# Shallow Neural Network

# What is a Neural Network?

- Logistic Regression

$x_1$
$x_2$  ○ → $\hat{y} = a$
$x_3$

– Computation Graph

$x$
$w$ → $\boxed{z = w^T x + b}$ → $\boxed{a = \sigma(z)}$ → $\boxed{\mathcal{L}(a, y)}$
$b$

# What is a Neural Network?

- Logistic Regression

$x_1$
$x_2$     ◯ → $\hat{y} = a$
$x_3$

   − Computation Graph

$x$

$w$ → | $z = w^T x + b$ | → | $a = \sigma(z)$ | → | $\mathcal{L}(a, y)$ |

$b$

# What is a Neural Network?

- Neural Network

$x_1$
$x_2$     → $\hat{y}$
$x_3$

   − Computation Graph

$x$

$W^{[1]}$ → | $z^{[1]} = W^{[1]}x + b^{[1]}$ | → | $a^{[1]} = \sigma(z^{[1]})$ | → | $z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$ | → | $a^{[2]} = \sigma(z^{[2]})$ | → | $\mathcal{L}(a^{[2]}, y)$ |

$b^{[1]}$

$W^{[2]}$

$b^{[2]}$

# What is a Neural Network?

- Neural Network

$x_1$

$x_2$

$x_3$

$\hat{y}$

[2]

[1]

[i]: represent i-th layer
(i): represent i-th example

- Computation Graph

$x$

$W^{[1]}$

$b^{[1]}$

$$z^{[1]} = W^{[1]}x + b^{[1]}$$ $$a^{[1]} = \sigma(z^{[1]})$$ $$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$ $$a^{[2]} = \sigma(z^{[2]})$$ $$\mathcal{L}(a^{[2]}, y)$$

$W^{[2]}$

$b^{[2]}$

---

# What is a Neural Network?

- Neural Network

$x_1$

$x_2$

$x_3$

$\hat{y}$

[2]

[1]

[i]: represent i-th layer
(i): represent i-th example

- Computation Graph

$x$

$W^{[1]}$

$b^{[1]}$

$$z^{[1]} = W^{[1]}x + b^{[1]}$$ $$a^{[1]} = \sigma(z^{[1]})$$ $$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$ $$a^{[2]} = \sigma(z^{[2]})$$ $$\mathcal{L}(a^{[2]}, y)$$

$W^{[2]}$

$b^{[2]}$

# What is a Neural Network?

- Neural Network



[i]: represent i-th layer
(i): represent i-th example

– Computation Graph



backpropagation

$x$

$W^{[1]}$   $z^{[1]} = W^{[1]}x + b^{[1]}$   $a^{[1]} = \sigma(z^{[1]})$   $z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$   $a^{[2]} = \sigma(z^{[2]})$   $\mathcal{L}(a^{[2]}, y)$

$dW^{[1]}$

$b^{[1]}$

$db^{[1]}$

$dz^{[1]}$   $da^{[1]}$

$dW^{[2]}$   $W^{[2]}$

$db^{[2]}$   $b^{[2]}$

$dz^{[2]}$   $da^{[2]}$

---

# Neural Network Representation



$x_1$

$x_2$

$x_3$

$\hat{y}$

4

# Neural Network Representation



input
layer

contains the inputs to the neural entwork

# Neural Network Representation



input
layer

hidden
layer

in supervised learning, a training set contains inputs as well as outputs
but, the true values for middle nodes are not observed (hidden!!)

# Neural Network Representation



input layer

hidden layer

output layer

responsible for generating the predicted value $\hat{y}$

# Neural Network Representation



input layer

hidden layer

output layer

$a^{[0]} = x$

Notations

# Neural Network Representation

$x_1$

$x_2$

$x_3$
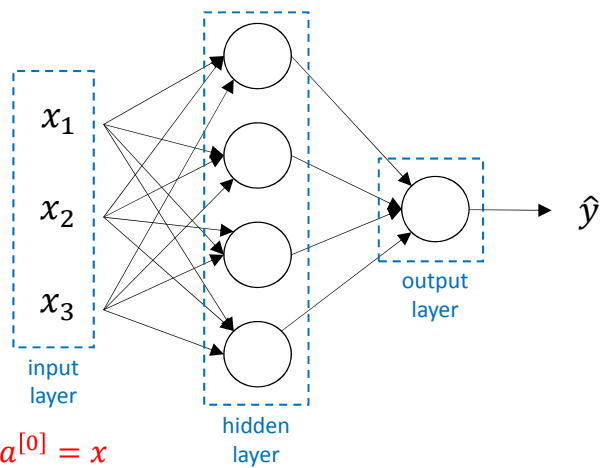
input
layer

$a^{[0]} = x$

Notations

$a_1^{[1]}$

$a_2^{[1]}$

$a_3^{[1]}$

$a_4^{[1]}$

hidden
layer

$a^{[1]}$

output
layer

$\hat{y}$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix}$$

# Neural Network Representation

$x_1$

$x_2$

$x_3$

input
layer

$a^{[0]} = x$

Notations

$a_1^{[1]}$

$a_2^{[1]}$

$a_3^{[1]}$

$a_4^{[1]}$

hidden
layer

$a^{[1]}$

$a^{[2]}$

output
layer

$\hat{y} = a^{[2]}$

# Neural Network Representation



$x_1$

$x_2$

$x_3$

input
layer

$a^{[0]} = x$

$a_1^{[1]}$

$a_2^{[1]}$

$a_3^{[1]}$

$a_4^{[1]}$

hidden
layer

$a^{[1]}$

$a^{[2]}$

output
layer

$\hat{y} = a^{[2]}$

2 Layer Neural Network

why? when we count layers,
we don't count the input layer
(just convention..)

---

# Neural Network Representation

Parameters $W^{[1]T}, b^{[1]T}$
$(4,3)$ $(4,1)$

Transpose why?
$3 \times 4 \rightarrow 4 \times 3$
$1 \times 4 \rightarrow 4 \times 1$
Just for matrix multiplication!

$x_0 = 1$

$b^{[1]T}$

$W^{[1]T}$

$x_1$

$x_2$

$x_3$

input
layer

$a^{[0]} = x$

$a_1^{[1]}$

$a_2^{[1]}$

$a_3^{[1]}$

$a_4^{[1]}$

hidden
layer

$a^{[1]}$

$a^{[2]}$

output
layer

$\hat{y} = a^{[2]}$

2 Layer Neural Network

why? when we count layers,
we don't count the input layer
(just convention..)

