

What you'll learn

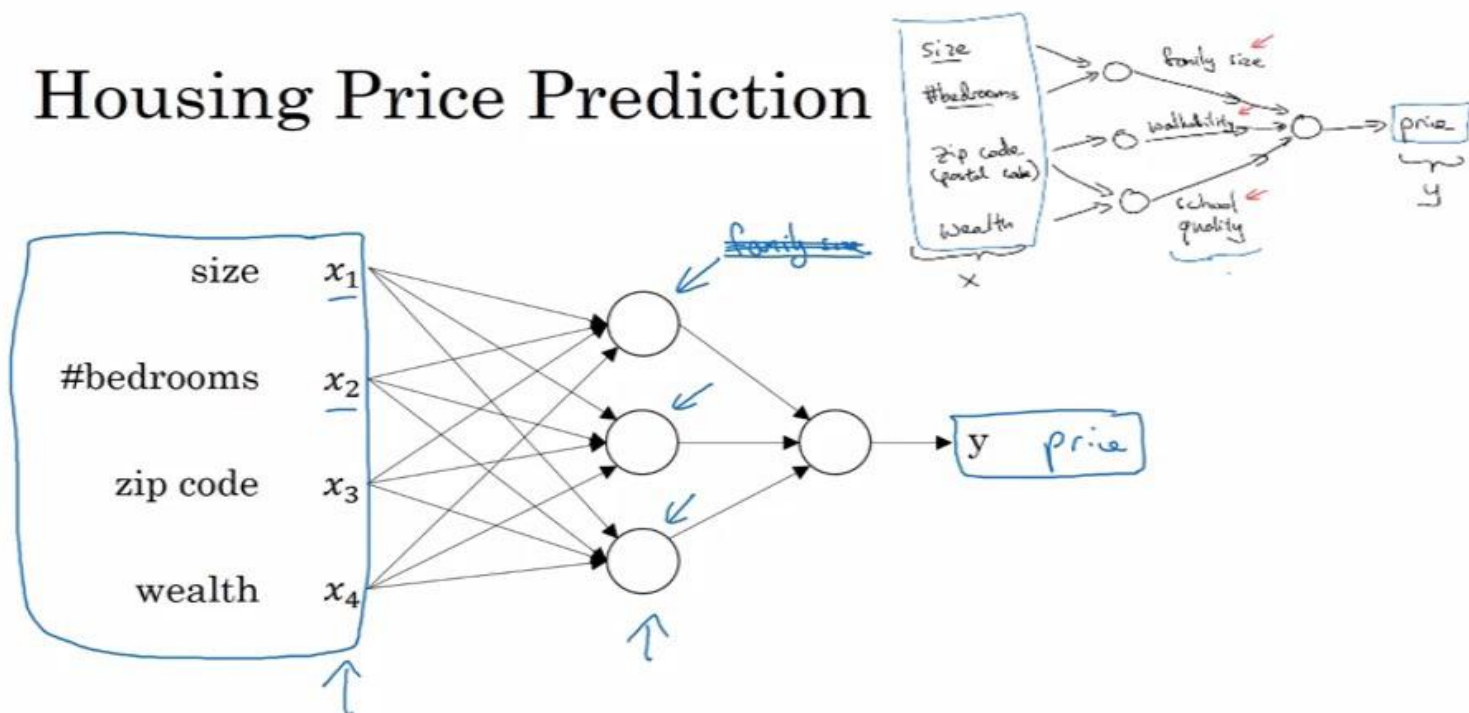


Courses in this sequence (Specialization):

1. Neural Networks and Deep Learning
2. Improving Deep Neural Networks: Hyperparameter tuning, Regularization and Optimization
3. Structuring your Machine Learning project *train/dev/test*
4. Convolutional Neural Networks *CNN* *end-to-end*
5. Natural Language Processing: Building sequence models *RNN, LSTM*

Andrew Ng

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

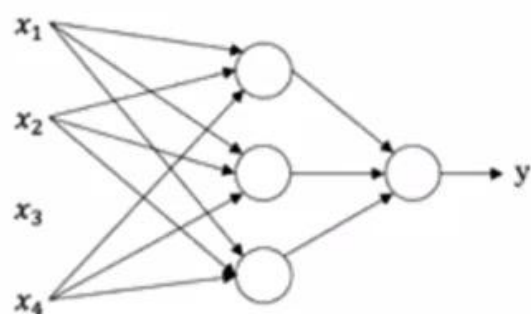
Standard
NN

CNN

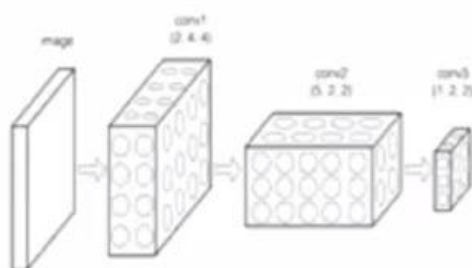
RNN

Custom/
Hybrid

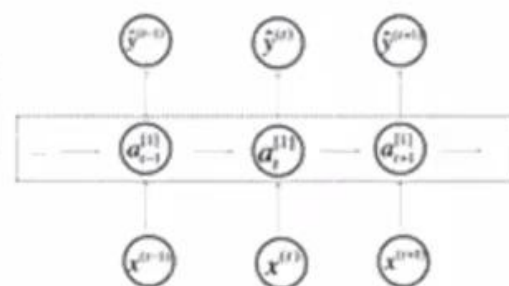
Neural Network examples



Standard NN



Convolutional NN



Recurrent NN

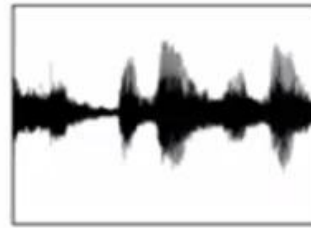
Supervised Learning

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

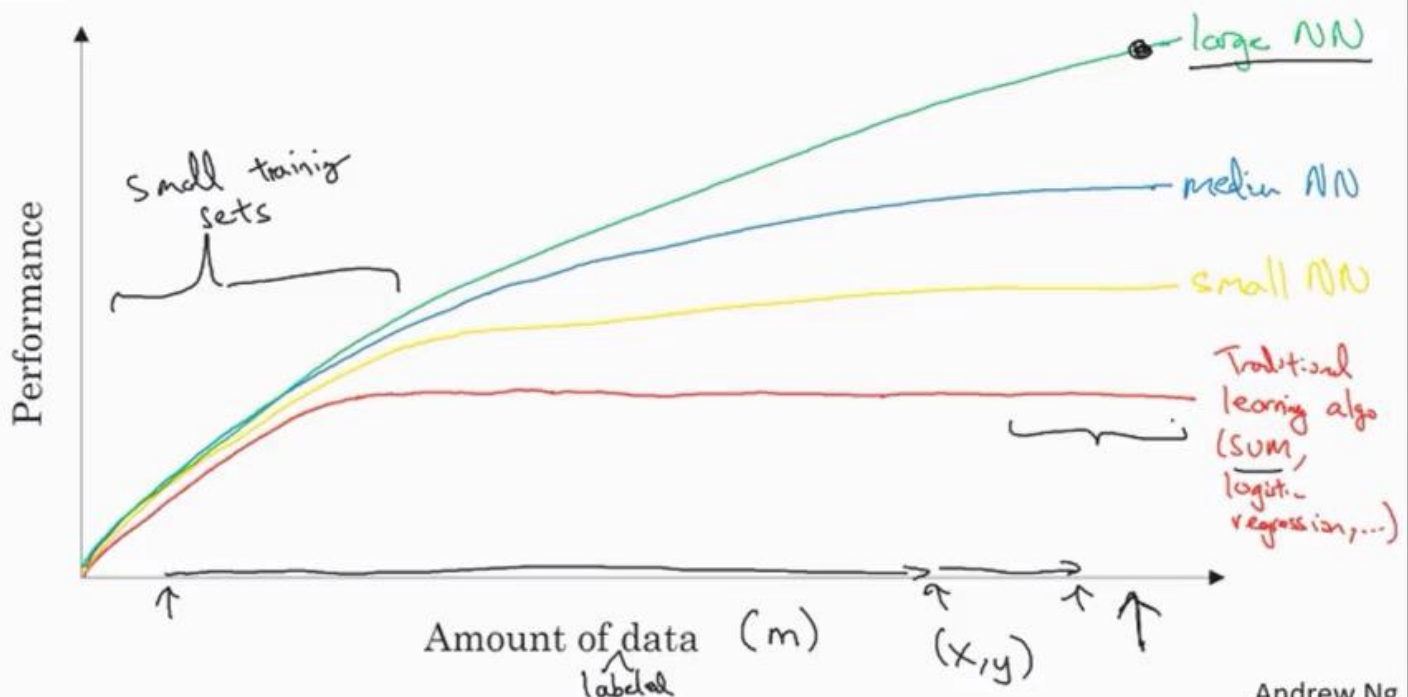


Image

Four scores and seven years ago...

Text

Scale drives deep learning progress



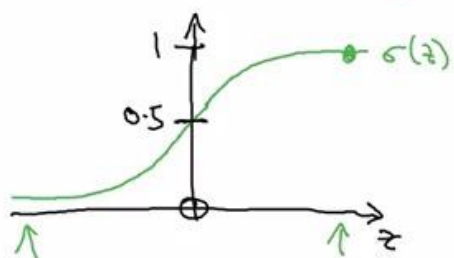
Andrew Ng

Logistic Regression

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

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

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



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

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

If z large negative number

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

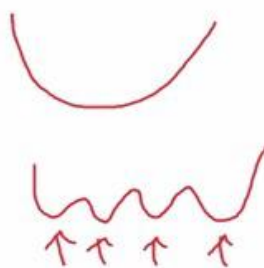
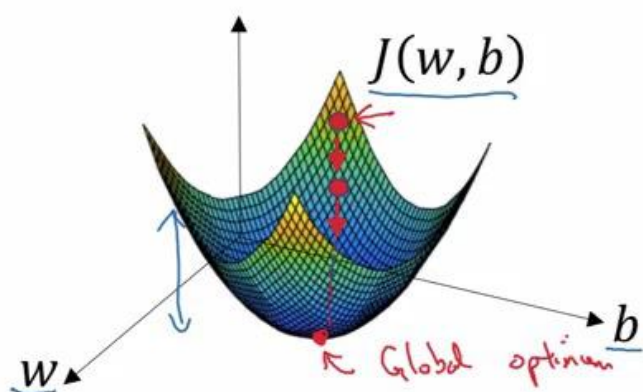
Andrew Ng

Gradient Descent

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

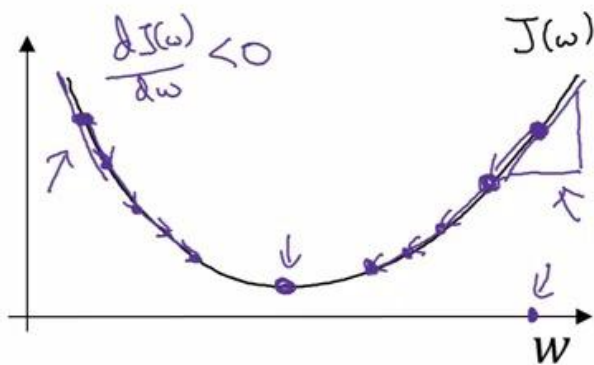
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

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



Andrew Ng

Gradient Descent



Repeat {

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

} $w := w - \alpha dw$

$\frac{dJ(w)}{dw} = ?$

learning rate

"dw"

$J(w, b)$

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

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

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

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

"partial derivative"

