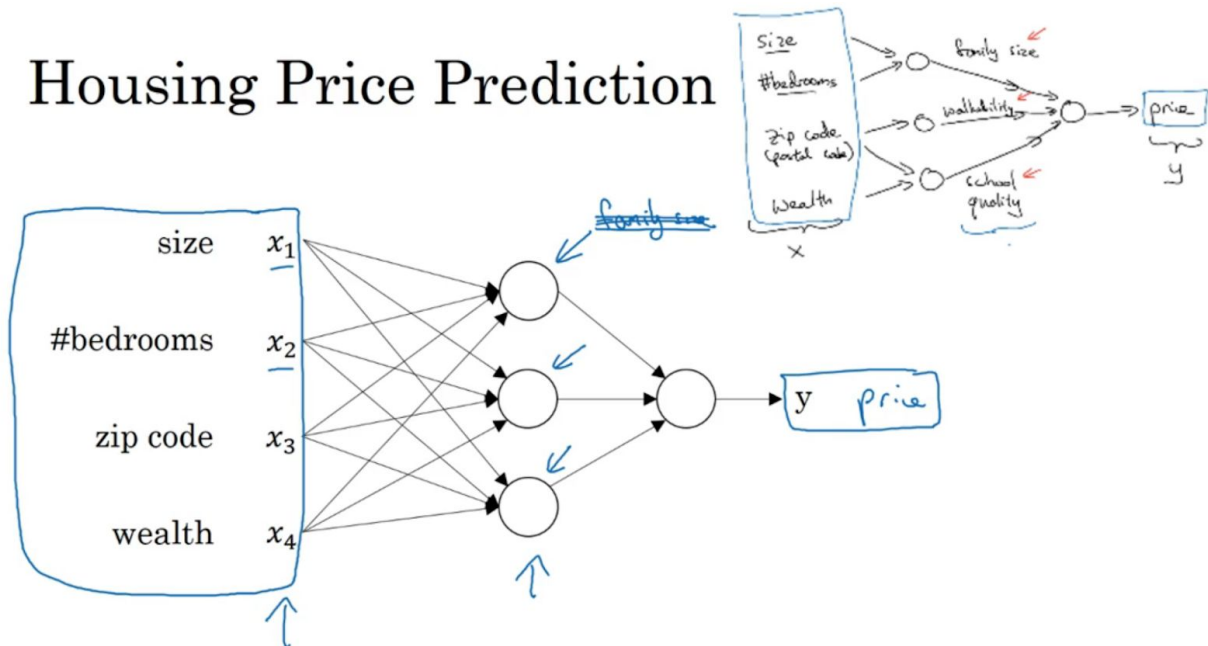


Course 1: NEURAL NETWORKS AND DEEP LEARNING

Housing Price Prediction



Supervised Learning

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
<u>Audio</u>	Text transcript	Speech recognition
<u>English</u>	Chinese	Machine translation
<u>Image, Radar info</u>	Position of other cars	Autonomous driving

} *Standard NN*
 } *CNN*
 } *RNN*
 } *Custom/Hybrid*

Supervised Learning

Structured Data

Size	#bedrooms	...	Price (1000\$)
2104	3		400
1600	3		330
2400	3		369
⋮	⋮		⋮
3000	4		540

User Age	Ad Id	...	Click
41	93242		1
80	93287		0
18	87312		1
⋮	⋮		⋮
27	71244		1

Unstructured Data



Audio



Image

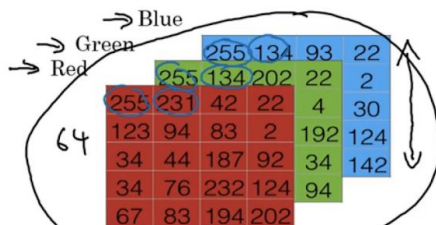
Four scores and seven
years ago...