8

# Neural Network Representation

Parameters  $W^{[1]T}, b^{[1]T}$

$a_0 = 1$  (4,3)  (4,1)

$b^{[2]T}$

$W^{[2]T}, b^{[2]T}$
(1,4)  (1,1)

$x_1$

$x_2$

$x_3$

$a_1^{[1]}$

$a_2^{[1]}$

$a_3^{[1]}$

$a_4^{[1]}$

$W^{[2]T}$

$a^{[2]}$

$\hat{y} = a^{[2]}$

output
layer

input
layer

hidden
layer

$a^{[0]} = x$

$a^{[1]}$

**2 Layer Neural Network**

why? when we count layers,
we don't count the input layer
(just convention..)

---

# Computing NN's Output

- Let's see details of exactly how NN computes outputs

- Logistic Regression

$x_1$

$x_2$

$x_3$

$w^T x + b$  $\sigma(z)$

$z$       $a$

$a = \hat{y}$

$z = w^T x + b$

$a = \sigma(z)$

# Computing NN's Output

- Let's see details of exactly how NN computes outputs
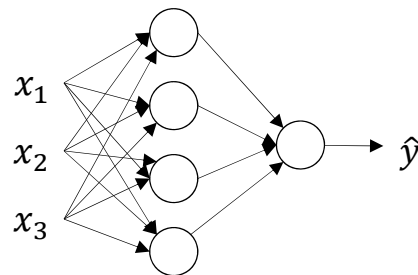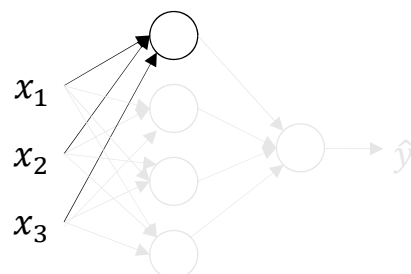
- Neural Network



# Computing NN's Output

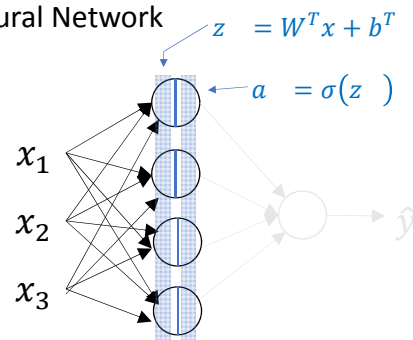- Let's see details of exactly how NN computes outputs

- Neural Network

# Computing NN's Output

- Let's see details of exactly how NN computes outputs

- Neural Network

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

$$a = \sigma(z)$$



$x_1$

$x_2$

$x_3$

$\hat{y}$

---

# Computing NN's Output

- Let's see details of exactly how NN computes outputs

- Neural Network

$$z_1^{[1]} = W_1^{[1]T} x + b_1^{[1]T}$$

$i$ th row

$$a_1^{[1]} = \sigma\left(z_1^{[1]}\right)$$



$x_1$

$x_2$

$x_3$

$\hat{y}$

# Computing NN's Output

- Let's see details of exactly how NN computes outputs

- Neural Network

$$z_1^{[1]} = W_1^{[1]T}x + b_1^{[1]T}$$

$$a_1^{[1]} = \sigma\left(z_1^{[1]}\right)$$

$x_1$

$x_2$

$x_3$

$\hat{y}$

layer

$a_i^{[l]}$

node in layer

---

# Computing NN's Output

- Let's see details of exactly how NN computes outputs

- Neural Network

$x_1$

$x_2$

$x_3$

$\hat{y}$

# Computing NN's Output

- Let's see details of exactly how NN computes outputs

- Neural Network

$$z_2^{[1]} = W_2^{[1]T} x + b_2^{[1]T}$$

$$a_2^{[1]} = \sigma\left(z_2^{[1]}\right)$$

$x_1$

$x_2$

$x_3$

$\hat{y}$

# Computing NN's Output

$x_1$
$x_2$
$x_3$

$a_1^{[1]}$
$a_2^{[1]}$
$a_3^{[1]}$
$a_4^{[1]}$

$\hat{y}$

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

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]T}, \ a_2^{[1]} = \sigma(z_2^{[1]})$$

# Computing NN's Output



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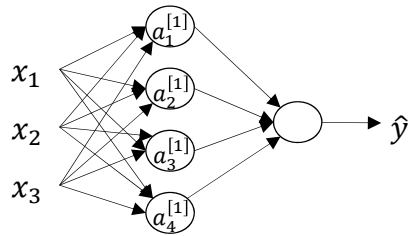$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]T}, \ a_4^{[1]} = \sigma(z_4^{[1]})$$

---

# Computing NN's Output



this is a vector

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

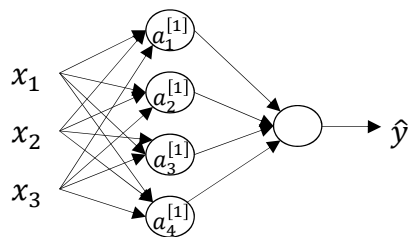for-loop is inefficient

# Computing NN's Output



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

$$W^{[1]}$$

$$\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_1^{[1]T} x \\ w_2^{[1]T} x \\ w_3^{[1]T} x \\ w_4^{[1]T} x \end{bmatrix}$$

(4,3)

---

# Computing NN's Output



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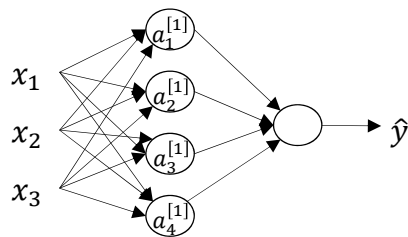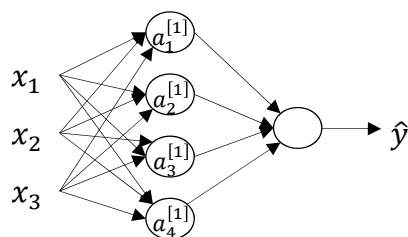$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]T}, \; a_2^{[1]} = \sigma(z_2^{[1]})$$
$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]T}, \; a_3^{[1]} = \sigma(z_3^{[1]})$$
$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]T}, \; a_4^{[1]} = \sigma(z_4^{[1]})$$