$\frac{\partial}{\partial w}$

$\frac{\partial}{\partial b}$

Andrew Ng

More derivative examples

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

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

$$a = 2$$

$$f(a) = 4$$

$$a = 2.001$$

$$f(a) \approx 4.004$$

$$f(a) = a^3$$

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

$$3 \times 2^2 = 12$$

$$a = 2$$

$$f(a) = 8$$

$$a = 2.001$$

$$f(a) \approx 8.012$$

$$f(a) = \log_e(a)$$

$$\ln(a)$$

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

$$\frac{d}{da} f(a) = \frac{1}{2}$$

$$a = 2$$

$$f(a) \approx 0.69315$$

$$a = 2.001$$

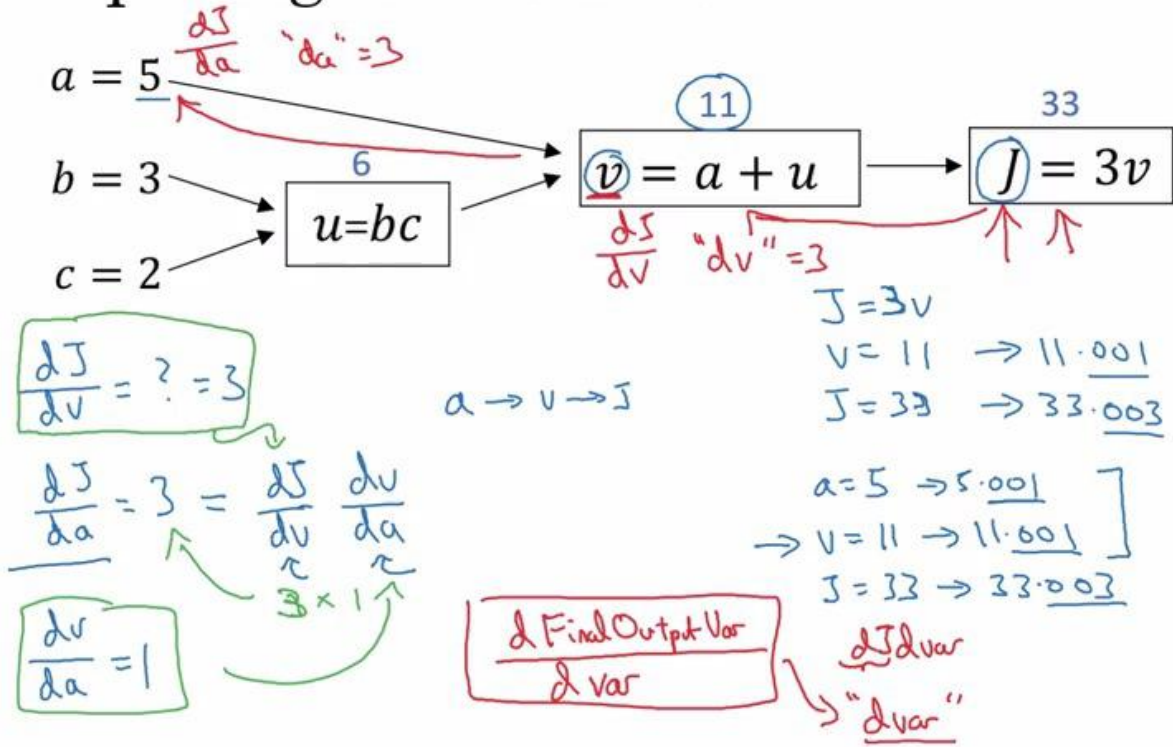
$$f(a) \approx 0.69365$$

$$0.0005$$

$$0.0005$$

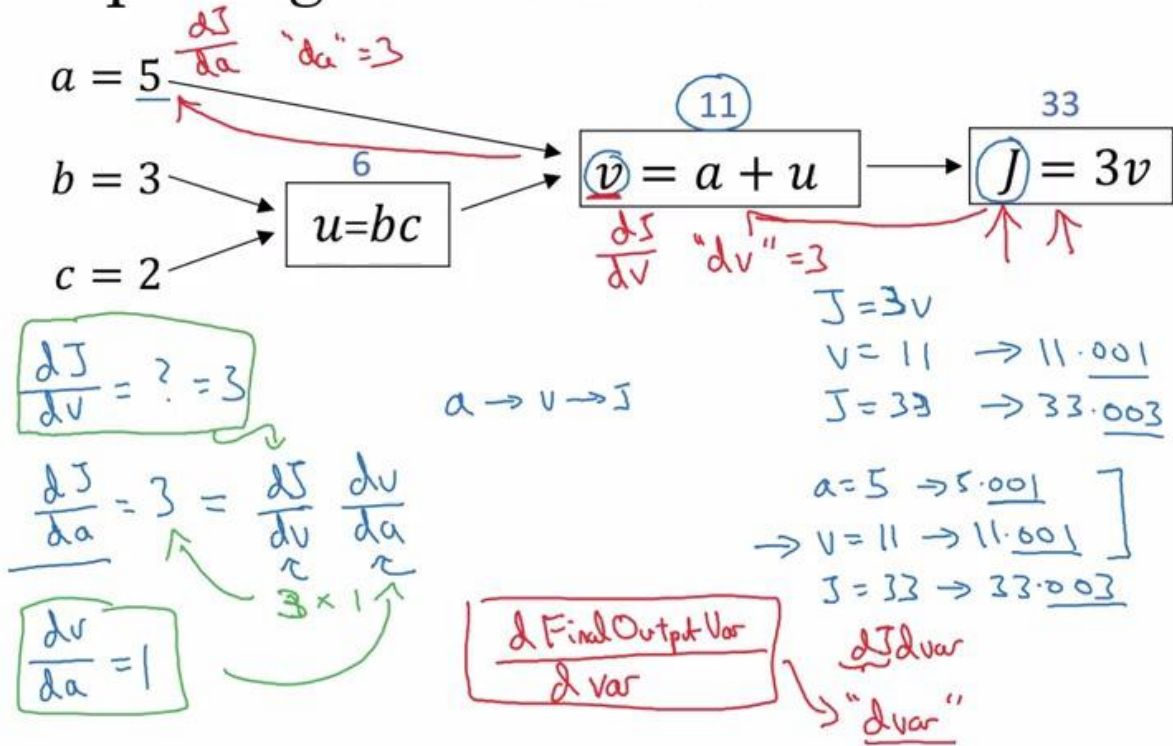
Andrew Ng

Computing derivatives



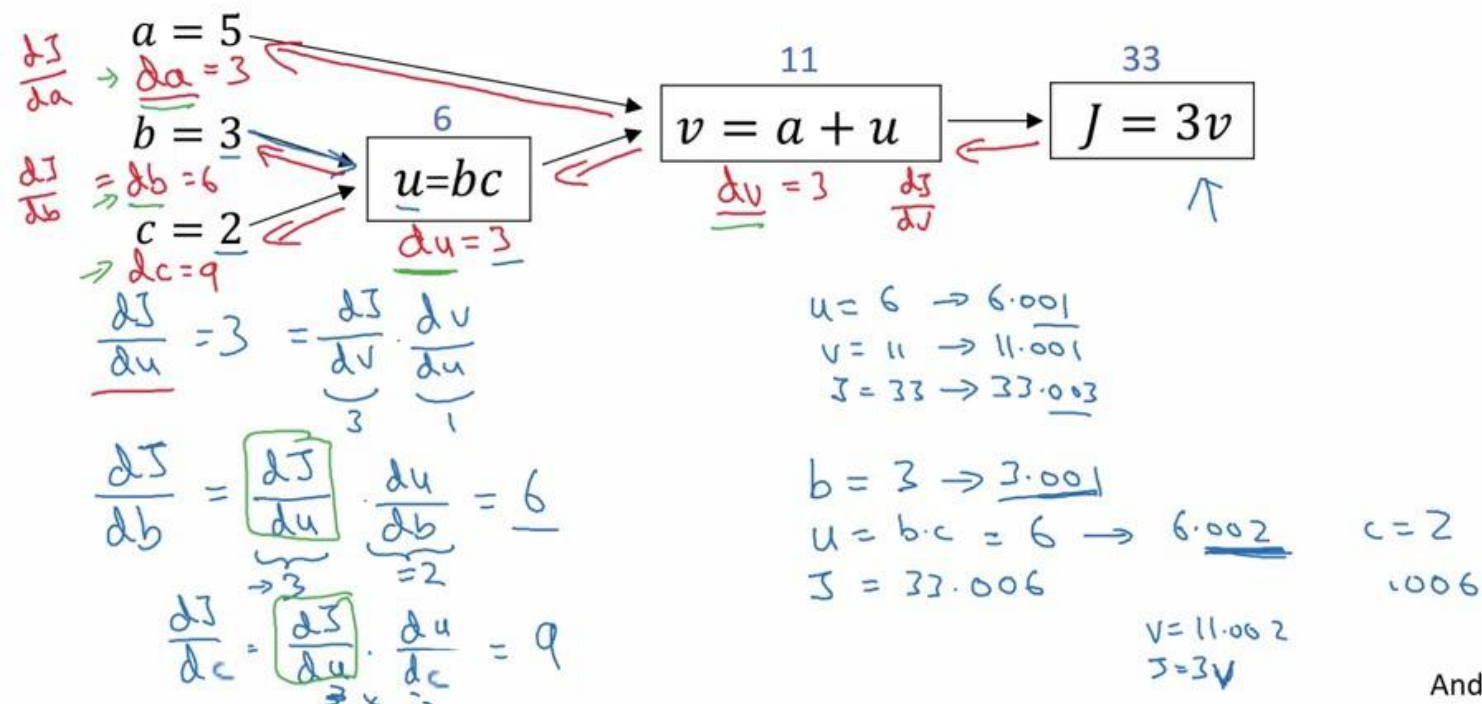
Andrew Ng

Computing derivatives



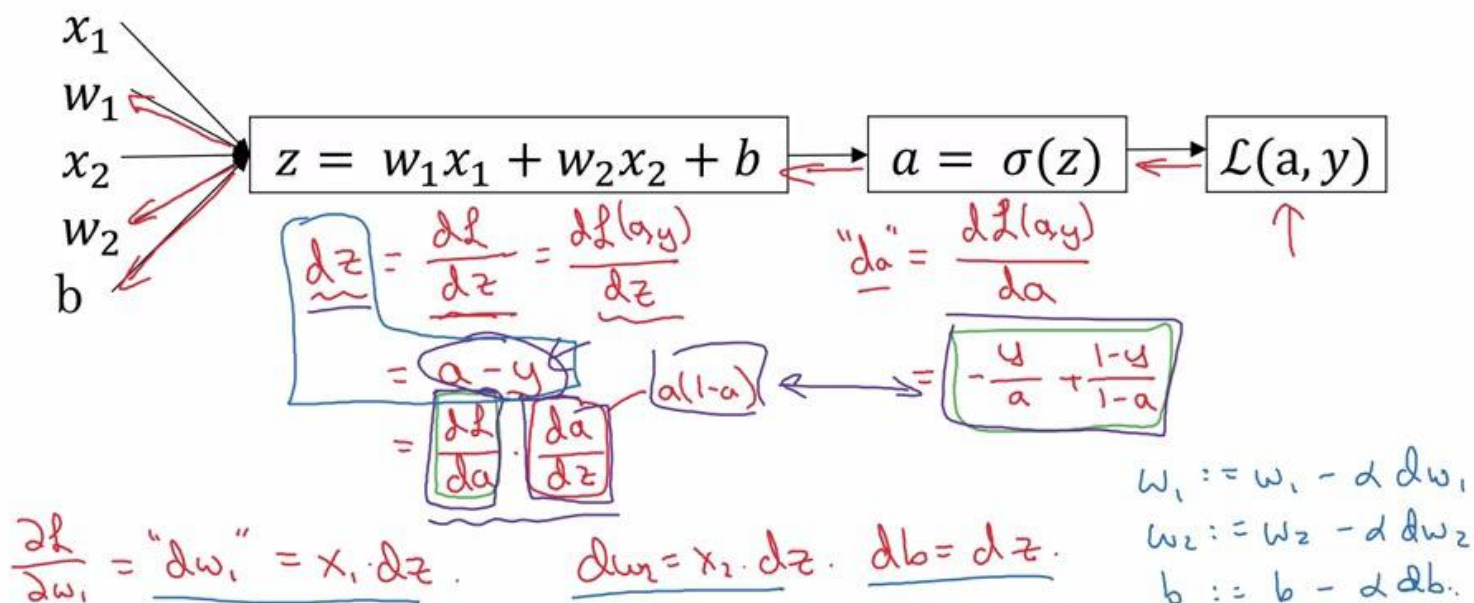
Andrew Ng

Computing derivatives



Andrew Ng

Logistic regression derivatives



Andrew Ng

Logistic regression on m examples

$$J=0; \underline{dw_1}=0; \underline{dw_2}=0; \underline{db}=0$$

→ For $i=1$ to m

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

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

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

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

$$\begin{cases} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{cases} \quad \begin{matrix} \uparrow \\ \downarrow \end{matrix} \quad \begin{matrix} n=2 \\ \end{matrix}$$

$$J /= m \leftarrow$$

$$\begin{matrix} \uparrow & \uparrow & \uparrow \\ dw_1 /= m; & dw_2 /= m; & db /= m. \end{matrix} \leftarrow$$

$$dw_1 = \frac{\partial J}{\partial w_1}$$

$$w_1 := w_1 - \alpha \underline{dw_1}$$

$$w_2 := w_2 - \alpha \underline{dw_2}$$

$$b := b - \alpha \underline{db}$$

Vectorization

Andrew Ng

Implementing Logistic Regression

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

for $i = 1$ to m :

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

$$a^{(i)} = \sigma(z^{(i)}) \leftarrow$$

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

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

$$\begin{cases} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{cases} \quad dw += x^{(i)} * dz^{(i)}$$

$$J = J/m, dw_1 = dw_1/m, dw_2 = dw_2/m$$

$$db = db/m$$

for iter in range(1000):

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

$$= np.dot(w.T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$dw = \frac{1}{m} X dZ^T$$

$$db = \frac{1}{m} np.sum(dZ)$$

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

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

Andrew Ng

Python/numpy vectors

`a = np.random.randn(5)`
`a.shape = (5,)`
"rank 1 array" } Don't use

`a = np.random.randn(5, 1)` → `a.shape = (5, 1)` column vector

`a = np.random.randn(1, 5)` → `a.shape = (1, 5)` row vec.

Andrew Ng

Logistic regression cost function

If $y = \underline{1}$: $p(y|x) = \hat{y}$
If $y = \underline{0}$: $p(y|x) = 1 - \hat{y}$ } $p(y|x)$

[

Andrew Ng

Logistic regression cost function

$$\begin{aligned} \rightarrow & \text{If } y = 1: p(y|x) = \hat{y} \\ \rightarrow & \text{If } y = 0: p(y|x) = 1 - \hat{y} \end{aligned} \quad \left. \vphantom{\begin{aligned} \rightarrow & \text{If } y = 1: p(y|x) = \hat{y} \\ \rightarrow & \text{If } y = 0: p(y|x) = 1 - \hat{y} \end{aligned}} \right\} p(y|x)$$

$$p(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)} \quad \leftarrow$$

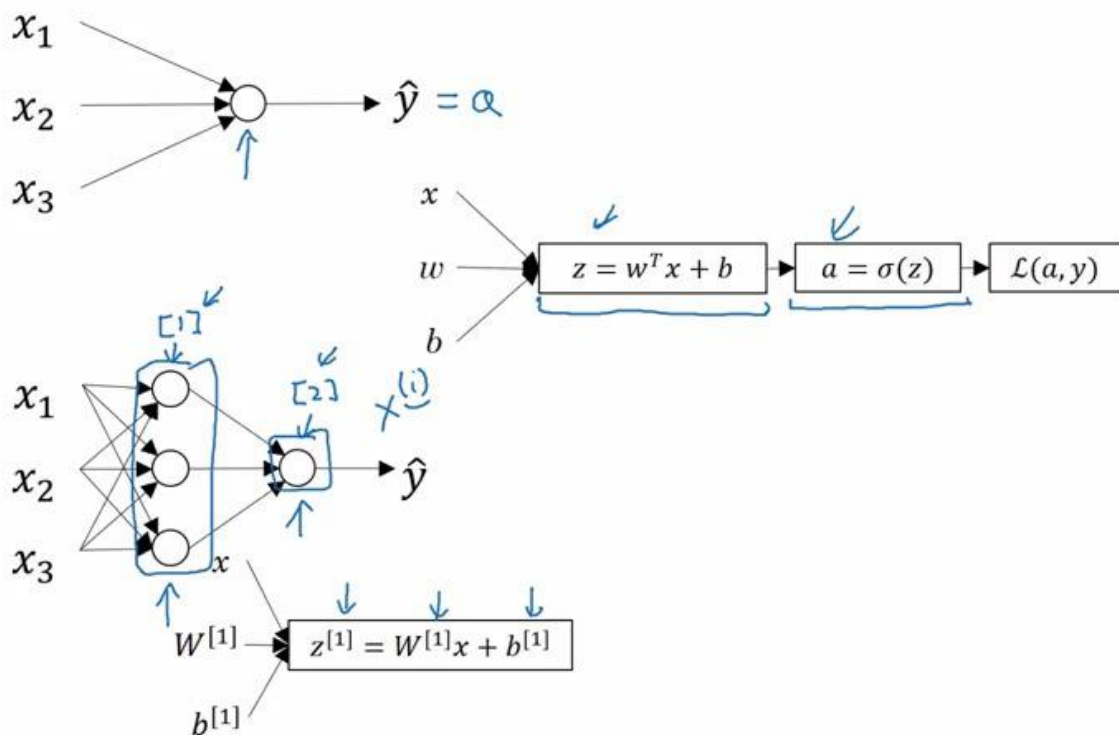
$$\text{If } y=1: p(y|x) = \hat{y} \underbrace{(1-\hat{y})^0}_{=1}$$

$$\text{If } y=0: p(y|x) = \hat{y}^0 \underbrace{(1-\hat{y})^1}_{=1} = 1 \times (1-\hat{y}) = 1-\hat{y}$$

