# Basics of neural networks
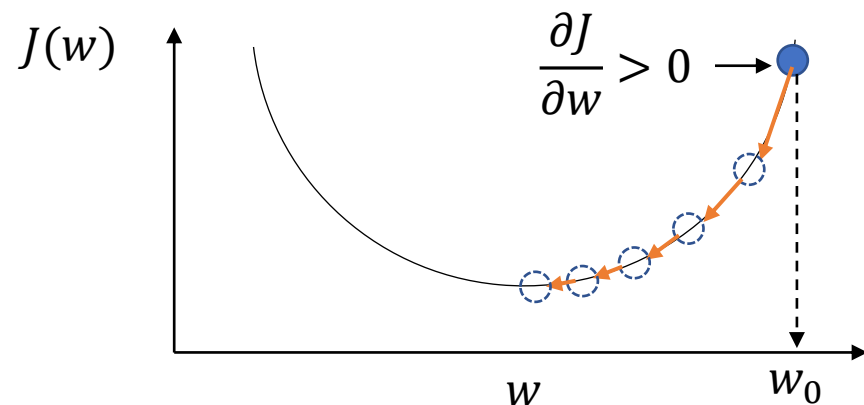
# How to find w and b? Gradient descent

- E.g., Model $y = \sigma(wx + b)$

- Cost function  $J(w, b) = \frac{1}{m} \sum_{i=1}^{m} \left( y(x_d^i) - y_d^i \right)^2$



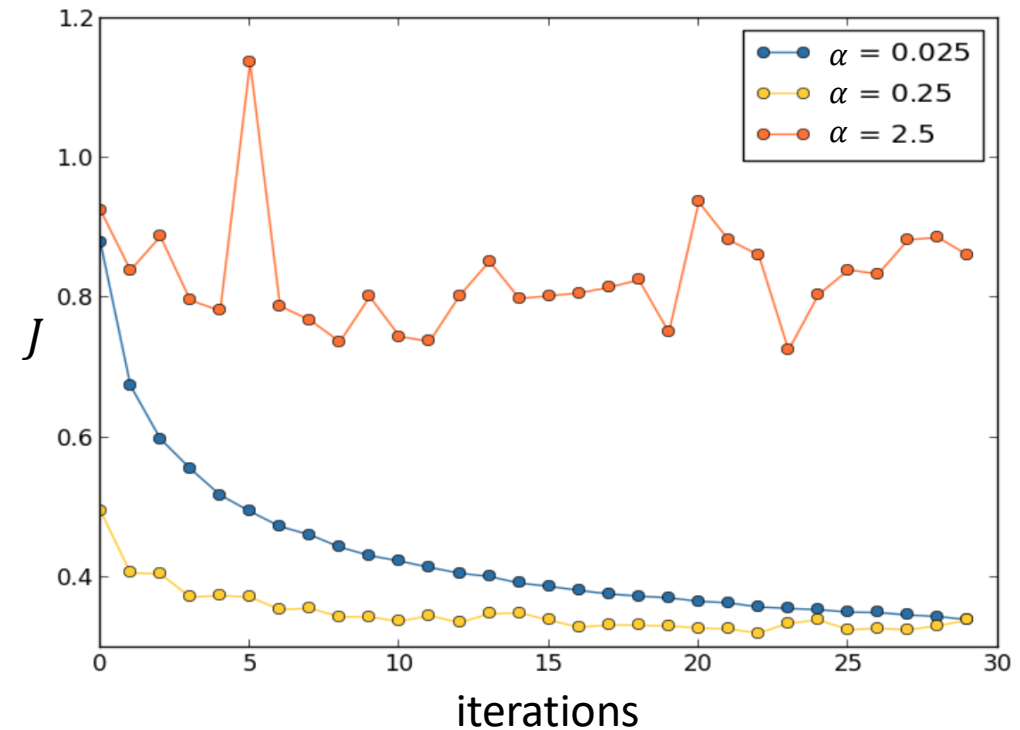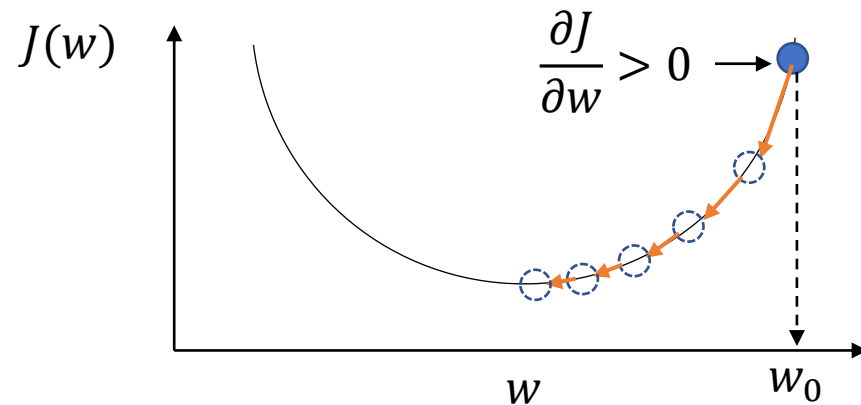I. calculate the slope $\frac{\partial J}{\partial w}$ corresponds to an initial w

II. Adjust $w$ according to the local slope

$w_{new} = w_{old} - \alpha \, \frac{\partial J}{\partial w}, \; \alpha > 0$ is the learning rate

III. Iterate until $\frac{\partial J}{\partial w} = 0$
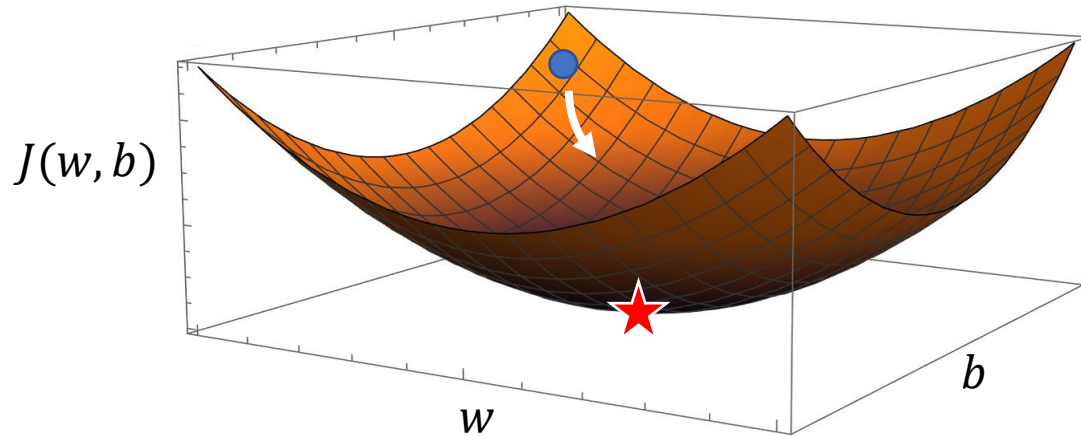
https://developers.google.com/machine-learning/crash-course/fitter/graph

# How to find w and b? Gradient descent

- E.g., Model $y = \sigma(wx + b)$

- Cost function $J(w, b) = \frac{1}{m}\sum_{i=1}^{m}\left(y(x_d^i) - y_d^i\right)^2$

# How to find w and b? Gradient descent

- E.g., Model $y = \sigma(wx + b)$

- Cost function $J(w, b) = \frac{1}{m} \sum_{i=1}^{m} \left( y(x_d^i) - y_d^i \right)^2$



$J(w, b)$

$b$

$w$

**Gradient on a surface**

- $-\nabla J$ gives the direction of the steepest decrease of J

$$-\nabla J(w, b) = -\left( \frac{\partial J}{\partial w}, \frac{\partial J}{\partial b} \right)$$

e.g. $-\nabla J(w_0, b_0) = -(10,1)$

Changing $w$ reduces J 10 times faster than changing $b$

- $-\nabla J$ tells you which weights and biases reduce cost function J the fastest!