$$W^{[1]} \qquad b^{[1]}$$

$$\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]T} \\ b_2^{[1]T} \\ b_3^{[1]T} \\ b_4^{[1]T} \end{bmatrix} = \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]T} \\ w_2^{[1]T} x + b_2^{[1]T} \\ w_3^{[1]T} x + b_3^{[1]T} \\ w_4^{[1]T} x + b_4^{[1]T} \end{bmatrix}$$
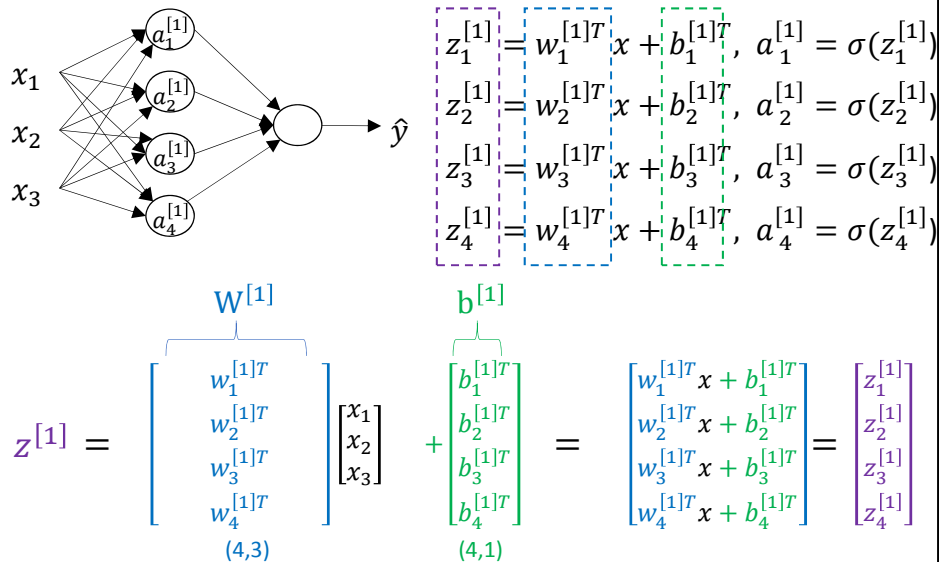
(4,3)      (4,1)

# Computing NN's Output



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

$$z^{[1]} = \overbrace{\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix}}^{W^{[1]}} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \overbrace{\begin{bmatrix} b_1^{[1]T} \\ b_2^{[1]T} \\ b_3^{[1]T} \\ b_4^{[1]T} \end{bmatrix}}^{b^{[1]}} = \begin{bmatrix} w_1^{[1]T}x + b_1^{[1]T} \\ w_2^{[1]T}x + b_2^{[1]T} \\ w_3^{[1]T}x + b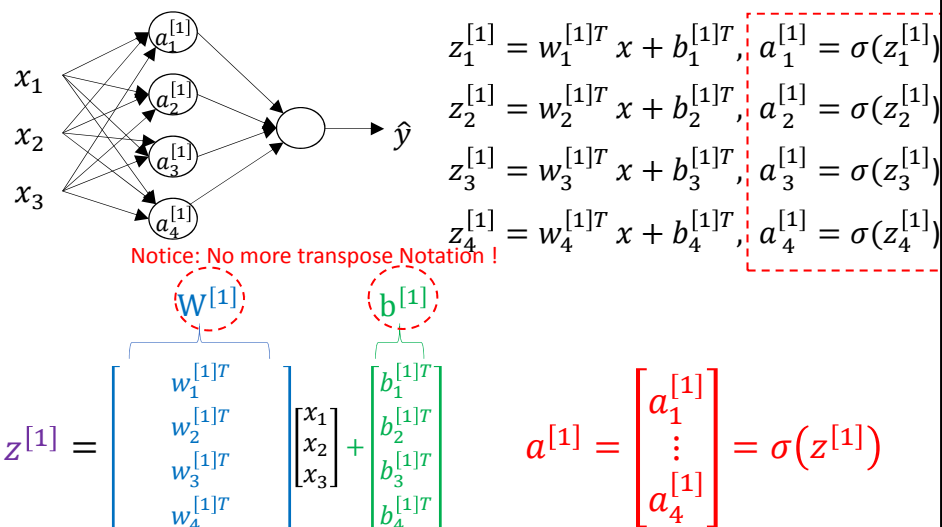_3^{[1]T} \\ w_4^{[1]T}x + b_4^{[1]T} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

(4,3)  (4,1)

# Computing NN's Output



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

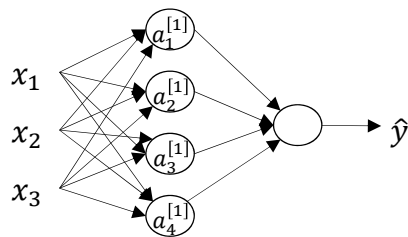Notice: No more transpose Notation !

$$z^{[1]} = \overbrace{\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix}}^{W^{[1]}} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \overbrace{\begin{bmatrix} b_1^{[1]T} \\ b_2^{[1]T} \\ b_3^{[1]T} \\ b_4^{[1]T} \end{bmatrix}}^{b^{[1]}} \qquad a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]})$$

# Computing NN's Output



Given input x:

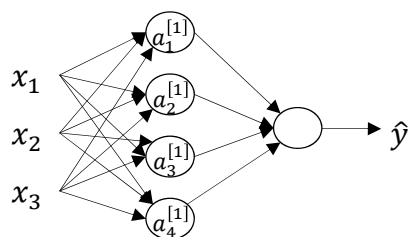$$z^{[1]} = W^{[1]}x + b^{[1]}$$

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

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

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

---

# Computing NN's Output



Given input x:

$$z^{[1]} = W^{[1]}x + b^{[1]}$$
(4,1)    (4,3)  (3,1)  (4,1)

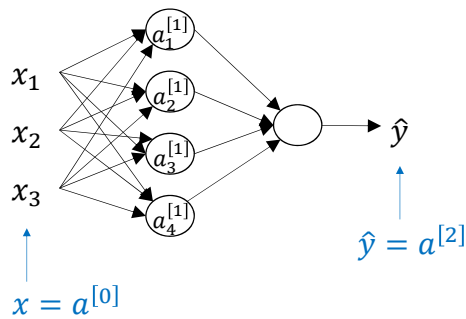$$a^{[1]} = \sigma(z^{[1]})$$
(4,1)    (4,1)

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
(1,1)    (1,4)  (4,1)    (1,1)

$$a^{[2]} = \sigma(z^{[2]})$$
(1,1)    (1,1)

