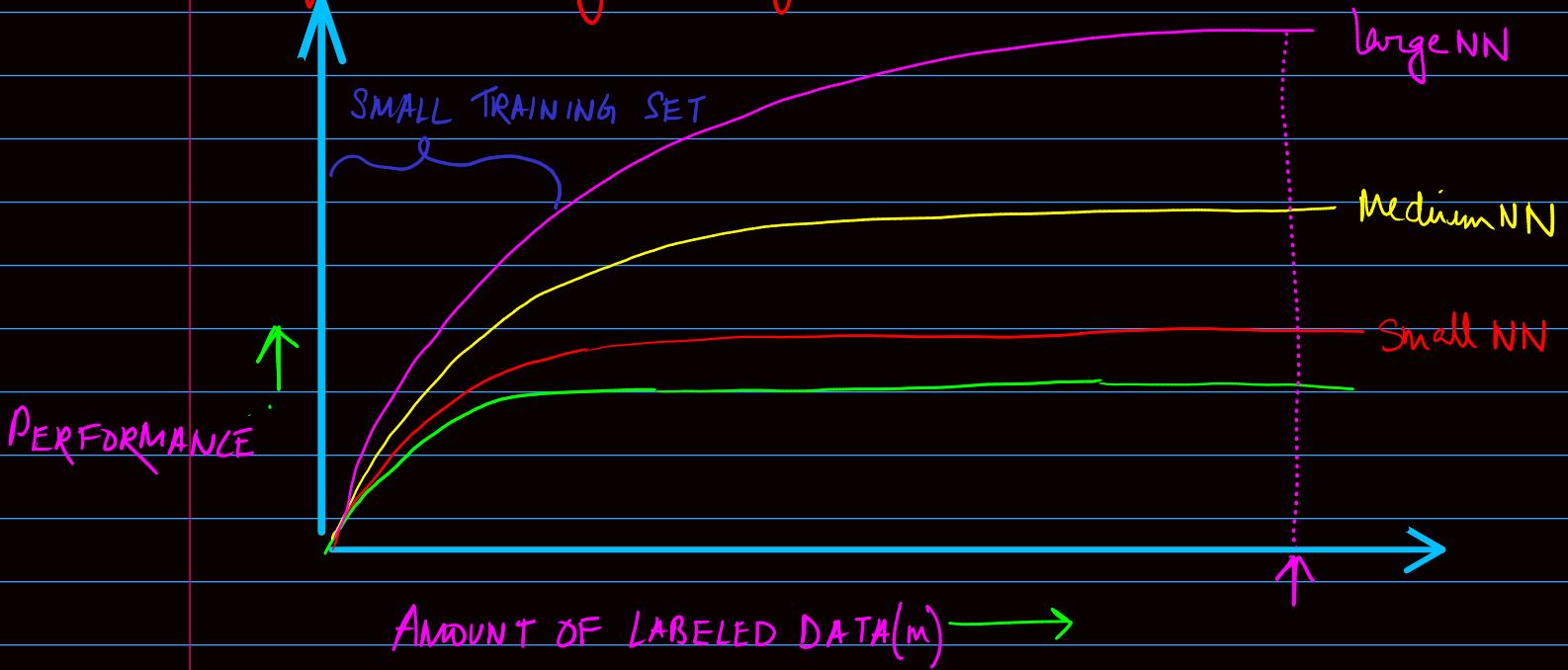


Other Classification

Structured Data				Unstructured Data																									
<table border="1"> <tr> <td>Size</td> <td>#bedrooms</td> <td>...</td> <td>Price (1000\$)</td> </tr> <tr> <td>2104</td> <td>3</td> <td></td> <td>400</td> </tr> <tr> <td>1600</td> <td>3</td> <td></td> <td>330</td> </tr> <tr> <td>2400</td> <td>3</td> <td></td> <td>369</td> </tr> <tr> <td>:</td> <td>:</td> <td></td> <td>:</td> </tr> <tr> <td>3000</td> <td>4</td> <td></td> <td>540</td> </tr> </table>				Size	#bedrooms	...	Price (1000\$)	2104	3		400	1600	3		330	2400	3		369	:	:		:	3000	4		540	 Audio	
Size	#bedrooms	...	Price (1000\$)																										
2104	3		400																										
1600	3		330																										
2400	3		369																										
:	:		:																										
3000	4		540																										
<table border="1"> <tr> <td>User Age</td> <td>Ad Id</td> <td>...</td> <td>Click</td> </tr> <tr> <td>41</td> <td>93242</td> <td></td> <td>1</td> </tr> <tr> <td>80</td> <td>93287</td> <td></td> <td>0</td> </tr> <tr> <td>18</td> <td>87312</td> <td></td> <td>1</td> </tr> <tr> <td>:</td> <td>:</td> <td></td> <td>:</td> </tr> <tr> <td>27</td> <td>71244</td> <td></td> <td>1</td> </tr> </table>				User Age	Ad Id	...	Click	41	93242		1	80	93287		0	18	87312		1	:	:		:	27	71244		1	 Image	
User Age	Ad Id	...	Click																										
41	93242		1																										
80	93287		0																										
18	87312		1																										
:	:		:																										
27	71244		1																										
				<div style="border: 1px solid black; padding: 5px; display: inline-block;"> Four scores and seven years ago... </div>																									
				Text																									

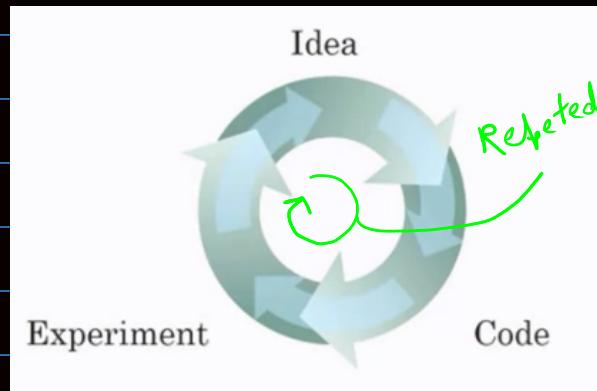
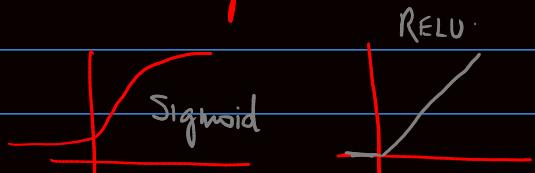
↑
 NN focuses more on
 Unstructured Data

Why deep learning is taking off now?