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

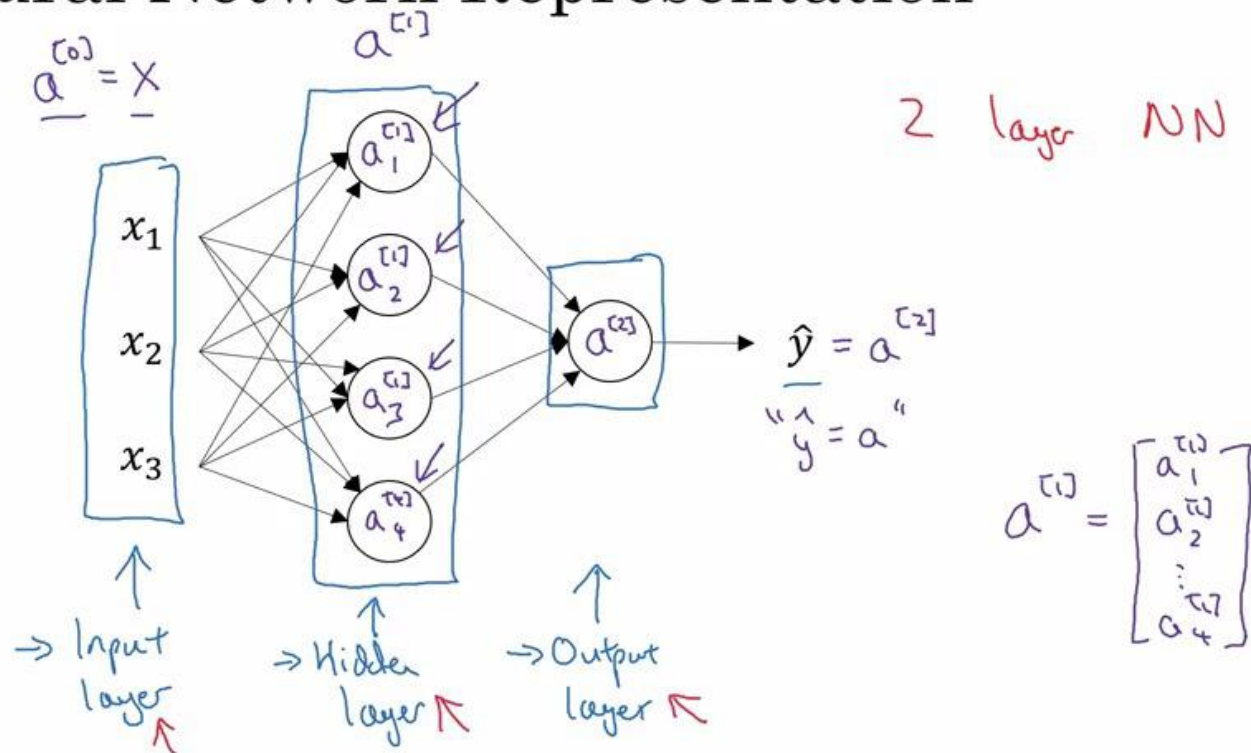
Andrew Ng

What is a Neural Network?



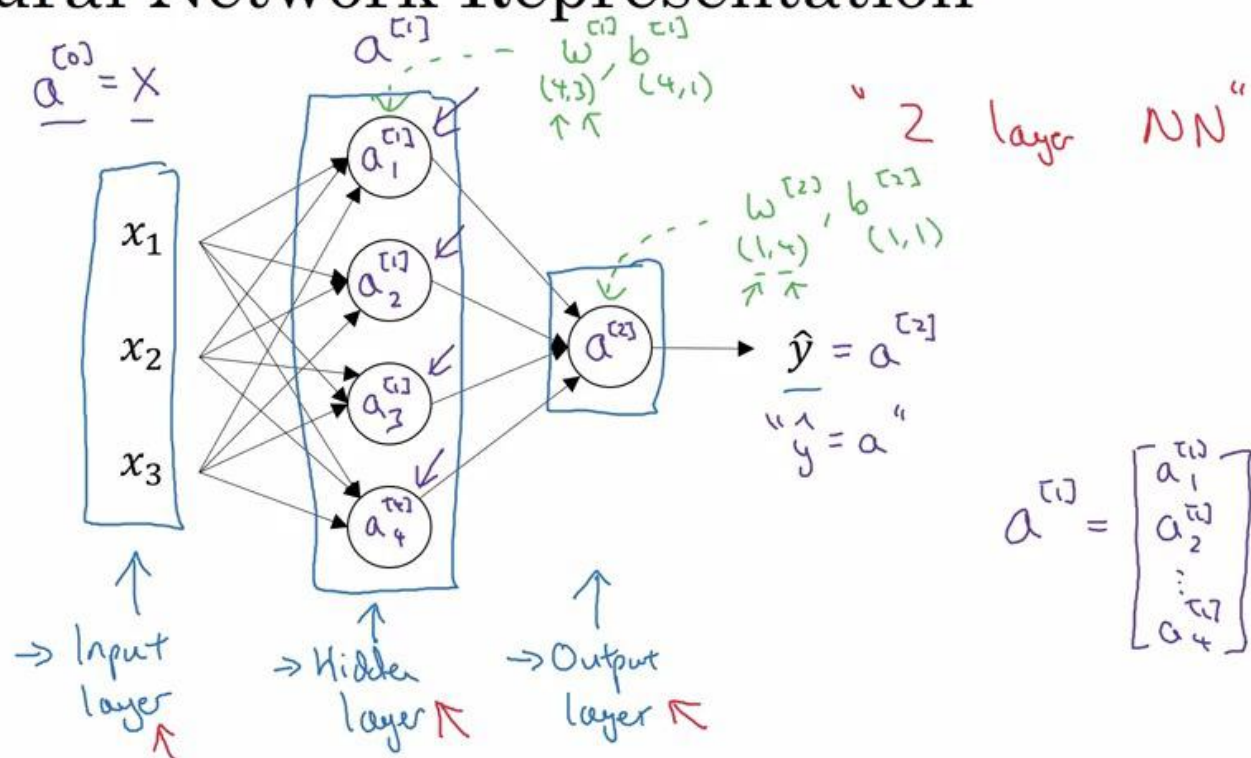
Andrew Ng

Neural Network Representation



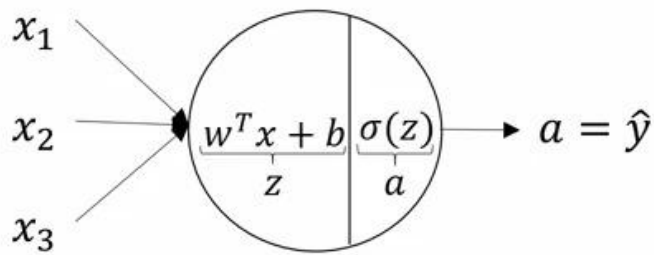
Andrew Ng

Neural Network Representation



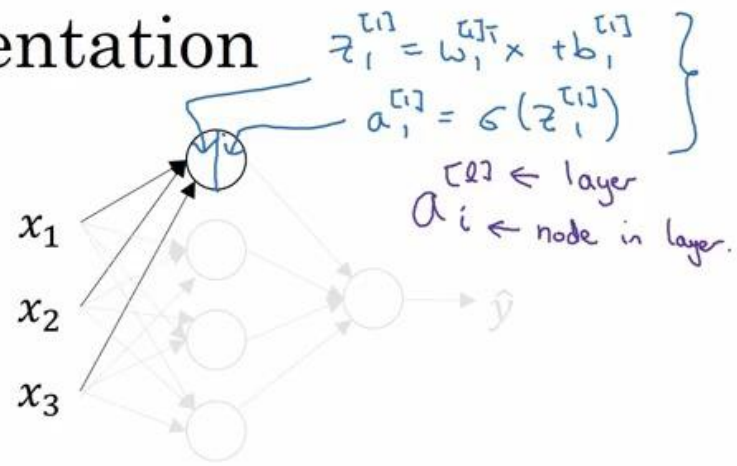
Andrew Ng

Neural Network Representation



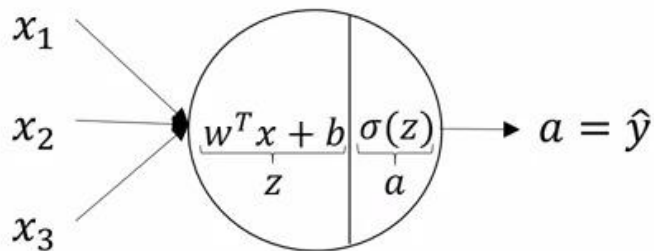
$$z = w^T x + b$$

$$a = \sigma(z)$$



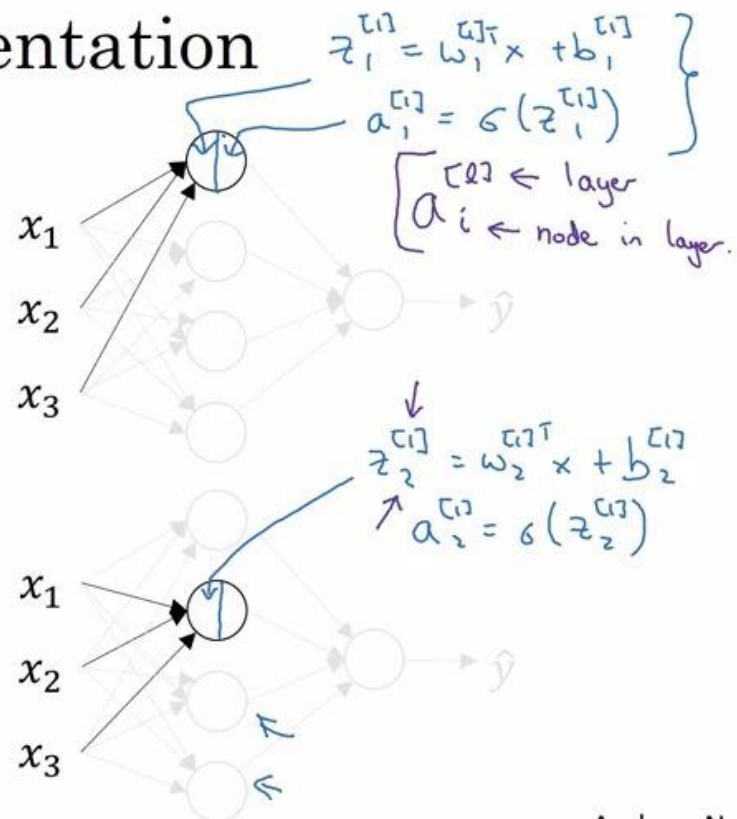
Andrew Ng

Neural Network Representation



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

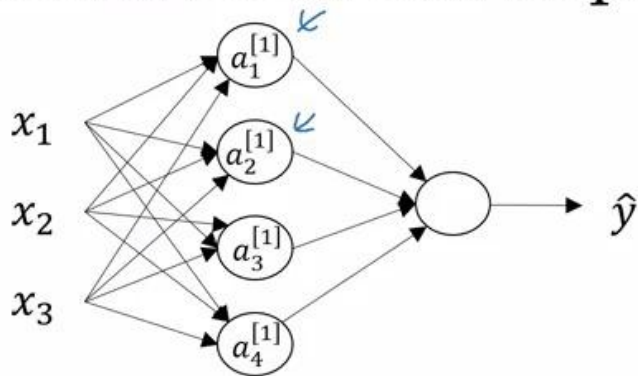
$$a = \sigma(z)$$



Andrew Ng

Neural Network Representation

$$(\omega_i^{[1]})^T x$$

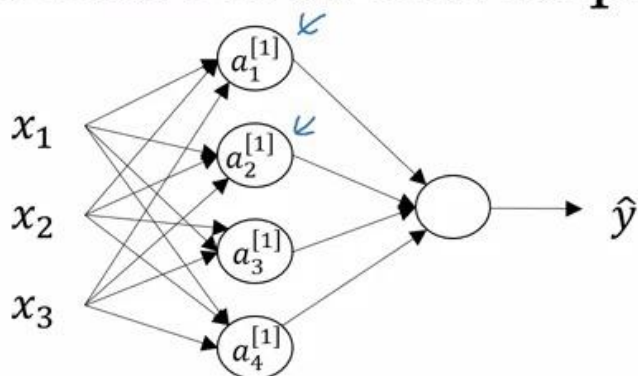


$$\begin{aligned} z_1^{[1]} &= \underline{w_1^{[1]}}^T x + b_1^{[1]}, & a_1^{[1]} &= \sigma(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, & a_2^{[1]} &= \sigma(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, & a_3^{[1]} &= \sigma(z_3^{[1]}) \\ z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, & a_4^{[1]}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

Andrew Ng

Neural Network Representation

$$(\omega_i^{[1]})^T x$$

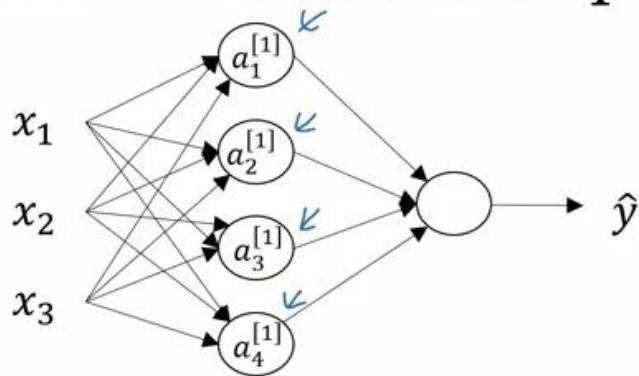


$$\begin{bmatrix} \omega_1^{[1]T} \\ \omega_2^{[1]T} \\ \omega_3^{[1]T} \\ \omega_4^{[1]T} \end{bmatrix}$$