# Computing NN's Output



Given input x:

$$z^{[1]} = W^{[1]}x + b^{[1]}$$
(4,1)   (4,3)   (3,1)   (4,1)

$$a^{[1]} = \sigma(z^{[1]})$$
(4,1)    (4,1)

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
(1,1)   (1,4)   (4,1)    (1,1)

$$a^{[2]} = \sigma(z^{[2]})$$
(1,1)    (1,1)

$\hat{y} = a^{[2]}$

$x = a^{[0]}$

---

# Computing NN's Output



Given input x:

$$z^{[1]} = W^{[1]}x + b^{[1]}$$
(4,1)   (4,3)   (3,1)   (4,1)

$$a^{[1]} = \sigma(z^{[1]})$$
(4,1)    (4,1)

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
(1,1)   (1,4)   (4,1)    (1,1)

$$a^{[2]} = \sigma(z^{[2]})$$
(1,1)    (1,1)

$\hat{y} = a^{[2]}$

$x = a^{[0]}$

# Computing NN's Output

1

$b^{[2]}$

$a^{[1]}$

$W^{[2]}$ → $\hat{y}$

$\hat{y} = a^{[2]}$

Logistic Regression
= 1 layer NN

Given input x:

$$z^{[1]} = W^{[1]}x + b^{[1]}$$
(4,1)   (4,3) (3,1)   (4,1)

$$a^{[1]} = \sigma(z^{[1]})$$
(4,1)          (4,1)

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
(1,1)     (1,4) (4,1)     (1,1)

$$a^{[2]} = \sigma(z^{[2]})$$
(1,1)          (1,1)

---

# Vectorizing across multiple examples

$x_1$

$x_2$

$x_3$

$a_1^{[1]}$

$a_2^{[1]}$

$a_3^{[1]}$

$a_4^{[1]}$

→ $\hat{y}$

$x \longrightarrow a^{[2]} = \hat{y}$

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

# Vectorizing across multiple examples



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

$x \longrightarrow a^{[2]} = \hat{y}$

$x^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)}$

$x^{(2)} \longrightarrow a^{[2](2)} = \hat{y}^{(2)}$

$\vdots \qquad\qquad\qquad \vdots$

number of training data

$x^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)}$

---

# Vectorizing across multiple examples



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

$x \longrightarrow a^{[2]} = \hat{y}$

$x^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)}$

$x^{(2)} \longrightarrow a^{[2](2)} = \hat{y}^{(2)}$

$\vdots \qquad\qquad\qquad \vdots$

$x^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)}$

non-vectorized implementation:

for i = 1 to m,
$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$
$a^{[1](i)} = \sigma(z^{[1](i)})$
$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$
$a^{[2](i)} = \sigma(z^{[2](i)})$

## Vectorizing across multiple examples

```
for i = 1 to m:
```
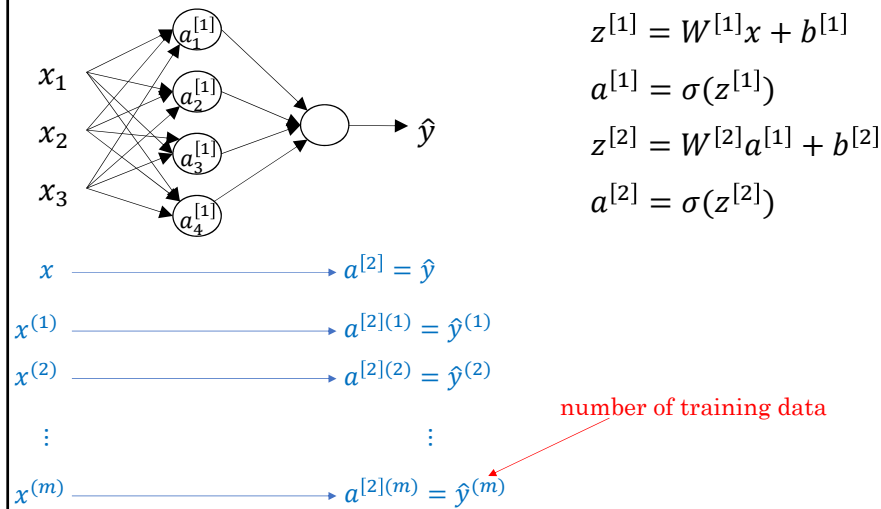$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$
$$a^{[1](i)} = \sigma(z^{[1](i)})$$
$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$
$$a^{[2](i)} = \sigma(z^{[2](i)})$$

## Vectorizing across multiple examples

```
for i = 1 to m:
```
$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$
$$a^{[1](i)} = \sigma(z^{[1](i)})$$
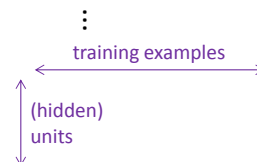$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$
$$a^{[2](i)} = \sigma(z^{[2](i)})$$

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \cdots & x^{(m)} \end{bmatrix}$$
(n,m)

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

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

$$\vdots$$

## Vectorizing across multiple examples

```
for i = 1 to m:
```
$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$
$$a^{[1](i)} = \sigma(z^{[1](i)})$$
$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$
$$a^{[2](i)} = \sigma(z^{[2](i)})$$

$$\mathrm{X} = \begin{bmatrix} x^{(1)} & x^{(2)} & \cdots & x^{(m)} \end{bmatrix}$$
(n,m)

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

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

⋮

training examples

(hidden) units

---

## Vectorizing across multiple examples

```
for i = 1 to m:
```
$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$
$$a^{[1](i)} = \sigma(z^{[1](i)})$$
$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$
$$a^{[2](i)} = \sigma(z^{[2](i)})$$

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$
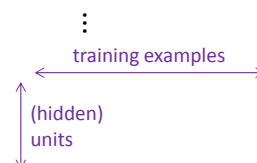$$A^{[1]} = \sigma(Z^{[1]})$$
$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$
$$A^{[2]} = \sigma(Z^{[2]})$$

$$\mathrm{X} = \begin{bmatrix} x^{(1)} & x^{(2)} & \cdots & x^{(m)} \end{bmatrix}$$
(n,m)

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

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

⋮

training examples

(hidden) units

# Vectorizing across multiple examples

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

# Vectorizing across multiple examples

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

$$W^{[1]} = \begin{bmatrix} \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \end{bmatrix}$$