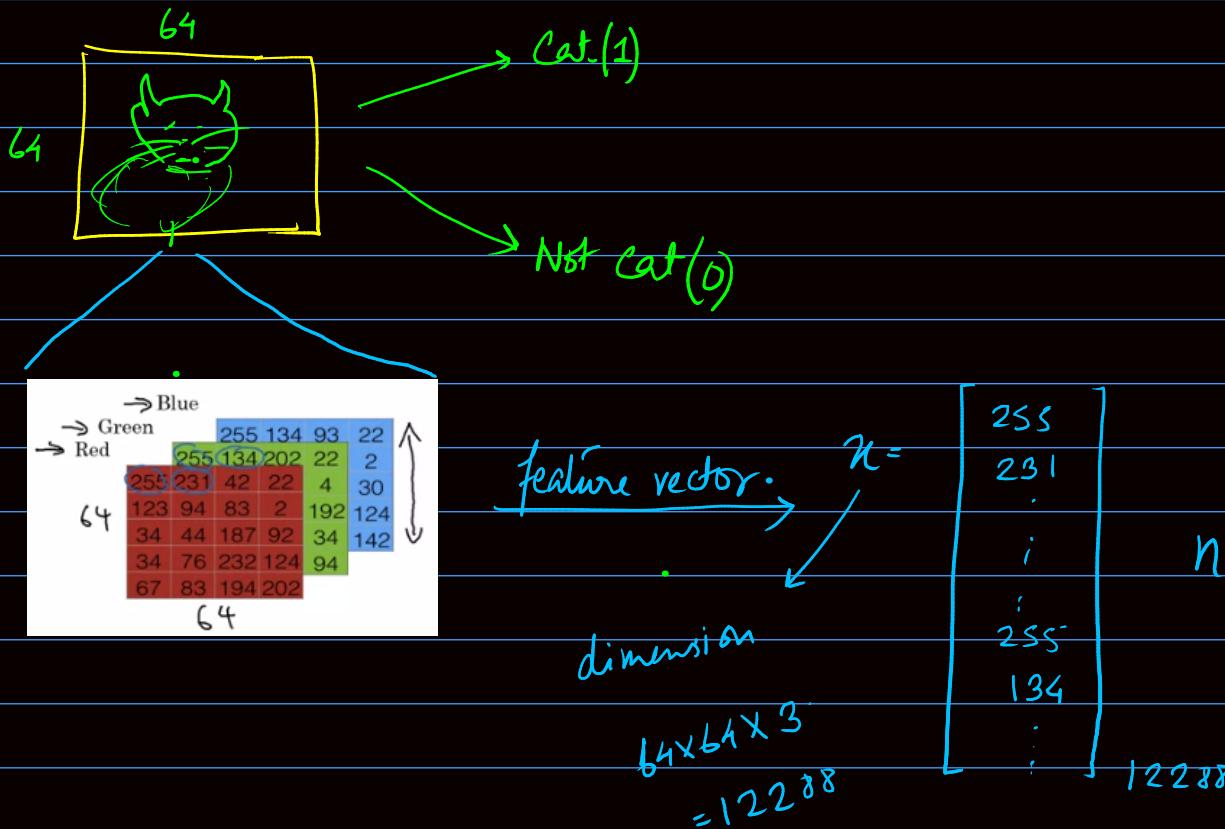
Driving forces:

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



Week 2

BINARY CLASSIFICATION



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

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

$\downarrow 1 \text{ feature vector (1st image)}$

$$\mathbf{X} = \left[\begin{array}{c|c|c|c}
 & & & \\
 \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\
 | & | & & | \\
 & & \ddots & \\
 & & & m
 \end{array} \right]_{n_x \times m}; \quad \mathbf{x} \in \mathbb{R}^{n_x \times m}$$

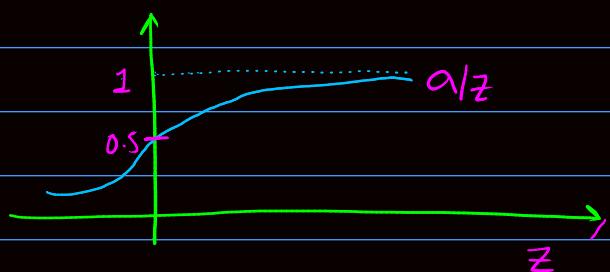
$$\mathbf{Y} = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]_{1 \times m}$$

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

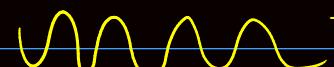
$y=0$; $L(\hat{y}, y) = -\log(1-\hat{y}) \rightarrow \hat{y}$ be be as small as possible

Loss funⁿ \rightarrow Single Training Example

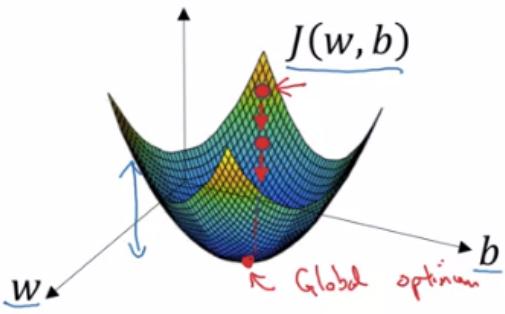
Cost funⁿ: (for the entire training set)

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

* Convex funⁿ  → One global minima

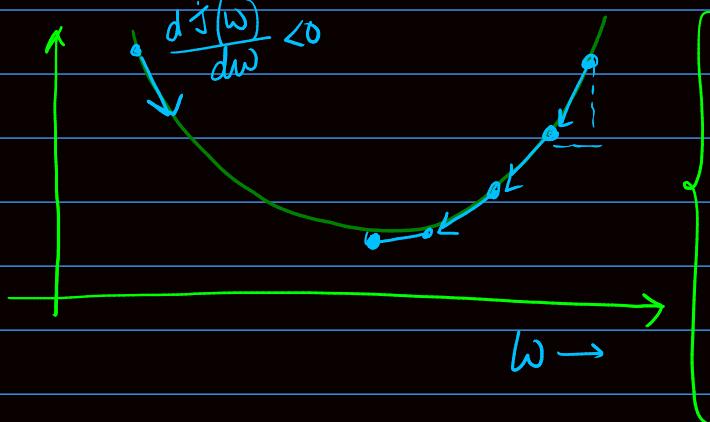
Non convex funⁿ  → Multiple global minima.

Want to find w, b that minimize $J(w, b)$



Gradient Descent

$w := w - \alpha \frac{dJ(w)}{dw}$ learning rate $\frac{dJ(w)}{dw}$

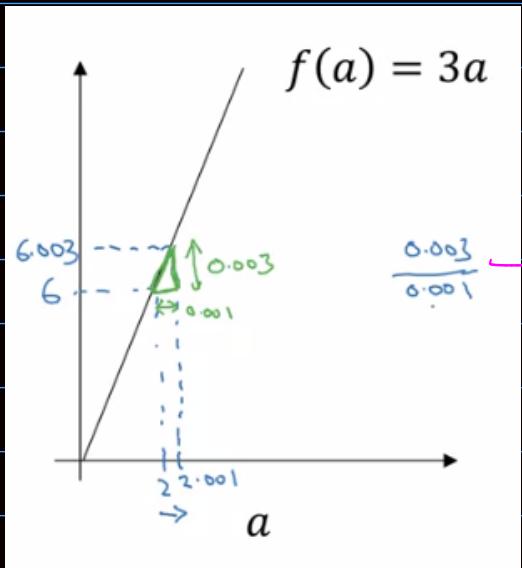


Updates

$w = w - \alpha \frac{dJ(w, b)}{dw}$ dw

$b = b - \alpha \frac{dJ(w, b)}{db}$ db

Derivatives



Slope

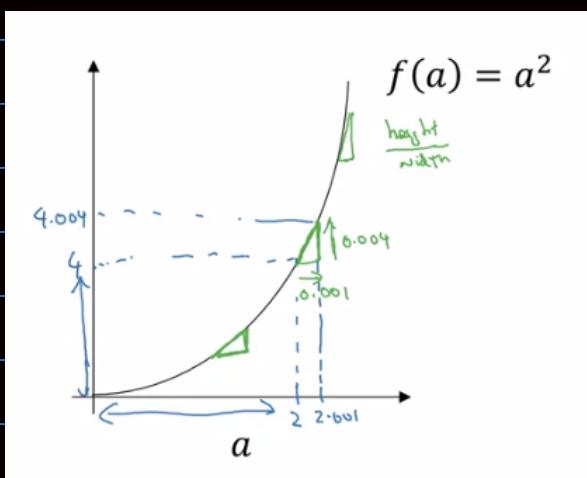
$$\frac{df(a)}{da} = 3 = \frac{d}{da} f(a)$$