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

Andrew Ng

Neural Network Representation



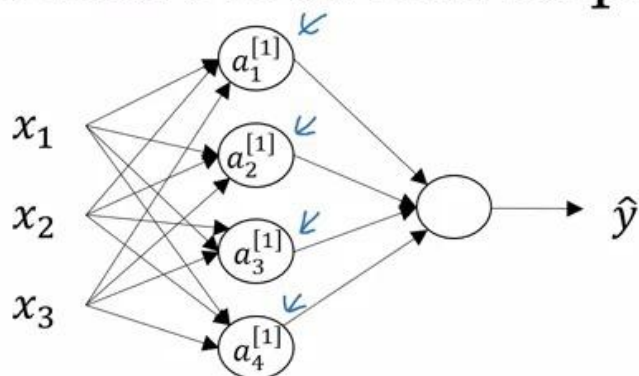
$$\begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \\ -w_3^{[1]T} \\ -w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} \rightarrow w_1^{[1]T} x + b_1^{[1]} \\ \rightarrow w_2^{[1]T} x + b_2^{[1]} \\ \rightarrow w_3^{[1]T} x + b_3^{[1]} \\ \rightarrow w_4^{[1]T} x + b_4^{[1]} \end{bmatrix}$$

(4, 3)

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

Andrew Ng

Neural Network Representation



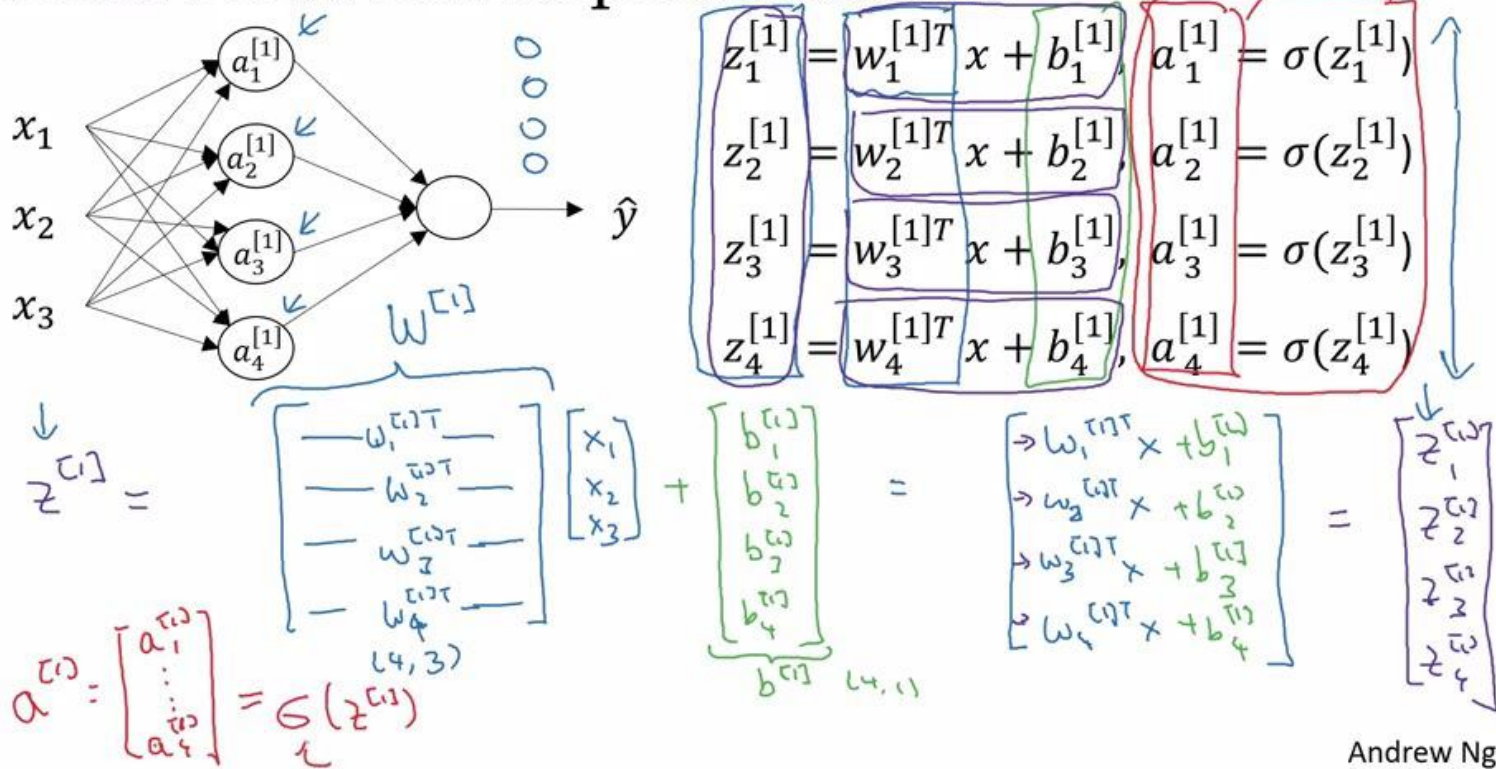
$$\begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \\ -w_3^{[1]T} \\ -w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} \rightarrow w_1^{[1]T} x + b_1^{[1]} \\ \rightarrow w_2^{[1]T} x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

(4, 3)

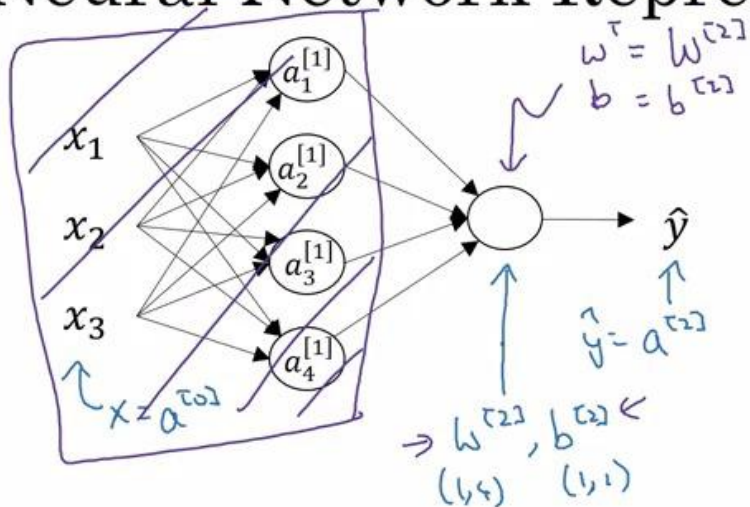
$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, & a_1^{[1]} &= \sigma(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, & a_2^{[1]} &= \sigma(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, & a_3^{[1]} &= \sigma(z_3^{[1]}) \\ z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, & a_4^{[1]} &= \sigma(z_4^{[1]})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

Andrew Ng

Neural Network Representation



Neural Network Representation learning



Given input x :

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

Dimensions: $(4,1)$, $(4,3)$, $(3,1)$, $(4,1)$

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

Dimensions: $(4,1)$, $(4,1)$

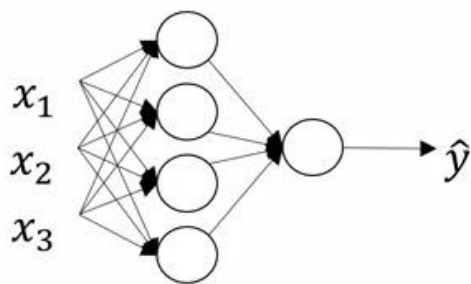
$$\rightarrow z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

Dimensions: $(1,1)$, $(1,4)$, $(4,1)$, $(1,1)$

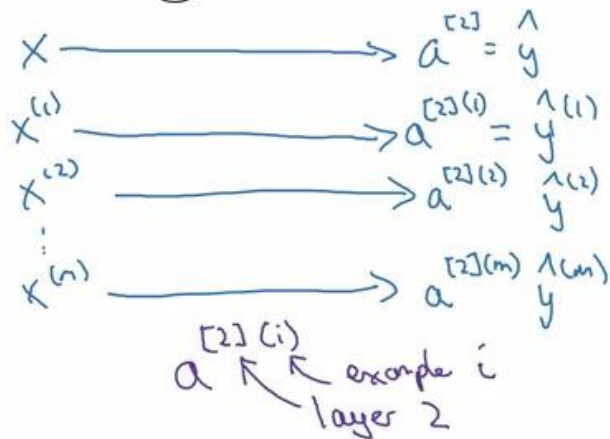
$$\rightarrow a^{[2]} = \sigma(z^{[2]})$$

Dimensions: $(1,1)$, $(1,1)$

Vectorizing across multiple examples

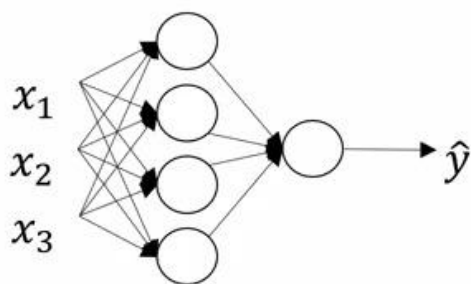


$$\begin{cases} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{cases}$$

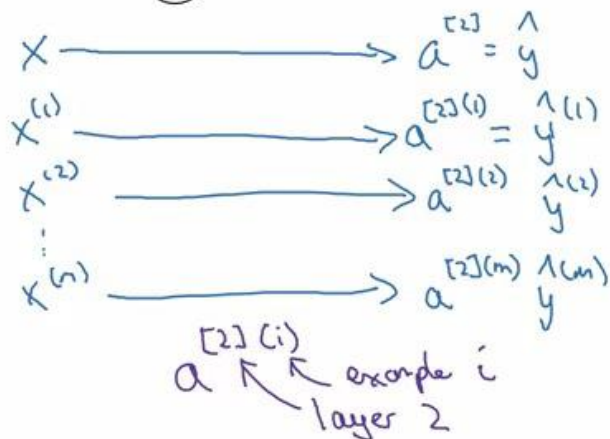


Andrew Ng

Vectorizing across multiple examples



$$\begin{cases} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{cases}$$



for $i = 1$ to m ,

$$\begin{aligned} z^{[1]}(i) &= W^{[1]}x^{(i)} + b^{[1]} \\ a^{[1]}(i) &= \sigma(z^{[1]}(i)) \\ z^{[2]}(i) &= W^{[2]}a^{[1]}(i) + b^{[2]} \\ a^{[2]}(i) &= \sigma(z^{[2]}(i)) \end{aligned}$$

Andrew Ng

Vectorizing across multiple examples

for $i = 1$ to m :

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

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

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

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

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

(n_x, m)

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

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

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

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

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

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

Andrew Ng

Vectorizing across multiple examples

for $i = 1$ to m :

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

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

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

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

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

(n_x, m)

training examples

hidden units

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

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

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

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

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

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

hidden units

Andrew Ng

Justification for vectorized implementation

$$z^{(1)} = w^{(1)} x^{(1)} + b^{(1)} \quad , \quad z^{(2)} = w^{(1)} x^{(2)} + b^{(1)} \quad , \quad z^{(3)} = w^{(1)} x^{(3)} + b^{(1)}$$

$$w^{(1)} = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix} \quad w^{(1)} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad w^{(1)} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad w^{(1)} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$w^{(1)} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

X

Andrew Ng

Justification for vectorized implementation

$$z^{(1)} = w^{(1)} x^{(1)} + b^{(1)} \quad , \quad z^{(2)} = w^{(1)} x^{(2)} + b^{(1)} \quad , \quad z^{(3)} = w^{(1)} x^{(3)} + b^{(1)}$$

$$w^{(1)} = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix} \quad w^{(1)} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad w^{(1)} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad w^{(1)} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$w^{(1)} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} | & | & | \\ z^{(1)} & z^{(2)} & z^{(3)} \\ | & | & | \end{bmatrix}$$

X

Andrew Ng

Justification for vectorized implementation

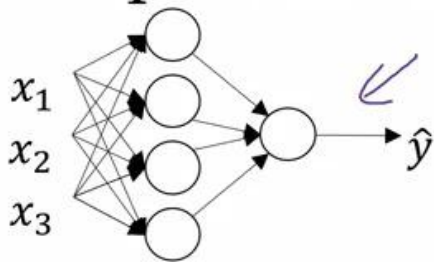
$$z^{1} = W^{[1]} x^{(1)} + b^{[1]}, \quad z^{[1](2)} = W^{[1]} x^{(2)} + b^{[1]}, \quad z^{[1](3)} = W^{[1]} x^{(3)} + b^{[1]}$$

$$W^{[1]} = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}, \quad W^{[1]} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, \quad W^{[1]} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, \quad W^{[1]} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$W^{[1]} \begin{bmatrix} | & | & | & \dots \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots \\ | & | & | & \dots \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot & \dots \end{bmatrix} = \begin{bmatrix} | & | & | & \dots \\ z^{1} & z^{[1](2)} & z^{[1](3)} & \dots \\ | & | & | & \dots \end{bmatrix} = z^{[1]}$$

Andrew Ng

Recap of vectorizing across multiple examples



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

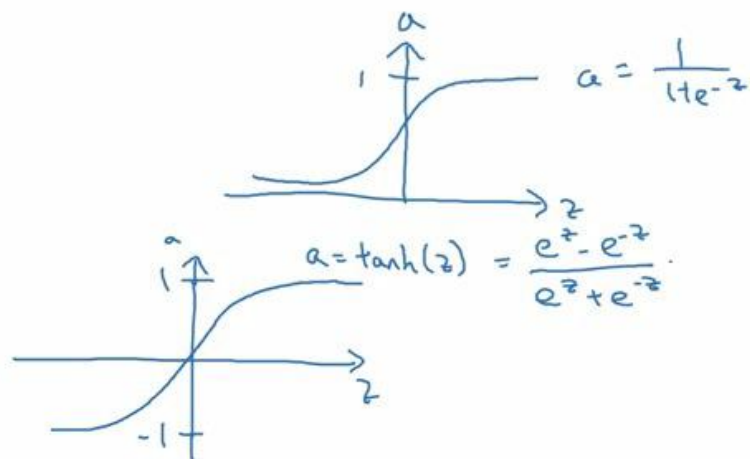
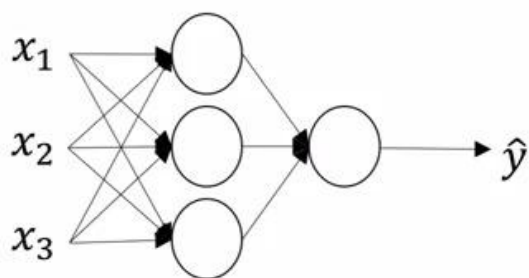
$$A^{[1]} = \begin{bmatrix} | & | & \dots & | \\ a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \\ | & | & \dots & | \end{bmatrix}$$

$$\left[\begin{array}{l} \text{for } i = 1 \text{ to } m \\ \rightarrow z^{[1](i)} = W^{[1]} x^{(i)} + b^{[1]} \\ \rightarrow a^{[1](i)} = \sigma(z^{[1](i)}) \\ \rightarrow z^{[2](i)} = W^{[2]} a^{[1](i)} + b^{[2]} \\ \rightarrow a^{[2](i)} = \sigma(z^{[2](i)}) \end{array} \right]$$

$$\left[\begin{array}{l} Z^{[1]} = W^{[1]} X + b^{[1]} \\ A^{[1]} = \sigma(Z^{[1]}) \\ Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} = \sigma(Z^{[2]}) \end{array} \right]$$