# Vectorizing across multiple examples

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

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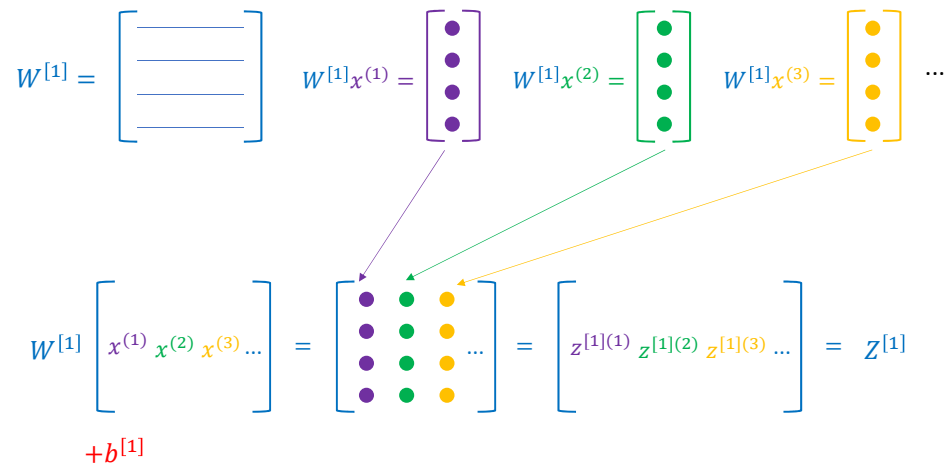
---

# Vectorizing across multiple examples

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

$$W^{[1]} = \begin{bmatrix} \\ \\ \end{bmatrix} \quad W^{[1]}x^{(1)} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix} \quad W^{[1]}x^{(2)} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix} \quad W^{[1]}x^{(3)} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ \bullet \end{bmatrix} \quad \cdots$$

$$W^{[1]} \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \cdots \end{bmatrix} = \begin{bmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{bmatrix} \cdots = \begin{bmatrix} z^{[1](1)} & z^{[1](2)} & z^{[1](3)} & \cdots \end{bmatrix} = Z^{[1]}$$

# Vectorizing across multiple examples

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

$W^{[1]} =$    $W^{[1]}x^{(1)} =$    $W^{[1]}x^{(2)} =$    $W^{[1]}x^{(3)} =$  ⋯

$W^{[1]} \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \cdots \end{bmatrix} = \begin{bmatrix} \cdots \end{bmatrix} = \begin{bmatrix} z^{[1](1)} & z^{[1](2)} & z^{[1](3)} & \cdots \end{bmatrix} = Z^{[1]}$
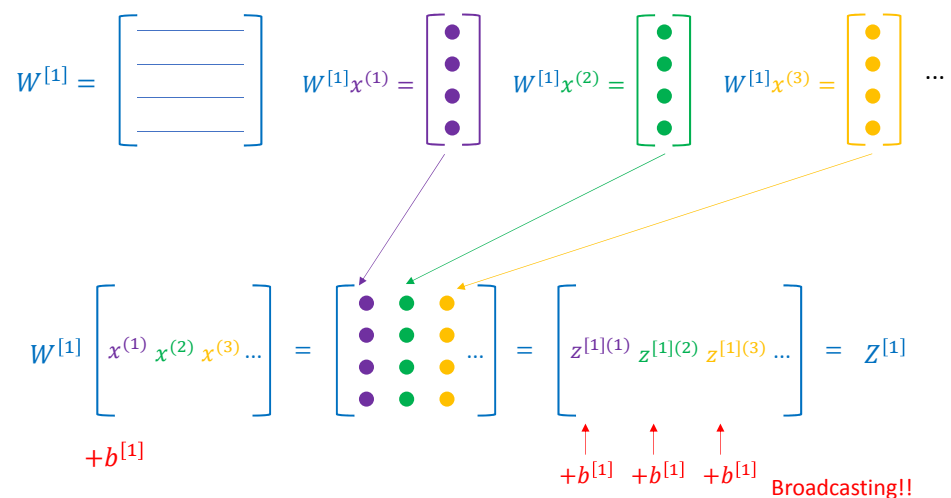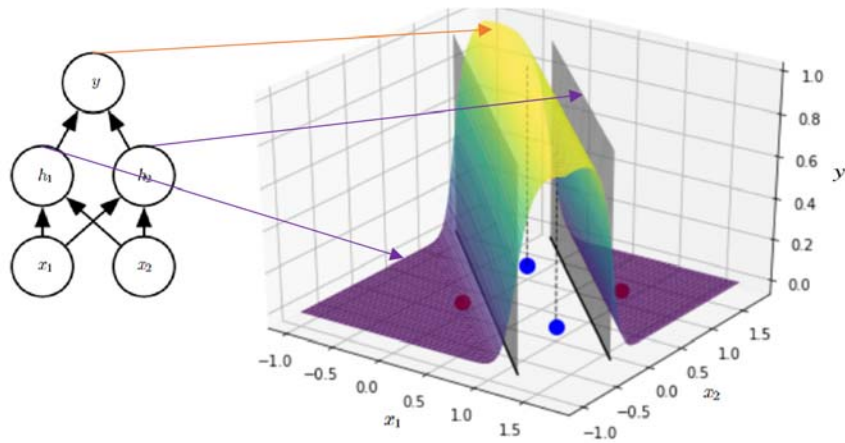
$+b^{[1]}$

---

# Vectorizing across multiple examples

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

$W^{[1]} =$    $W^{[1]}x^{(1)} =$    $W^{[1]}x^{(2)} =$    $W^{[1]}x^{(3)} =$  ⋯

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

$+b^{[1]}$

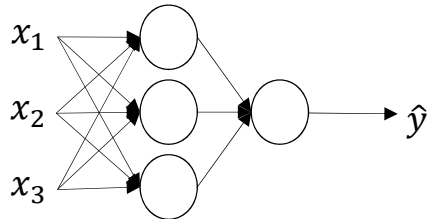$+b^{[1]}$  $+b^{[1]}$  $+b^{[1]}$

Broadcasting!!