# How to find w and b? Gradient descent

- E.g., Model $y = \sigma(wx + b)$

- Cost function $J(w, b) = \frac{1}{m} \sum_{i=1}^{m} \left(y(x_d^i) - y_d^i\right)^2$



$J(w, b)$

$w$

$b$

**Gradient on a surface**
- $-\nabla J$ gives the direction of the steepest decrease of J

$$-\nabla J(w, b) = -\left(\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}\right)$$

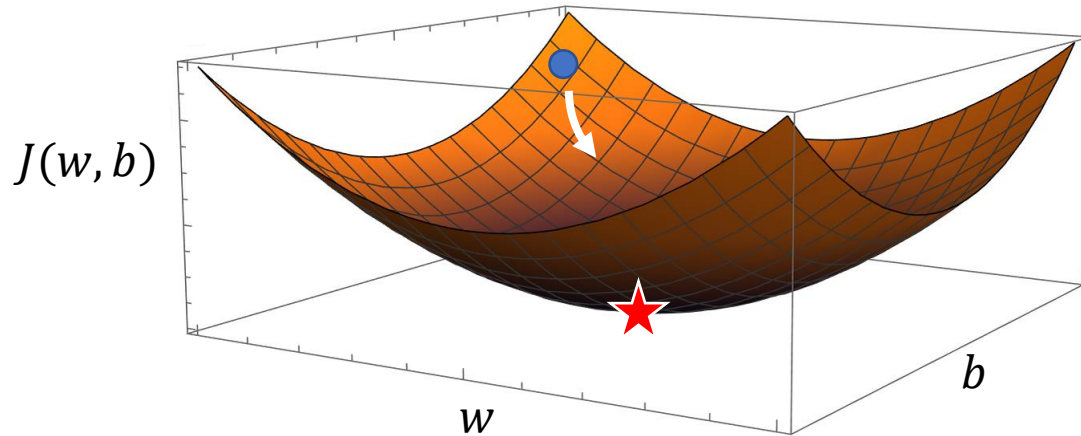- $(w_{new}, b_{new}) = (w_{old}, b_{old}) - \alpha \, \nabla J(w, b)$, $\alpha$ is the learning rate

- Iterate until $\nabla J = 0$

# How to find w and b? Gradient descent

- E.g., Model $y = \sigma(wx + b)$

- Cost function $J(w, b) = \frac{1}{m} \sum_{i=1}^{m} \left( y\left(x_d^i\right) - y_d^i \right)^2$
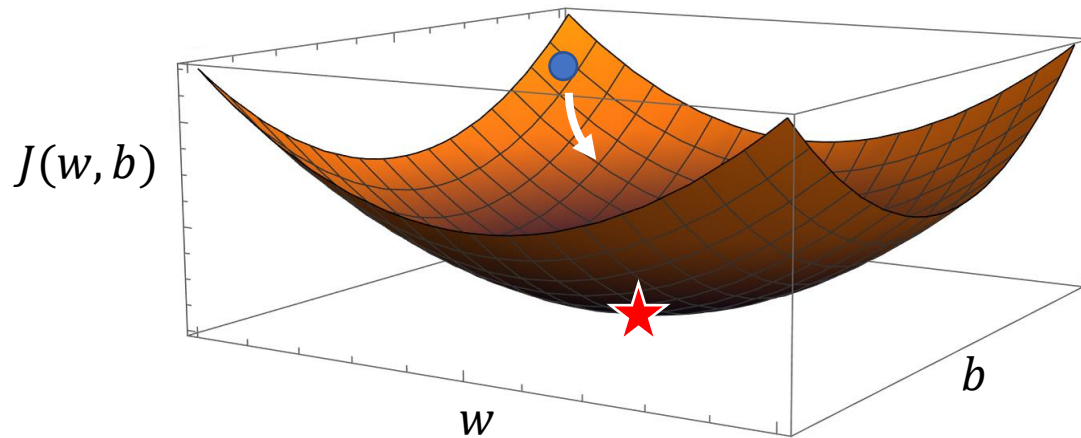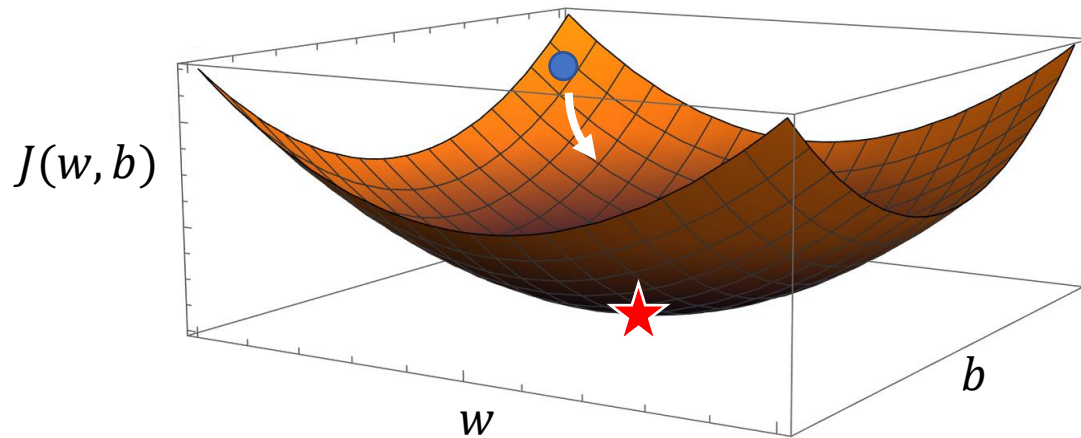


**Gradient on a surface**

- $-\nabla J$ gives the direction of the steepest decrease of J

$$-\nabla J(w, b) = -\left( \frac{\partial J}{\partial w}, \frac{\partial J}{\partial b} \right)$$

But how are $\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}$ calculated?

# How to find w and b? Gradient descent

- E.g., Model $y = \sigma(wx + b)$

- Cost function $J(w, b) = \frac{1}{m}\sum_{i=1}^{m}\left(y(x_d^i) - y_d^i\right)^2 = \frac{1}{m}\sum_{i=1}^{m} L(y^i, y_d^i)$

$$L \equiv (y(x) - y_d)^2, \quad J = \frac{1}{m}\sum_{i=1}^{m} L$$

Back propagation (chain rule)

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y}\frac{\partial y}{\partial(wx+b)}\frac{\partial(wx+b)}{\partial w} = 2(y - y_d)\sigma' x$$

1) For one example:

$$\frac{\partial L}{\partial w}(w, b, x_d^i, y_d^i) = 2(\sigma(wx_d^i + b) - y_d^i)\sigma' x_d^i$$

2) For the full data set:

$$\frac{\partial J}{\partial w} = \frac{1}{m}\sum_{i=1}^{m}\frac{\partial L}{\partial w} = \frac{1}{m}\sum_{i=1}^{m} 2(\sigma(wx_d^i + b) - y_d^i)\sigma' x_d^i$$



$J(w, b)$

$w$

$b$

# Example:

$$x \xrightarrow{\ w\ } \bigcirc \longrightarrow y$$

- Linear Regression Model $y = wx$

- Cost function $J(w, b) = \frac{1}{m}\sum_{i=1}^{m}\left(y(x_d^i) - y_d^i\right)^2 = \frac{1}{m}\sum_{i=1}^{m} L(y^i, y_d^i)$

1) For one example:

$$L = \left(y(x_d^i) - y_d^i\right)^2$$

2) For the full data set:

$$J(w) = \frac{1}{m}\sum_{i=1}^{m} L$$

$|x_d^i|$

$y_d^i$

$$x_d^i, w \longrightarrow y \longrightarrow J(w)$$

$$\text{Update } w \longleftarrow \frac{\partial J}{\partial w}(w)$$

$$\boxed{L \equiv (y(x) - y_d)^2, \quad J = \frac{1}{m}\sum_{i=1}^{m} L \quad \rightarrow \quad \frac{\partial J}{\partial w} = \frac{1}{m}\sum_{i=1}^{m}\frac{\partial L}{\partial w}}$$

Back propagation (chain rule)
$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y}\frac{\partial y}{\partial (wx)}\frac{\partial (wx)}{\partial w} = 2(y - y_d)x$$

1) For one example:
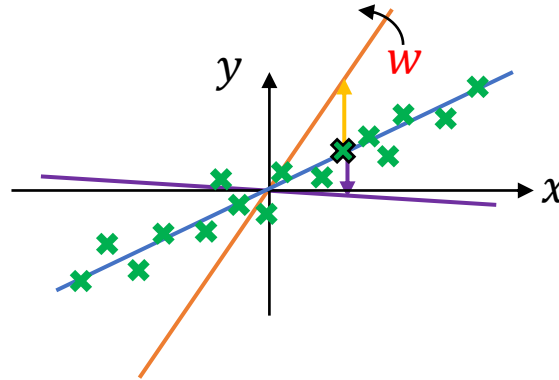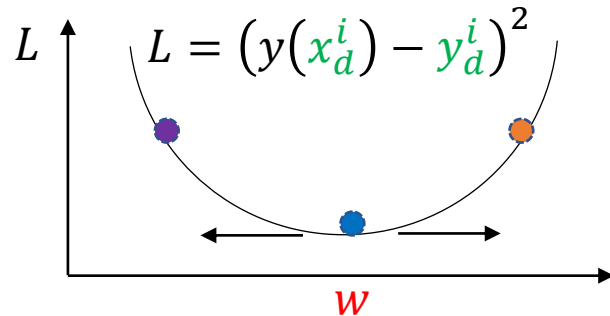$$\frac{\partial L}{\partial w}(w, b, x_d^i, y_d^i) = 2\left(y(x_d^i) - y_d^i\right)x_d^i$$