Andrew Ng

Activation functions



Given x :

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

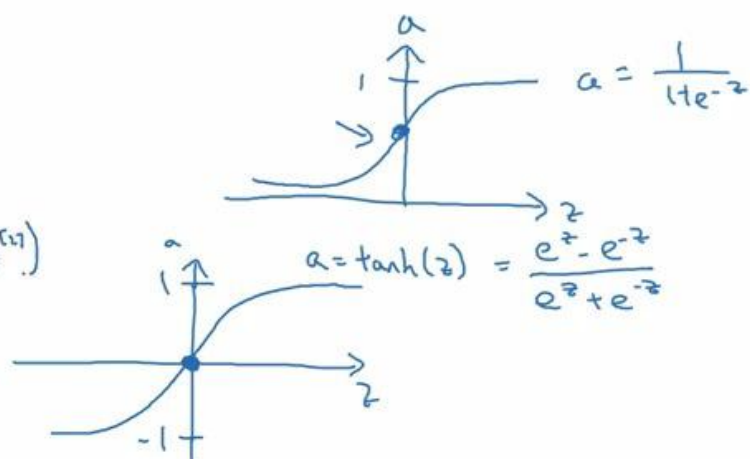
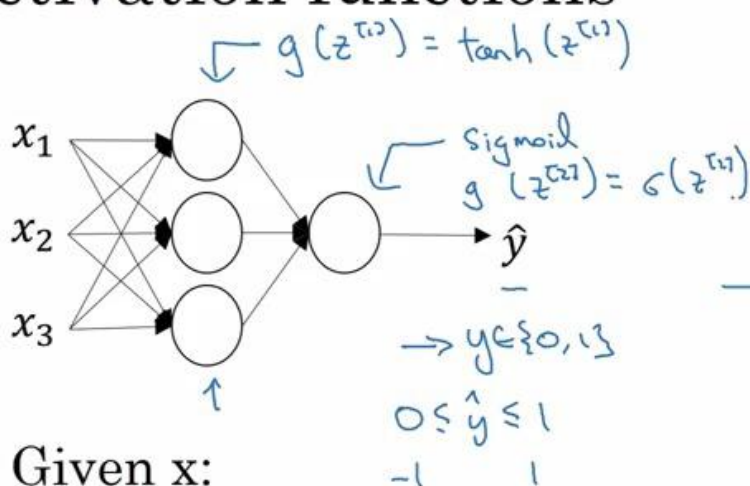
$$\rightarrow a^{[1]} = \cancel{\sigma(z^{[1]})} g(z^{[1]})$$

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

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

Andrew Ng

Activation functions



Given x :

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

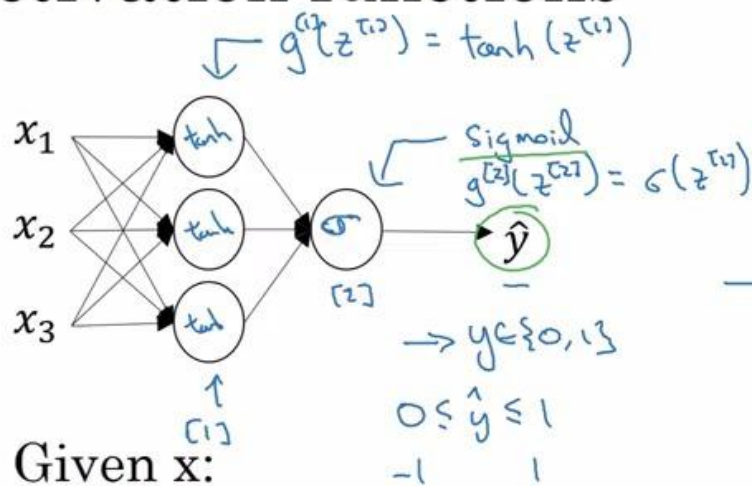
$$\rightarrow a^{[1]} = \cancel{\sigma(z^{[1]})} g(z^{[1]})$$

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

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

Andrew Ng

Activation functions

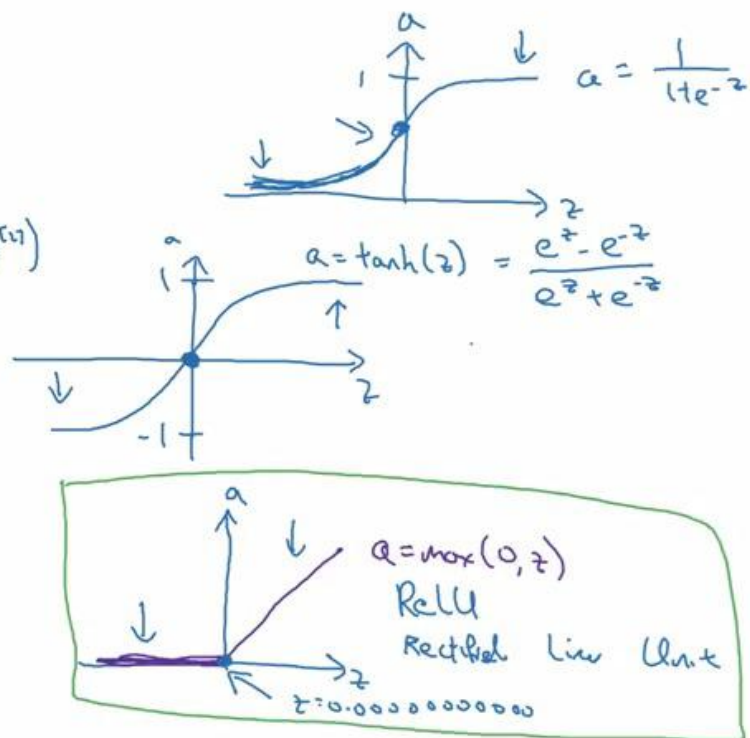


$$z^{[1]} = W^{[1]}x + b^{[1]}$$

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

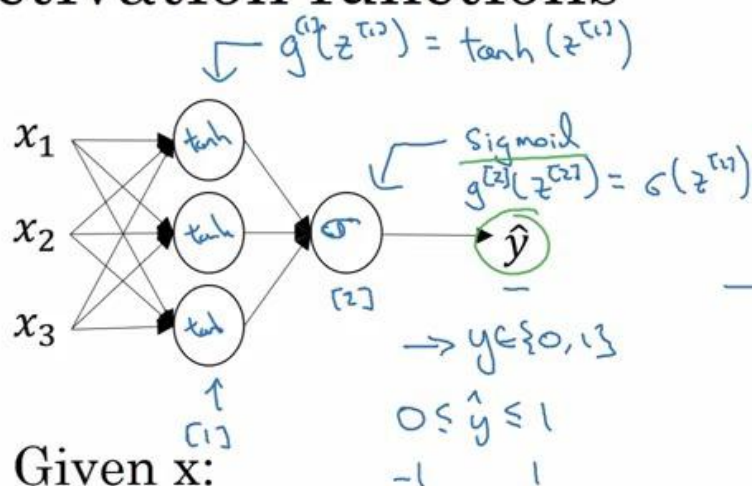
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

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



Andrew Ng

Activation functions

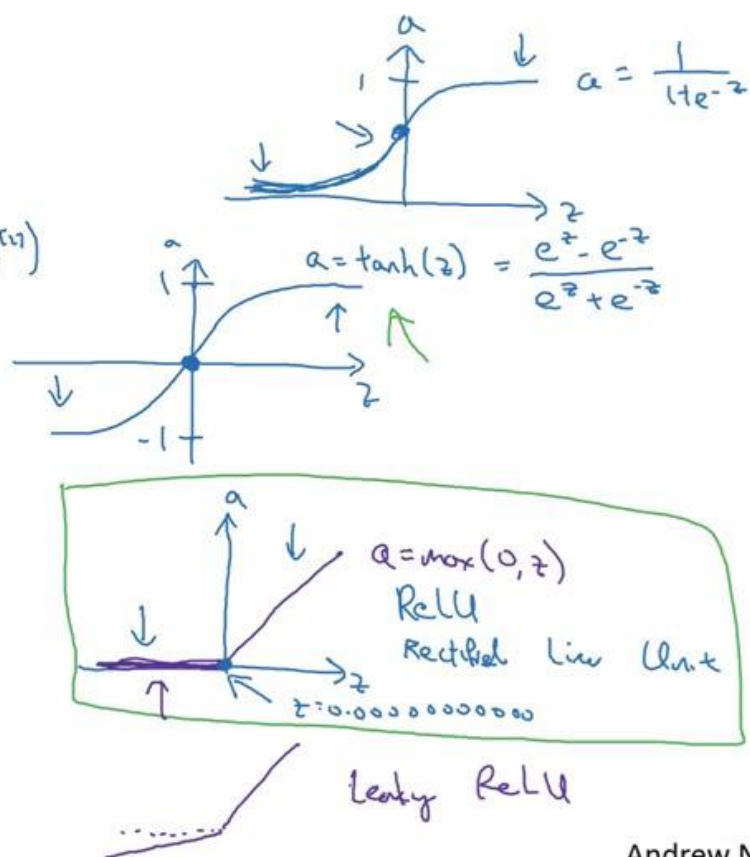


$$z^{[1]} = W^{[1]}x + b^{[1]}$$

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

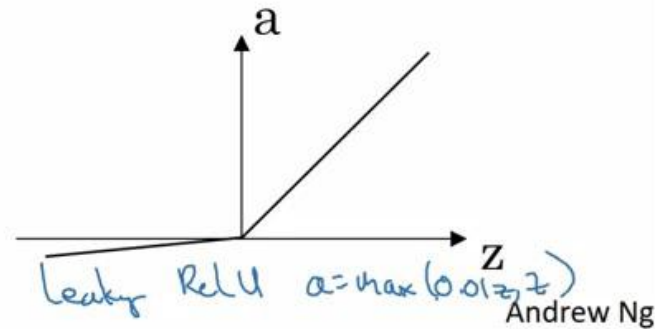
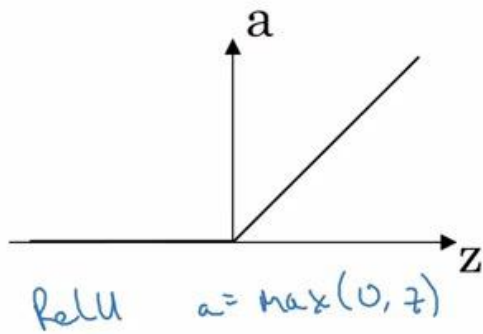
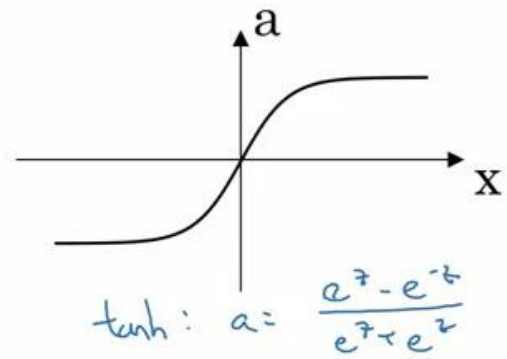
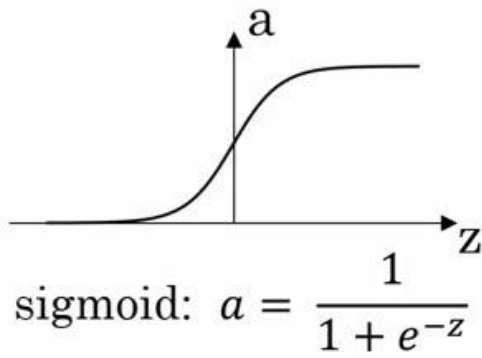
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

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



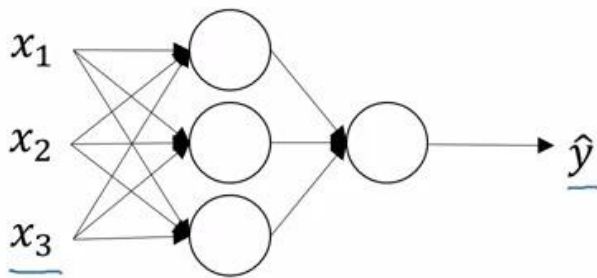
Andrew Ng

Pros and cons of activation functions



Andrew Ng

Activation function



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

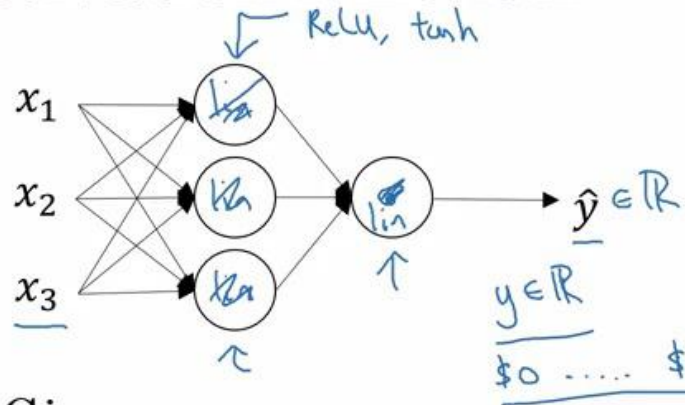
Given x :

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

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

Andrew Ng

Activation function



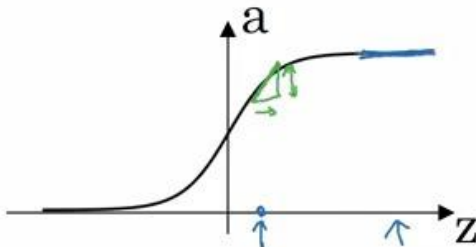
$$a^{[1]} = z^{[1]} = W^{[1]}x + b^{[1]}$$
$$a^{[2]} = z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
$$a^{[2]} = W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]}$$
$$= (W^{[2]}W^{[1]})x + (W^{[2]}b^{[1]} + b^{[2]})$$
$$= W'x + b'$$
$$g(z) = z$$

- Given x:
- $z^{[1]} = W^{[1]}x + b^{[1]}$
 - $a^{[1]} = g^{[1]}(z^{[1]})$
 - $z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$
 - $a^{[2]} = g^{[2]}(z^{[2]})$