25

# A Quick Example:
## 2 Layered (shallow) Neural Network by Gradient Descent



# Another Example:
## Deep (Layered) Neural Network by Gradient Descent

# Activation functions



Given x:

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

# Activation functions



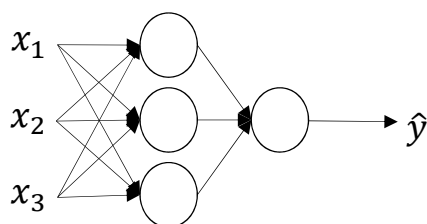In general, sigmoid function is replaced with other non-linear activation functions
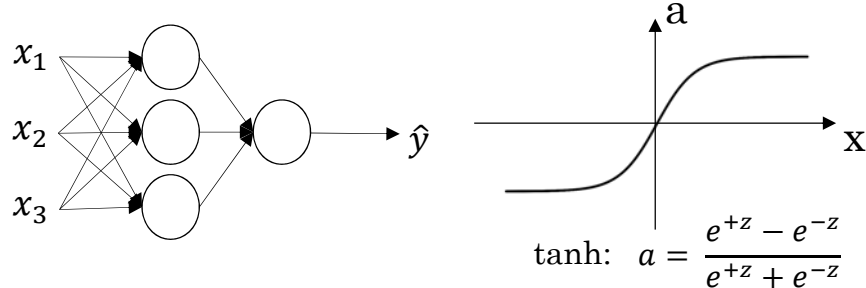
Given x:

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

# Activation functions



Given x:

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

tanh: $a = \dfrac{e^{+z} - e^{-z}}{e^{+z} + e^{-z}}$

---

# Activation functions



Given x:

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

tanh: $a = \dfrac{e^{+z} - e^{-z}}{e^{+z} + e^{-z}}$

**tanh (hyperboic tangent)**
- mathematically a sifted version of the sigmoid function
- almost always works better than the sigmoid function (make learning easier)

# Activation functions



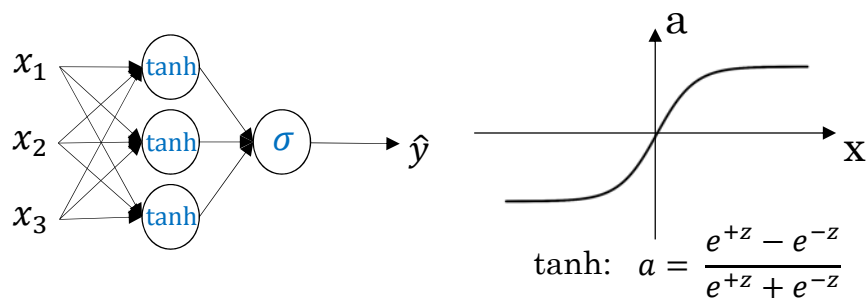$$\text{tanh:} \quad a = \frac{e^{+z} - e^{-z}}{e^{+z} + e^{-z}}$$

Given x:

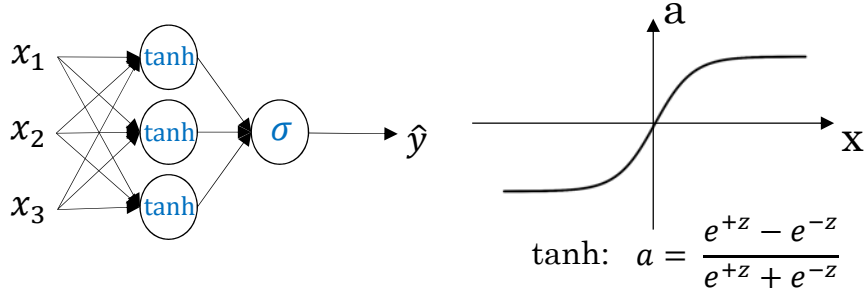$$z^{[1]} = W^{[1]}x + b^{[1]}$$
$$a^{[1]} = g^{[1]}(z^{[1]})$$
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
$$a^{[2]} = g^{[2]}(z^{[2]})$$

**tanh (hyperboic tangent)**
- mathematically a sifted version of the sigmoid function
- almost always works better than the sigmoid function (make learning easier)
- if z is either very large or very small, then the gradient (derivative) becomes very small, and this slows down gradient descent → vanishing gradient

---

# Activation functions



$$\text{ReLU:} \quad a = \max(0, z)$$

Given x:

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

**ReLU (rectified linear unit)**
- solves vanishing gradient problem
- don't need to worry about derivative at origin (just select 1 or 0)

# Activation functions



ReLU:  $a = \max(0, z)$

Given x:

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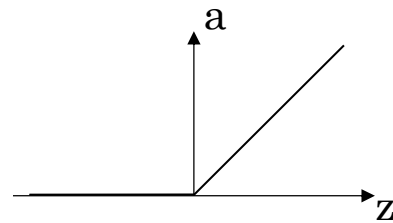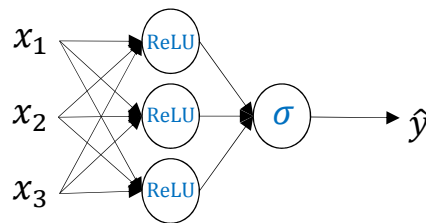$$a^{[1]} = g^{[1]}(z^{[1]})$$
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
$$a^{[2]} = g^{[2]}(z^{[2]})$$

**Rule of Thumb**
- if your output is 0/1 (binary classification), then the sigmoid activation function is very natural for the output layer
- for all other units, ReLU is increasingly the default choice of activation functions

---

# Activation functions



Leaky ReLU:  $a = \max(0.01z, z)$

Given x:

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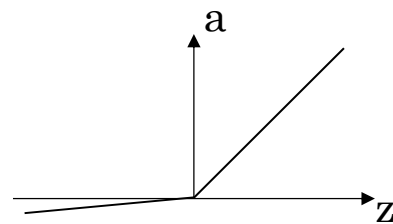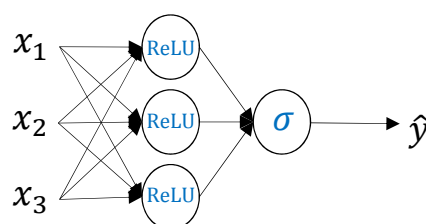$$a^{[1]} = g^{[1]}(z^{[1]})$$
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
$$a^{[2]} = g^{[2]}(z^{[2]})$$

**Rule of Thumb**
- if your output is 0/1 (binary classification), then the sigmoid activation function is very natural for the output layer
- for all other units, ReLU is increasingly the default choice of activation functions
- derivative is equal to 0.01 when z is negative → leaky ReLU

# Recap: Activation functions

sigmoid:  $a = \dfrac{1}{1 + e^{-z}}$

tanh:  $a = \dfrac{e^{+z} - e^{-z}}{e^{+z} + e^{-z}}$

ReLU:  $a = \max(0, z)$

Leaky ReLU:  $a = \max(0.01z, z)$

# Why does neural network need non-linear activation function?

- For your neural network to compute interesting functions, you do need to take a non-linear function
- Try simple linear activation function (identity function)