Text

Binary Classification



64



→ 1 (cat) vs 0 (non cat)
y

$$X = \begin{bmatrix} 255 \\ 134 \\ 93 \\ 22 \\ \vdots \\ 255 \\ 134 \end{bmatrix}$$

$$64 \times 64 \times 3 = 12288$$

$$n = n_x = 12288$$

$$X \rightarrow y$$

Notation

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

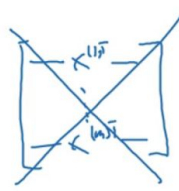
$$m \text{ training examples: } \{(\underline{x}^{(1)}, \underline{y}^{(1)}), (\underline{x}^{(2)}, \underline{y}^{(2)}), \dots, (\underline{x}^{(m)}, \underline{y}^{(m)})\}$$

$$M = M_{\text{train}}$$

$$M_{\text{test}} = \# \text{test examples.}$$

$$X = \begin{bmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix} \quad \begin{matrix} \uparrow \\ n_x \\ \downarrow \end{matrix}$$

\xleftarrow{m}

$$X \in \mathbb{R}^{n_x \times m} \quad X.\text{shape} = (n_x, m)$$


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

$$Y \in \mathbb{R}^{1 \times m}$$

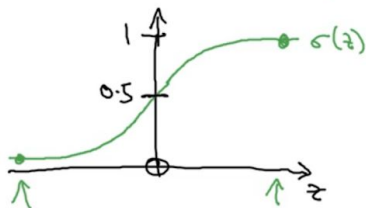
$$Y.\text{shape} = (1, m)$$

Logistic Regression

Given x , want $\hat{y} = P(y=1|x)$
 $x \in \mathbb{R}^{n_x}$ $0 \leq \hat{y} \leq 1$

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

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



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

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

If z large $\sigma(z) \approx \frac{1}{1+0} = 1$

If z large negative number

$$\sigma(z) = \frac{1}{1 + e^{-z}} \approx \frac{1}{1 + \text{Big num}} \approx 0$$

Andrew I

Logistic Regression cost function

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

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$.

$x^{(i)}$
 $y^{(i)}$
 $z^{(i)}$

\bar{c} - the
example.

Loss (error) function:

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

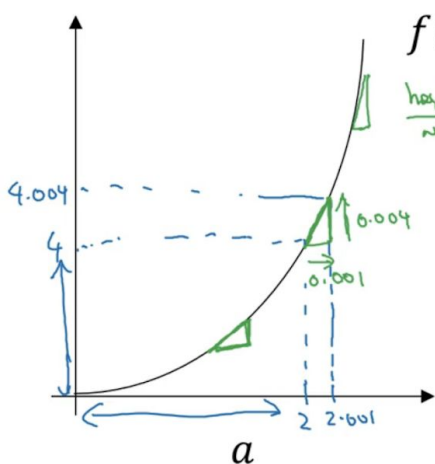
$$\mathcal{L}(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y})) \leftarrow$$

If $y=1$: $\mathcal{L}(\hat{y}, y) = -\log \hat{y} \leftarrow$ want $\log \hat{y}$ large, want \hat{y} large.

If $\frac{y}{y=0}$: $L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow$ want $\log(1-\hat{y})$ large ... want \hat{y} small

Cost function: $J(a, b) = \frac{1}{n} \sum_{i=1}^n \ell(y^{(i)}, \hat{y}^{(i)}) = \frac{1}{n} \sum_{i=1}^n [y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})]$

Intuition about derivatives



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

$$\frac{d}{da} a^2 = 2a$$

0.001
(2a).

$$a = 2$$

$$Q = 2.001$$

$$f(a) = 4$$

$$f(a) \approx 4.004$$

 (4.004004)

slope (derivative) of $f(x)$ at $x=2$ is 4.