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

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

$$a = 5.001, f(a) = 15.003$$

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

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

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

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

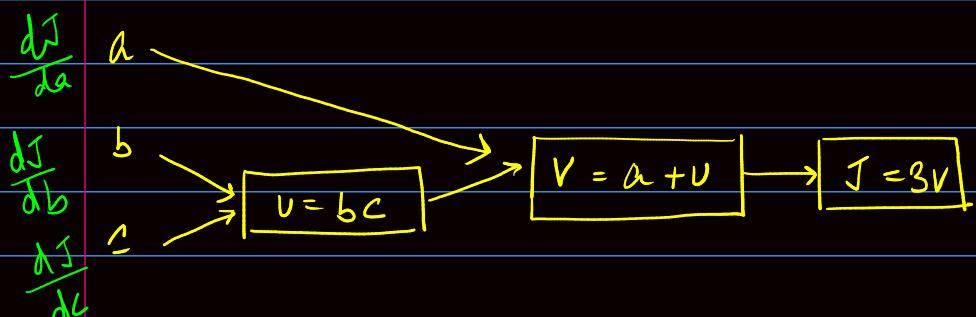
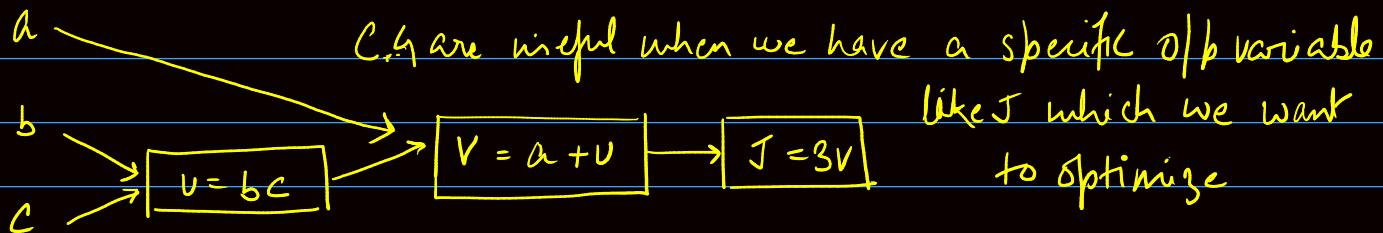
$$\frac{df(a)}{da} = \frac{1}{a}$$

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

COMPUTATION GRAPH

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

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



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

$$V = 11 \rightarrow 11.001 \{ \text{xyz} \}$$

$$J = 33 \rightarrow 33.003 \{ \text{3xyz} \}$$

$$\therefore \frac{dJ}{dV} = \frac{d(3V)}{dV} = 3$$

(Chain Rule)

$$\frac{dJ}{da} = a = 5 \rightarrow 5.001$$

$$V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.001$$

a effect \rightarrow V effect \rightarrow J

$$\frac{dJ}{da} = \frac{dJ}{dV} \times \frac{dV}{da}$$

$$= 3 \times 1$$

Final O/p Variable to optimize = J

$$\frac{d(\text{Final O/p Variable})}{d(Van)} = dJ_{dvan} \leftarrow "dvan"$$

Smart Notation for o/p

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

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

$$b = 3 \rightarrow 3.001$$

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

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

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

$$= 6$$

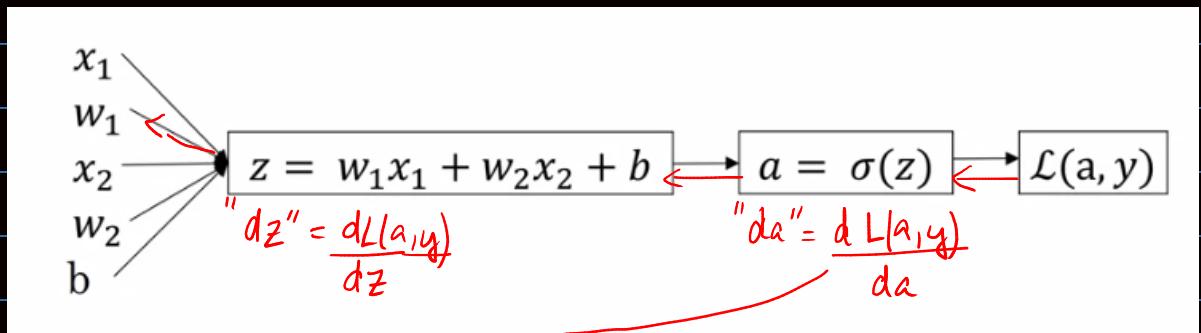
$$\frac{\partial J}{\partial c} = 6 \quad (\text{Same})$$

* 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

Recap $\hat{y} = a = \sigma(z)$ $L(a, y) = -(y \log a + (1-y) \log(1-a))$



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

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

$$= (a - y)$$

$$*\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w_1}$$

$(a - y) \times n_1$
 $\underbrace{dz}_{dZ} \times n_1$

$$\therefore \frac{\partial L}{\partial w_1} = "dw_1" = n_1 dz, \frac{\partial L}{\partial w_2} = n_2 dw_2, \frac{dL}{db} = dz$$

Updates learning rate

$$w_1 = w_1 - \alpha dw_1; \quad w_2 = w_2 - \alpha dw_2; \quad b = b - \alpha db$$

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

Steps

① Initialization

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

For $i=1$ to m

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

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

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

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

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

This increases in only 2
features w_1 & w_2

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

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

$$\int / = m$$

$$d\omega_1 / = m$$

$$d\omega_2 / = m$$

$$db / = m$$

This technique is very tedious

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

VECTORIZATION

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

In python

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

```
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

for j

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

} Not Vectorized

\textcircled{2}

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

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

$\rightarrow u = np.zeros((n, 1))$
for i in range(n):
 $u[i] = \text{math.exp}(v[i])$

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

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

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

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

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

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

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

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

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

$$\frac{J}{m}$$

$$\frac{d\omega_1}{m} = m$$

$$\frac{d\omega_2}{m} = m$$

$$\frac{db}{m} = m$$

$$\frac{d\omega}{m} = m$$

VECTORIZING LOGISTIC REGRESSION

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

$$X = \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \\ | & | & | \end{bmatrix}_{(n_x \times m)} \quad w = \begin{bmatrix} \vdots \\ w \\ \vdots \end{bmatrix}$$

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

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

$$Z = np.\text{dot}(w.T, X) + b \quad \begin{array}{l} \xrightarrow{\text{Result broadcasted to } (1 \times m) \text{ matrix}} \\ \text{to add to the previous term} \end{array}$$

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

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

$$dz = [\delta z^{(1)} \ \delta z^{(2)} \ \dots \ \delta z^{(m)}]$$

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(n)}] \quad Y = [y^{(1)} \ \dots \ y^{(n)}]$$

$$\delta Z = A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots \ a^{(n)} - y^{(n)}]$$

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

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

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

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

* Code for $\overrightarrow{\text{iter}} \in \text{range}(1000)$:

$$\left. \begin{array}{l} Z = \text{np.dot}(w.T, X) + b \\ A = \sigma(Z) \end{array} \right\}$$

$$\delta Z = A - Y$$

$$dW = \frac{1}{m} X \times \delta Z^T$$

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

$$w = w - \alpha dW$$

$$b = b - \alpha db$$

Vectorized imp of Gradient Descent

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

↓ 94%.

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)           ↗ reshape command is redundant here
print(percentage)                                     but good to be safe
[[94.91525424  0.          2.83140283  88.42652796]
 [ 2.03389831 43.51464435 33.46203346 10.40312094]
 [ 3.05084746 56.48535565 63.70656371  1.17035111]]
```

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

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

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

etc..

$$\begin{array}{ccc} * & (m, n) & + \quad (1, n) \rightsquigarrow (m, n) \\ & \text{matrix} & * \quad (n, 1) \rightsquigarrow (m, n) \\ & & / \end{array}$$

* 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 = 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

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

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

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

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

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

$$\log(p/y|x) = \underset{\text{to minimize the loss}}{\underset{\uparrow}{L(\hat{y}, y)}}$$