2) For the full data set:
$$\frac{\partial J}{\partial w} = \frac{1}{m}\sum_{i=1}^{m}\frac{\partial L}{\partial w} = \frac{1}{m}\sum_{i=1}^{m} 2\left(y(x_d^i) - y_d^i\right)x_d^i$$

8

# How to find w and b? Gradient descent

- E.g., Model $y = \sigma(wx + b)$

- Cost function $J(w, b) = \frac{1}{m} \sum_{i=1}^{m} \left( y(x_d^i) - y_d^i \right)^2 = \frac{1}{m} \sum_{i=1}^{m} L(y^i, y_d^i)$



- $-\boldsymbol{\nabla} J(w, b) = -\left( \frac{\partial J}{\partial w}, \frac{\partial J}{\partial b} \right)$ **for a given data set at a given $\boldsymbol{w}, \boldsymbol{b}$ is known analytically**

- $(w_{new}, b_{new}) = (w_{old}, b_{old}) - \alpha \nabla J(w, b)$, $\alpha$ is the learning rate

- Iterate until $\nabla J = 0$

9

# So far we have discussed…

- Universal function approximator

- Gradient descent: a method to find weights and biases that minimize J

- Calculate $-\boldsymbol{\nabla} J(w, b) = -\left(\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}\right)$ for a simple model $y = \sigma(wx + b)$
  - One example
  - A full dataset

We know how to find w and b for a simple model $y = \sigma(wx + b)$.

What about finding w and b in a neural network?

# Writing a neural network in terms of vectors

$$y = a_1^{(2)} = \sigma\left(\sum_{k=1}^{3} w_{1k}^{(2)} \sigma\left(\sum_{l=1}^{2} w_{kl}^{(1)} x_l + b_k^{(1)}\right) + b_1^{(2)}\right)$$

$w_{jk}^{(1)}$

$b_1^{(1)}$

$a_1^{(1)}$

$w_{jk}^{(2)}$

$x_1$

$b_2^{(1)}$

$a_2^{(1)}$

$b_1^{(2)}$

$a_1^{(2)}$

$y$

$x_2$

$b_3^{(1)}$

$a_3^{(1)}$

Input layer

Hidden layer

Output layer

| # of units | $j=1,2$ | $j=1,2,3$ | $j=1$ |
|---|---|---|---|
| Forward propagation | $a_j^{(0)} = x_j$ | $z_j^{(1)} = \sum_{k=1}^{2} w_{jk}^{(1)} a_k^{(0)} + b_j^{(1)}$ <br> $a_j^{(1)} = \sigma\left(z_j^{(1)}\right)$ | $z_j^{(2)} = \sum_{k=1}^{3} w_{jk}^{(2)} a_k^{(1)} + b_j^{(2)}$ <br> $a_j^{(2)} = \sigma\left(z_j^{(2)}\right) = z_j^{(2)}$ |

# Writing a neural network in terms of vectors

$$b_1^{(1)}$$

$$\boldsymbol{w}_1^{(1)} \quad a_1^{(1)}$$

$$b_2^{(1)} \qquad \boldsymbol{w}_1^{(2)} \qquad b_1^{(2)}$$

$$x_1 \qquad \boldsymbol{w}_2^{(1)} \quad a_2^{(1)} \qquad a_1^{(2)} \longrightarrow y$$

$$x_2$$

$$b_3^{(1)}$$

$$\boldsymbol{w}_3^{(1)} \quad a_3^{(1)}$$

Input layer

Hidden layer

Output layer

$$\boldsymbol{a}^{(0)} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \boldsymbol{w}_j^{(1)} = \begin{bmatrix} w_{j1}^{(1)} \\ w_{j2}^{(1)} \end{bmatrix}$$

$$\boldsymbol{a}^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \end{bmatrix}, \boldsymbol{w}_j^{(2)} = \begin{bmatrix} w_{j1}^{(2)} \\ w_{j2}^{(2)} \\ w_{j3}^{(2)} \end{bmatrix}$$

| # of units | $j=1,2$ | $j=1,2,3$ | $j=1$ |
|---|---|---|---|
| Forward propagation | $\boldsymbol{a}^{(0)} = \boldsymbol{x}$ | $z_j^{(1)} = \mathbf{w}_j^{(1)T} \boldsymbol{a}^{(0)} + b_j^{(1)}$ <br> $a_j^{(1)} = \sigma(z_j^{(1)})$ | $z_j^{(2)} = \mathbf{w}_j^{(2)T} \boldsymbol{a}^{(1)} + b_j^{(2)}$ <br> $a_j^{(2)} = \sigma(z_j^{(2)})$ |

13

# Writing a neural network in terms of vectors



3 x 2  3 x 1
$\boldsymbol{W}^{(1)}, \boldsymbol{b}^{(1)}$

1 x 3  1 x 1
$\boldsymbol{W}^{(2)}, \boldsymbol{b}^{(2)}$

$a_1^{(1)}$

$a_2^{(1)}$

$a_1^{(2)}$

$a_3^{(1)}$

$y$

Input
layer

Hidden
layer

Output
layer

$$\boldsymbol{a}^{(0)} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \boldsymbol{W}^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} \\ w_{31}^{(1)} & w_{32}^{(1)} \end{bmatrix}$$
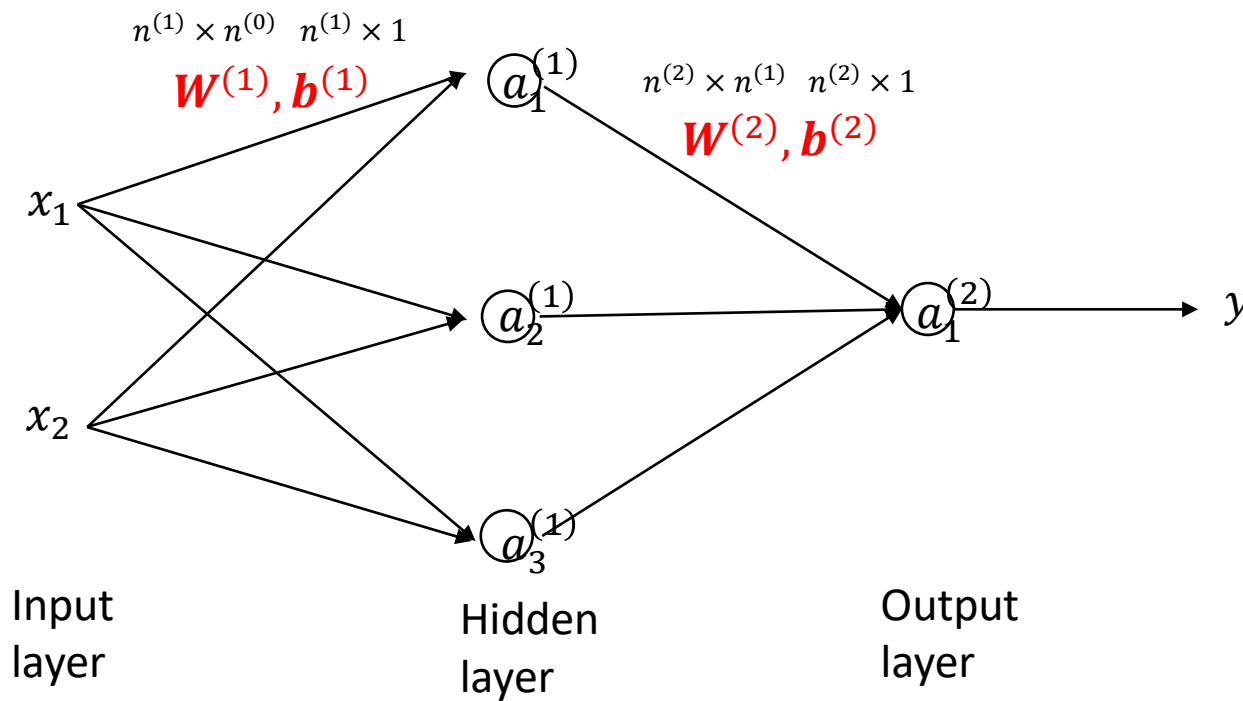
$$\boldsymbol{b}^{(1)} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{bmatrix}, \boldsymbol{z}^{(1)} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \end{bmatrix},$$

$$\boldsymbol{a}^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \end{bmatrix},$$

$$\boldsymbol{W}^{(2)} = \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} & w_{13}^{(2)} \end{bmatrix}$$

| # of units | $j=1,2$ | $j=1,2,3$ | $j=1$ |
|---|---|---|---|
| Forward propagation | $\boldsymbol{a}^{(0)} = \boldsymbol{x}$ | $\boldsymbol{z}^{(1)} = \boldsymbol{W}^{(1)}\boldsymbol{a}^{(0)} + \boldsymbol{b}^{(1)}$ <br> $\boldsymbol{a}^{(1)} = \sigma(\boldsymbol{z}^{(1)})$ | $\boldsymbol{z}^{(2)} = \boldsymbol{W}^{(2)}\boldsymbol{a}^{(1)} + \boldsymbol{b}^{(2)}$ <br> $\boldsymbol{a}^{(2)} = \sigma(\boldsymbol{z}^{(2)})$ |