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

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

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

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

# Why does neural network need non-linear activation function?

- For your neural network to compute interesting functions, you do need to take a non-linear function

- Try simple linear activation function (identity function)
  → neural network just outputs a linear function of the input

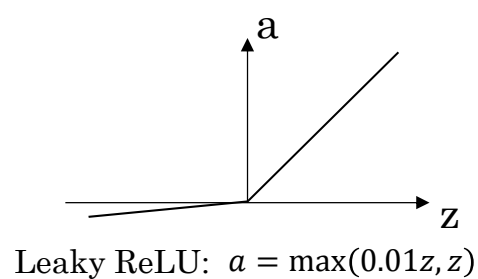$$z^{[1]} = W^{[1]}x + b^{[1]}$$
$$a^{[1]} = g^{[1]}(z^{[1]}) = z^{[1]}$$
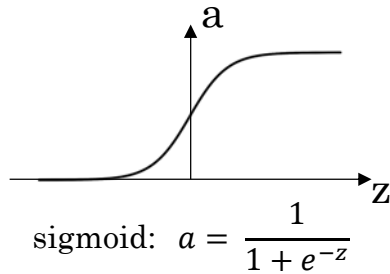$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
$$a^{[2]} = g^{[2]}(z^{[2]}) = z^{[2]}$$

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

---

# Why does neural network need non-linear activation function?

- For your neural network to compute interesting functions, you do need to take a non-linear function

- Try simple linear activation function (identity function)
  → neural network just outputs a linear function of the input



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

# Why does neural network need non-linear activation function?



### how about this?

this model may be no much more expressive than standard logistic regression

---

# Derivatives of activation functions

- When you implement back-propagation for neural network, you need to compute the slope or the derivative of the activation functions

# Sigmoid activation function



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

---

# Sigmoid activation function



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

$$g'(z) = \frac{d}{dz} g(z) = \text{ slope of } g(x) \text{ at } z$$

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

$$= g(z)\big(1 - g(z)\big)$$

# Tanh activation function



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

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

$$g'(z) = \frac{d}{dz} g(z) = \text{ slope of } g(x) \text{ at } z$$

$$= 1 - \big(\tanh(z)\big)^2$$

$$= 1 - g(z)^2 = (1 + g(z))(1 - g(z))$$

# ReLU and Leaky ReLU



ReLU

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

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

Leaky ReLU

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

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

# Gradient descent for neural networks

For neural networks with a single hidden layer

# of units: $n^{[0]}, n^{[1]}, n^{[2]}$

dimension of training data
= input dimension

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
$(n^{[1]}, n^{[0]}) \quad (n^{[1]}, 1) \quad (n^{[2]}, n^{[1]}) \quad (n^{[2]}, 1)$

# Gradient descent for neural networks

For neural networks with a single hidden layer

# of units:   $n^{[0]}, n^{[1]}, n^{[2]}$

Parameters:   $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
      $(n^{[1]}, n^{[0]})$   $(n^{[1]}, 1)$   $(n^{[2]}, n^{[1]})$   $(n^{[2]}, 1)$

Cost function:   $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \dfrac{1}{m} \sum\limits_{i=1}^{m} \mathcal{L}(\hat{y}, y)$   $\nearrow a^{[2]}$

---

# Gradient descent for neural networks

For neural networks with a single hidden layer

# of units:   $n^{[0]}, n^{[1]}, n^{[2]}$

Parameters:   $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
      $(n^{[1]}, n^{[0]})$   $(n^{[1]}, 1)$   $(n^{[2]}, n^{[1]})$   $(n^{[2]}, 1)$

Cost function:   $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \dfrac{1}{m} \sum\limits_{i=1}^{m} \mathcal{L}(\hat{y}, y)$   $\nearrow a^{[2]}$

Gradient descent:      Repeat {

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

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

     Update parameters   $W^{[1]} := W^{[1]} - \alpha \cdot dW^{[1]}$
               $b^{[1]} := b^{[1]} - \alpha \cdot db^{[1]}$
               $W^{[2]} := W^{[2]} - \alpha \cdot dW^{[2]}$
               $b^{[2]} := b^{[2]} - \alpha \cdot db^{[2]}$

       } parameters converged;

# Gradient descent for neural networks

For neural networks with a single hidden layer

# of units:  $n^{[0]}, n^{[1]}, n^{[2]}$

Parameters:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
$\quad\quad\quad\quad (n^{[1]}, n^{[0]}) \quad (n^{[1]}, 1) \quad (n^{[2]}, n^{[1]}) \quad (n^{[2]}, 1)$

$a^{[2]}$

Cost function:  $J\left(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}\right) = \dfrac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}, y)$

Gradient descent:  Repeat {

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

$\quad\quad\quad$ back-propagation $\quad$ Compute gradients $\quad dW^{[1]} = \dfrac{\partial J}{\partial W^{[1]}}, \ b^{[1]} = \dfrac{\partial J}{\partial b^{[1]}}, \dots$

$\quad\quad\quad\quad$ Update parameters $\quad W^{[1]} := W^{[1]} - \alpha \cdot dW^{[1]}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad b^{[1]} := b^{[1]} - \alpha \cdot db^{[1]}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad W^{[2]} := W^{[2]} - \alpha \cdot dW^{[2]}$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad b^{[2]} := b^{[2]} - \alpha \cdot db^{[2]}$

$\quad\quad\quad\quad$ } parameters converged;

---

# Computing gradients

- Logistic regression

$x$

forward propagation

$w \rightarrow$

$b$

$\boxed{z = w^T x + b} \rightarrow \boxed{a = \sigma(z)} \rightarrow \boxed{\mathcal{L}(a, y)}$

# Computing gradients

- Logistic regression

$x$

$w$
$dw$

$z = w^T x + b$    $a = \sigma(z)$    $\mathcal{L}(a, y)$

$dz$    $da$

$b$
$db$

---

# Computing gradients

- Logistic regression

back-propagation

$x$

$$\mathcal{L}(a, y) = -y \log a - (1 - y) \log(1 - a)$$

$w$
$dw$
$= dz \cdot x$

$z = w^T x + b$    $a = \sigma(z)$    $\mathcal{L}(a, y)$

$b$

$db$
$= dz$

$dz$    $da$

$$= \frac{d}{dz} \mathcal{L}(a, y)$$

$$= \frac{d\mathcal{L}}{da} \cdot \frac{da}{dz}$$

$$= da \cdot \sigma'(z)$$

$$= da \cdot \sigma(z) \cdot \big(1 - \sigma(z)\big)$$

$$= da \cdot a \cdot (1 - a)$$

$$= a - y$$

$$= \frac{d}{da} \mathcal{L}(a, y) = -\frac{y}{a} + \frac{1 - y}{1 - a}$$

# Computing gradients

- Two-layer neural network

$x$

$W^{[1]}$

$b^{[1]}$

$W^{[2]}$

$b^{[2]}$

$\boxed{z^{[1]} = W^{[1]}x + b^{[1]}}$ → $\boxed{a^{[1]} = g(z^{[1]})}$ → $\boxed{z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}}$ → $\boxed{a^{[2]} = \sigma(z^{[2]})}$ → $\boxed{\mathcal{L}(a^{[2]}, y)}$