Cost on m examples

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

By maximum likelihood estimation

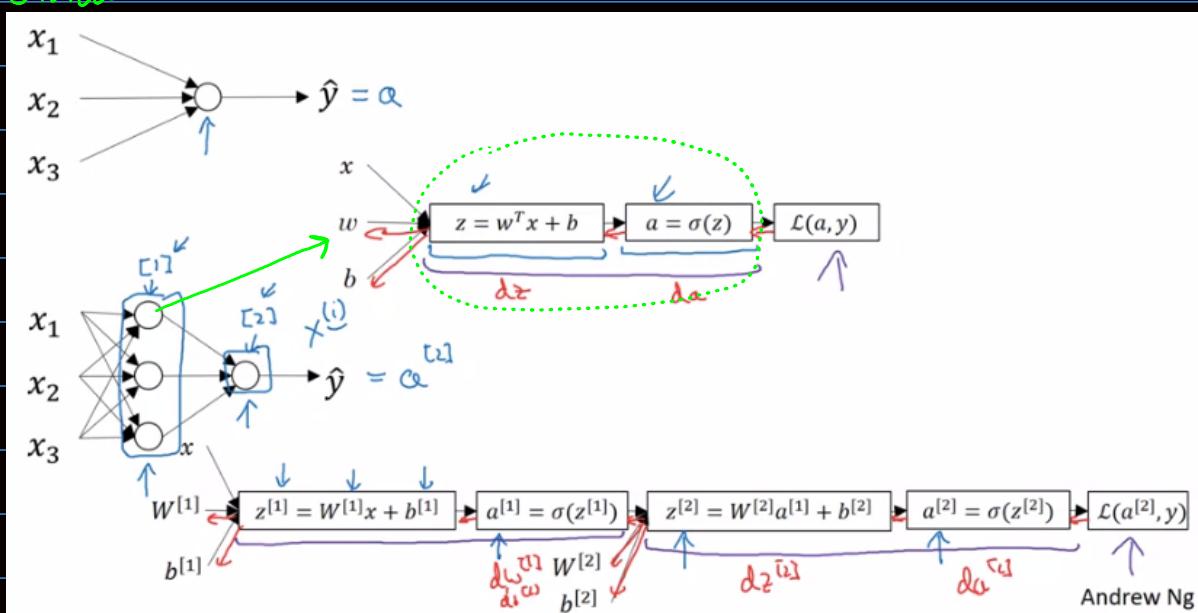
$$= -\sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

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

minimize //

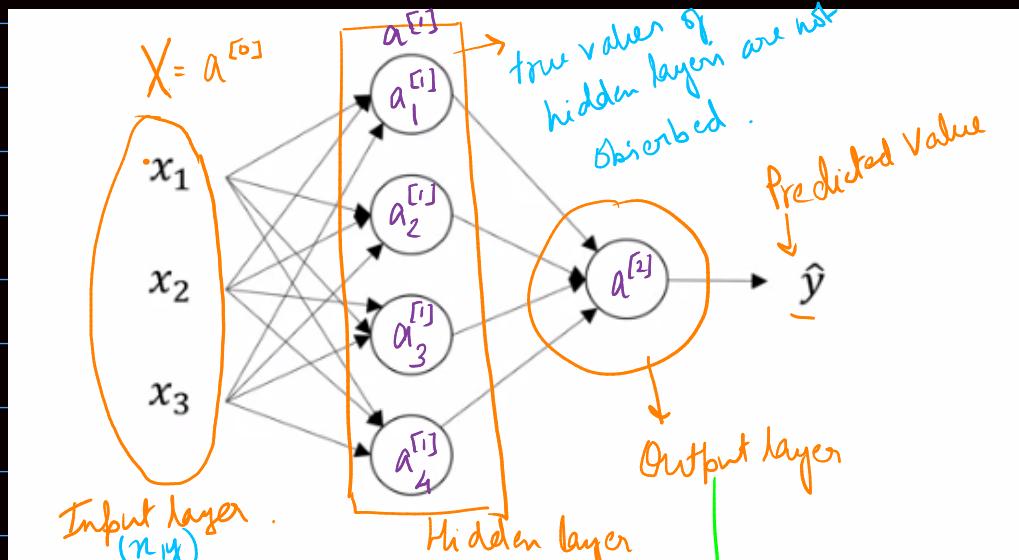
WEEK 3

Overview



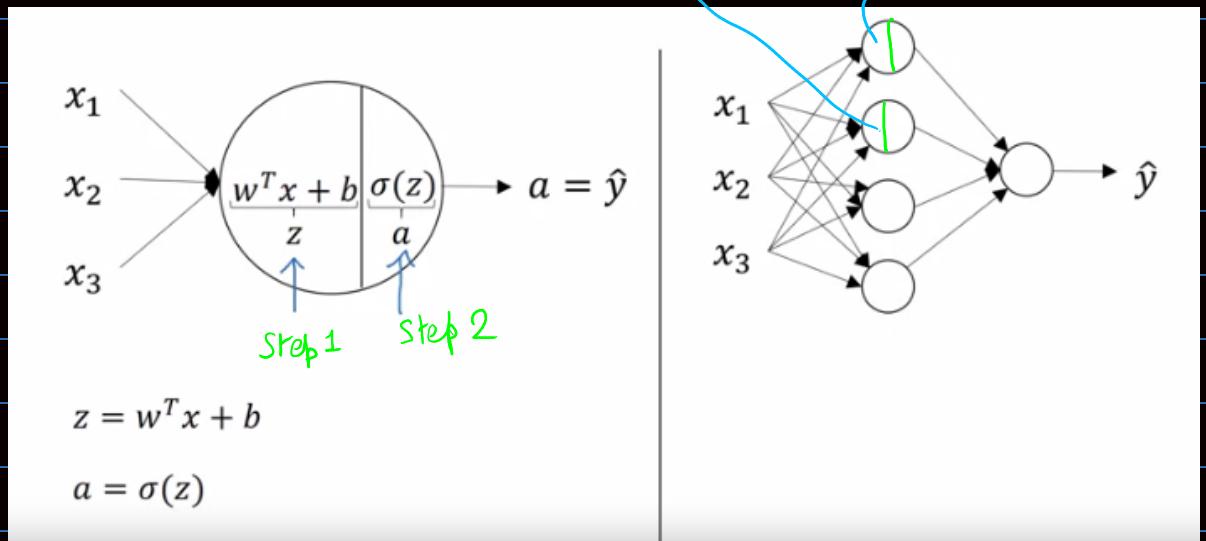
Andrew Ng

2-layer Network [if p excluded]



$$\begin{aligned} & \text{Input} \rightarrow \text{label} \\ & \text{if } p \text{ is not } \\ & \text{label} \quad \downarrow \\ & w^1, b^1 \quad \downarrow \\ & (4, 3) \quad (3, 1) \\ & \downarrow \\ & z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]} \\ & a_1^{[1]} = \sigma(z_1^{[1]}) \\ & \vdots \\ & z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]} \\ & a_4^{[1]} = \sigma(z_4^{[1]}) \end{aligned}$$

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



$$Z = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} \quad Z^{[1]}$$

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

$$Z = \begin{bmatrix} -w_1^{[1]T} - \\ -w_2^{[1]T} - \\ -w_3^{[1]T} - \\ -w_4^{[1]T} - \end{bmatrix} \begin{bmatrix} n_1 \\ n_2 \\ n_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}_{3 \times 1} \quad W^{[1]}_{4 \times 3}$$

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

∴

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.

$a^{[2]}(i) \rightarrow$ i^{th} Training example
layer 2

Vectorized Solution

training example

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

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

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

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

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

where, no. of training examples

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

No. of training examples

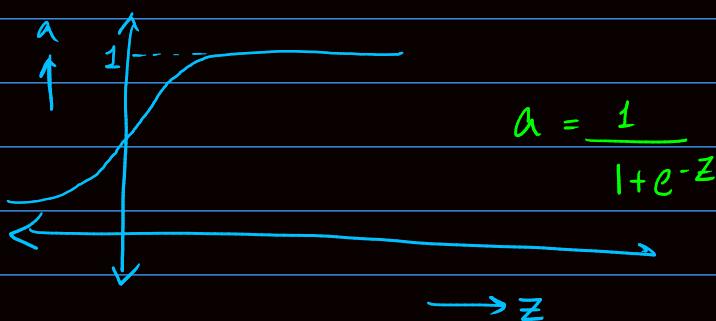
$$A^{[1]} = \begin{bmatrix} | & | & | & | \\ a_1^{1} & a^{[1](2)} & \dots & a^{[1](m)} \\ a_2^{1} \\ \vdots \end{bmatrix}$$

No. of nodes in that particular layer.