# Writing a neural network in terms of vectors



$n^{(1)} \times n^{(0)} \quad n^{(1)} \times 1$

$\boldsymbol{W}^{(1)}, \boldsymbol{b}^{(1)}$

$n^{(2)} \times n^{(1)} \quad n^{(2)} \times 1$

$\boldsymbol{W}^{(2)}, \boldsymbol{b}^{(2)}$

$a_1^{(1)}$

$a_2^{(1)}$

$a_3^{(1)}$

$a_1^{(2)}$

$y$

$x_1$

$x_2$

Input layer

Hidden layer

Output layer

$$\boldsymbol{a}^{(0)} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \boldsymbol{W}^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} \\ w_{31}^{(1)} & w_{32}^{(1)} \end{bmatrix}$$
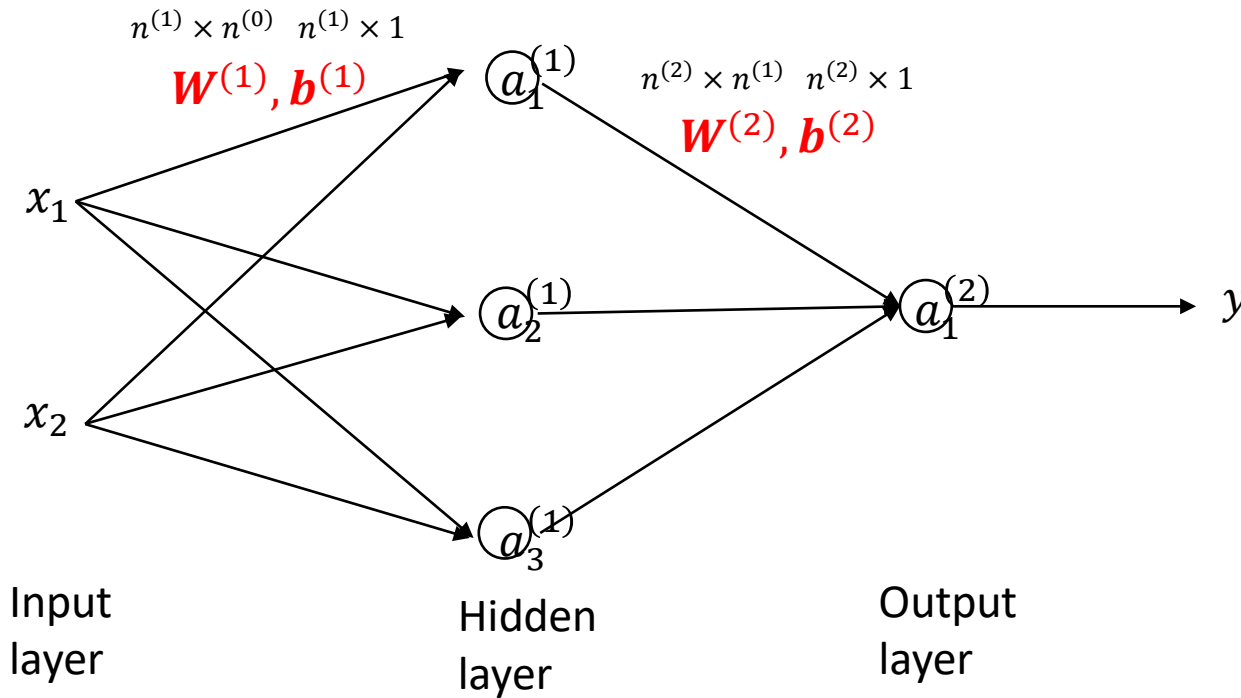
$$\boldsymbol{b}^{(1)} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{bmatrix}, \boldsymbol{z}^{(1)} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \end{bmatrix},$$

$$\boldsymbol{a}^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \end{bmatrix},$$

$$\boldsymbol{W}^{(2)} = \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} & w_{13}^{(2)} \end{bmatrix}$$

| | # of units | | |
|---|---|---|---|
| # of units | $n^{(0)} = 2$ | $n^{(1)} = 3$ | $n^{(2)} = 1$ |
| Forward propagation | $\boldsymbol{a}^{(0)} = \boldsymbol{x}$ | $\boldsymbol{z}^{(1)} = \boldsymbol{W}^{(1)}\boldsymbol{a}^{(0)} + \boldsymbol{b}^{(1)}$ $\boldsymbol{a}^{(1)} = \sigma(\boldsymbol{z}^{(1)})$ | $\boldsymbol{z}^{(2)} = \boldsymbol{W}^{(2)}\boldsymbol{a}^{(1)} + \boldsymbol{b}^{(2)}$ $\boldsymbol{a}^{(2)} = \sigma(\boldsymbol{z}^{(2)})$ |

# Writing a neural network in terms of vectors

$$n^{(1)} \times n^{(0)} \quad n^{(1)} \times 1$$

$$\boldsymbol{W}^{(1)}, \boldsymbol{b}^{(1)}$$

$$a_1^{(1)}$$

$$n^{(2)} \times n^{(1)} \quad n^{(2)} \times 1$$

$$\boldsymbol{W}^{(2)}, \boldsymbol{b}^{(2)}$$

$$x_1$$

$$a_2^{(1)}$$

$$a_1^{(2)} \longrightarrow y$$

$$x_2$$

$$a_3^{(1)}$$

Input
layer

Hidden
layer

Output
layer

$$\boldsymbol{a}^{(0)} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \boldsymbol{W}^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} \\ w_{31}^{(1)} & w_{32}^{(1)} \end{bmatrix}$$

$$\boldsymbol{b}^{(1)} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{bmatrix}, \boldsymbol{z}^{(1)} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \end{bmatrix},$$

$$\boldsymbol{a}^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \end{bmatrix},$$

$$\boldsymbol{W}^{(2)} = \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} & w_{13}^{(2)} \end{bmatrix}$$

In short, $\quad y = a_1^{(2)} = \sigma \left( \sum_{k=1}^{3} w_{1k}^{(2)} \sigma \left( \sum_{l=1}^{2} w_{kl}^{(1)} x_l + b_k^{(1)} \right) + b_1^{(2)} \right)$
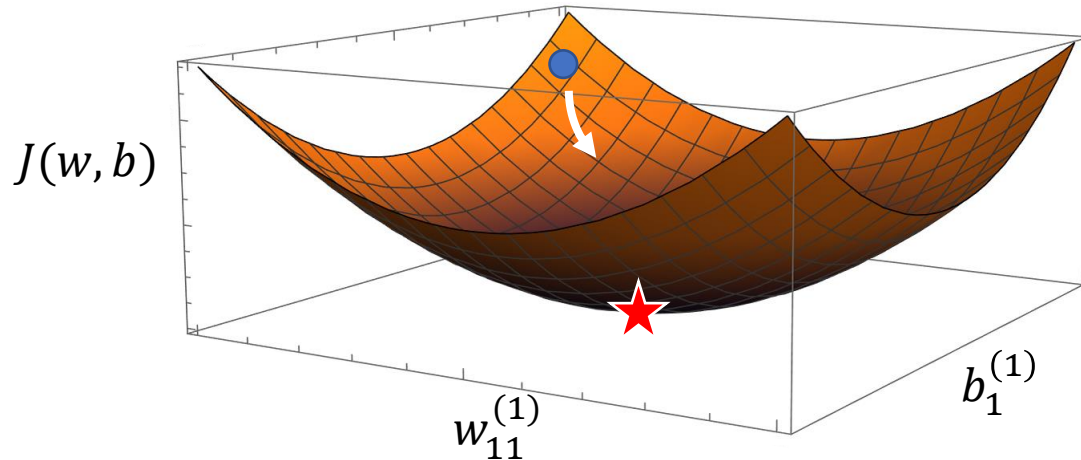
$$\longrightarrow \quad y = \boldsymbol{a}^{(2)} = \sigma \left( \boldsymbol{W}^{(2)} \sigma \left( \boldsymbol{W}^{(1)} \boldsymbol{x} + \boldsymbol{b}^{(1)} \right) + \boldsymbol{b}^{(2)} \right)$$

We know how to find w and b for a simple model
$y = \sigma(wx + b)$.