$$\frac{d}{da} f(a) = 4 \quad \text{when } a=2.$$

$$a = 5$$

$$a = 5.601$$

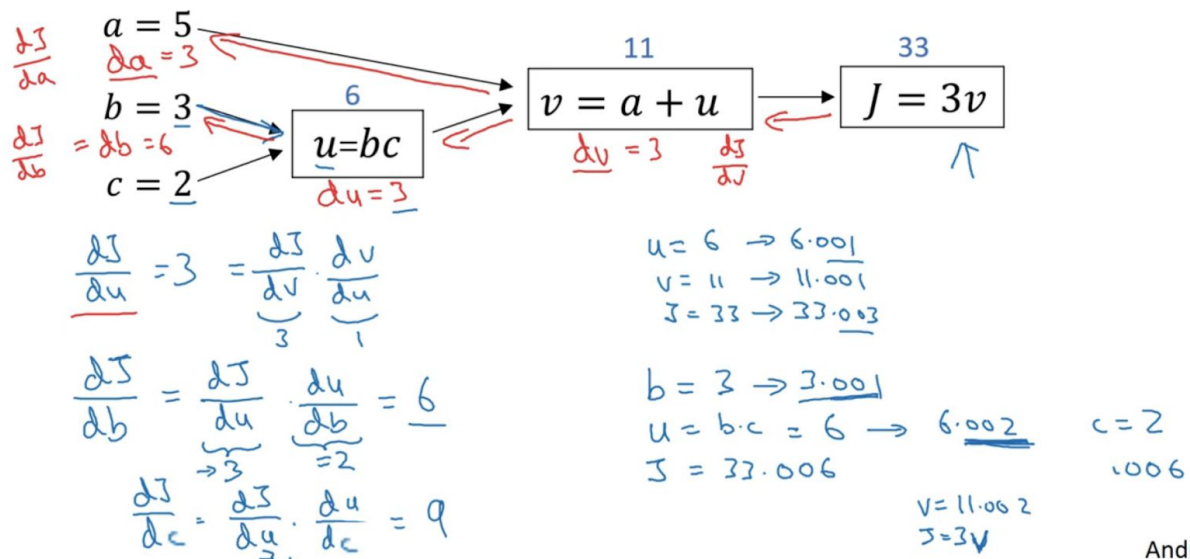
$$f(w) = 25$$

$$f(u) \approx \underline{25.010}$$

$$\frac{d}{da} f(a) = 10 \quad \text{when} \quad a = 5$$

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

Computing derivatives



Andrew N

VECTORIZATION:

Avoid for loops and just use dot function -- this will allow better parallel use of MPU power.

What is vectorization?

$$z = \omega^T x + b$$

Non-vectorized:

$$z = 0$$

for i in $\text{range}(n-x)$:

$$z += \omega[i] * x[i]$$

$$z += b$$

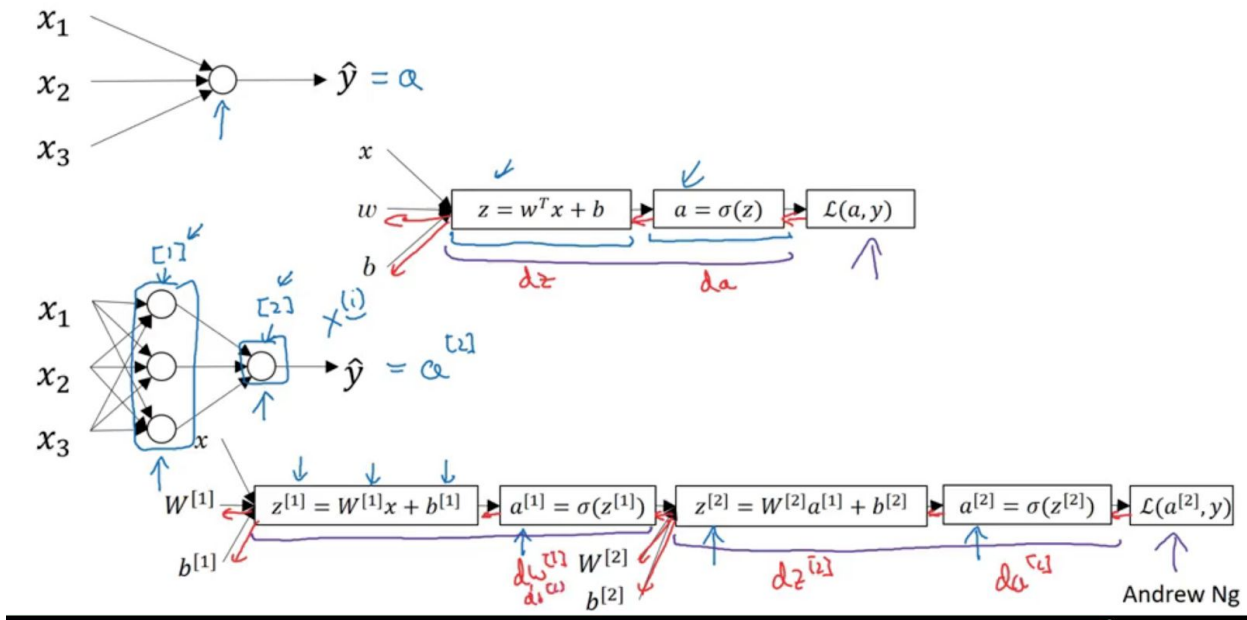
$$\omega = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad \omega \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x}$$