ACTIVATION FUNCTIONS

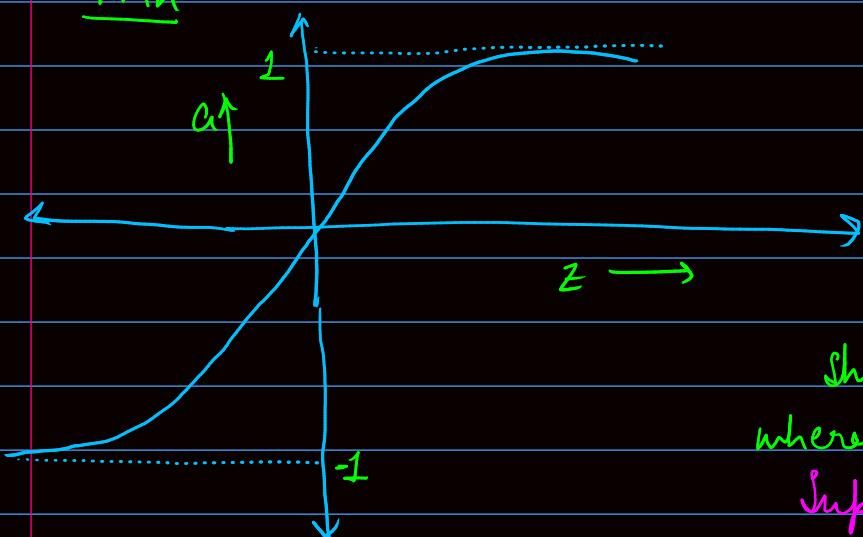
Sigmoid

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

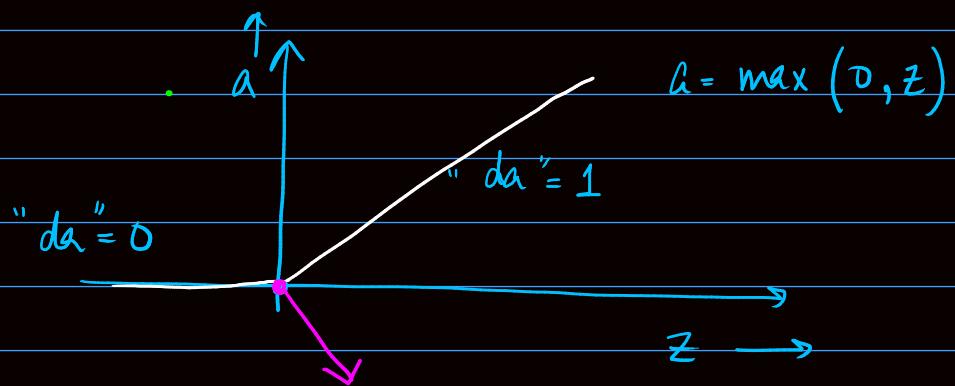
Tanh



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

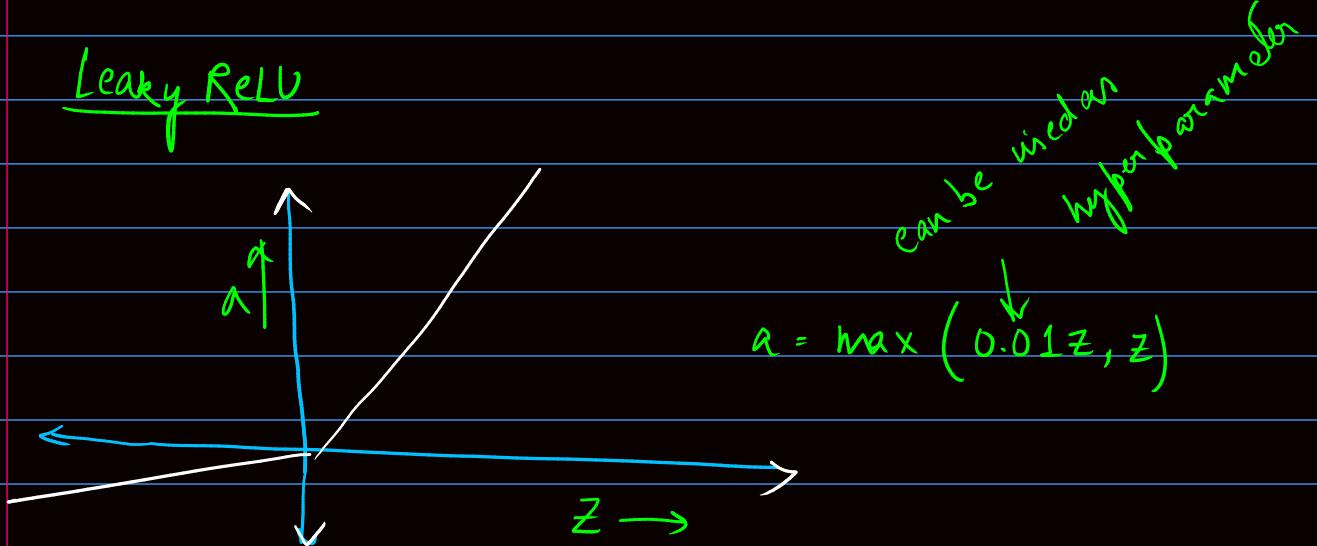
ReLU (Rectified Linear Unit)



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

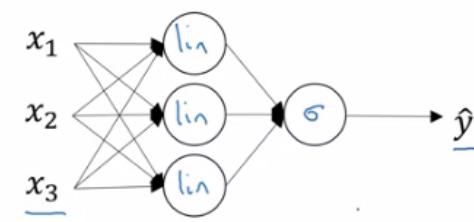
- * If our model is for binary classification, let the opp layer activation fun by sigmoid & all the rest layers, ReLU must be the default choice.

Leaky ReLU



Why non-linear activation fun?

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



Given x :

$$\begin{aligned} \rightarrow z^{[1]} &= W^{[1]}x + b^{[1]} \\ \rightarrow a^{[1]} &= g^{[1]}(z^{[1]}) \geq^{c^{[1]}} \\ \rightarrow z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\ \rightarrow a^{[2]} &= g^{[2]}(z^{[2]}) \geq^{c^{[2]}} \end{aligned}$$

$g(z) = z$
"linear activation function"

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

Andrew Ng

One place where we use linear activation funⁿ age regression problem (ML)

Eg Predict movie prices

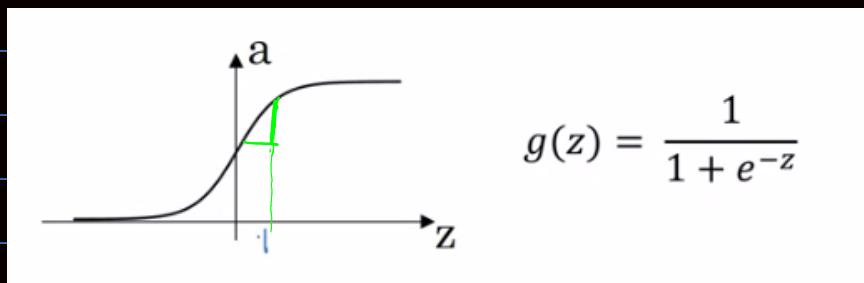
$$y \in \mathbb{R} \{ \$0, \dots, \$100000 \}$$