What about finding w and b in a neural network?
$y = a^{(2)} = \sigma(W^{(2)}\sigma(W^{(1)}x + b^{(1)}) + b^{(2)})$

# How to find w and b in a neural network?

- E. g., Model $y = a^{(2)} = \sigma(\boldsymbol{W}^{(2)}\sigma(\boldsymbol{W}^{(1)}\boldsymbol{x} + \boldsymbol{b}^{(1)}) + \boldsymbol{b}^{(2)})$

- Cost function $J(\boldsymbol{W}^{(1)}, \boldsymbol{W}^{(2)}, \boldsymbol{b}^{(1)}, \boldsymbol{b}^{(2)}) = \frac{1}{m}\sum_{i=1}^{m} L(y^i, y_d^i)$



- For the full data set (i=1…m), compute

$$-\boldsymbol{\nabla}J = -\left(\frac{\partial J}{\partial w_{jk}^{(l)}}, \dots, \frac{\partial J}{\partial b_j^{(l)}}\right)$$
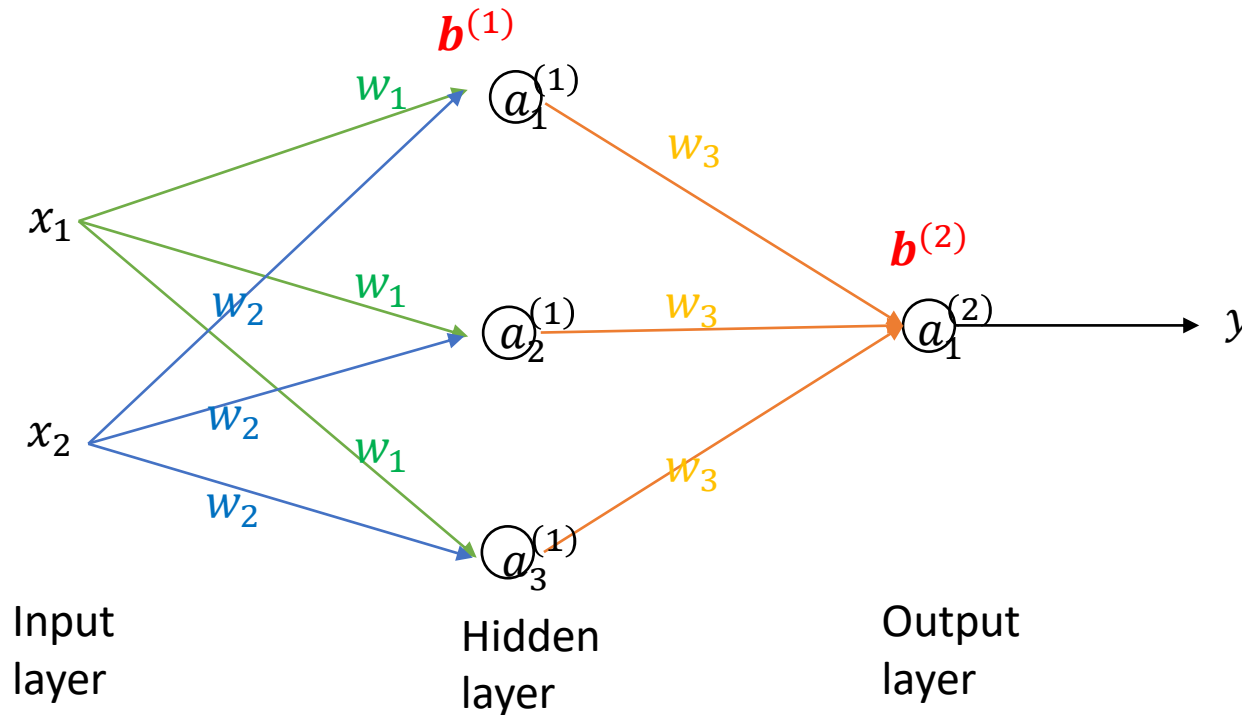
- $w_{jk}{}^{(l)}_{new} = w_{jk}{}^{(l)}_{old} - \alpha\, \frac{\partial J}{\partial w_{jk}^{(l)}},\ b_j{}^{(l)}_{new} = b_j{}^{(l)}_{old} - \alpha\, \frac{\partial J}{\partial b_j^{(l)}}$

  $\alpha$ is the learning rate

- Iterate until $\boldsymbol{\nabla}J = 0$

# NN weight Initialization



Input
layer

Hidden
layer

Output
layer

Q: What would happen if all initial weights are
chosen to be the same (e.g., all zeros)?

# NN weight Initialization

More specifically, if weights are symmetric:

$$W^{(1)} = \begin{bmatrix} w_1 & w_2 \\ w_1 & w_2 \\ w_1 & w_2 \end{bmatrix}, W^{(2)} = \begin{bmatrix} w_3 & w_3 & w_3 \end{bmatrix}$$

(e.g. all zero weights, all constant weights)

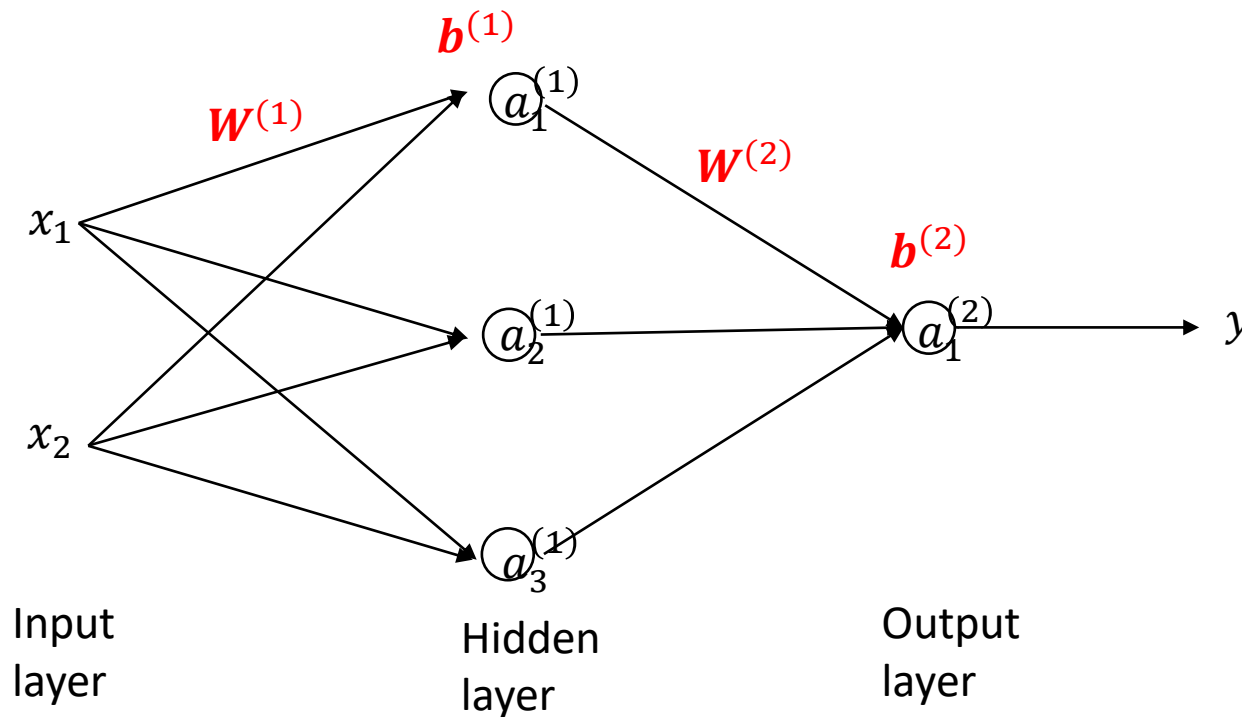Then $a_1^{(1)} = a_2^{(1)} = a_3^{(1)}$

and

$$\begin{bmatrix} w_{1new} & w_{2new} \\ w_{1new} & w_{2new} \\ w_{1new} & w_{2new} \end{bmatrix} = \begin{bmatrix} w_1 & w_2 \\ w_1 & w_2 \\ w_1 & w_2 \end{bmatrix} - \alpha \begin{bmatrix} \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_2} \\ \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_2} \\ \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_2} \end{bmatrix}$$

The NN never learns to have non-symmetric weights -> $a_1^{(1)} = a_2^{(1)} = a_3^{(1)}$ during training



Input layer

Hidden layer

Output layer

If the NN weights are **symmetric** → Gradient descent will update these hidden units in the same way → All hidden units in the same layer will be identical throughout training iterations → Equivalent to NN with just 1 hidden unit.

# NN weight Initialization

$b^{(1)}$

$W^{(1)}$

$W^{(2)}$

$b^{(2)}$

$x_1$

$a_1^{(1)}$

$a_2^{(1)}$

$a_1^{(2)}$

$y$

$x_2$

$a_3^{(1)}$

Input
layer

Hidden
layer

Output
layer

To make the different hidden units approximate different functions

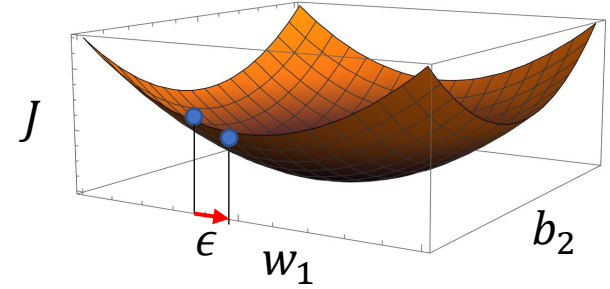$\rightarrow$ Initialized weights $w_{jk}^{(l)}$ to random values

What about biases?