Vectorized

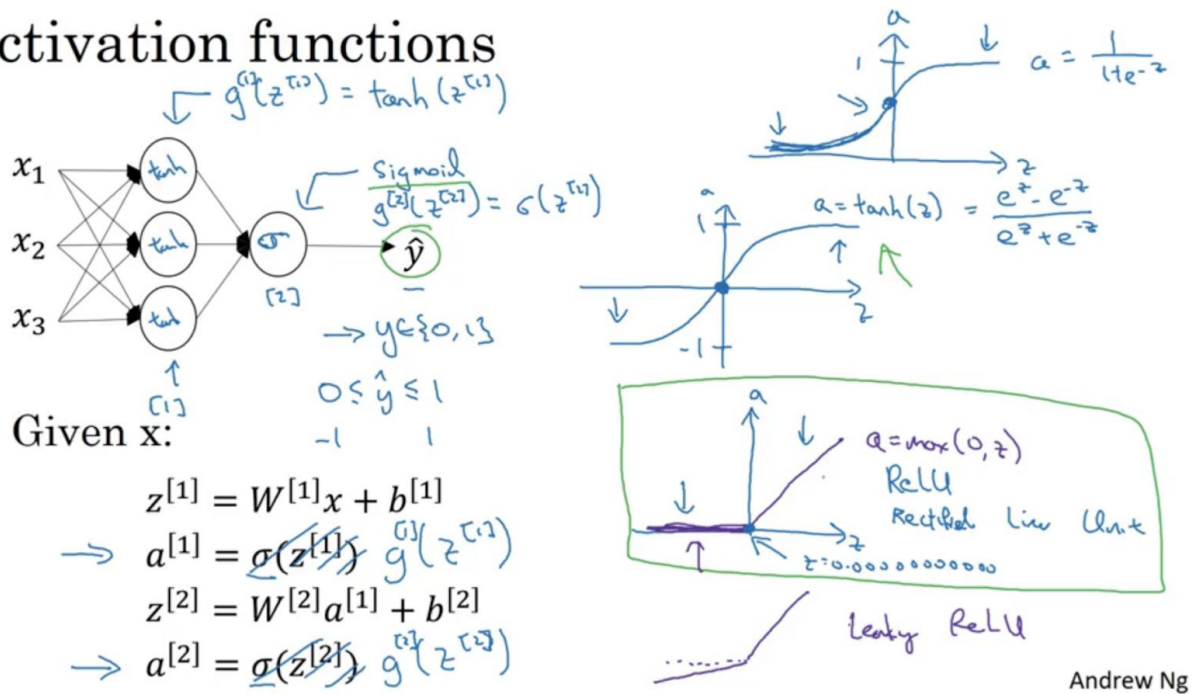
$$z = \text{np.dot}(\omega, x) + b$$

GPU } SIMD - single instruction
CPU } multiple data.