Use linear here

But all other hidden layers must use ReLU

Derivative of Activation funⁿ

① sigmoid

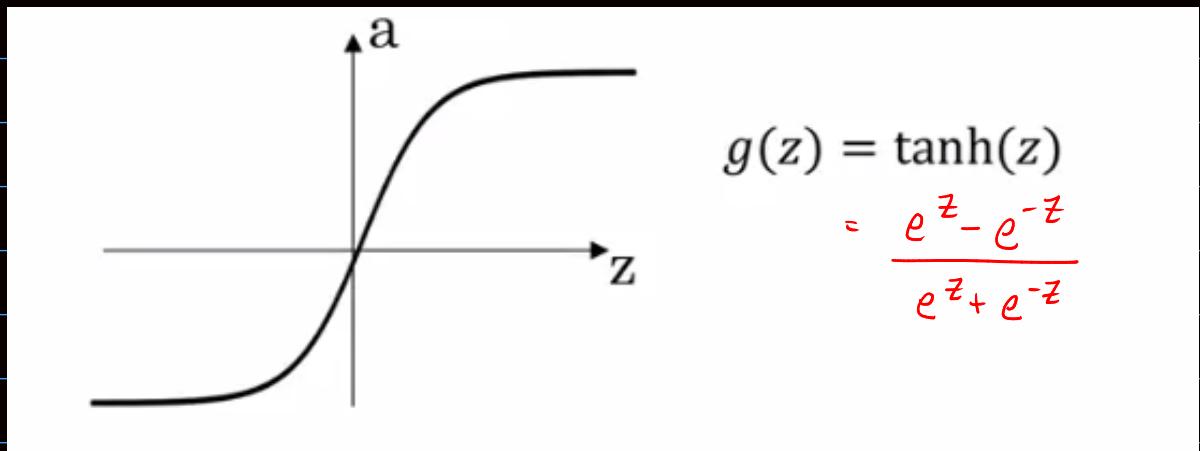


$$g(z) = \frac{1}{1 + e^{-z}}$$

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

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

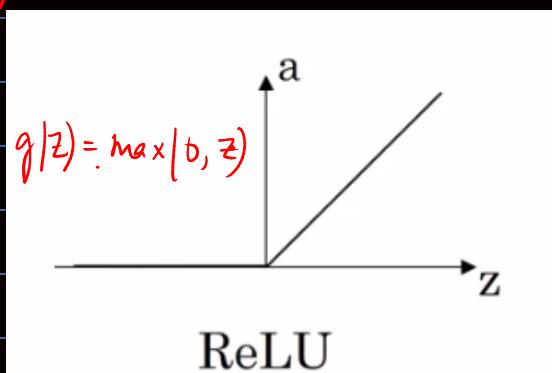
② Tanh



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

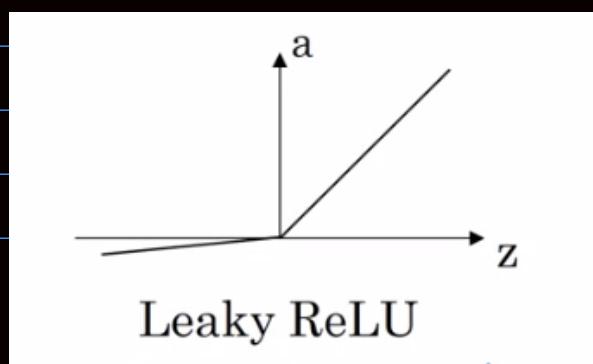
$$a = g(z) = g'(z) = 1 - a^2$$

③ ReLU



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

④ Leaky ReLU



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

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

Repeat

Till parameter convergence
 ① Compute predictions = $(\hat{y}^{(i)}, i = 1, 2, \dots, m)$

$$② \frac{dJ}{dw^{[1]}} = d\omega^{[1]}, d\omega^{[2]}, db^{[1]}, db^{[2]}$$

$$③ w^{[1]} = w^{[1]} - \alpha d\omega^{[1]} ; b^{[1]} = b^{[1]} - \alpha db^{[1]} \\ w^{[2]} = w^{[2]} - \alpha d\omega^{[2]} ; b^{[2]} = b^{[2]} - \alpha db^{[2]}$$

}

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

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

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

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

Automatic casting changes
any rank array to
a vector

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

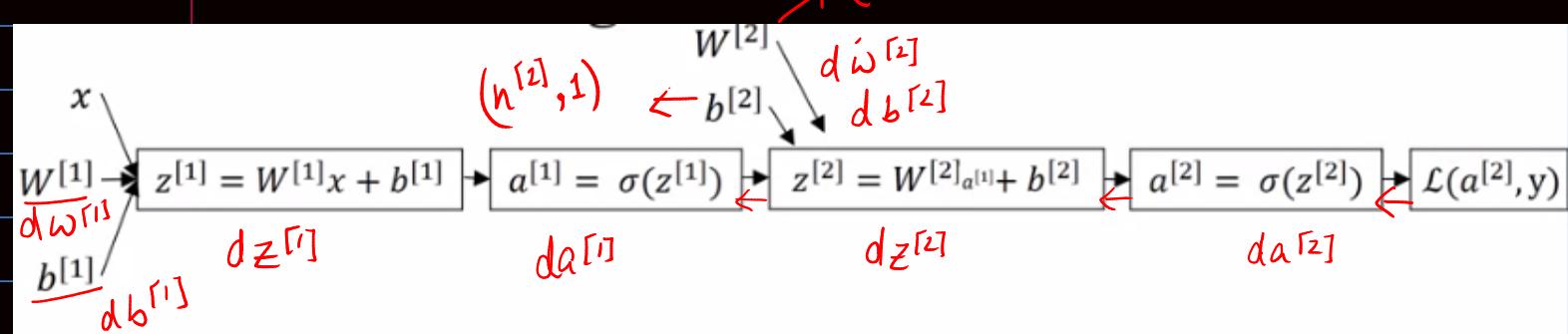
elementwise product

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

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

Intuition behind the
Back Prop formulae

ω & $d\omega \rightarrow$ same dimension



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

$$d\omega^{[2]} = dz^{[2]} a^{[1]T} \approx \begin{cases} d\omega = dz \cdot c \end{cases}$$

Logistic Regm.

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

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

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

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

Summary of gradient descent

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

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

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

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

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

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

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

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

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

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

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

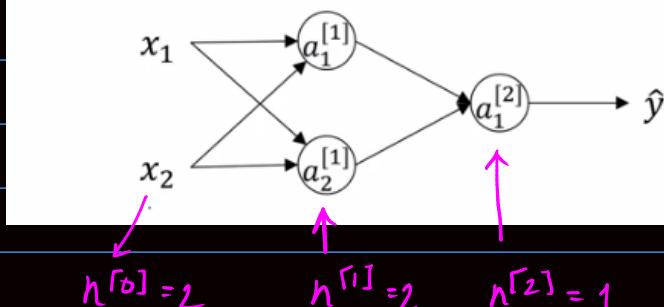
$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$$

Vectorized Implementation

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

RANDOM INITIALIZATION

What happens if you initialize weights to zero?



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

Then $a_1^{[1]} = a_2^{[1]}$ and $d z_1^{[1]} = d z_2^{[1]}$

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

$$\partial w_1 = \partial w_2$$

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:

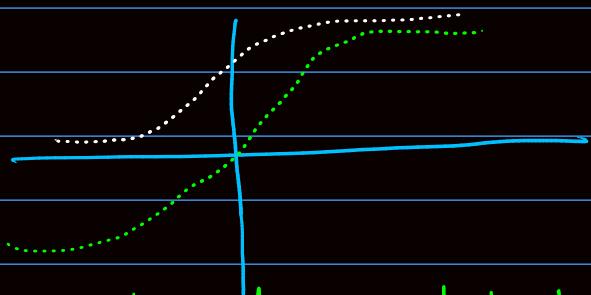
$$w^{[1]} = np.random.randn(2, 2) * 0.01 \quad \text{to initialize } w \text{ by very small values}$$