$\rightarrow$ Biases can just be zeros

$$b^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \quad b^{(2)} = [0]$$

If the NN weights are **symmetric** $\rightarrow$ Gradient descent will update these hidden units in the same way $\rightarrow$ All hidden units in the same layer will be identical throughout training iterations $\rightarrow$ Equivalent to NN with just 1 hidden unit.

# How are $\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial b_1}$ ... calculated?



i. Slow: Calculate $J$ at two adjacent locations, evaluate their differences

$\frac{\partial J}{\partial w_j} = \frac{J(w_0 + \epsilon e_j) - J(w_0)}{\epsilon}$    For 1M weights, have to do 1M forward passes
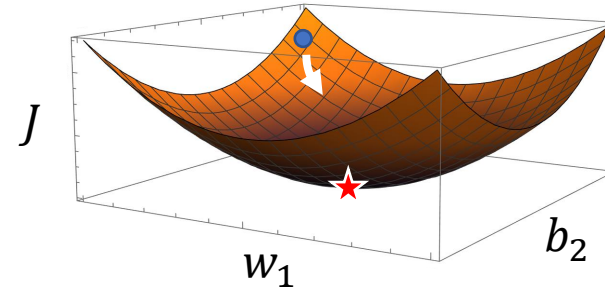
ii. Fast: Back propagation
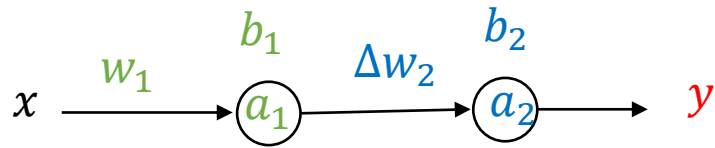Compute derivatives w.r.t. all weights analytically with chain rules

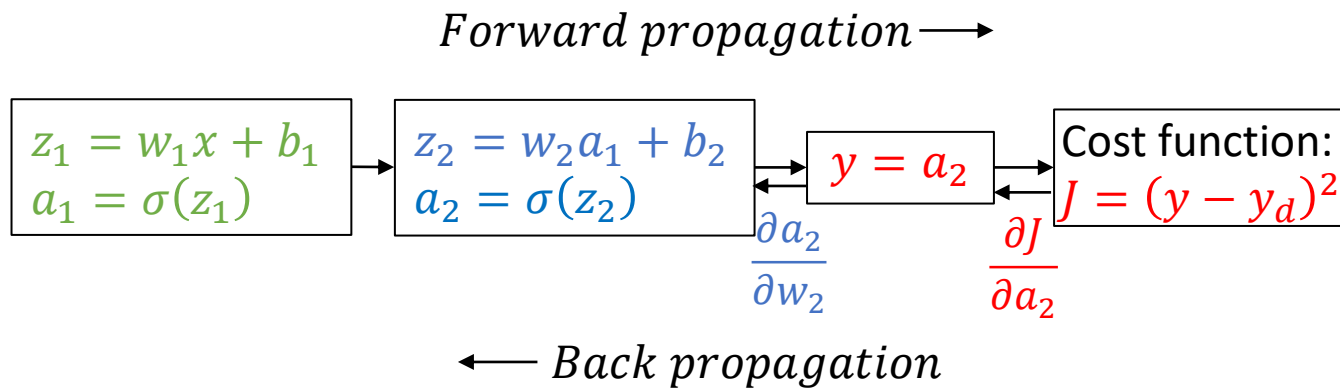# Forward & Back Propagation: 1. single neuron

Exercise:
Model: $y = \sigma(w_2\sigma(w_1 x + b_1) + b_2)$



$\nabla J$ vary with
$w_1, w_2, b_1, b_2$

$x \xrightarrow{\;w_1\;} \underset{a_1}{\bigcirc} \xrightarrow{\;\Delta w_2\;} \underset{a_2}{\bigcirc} \longrightarrow y$

with $b_1$ above $a_1$ and $b_2$ above $a_2$

How important is $w_2$ for
changing the cost function $J$?

*Forward propagation* $\longrightarrow$

| $z_1 = w_1 x + b_1$ <br> $a_1 = \sigma(z_1)$ | $\rightarrow$ | $z_2 = w_2 a_1 + b_2$ <br> $a_2 = \sigma(z_2)$ | $y = a_2$ | Cost function: <br> $J = (y - y_d)^2$ |

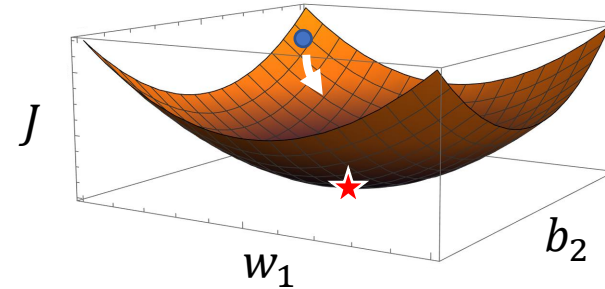$\dfrac{\partial a_2}{\partial w_2}$  $\dfrac{\partial J}{\partial a_2}$

$\longleftarrow$ *Back propagation*

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial a_2}\frac{\partial a_2}{\partial w_2}$$
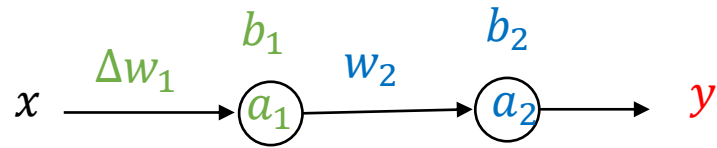$$= 2(y - y_d)\sigma'(z_2)a_1$$
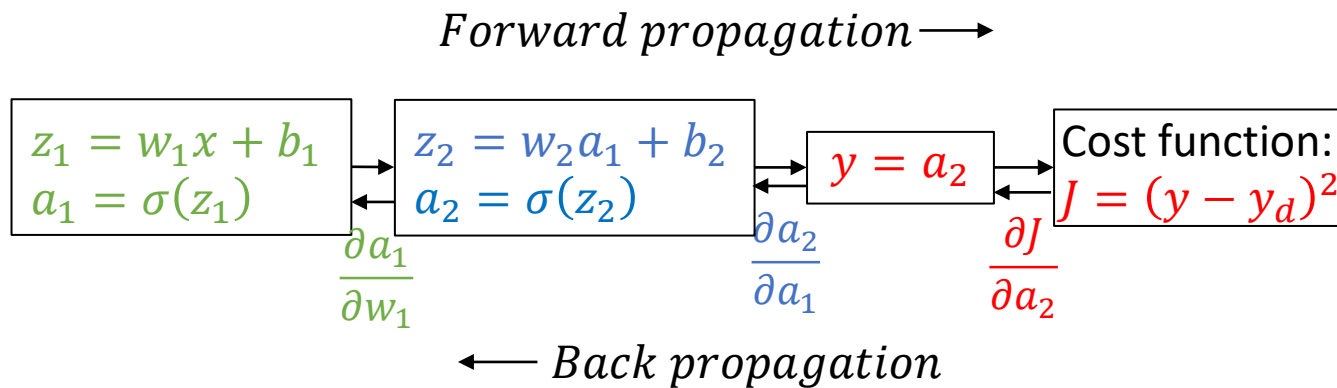
# Forward & Back Propagation: 1. single neuron

Exercise:
Model: $y = \sigma(w_2 \sigma(w_1 x + b_1) + b_2)$
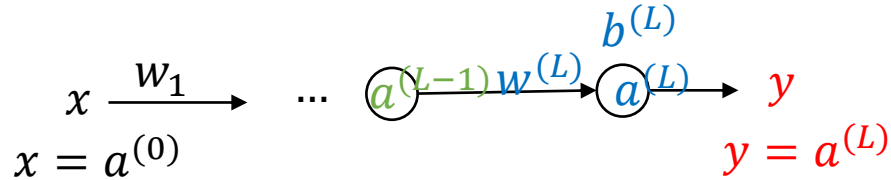
$\nabla J$ vary with $w_1, w_2, b_1, b_2$



$x \xrightarrow{\Delta w_1} a_1 \xrightarrow{w_2} a_2 \longrightarrow y$

$b_1$    $b_2$

How important is $w_1$ for changing the cost function $J$?