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

Andrew Ng

Sigmoid activation function



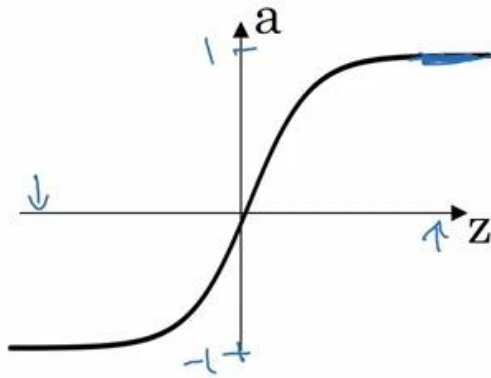
$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{d}{dz}g(z) = \text{slope of } g(z) \text{ at } z$$
$$= \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}}\right)$$
$$= g(z) (1 - g(z)) \leftarrow$$

$$z = 10, \quad g(z) \approx 1$$
$$\frac{d}{dz}g(z) \approx 1(1-1) \approx 0$$

Andrew Ng

Tanh activation function



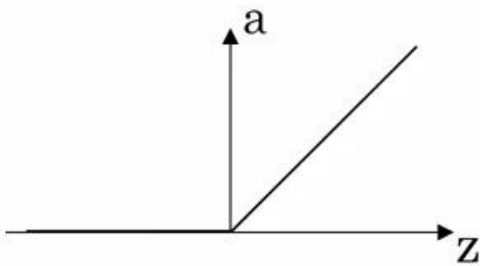
$$g(z) = \tanh(z) \\ = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z \\ = 1 - (\tanh(z))^2 \leftarrow$$

$$\left| \begin{array}{l} z=10 \quad \tanh(z) \approx 1 \\ \quad \quad g'(z) \approx 0 \end{array} \right|$$

Andrew Ng

ReLU and Leaky ReLU



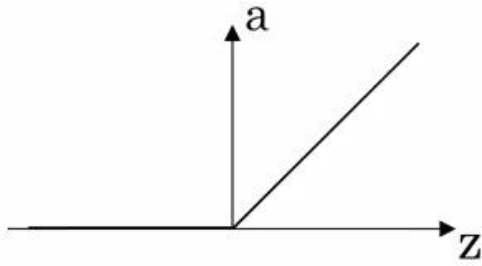
ReLU

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

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

Andrew Ng

ReLU and Leaky ReLU

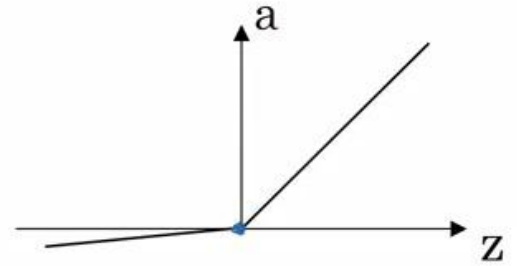


ReLU

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

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

$z = 0.0000000000$



Leaky ReLU

$$g(z) = \max(0.01z, z)$$

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

Andrew Ng

Gradient descent for neural networks

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
 $(n^{[1]}, n^{[2]})$ $(n^{[1]}, 1)$ $(n^{[2]}, n^{[3]})$ $(n^{[2]}, 1)$

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

$$\text{Cost function: } J(W^{[1]}, b^{[1]}, \underline{W^{[2]}}, \underline{b^{[2]}}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i)$$

$\uparrow a^{[2]}$

Gradient descent:

Repeat {

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

$$dW^{[1]} = \frac{\partial J}{\partial W^{[1]}}, db^{[1]} = \frac{\partial J}{\partial b^{[1]}}, \dots$$

$$W^{[1]} := W^{[1]} - \alpha dW^{[1]}$$

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

$W^{[2]}, b^{[2]}$

Formulas for computing derivatives

Forward propagation:

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

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

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

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

Formulas for computing derivatives

Forward propagation:

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

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

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

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

Back propagation:

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

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

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

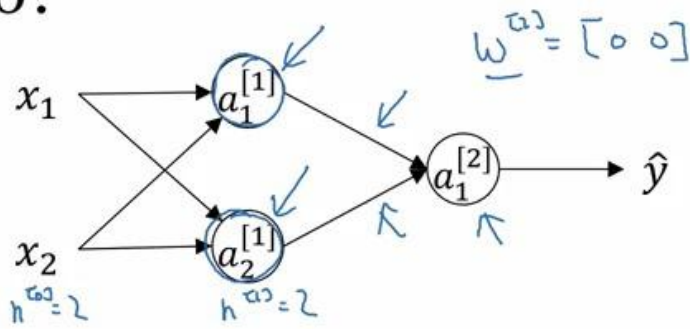
$$db^{[2]} = \frac{1}{n} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims}=\text{True})$$

$$dz^{[1]} = \underbrace{w^{[2]T} dz^{[2]}}_{(n^{[1]}, m)} \times \underbrace{g^{[1]'}(z^{[1]})}_{\text{element-wise product}} \quad (n^{[1]}, m)$$

$$dw^{[1]} = \frac{1}{n} dz^{[1]} x^T$$

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

What happens if you initialize weights to zero?



$$w_{\kappa}^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

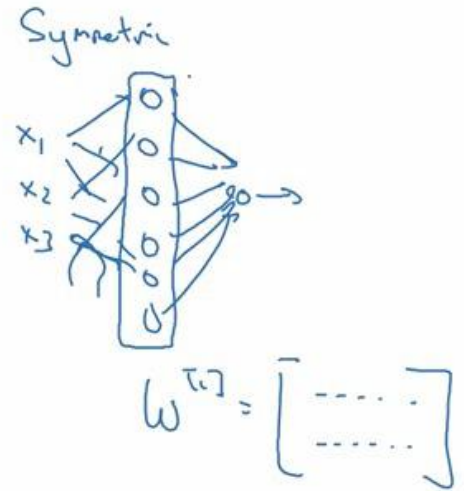
$$b_{\kappa}^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]}$$

$$\Delta z_1 = \Delta z_2$$

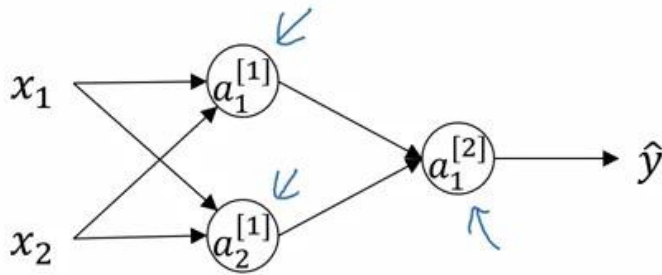
$$\Delta w = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

$$w^{[1]} = w^{[1]} - \Delta \Delta w$$



Andrew Ng

Random initialization



$$\rightarrow w^{[1]} = \text{np.random.randn}(2,2) * \frac{0.01}{100?}$$

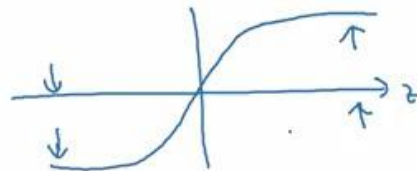
$$b^{[1]} = \text{np.zeros}(2,1)$$

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

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

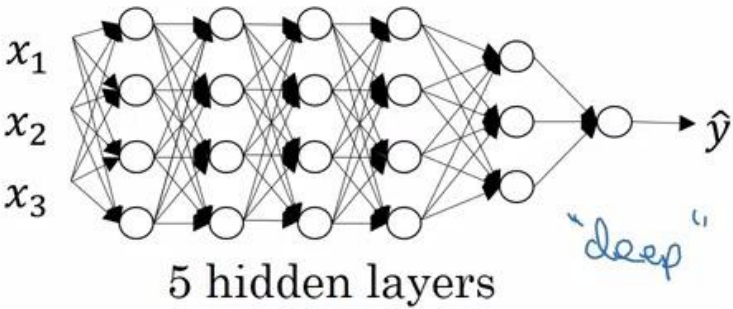
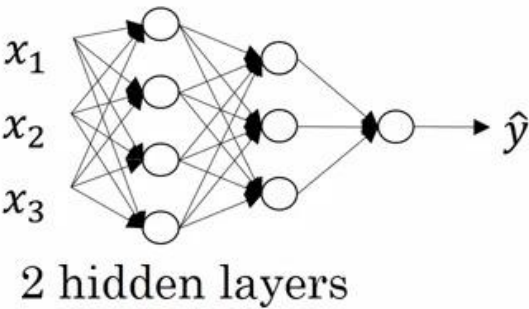
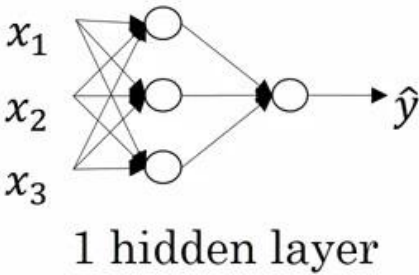
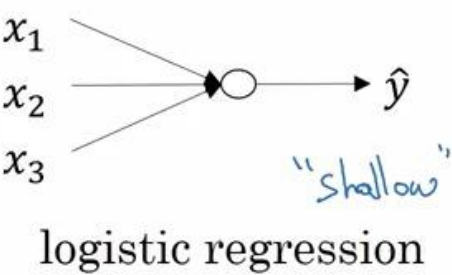
$$\rightarrow z^{[1]} = w^{[1]}x + b^{[1]}$$

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



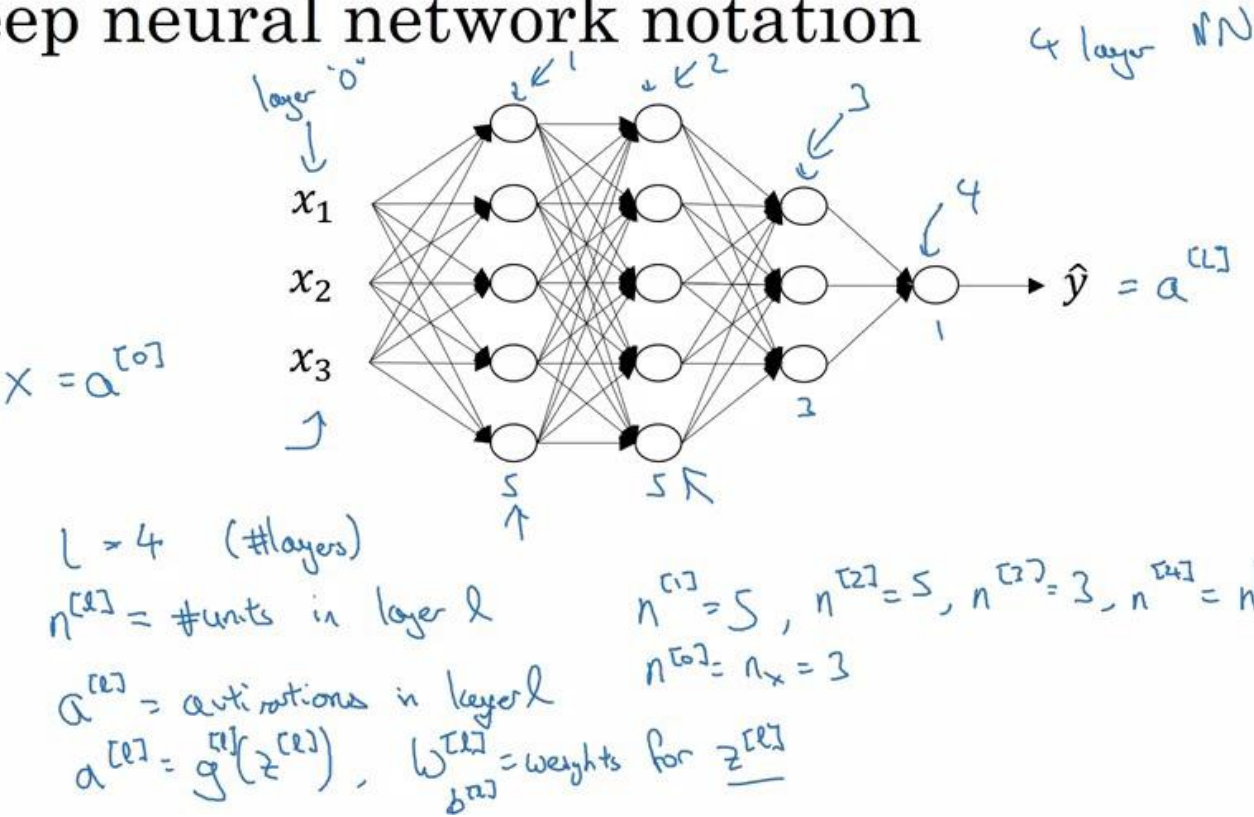
Andrew Ng

What is a deep neural network?