What is a Neural Network?

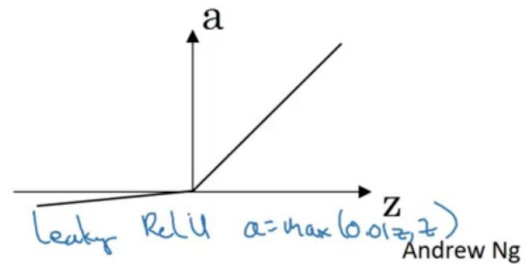
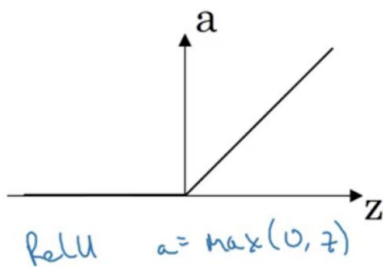
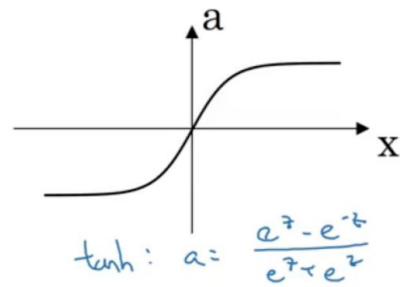
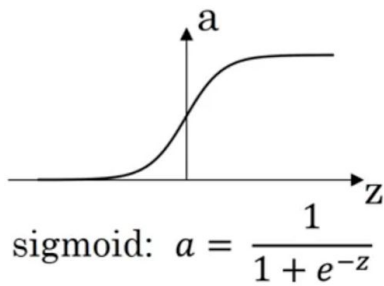


Activation functions



Relu is really good - no vanishing gradient.

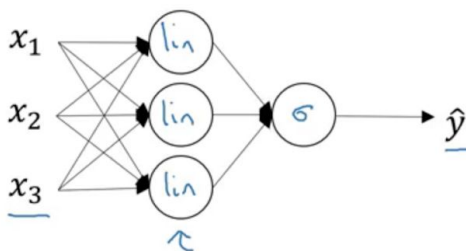
Pros and cons of activation functions



if only linear activations then its no better to use hidden layers, all linear and last sigmoid is just logistic regression. To learn interesting functions you need non linearity.

Exception: Use linear output neuron for regression problems but inner hidden layers are all sigmoid or tanh or relu

Activation function



Given x :

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

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

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

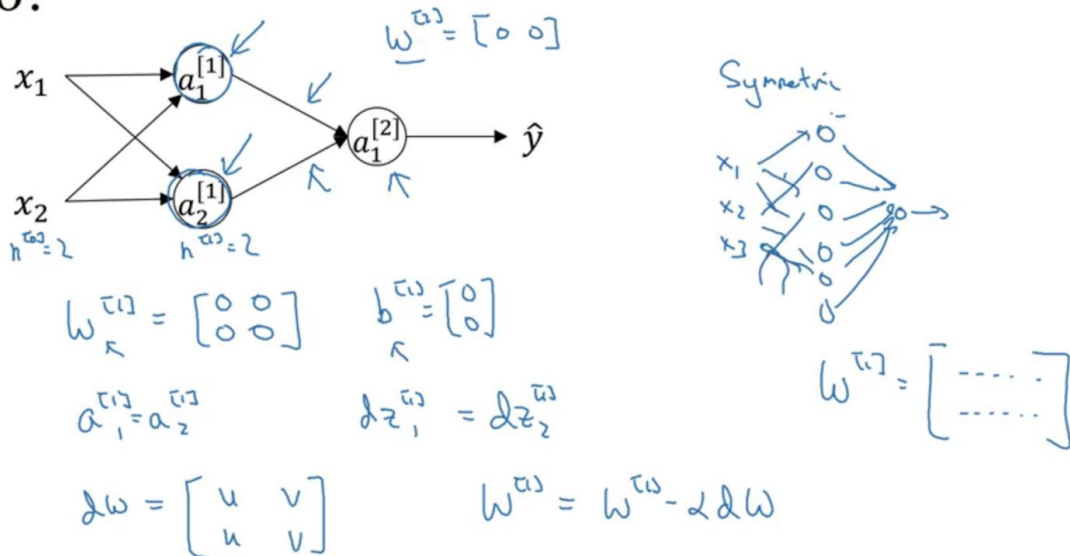
$$\rightarrow a^{[2]} = g^{[2]}(z^{[2]}) = z^{[2]}$$