*Forward propagation* $\longrightarrow$

| $z_1 = w_1 x + b_1$ $a_1 = \sigma(z_1)$ | $z_2 = w_2 a_1 + b_2$ $a_2 = \sigma(z_2)$ | $y = a_2$ | Cost function: $J = (y - y_d)^2$ |

$\dfrac{\partial a_1}{\partial w_1}$     $\dfrac{\partial a_2}{\partial a_1}$     $\dfrac{\partial J}{\partial a_2}$

$\longleftarrow$ *Back propagation*

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial a_2} \frac{\partial a_2}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$
$$= 2(y - y_d)\sigma'(z_2)w_2\sigma'(z_1)x$$

# Forward & Back Propagation: 1. single neuron

$$x \xrightarrow{w_1} \quad \cdots \quad a^{(L-1)} \xrightarrow[\;]{\;w^{(L)}\;} a^{(L)} \xrightarrow{\;} y$$

$$b^{(L)}$$

$$x = a^{(0)}$$

$$y = a^{(L)}$$

## Forward propagation

Input: $a^{(0)} = x$

$\cdots$

$z^{(L-1)} = w^{(L-1)}a^{(L-2)} + b^{(L-1)}$

$a^{(L-1)} = \sigma(z^{(L-1)})$

$z^{(L)} = w^{(L)}a^{(L-1)} + b^{(L)}$

$a^{(L)} = \sigma(z^{(L)})$
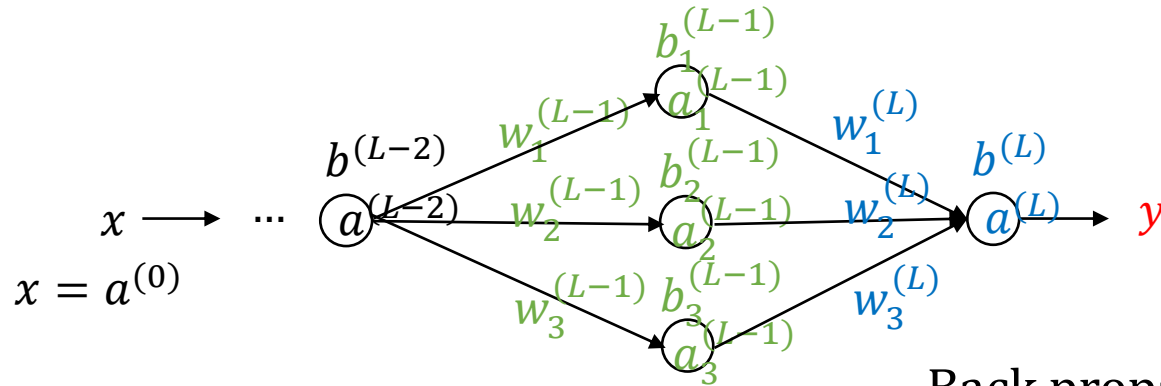
Output: $y = a^{(L)}$

Cost function $J = (y - y_d)^2$

## Back propagation (use chain rules)

$$\frac{\partial J}{\partial w^{(L)}} = \frac{\partial J}{\partial a^{(L)}}\frac{\partial a^{(L)}}{\partial w^{(L)}} = 2(a^{(L)} - y_d)\,\sigma'(z^{(L)})\,a^{(L-1)}$$

$$\frac{\partial J}{\partial b^{(L)}} = \frac{\partial J}{\partial a^{(L)}}\frac{\partial a^{(L)}}{\partial b^{(L)}} = 2(a^{(L)} - y_d)\,\sigma'(z^{(L)})$$

$$\frac{\partial J}{\partial w^{(L-1)}} = \frac{\partial J}{\partial a^{(L)}}\frac{\partial a^{(L)}}{\partial a^{(L-1)}}\frac{\partial a^{(L-1)}}{\partial w^{(L-1)}}$$

$$= 2(a^{(L)} - y_d)\,\sigma'(z^{(L)})w^{(L)}\sigma'(z^{(L-1)})a^{(L-2)}$$

$$\frac{\partial J}{\partial b^{(L-1)}} = \frac{\partial J}{\partial a^{(L)}}\frac{\partial a^{(L)}}{\partial a^{(L-1)}}\frac{\partial a^{(L-1)}}{\partial b^{(L-1)}}$$

$$= 2(a^{(L)} - y_d)\,\sigma'(z^{(L)})w^{(L)}\sigma'(z^{(L-1)})$$

# Forward & Back Propagation: 2. multiple hidden units



Back propagation (use chain rules)

Forward propagation

Input: $a^{(0)} = x$

...

$z_k^{(L-1)} = w_k^{(L-1)} a^{(L-2)} + b_k^{(L-1)}$

$a_k^{(L-1)} = \sigma(z_k^{(L-1)})$

$z^{(L)} = \sum_{k=1}^{3} w_k^{(L)} a_k^{(L-1)} + b^{(L)}$

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

Output: $y = a^{(L)}$

Cost function $\quad J = (y - y_d)^2$

$$\frac{\partial J}{\partial w_k^{(L)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial w_k^{(L)}} = 2(a^{(L)} - y_d)\, \sigma'(z^{(L)})\, a_k^{(L-1)}$$

$$\frac{\partial J}{\partial b^{(L)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial b^{(L)}} = 2(a^{(L)} - y_d)\, \sigma'(z^{(L)})$$
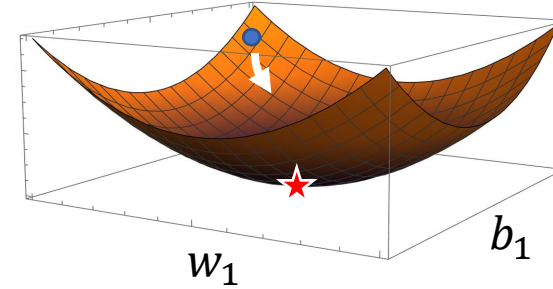
$$\frac{\partial J}{\partial w_k^{(L-1)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial a_k^{(L-1)}} \frac{\partial a_k^{(L-1)}}{\partial w_k^{(L-1)}}$$

$$= 2(a^{(L)} - y_d)\, \sigma'(z^{(L)}) w_k^{(L)} \sigma'(z_k^{(L-1)}) a^{(L-2)}$$

$$\frac{\partial J}{\partial b_k^{(L-1)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial a_k^{(L-1)}} \frac{\partial a_k^{(L-1)}}{\partial b_k^{(L-1)}}$$

$$= 2(a^{(L)} - y_d)\, \sigma'(z^{(L)}) w_k^{(L)} \sigma'(z_k^{(L-1)})$$

# Forward & Back Propagation: 2. multiple hidden units

Exercise:

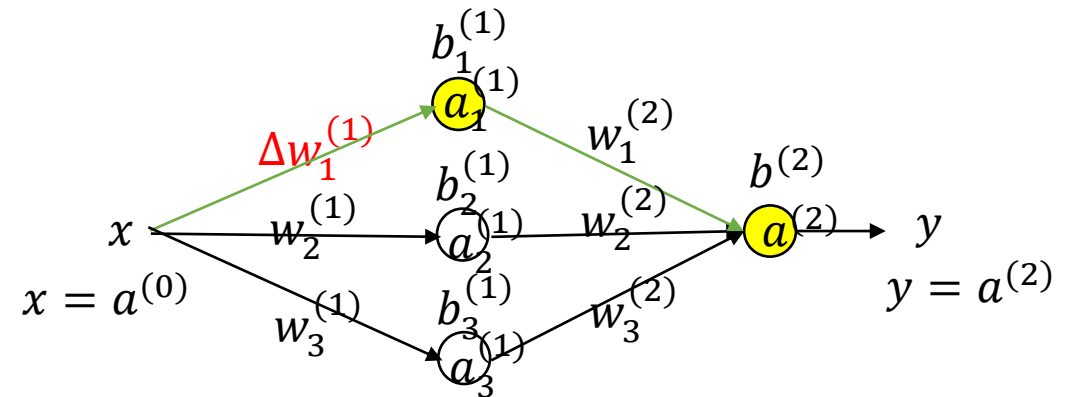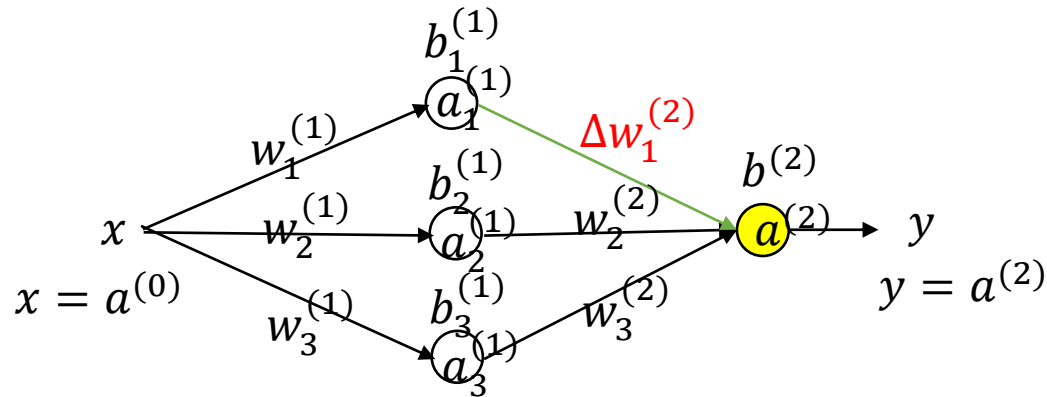Model: $y = \sigma\left(\sum_{k=1}^{3} w_k^{(2)}\, \sigma\left(w_k^{(1)} x + b_k^{(1)}\right) + b^{(2)}\right)$