$$h^{[1]} = np.zeros((2, 1)) \quad \{ \text{No symmetry issue} \}$$

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

$$b^{[1]} = \dots$$

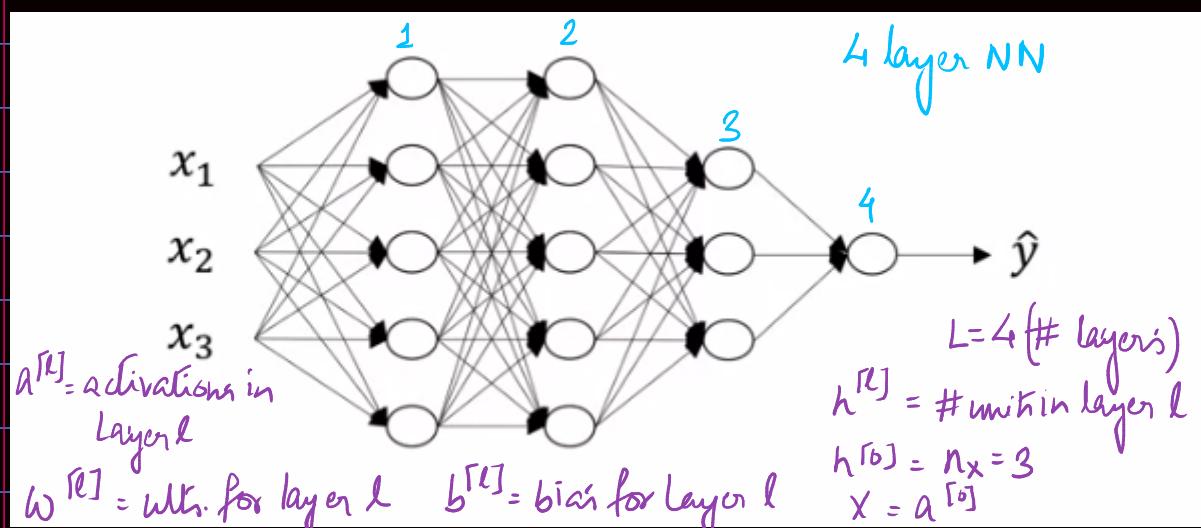
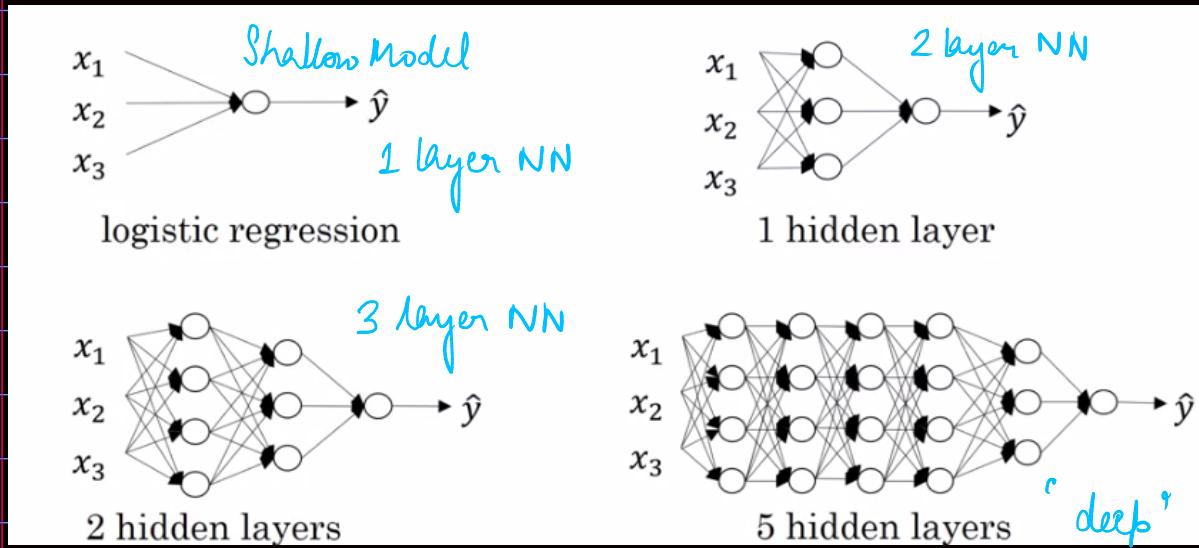
* If w are large $g(x)$ will be big & we will end up in a flat slope if we use tanh or sigmoid as $g(x)$ fun



* 0.01 can be used as a hyperparameter

Week 4

Deep L-Layer neural network



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

Forward Propagation

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

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

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

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

$$\vdots \quad \vdots$$

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

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

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

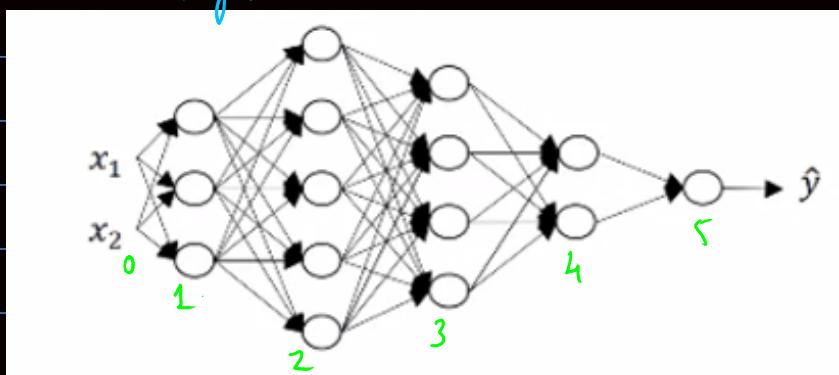
Vectorized

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



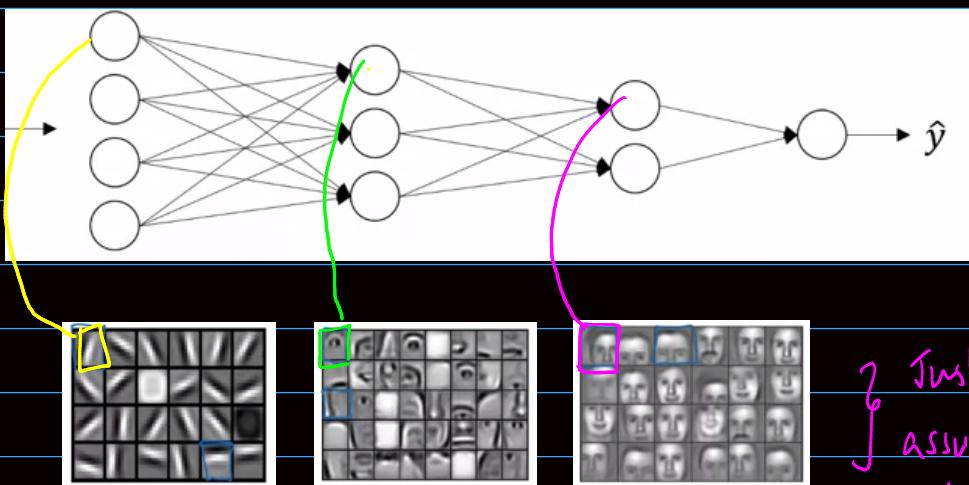
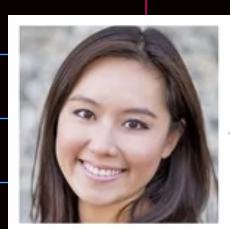
Eg. $Z^{[1]} = \omega^{[1]} X + b^{[1]}$

3×1	(3×2)	(2×1)	\downarrow	(3×1)	{
\vdots					

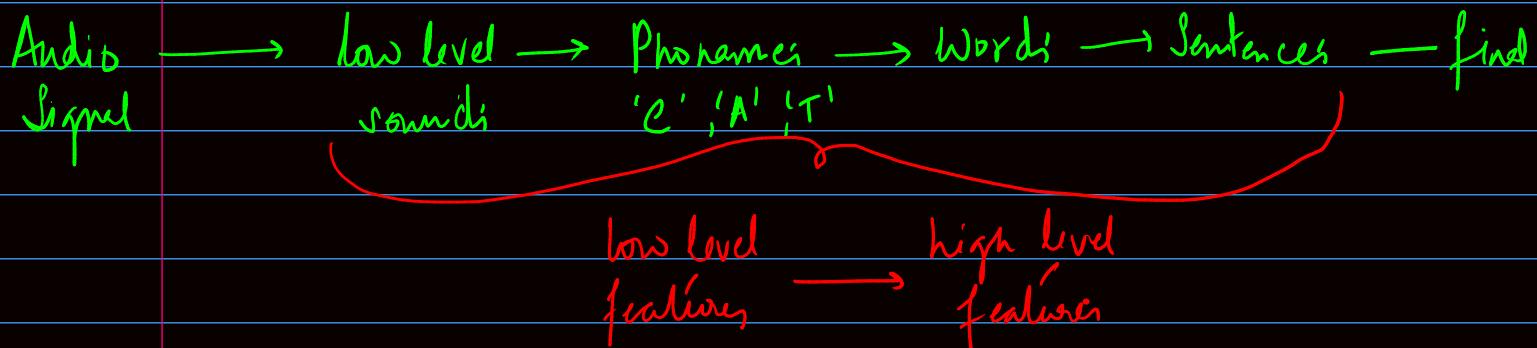
General

$\omega^{[\ell]} : (n^{[\ell]}, n^{[\ell-1]})$
 $b^{[\ell]} : (n^{[\ell]}, 1)$
 $d\omega^{[\ell]} : \text{same as } \omega^{[\ell]}$
 $db^{[\ell]} : \dots \dots b^{[\ell]}$
 $Z^{[\ell]}, a^{[\ell]} : (n^{[\ell]}, 1)$

Q Why deep representation?