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

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

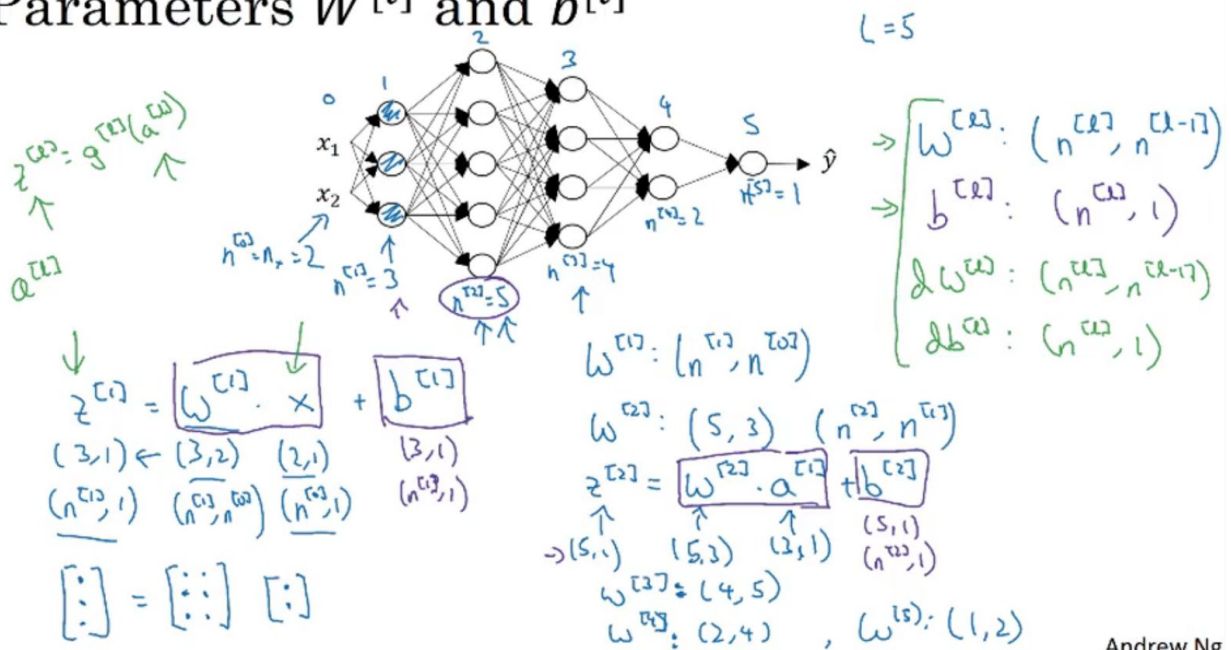
Andrew Ng

What happens if you initialize weights to zero?



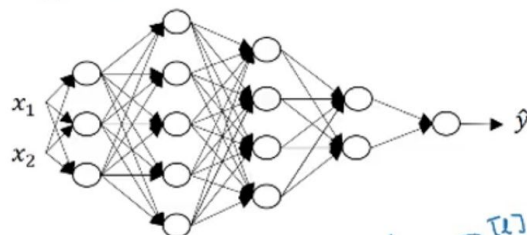
Very good overview of dimensions of W and b .