Andrew Ng

Deep neural network notation



Andrew Ng

Forward propagation in a deep network

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

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

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

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

$$z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}, a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$$

$$\begin{cases} z^{[1]} = W^{[1]} a^{[0]} + b^{[1]} \\ a^{[1]} = g^{[1]}(z^{[1]}) \end{cases}$$

$$\begin{cases} z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} = g^{[2]}(z^{[2]}) \end{cases}$$

$$z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}, a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$$

Andrew Ng

Forward propagation in a deep network

$$X : \begin{cases} z^{[1]} = W^{[1]} a^{[0]} + b^{[1]} \\ a^{[1]} = g^{[1]}(z^{[1]}) \end{cases}$$

$$\begin{cases} z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} = g^{[2]}(z^{[2]}) \end{cases}$$

$$z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}, a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$$

$$\begin{cases} z^{[1]} = W^{[1]} A^{[0]} + b^{[1]} \\ A^{[1]} = g^{[1]}(z^{[1]}) \end{cases}$$

$$\begin{cases} z^{[2]} = W^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} = g^{[2]}(z^{[2]}) \end{cases}$$

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

Vertical:

$$\begin{bmatrix} z^{1} & z^{[1](2)} & \dots & z^{[1](4)} \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

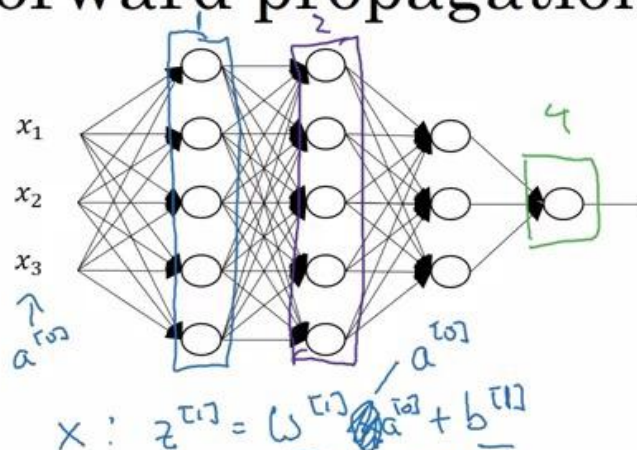
$$\begin{cases} z^{[1]} = W^{[1]} A^{[0]} + b^{[1]} \\ A^{[1]} = g^{[1]}(z^{[1]}) \end{cases}$$

$$\begin{cases} z^{[2]} = W^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} = g^{[2]}(z^{[2]}) \end{cases}$$

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

Andrew Ng

Forward propagation in a deep network



$x: z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$
 $a^{[1]} = g^{[1]}(z^{[1]})$
 $z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$
 $a^{[2]} = g^{[2]}(z^{[2]})$
 $z^{[3]} = W^{[3]} a^{[2]} + b^{[3]}$
 $a^{[3]} = g^{[3]}(z^{[3]}) = \hat{y}$

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

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

Vertical:

$$z^{[1]} = W^{[1]} A^{[0]} + b^{[1]} \rightarrow X = A^{[0]}$$

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

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

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

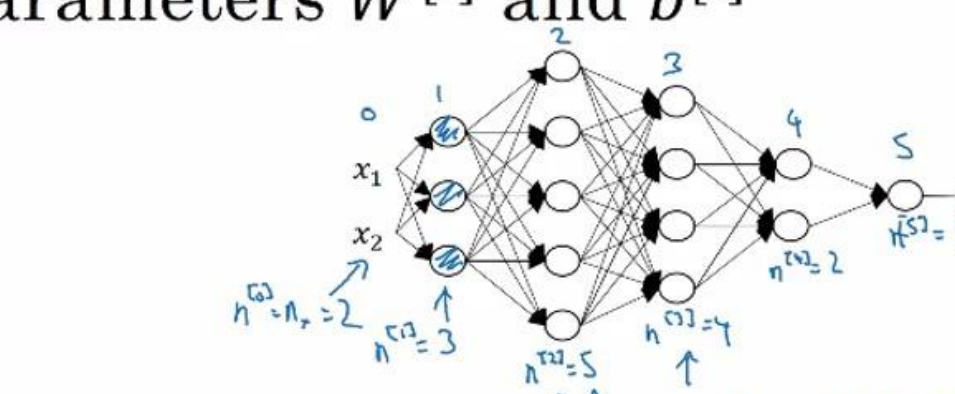
$$z^{[3]} = W^{[3]} A^{[2]} + b^{[3]}$$

$$\hat{y} = g^{[3]}(z^{[3]}) = A^{[3]}$$

for $l=1 \dots 4$

Andrew Ng

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



$z^{[1]} = W^{[1]} \cdot x + b^{[1]}$
 $(3,1) \leftarrow (3,2) \quad (2,1)$
 $(n^{[1]},1) \quad (n^{[1]},n^{[0]}) \quad (n^{[0]},1)$
 $\begin{bmatrix} \vdots \end{bmatrix} = \begin{bmatrix} \vdots \end{bmatrix} \begin{bmatrix} \vdots \end{bmatrix}$

$L=5$

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

$W^{[2]}: (5, 3) \quad (n^{[2]}, n^{[1]})$

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

$\uparrow \quad \uparrow \quad \uparrow$

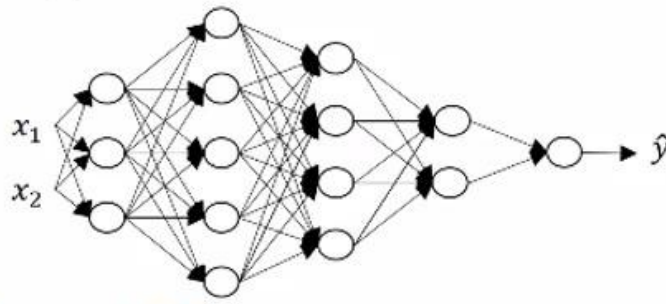
$(5,1) \quad (5,3) \quad (3,1)$

$W^{[3]}: (4, 5)$

$W^{[4]}: (2, 4) \quad , \quad W^{[5]}: (1, 2)$

Andrew Ng

Vectorized implementation



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

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

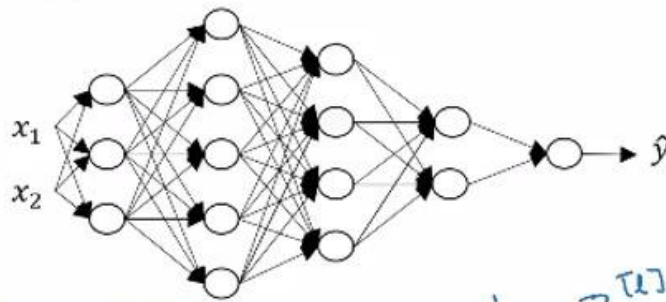
$$[z^{[1]}, z^{[2]}, \dots, z^{[L]}]$$

$$\rightarrow \hat{z}^{[1]} = W^{[1]} \cdot X + b^{[1]}$$

$(n^{[1]}, m)$ $(n^{[1]}, n)$ $(n^{[1]}, m)$ $(n^{[1]}, 1)$
 $n^{[1]}, m$

Andrew Ng

Vectorized implementation



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

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

$$[z^{[1]}, z^{[2]}, \dots, z^{[L]}]$$

$$\rightarrow \hat{z}^{[1]} = W^{[1]} \cdot X + b^{[1]}$$

$(n^{[1]}, m)$ $(n^{[1]}, n)$ $(n^{[1]}, m)$ $(n^{[1]}, 1)$
 $(n^{[1]}, m)$

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

$$z^{[2]}, A^{[2]} : (n^{[2]}, m)$$

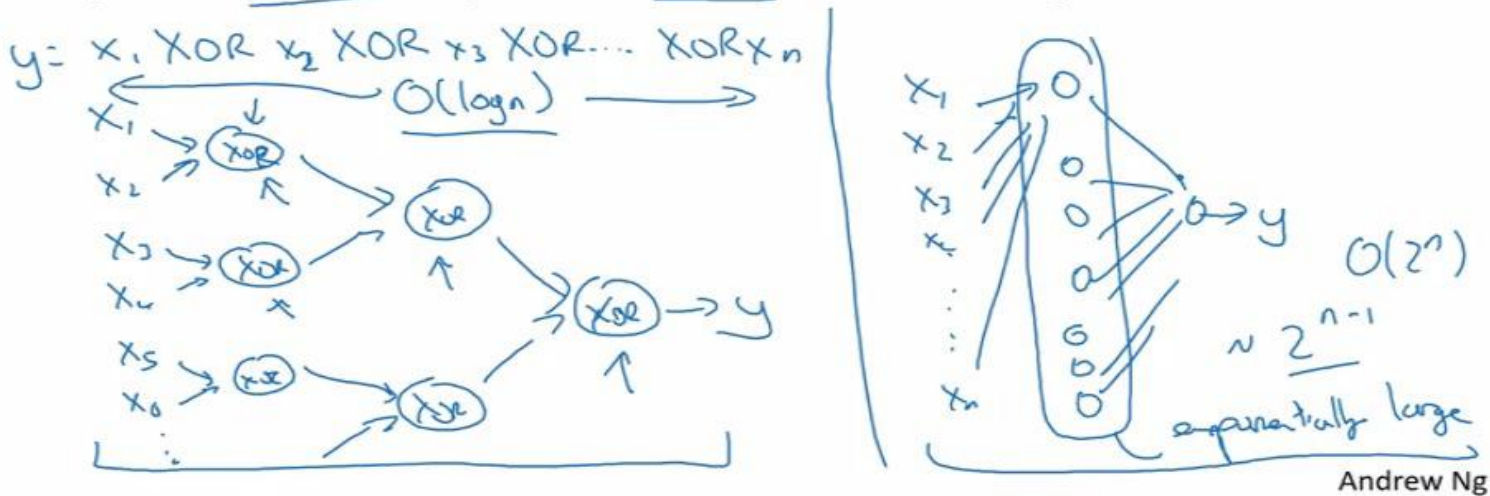
$l=0 \quad A^{[0]} = X = (n^{[0]}, m)$

$$dz^{[2]}, dA^{[2]} : (n^{[2]}, m)$$

Andrew Ng

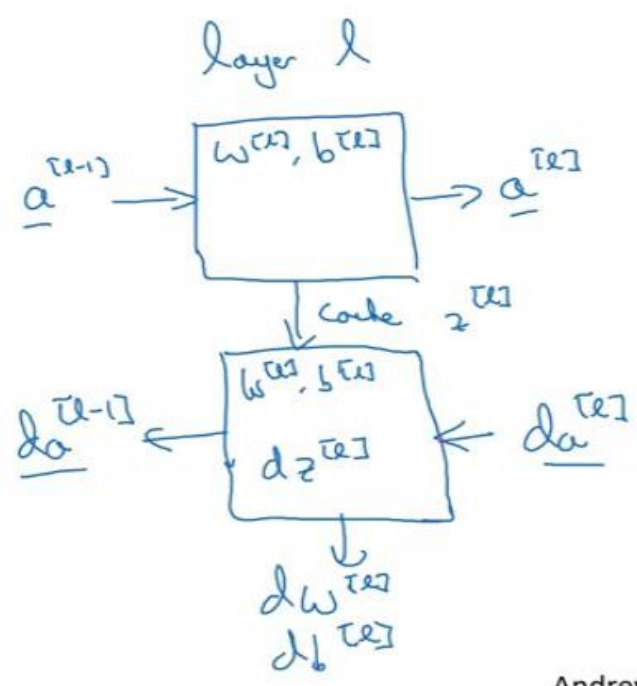
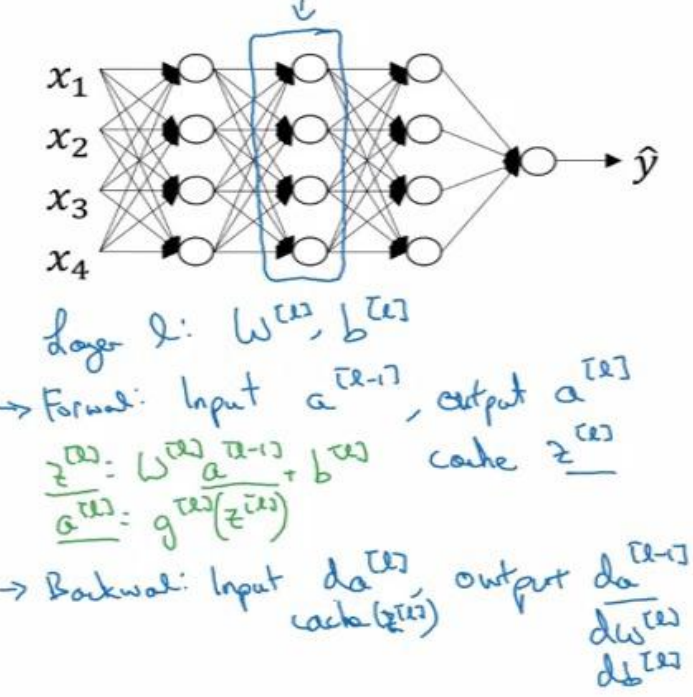
Circuit theory and deep learning

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.



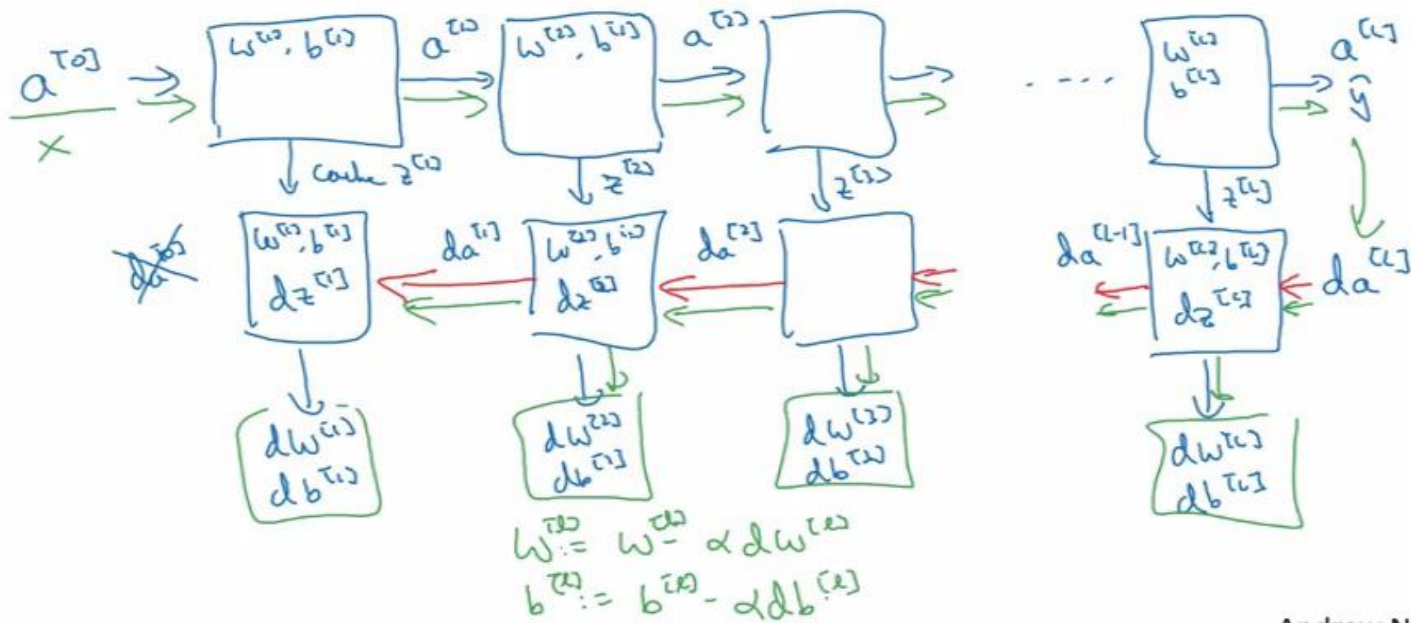
Andrew Ng

Forward and backward functions



Andrew Ng

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

Andrew Ng

What are hyperparameters?