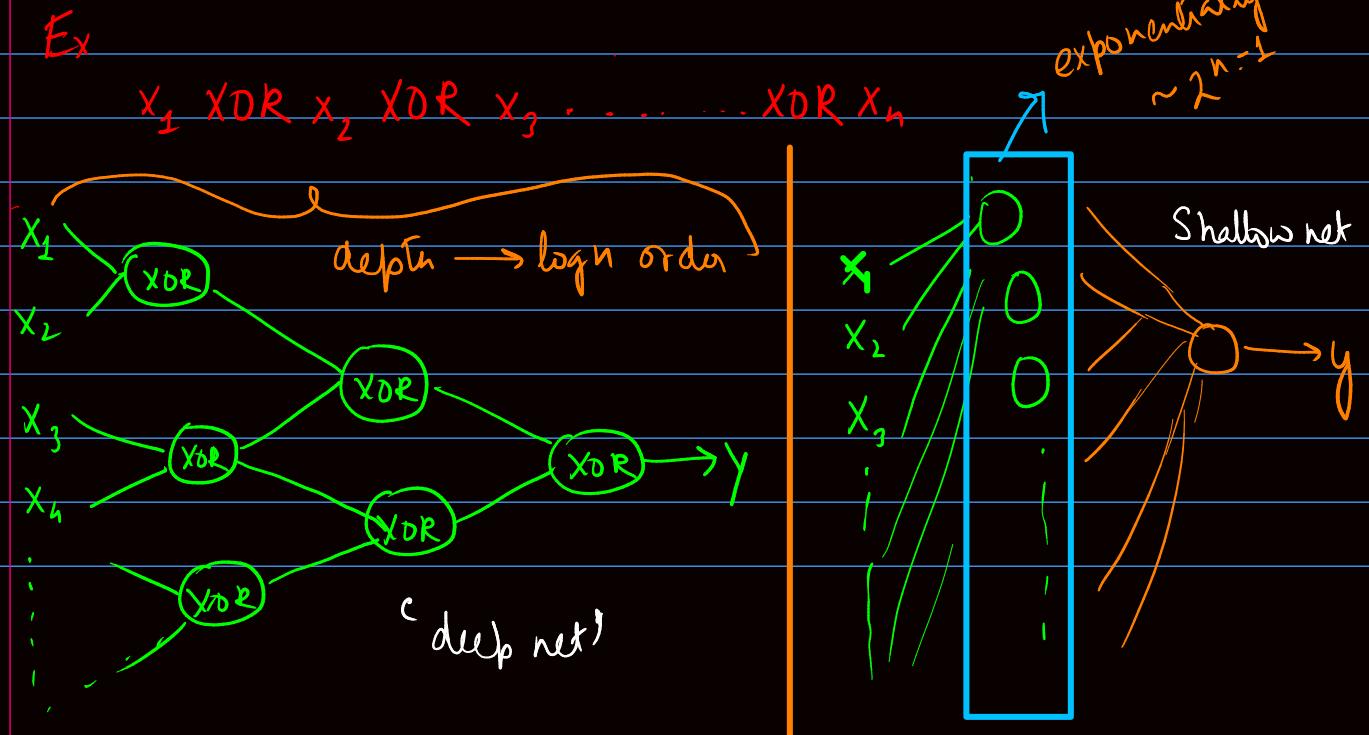
? Just an arbitrary assumption about what a layer does



Circuit theory and deep learning

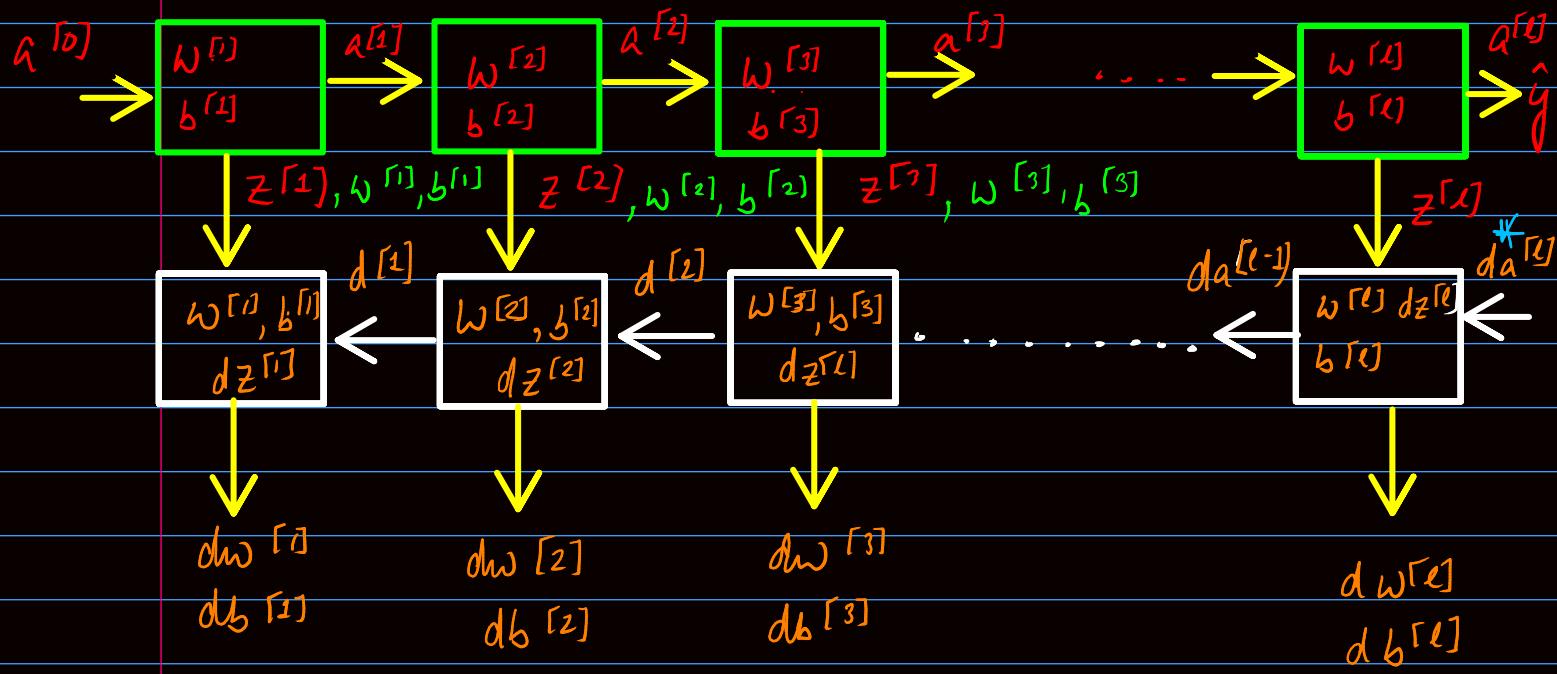
Informally: There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

Ex



Building Blocks of Deep Neural Network

$$\text{for logistic regression} \quad * \quad d\alpha[t] = \frac{-y}{\alpha} + \frac{(1-y)}{1-\alpha}$$



Formulas:

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

$$dw^{[l]} = dz^{[l]} * a^{[l-1]}' (z^{[l]})$$

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

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

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

Parameters & Hyperparameters

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

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

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