Parameters $W^{[l]}$ and $b^{[l]}$



Andrew Ng

Vectorized implementation



$$z^{[l]} = W^{[l]} \cdot X + b^{[l]}$$

$(n^{[l]}, 1)$ $(n^{[l]}, n)$ $(n^{[l]}, 1)$ $(n^{[l]}, 1)$
 $[z^{[0]}, z^{[1]}, \dots, z^{[L-1]}]$
 $(n^{[l]}, m)$ $(n^{[l]}, n)$ $(n^{[l]}, m)$ $(n^{[l]}, 1)$
 $(n^{[l]}, m)$ $(n^{[l]}, m)$

$$z^{[L]}, a^{[L]} : (n^{[L]}, 1)$$

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

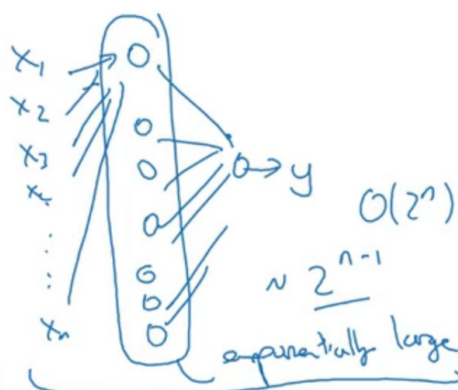
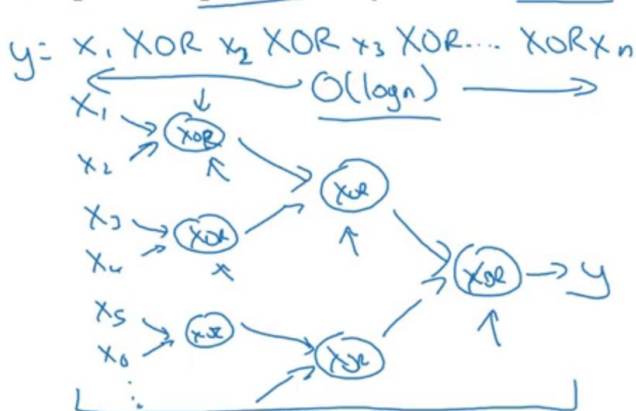
$l=0 \quad A^{[0]} = X = (n^{[0]}, m)$

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

Andrew Ng

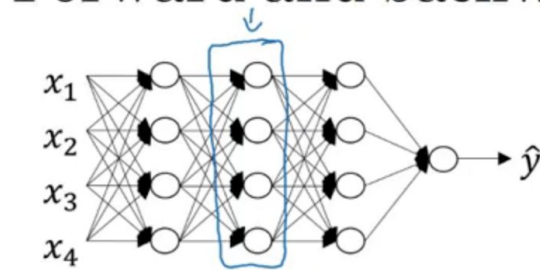
Circuit theory and deep learning

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.