---

# Computing gradients

- Two-layer neural network

$x$

$W^{[1]}$

$dW^{[1]}$

$b^{[1]}$

$db^{[1]}$

$W^{[2]}$

$dW^{[2]}$

$b^{[2]}$

$db^{[2]}$

$\boxed{z^{[1]} = W^{[1]}x + b^{[1]}}$ → $\boxed{a^{[1]} = g(z^{[1]})}$ → $\boxed{z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}}$ → $\boxed{a^{[2]} = \sigma(z^{[2]})}$ → $\boxed{\mathcal{L}(a^{[2]}, y)}$

$dz^{[1]}$

$da^{[1]}$

$dz^{[2]}$

$da^{[2]}$

# Computing gradients

- Two-layer neural network

$W^{[2]}$

$b^{[2]}$

$dW^{[2]} = dz^{[2]}a^{[1]T}$  $(n^{[2]}, n^{[1]}) = (n^{[2]}, 1) \times (1, n^{[1]})$

$db^{[2]} = dz^{[2]}$  $(n^{[2]}, 1)$

$x$

$W^{[1]}$

| $z^{[1]} = W^{[1]}x + b^{[1]}$ | $a^{[1]} = g(z^{[1]})$ | $z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$ | $a^{[2]} = \sigma(z^{[2]})$ | $\mathcal{L}(a^{[2]}, y)$ |

$dW^{[1]}$

$dz^{[1]}$  $da^{[1]}$  $dz^{[2]} = a^{[2]} - y$  $da^{[2]} = -\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}}Z$

$b^{[1]}$

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

element-wise product  $(n^{[1]}, 1) = (n^{[1]}, n^{[2]}) \times (n^{[2]}, 1)$

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

$(n^{[1]}, 1) = (n^{[1]}, n^{[2]}) \times (n^{[2]}, 1) * (n^{[2]}, 1)$

$dW^{[1]} = dz^{[1]} \cdot x^T$  $(n^{[1]}, n^{[0]}) = (n^{[1]}, 1) \times (1, n^{[0]})$

$db^{[1]} = dz^{[1]}$  $(n^{[1]}, 1)$

# Computing gradients

- Two-layer neural network

$W^{[2]}$

$b^{[2]}$

$dW^{[2]} = dz^{[2]}a^{[1]T}$  $(n^{[2]}, n^{[1]}) = (n^{[2]}, 1) \times (1, n^{[1]})$

$db^{[2]} = dz^{[2]}$  $(n^{[2]}, 1)$

$x$

$W^{[1]}$

| $z^{[1]} = W^{[1]}x + b^{[1]}$ | $a^{[1]} = g(z^{[1]})$ | $z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$ | $a^{[2]} = \sigma(z^{[2]})$ | $\mathcal{L}(a^{[2]}, y)$ |

$dW^{[1]}$

$dz^{[1]}$  $da^{[1]}$  $dz^{[2]} = a^{[2]} - y$  $da^{[2]} = -\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}}Z$

$b^{[1]}$

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

element-wise product  $(n^{[1]}, 1) = (n^{[1]}, n^{[2]}) \times (n^{[2]}, 1)$

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

$(n^{[1]}, 1) = (n^{[1]}, n^{[2]}) \times (n^{[2]}, 1) * (n^{[1]}, 1)$

$dW^{[1]} = dz^{[1]} \cdot x^T$  $(n^{[1]}, n^{[0]}) = (n^{[1]}, 1) \times (1, n^{[0]})$

$db^{[1]} = dz^{[1]}$  $(n^{[1]}, 1)$

## Summary of gradient descent

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

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

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

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

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

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

## Summary of gradient descent

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

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

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

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

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

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

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

$dW^{[2]} = \dfrac{1}{m}(A^{[2]} - Y)A^{[1]^T}$

$db^{[2]} = \dfrac{1}{m}(A^{[2]} - Y)I$

$dZ^{[1]} = W^{[2]T}(A^{[2]} - Y) * g^{[1]'}(Z^{[1]})$

$dW^{[1]} = \dfrac{1}{m}(W^{[2]T}(A^{[2]} - Y) * g^{[1]'}(Z^{[1]}))X^T$

$db^{[1]} = \dfrac{1}{m}(W^{[2]T}(A^{[2]} - Y) * g^{[1]'}(Z^{[1]}))I$

# Summary of gradient descent

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

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

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

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]\prime}(z^{[1]})$$

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

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

$$dZ^{[2]} = \underset{(n^{[2]},m)}{A^{[2]}} - \underset{(n^{[2]},m)}{Y}$$

$$dW^{[2]} = \underset{(n^{[2]},n^{[1]})}{\frac{1}{m}} (A^{[2]} - Y) A^{[1]T} = \underset{(n^{[2]},m)}{} \; \underset{(m,n^{[1]})}{}$$

$$db^{[2]} = \underset{(n^{[2]},1)}{\frac{1}{m}} \underset{(n^{[2]},m)}{(A^{[2]} - Y)} \underset{(m,1)}{I}$$

$$dZ^{[1]} = \underset{(n^{[1]},m)}{W^{[2]T}(A^{[2]} - Y)} * \underset{(n^{[1]},m)}{g^{[1]\prime}(Z^{[1]})}$$

$$dW^{[1]} = \underset{(n^{[1]},n^{[0]})}{\frac{1}{m}} (W^{[2]T}(A^{[2]} - Y) * g^{[1]\prime}(Z^{[1]})) X^T = \underset{(n^{[1]},m)}{} \; \underset{(m,n^{[0]})}{}$$

$$db^{[1]} = \underset{(n^{[1]},1)}{\frac{1}{m}} (W^{[2]T}(A^{[2]} - Y) * g^{[1]\prime}(Z^{[1]})) I = \underset{(n^{[1]},m)}{} \; \underset{(m,1)}{}$$