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

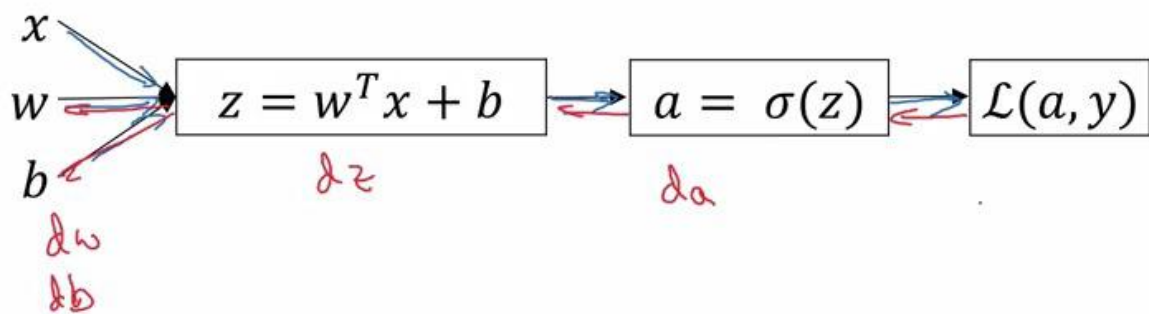
Hyperparameters: $\left. \begin{array}{l} \text{learning rate } \alpha \\ \text{\#iterations} \\ \text{\#hidden layers } L \\ \text{\#hidden units } n^{[1]}, n^{[2]}, \dots \\ \text{choice of activation function} \end{array} \right\}$

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

Andrew Ng

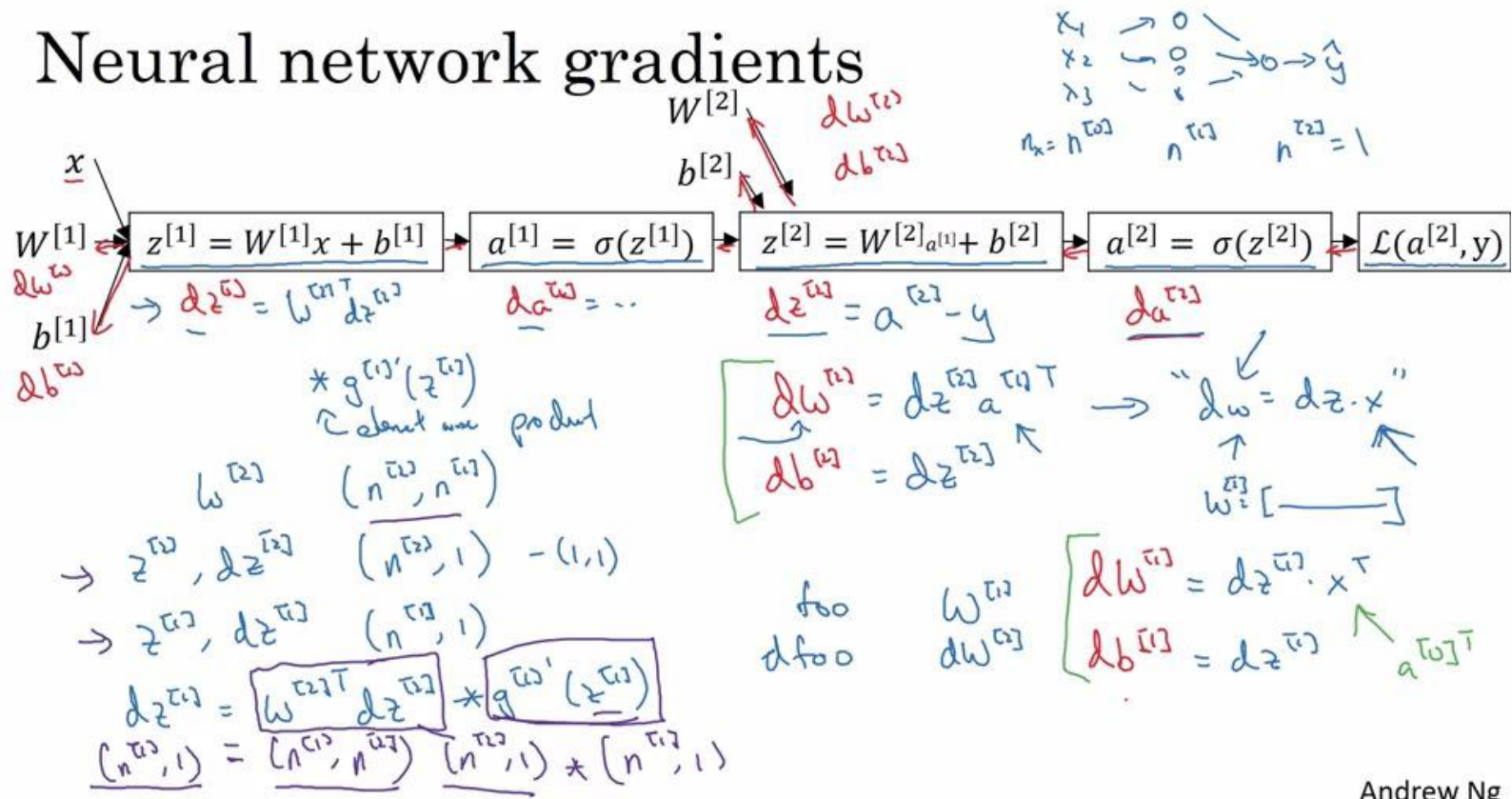
Computing gradients

Logistic regression



Andrew Ng

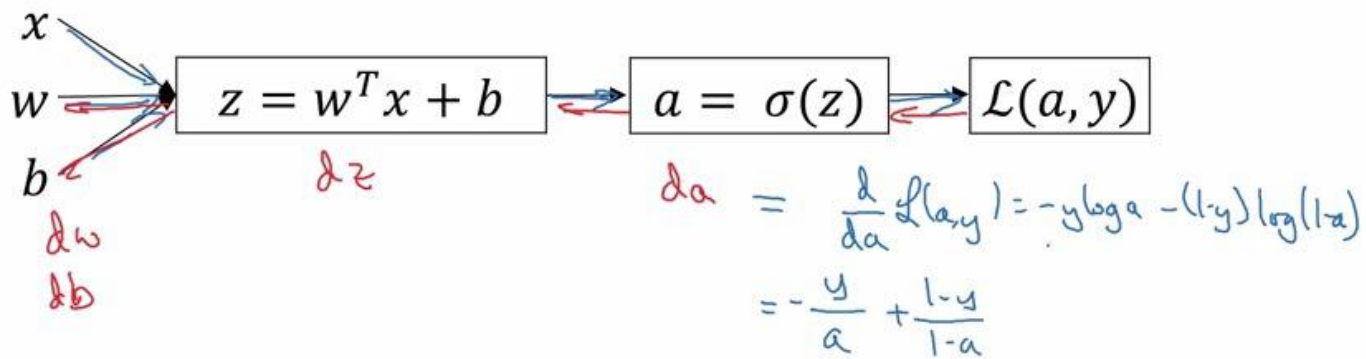
Neural network gradients



Andrew Ng

Computing gradients

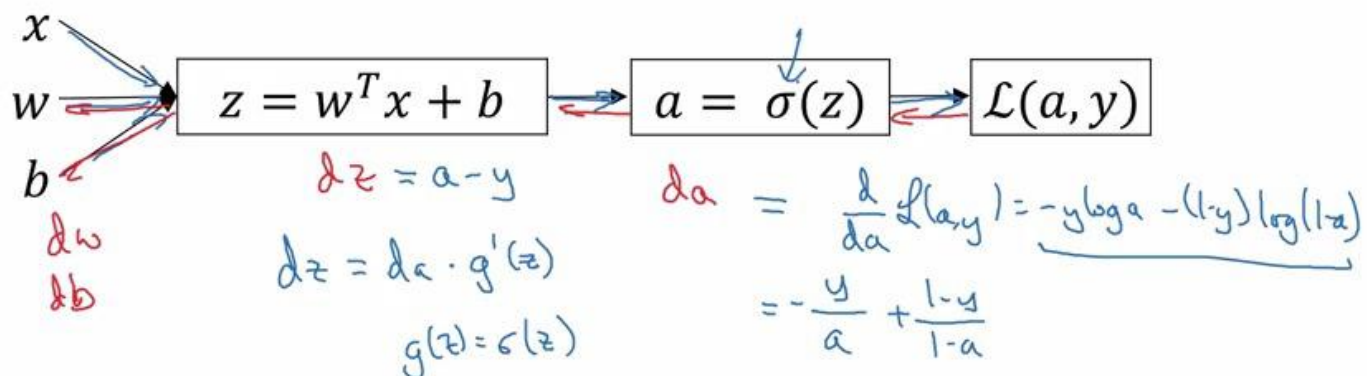
Logistic regression



Andrew Ng

Computing gradients

Logistic regression



Andrew Ng

Summary of gradient descent

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

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

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

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

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

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

Vectorized Implementation:

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

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

$$z^{[2]} = \begin{bmatrix} z^{[2](1)} \\ z^{2} \\ \vdots \\ z^{[2](n)} \end{bmatrix}$$

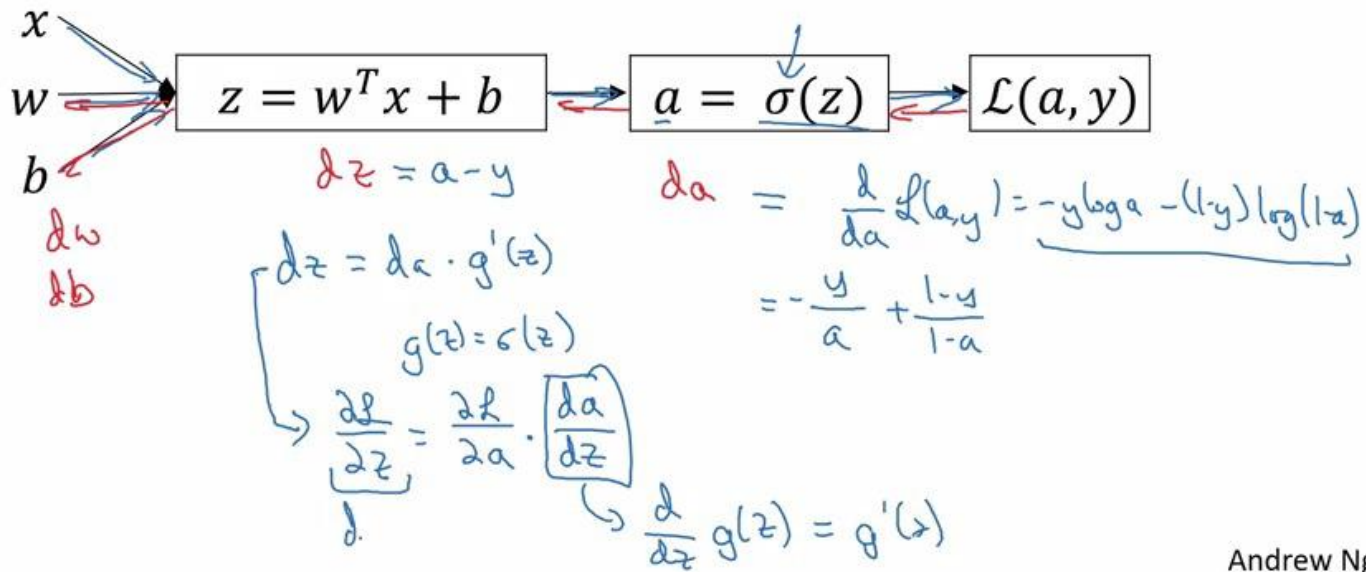
$$z^{[2]} = W^{[2]} X + b^{[2]}$$

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

Andrew Ng

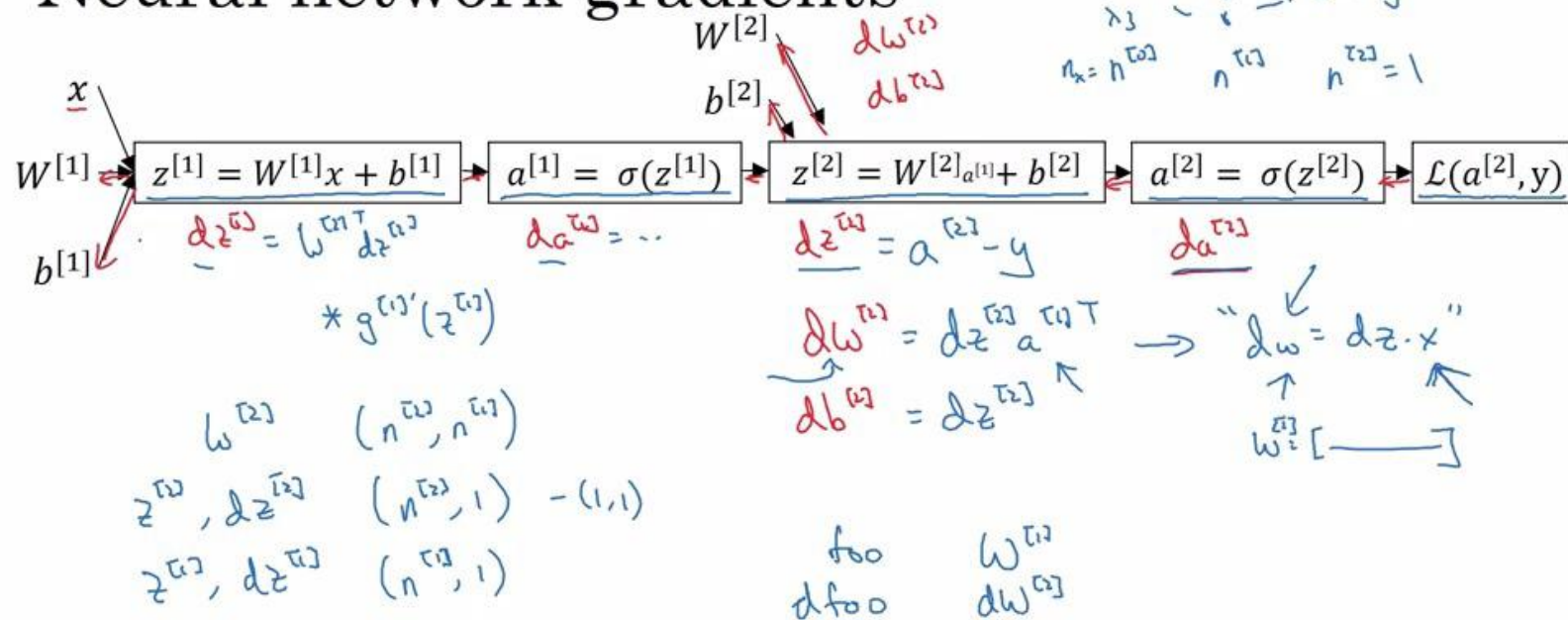
Computing gradients

Logistic regression



Andrew Ng

Neural network gradients



Andrew Ng