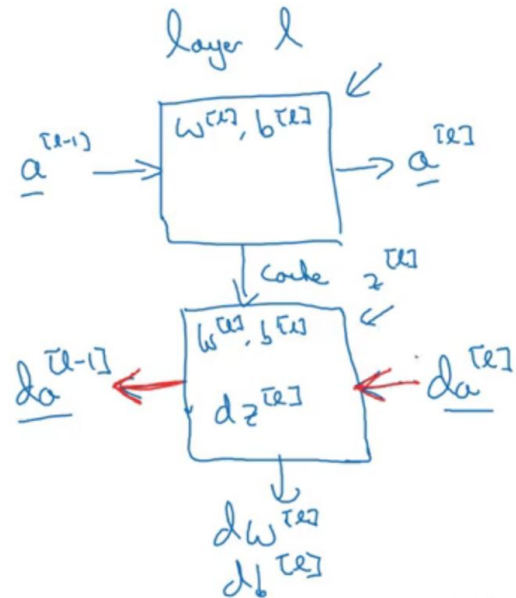


Andrew Ng

Forward and backward functions

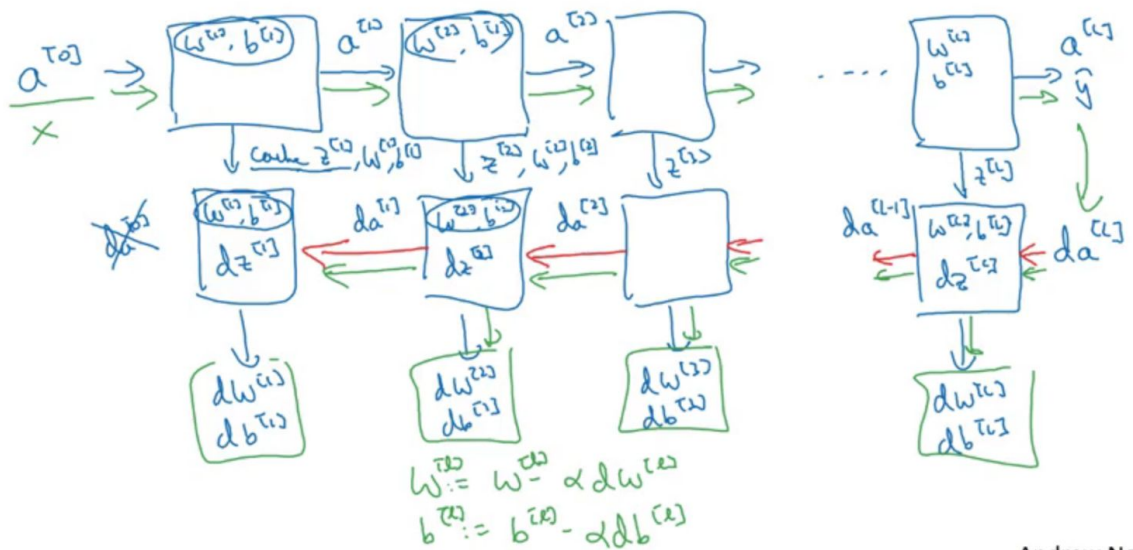


layer l : $W^{[l]}, b^{[l]}$
 \rightarrow Forward: Input $a^{[l-1]}$, output $a^{[l]}$
 $z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$ cache $z^{[l]}$
 $a^{[l]} = g(z^{[l]})$
 \rightarrow Backward: Input $da^{[l]}$ output $da^{[l-1]}$
 cache $(z^{[l]})$
 $\frac{dw^{[l]}}{da^{[l-1]}}$
 $\frac{db^{[l]}}{da^{[l-1]}}$



Andrew I

Forward and backward functions



Andrew Ng

What are hyperparameters?

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]} \dots$

Hyperparameters: α
#iterations
#hidden layers L
#hidden units $n^{[1]}, n^{[2]}, \dots$
choice of activation function

Loss: Momentum, mini-batch size, regularizations, ...

Forward and backward propagation

$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ Z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(Z^{[2]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) = \hat{Y} \end{aligned}$$

$$\begin{aligned} dZ^{[L]} &= A^{[L]} - Y \\ dW^{[L]} &= \frac{1}{m} dZ^{[L]} A^{[L]T} \\ db^{[L]} &= \frac{1}{m} \text{np.sum}(dZ^{[L]}, \text{axis} = 1, \text{keepdims} = \text{True}) \\ dZ^{[L-1]} &= dW^{[L]T} dZ^{[L]} g'^{[L]}(Z^{[L-1]}) \\ &\vdots \\ dZ^{[1]} &= dW^{[L]T} dZ^{[2]} g'^{[1]}(Z^{[1]}) \\ dW^{[1]} &= \frac{1}{m} dZ^{[1]} A^{[1]T} \\ db^{[1]} &= \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis} = 1, \text{keepdims} = \text{True}) \end{aligned}$$