$\nabla J$ vary with
$w_k^{(1)}, w_k^{(2)}, b_k^{(1)}, \mathrm{b}^{(2)}$

$J$

$b_1$

$w_1$

$$\frac{\partial J}{\partial w_1^{(2)}} = 2\left(a^{(2)} - y_d\right)\sigma'\left(z^{(2)}\right)a_1^{(1)}$$

$$= \mathrm{fn}(\mathrm{x}, \mathrm{y_d}, w_k^{(1)}, w_k^{(2)}, b_k^{(1)}, \mathrm{b}^{(2)})$$

$$\frac{\partial J}{\partial w_1^{(1)}} = 2\left(a^{(2)} - y_d\right)\sigma'\left(z^{(2)}\right)w_1^{(2)}\sigma'\left(z_1^{(1)}\right)x$$

$$= \mathrm{fn}(\mathrm{x}, \mathrm{y_d}, w_k^{(1)}, w_k^{(2)}, b_k^{(1)}, \mathrm{b}^{(2)})$$

# Forward & Back Propagation: 3. multiple outputs

Exercise:

Model: $y_j = \sigma\left(\sum_{k=1}^{3} w_{jk}^{(2)} \sigma\left(w_k^{(1)} x + b_k^{(1)}\right) + b_j^{(2)}\right)$
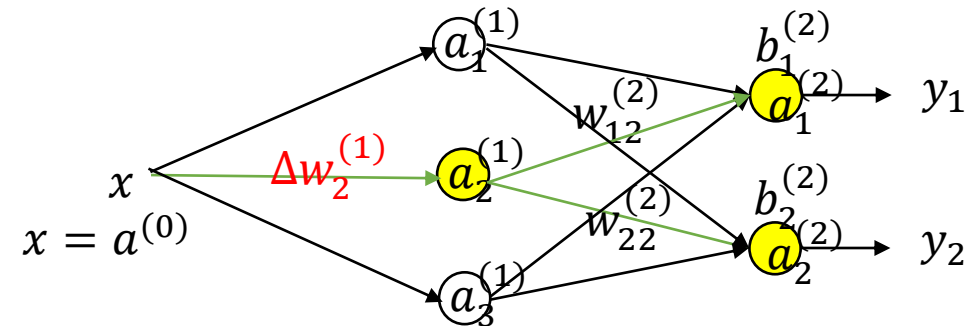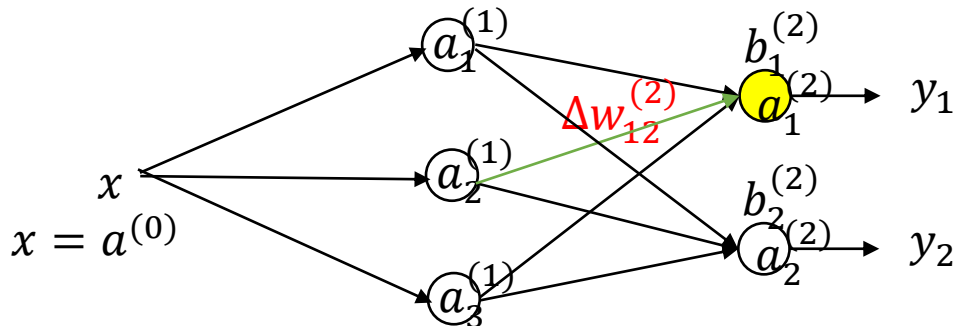
$J$



$\nabla J$ vary with
$w_k^{(1)}, w_{jk}^{(2)}, b_k^{(1)}, b_j^{(2)}$

$b_1$

$w_1$

$$\frac{\partial J}{\partial w_{12}^{(2)}} = 2\left(a_1^{(2)} - y_{d,1}\right)\sigma'\left(z_1^{(2)}\right)a_2^{(1)}$$
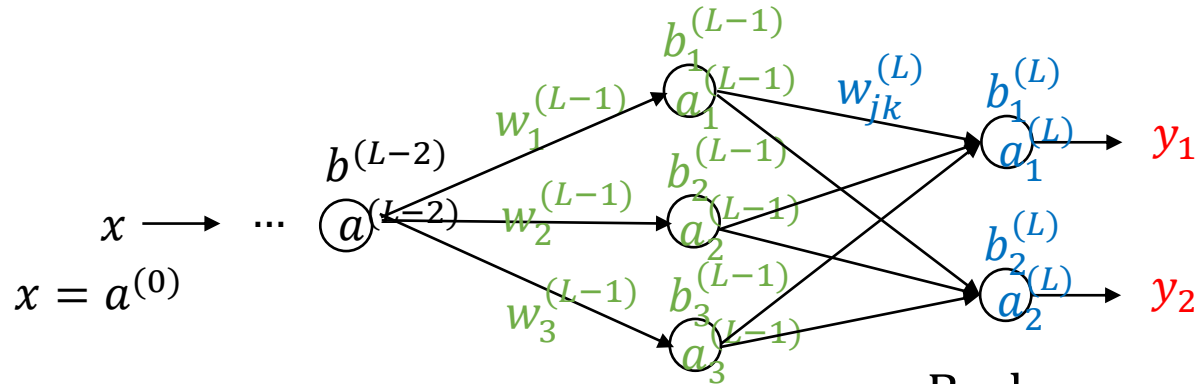
$$= \mathrm{fn}(x, y_{d,1}, w_k^{(1)}, w_{1k}^{(2)}, b_k^{(1)}, b_1^{(2)})$$

$$\frac{\partial J}{\partial w_2^{(1)}} = \sum_{j=1}^{2} 2(a_j^{(2)} - y_{d,j})\,\sigma'(z_j^{(2)}) w_{j2}^{(2)}\,\sigma'\left(z_2^{(1)}\right)x$$

$$= \mathrm{fn}(x, y_{d,j}, w_k^{(1)}, w_{jk}^{(2)}, b_k^{(1)}, b_j^{(2)})$$

## Forward propagation

Input: $a^{(0)} = x$

$\dots$

$z_k^{(L-1)} = w_k^{(L-1)} a^{(L-2)} + b_k^{(L-1)}$

$a_k^{(L-1)} = \sigma(z_k^{(L-1)})$

$z_j^{(L)} = \sum_{k=1}^3 w_{jk}^{(L)} a_k^{(L-1)} + b_j^{(L)}$

$a_j^{(L)} = \sigma(z_j^{(L)})$

Output: $y_j = a_j^{(L)}$

Cost function $J = \sum_{j=1}^2 (y_j - y_{d,j})^2$

## Back propagation (use chain rules)

$\dfrac{\partial J}{\partial w_{jk}^{(L)}} = \dfrac{\partial J}{\partial a_j^{(L)}} \dfrac{\partial a_j^{(L)}}{\partial w_{jk}^{(L)}} = 2(a_j^{(L)} - y_{d,j})\, \sigma'(z_j^{(L)}) a_k^{(L-1)}$

$\dfrac{\partial J}{\partial b_j^{(L)}} = \dfrac{\partial J}{\partial a_j^{(L)}} \dfrac{\partial a_j^{(L)}}{\partial b_j^{(L)}} = 2(a_j^{(L)} - y_{d,j})\, \sigma'(z_j^{(L)})$

$\dfrac{\partial J}{\partial w_k^{(L-1)}} = \sum_{j=1}^2 \dfrac{\partial J}{\partial a_j^{(L)}} \dfrac{\partial a_j^{(L)}}{\partial a_k^{(L-1)}} \dfrac{\partial a_k^{(L-1)}}{\partial w_k^{(L-1)}}$

$\qquad = \sum_{j=1}^2 2(a_j^{(L)} - y_{d,j})\, \sigma'(z_j^{(L)}) w_{jk}^{(L)}\ \sigma'(z_k^{(L-1)}) a^{(L-2)}$

$\dfrac{\partial J}{\partial b_k^{(L-1)}} = \sum_{j=1}^2 \dfrac{\partial J}{\partial a_j^{(L)}} \dfrac{\partial a_j^{(L)}}{\partial a_k^{(L-1)}} \dfrac{\partial a_k^{(L-1)}}{\partial b_k^{(L-1)}}$

$\qquad = \sum_{j=1}^2 2(a_j^{(L)} - y_{d,j})\, \sigma'(z_j^{(L)}) w_{jk}^{(L)}\ \sigma'(z_k^{(L-1)})$

# When training a neural network

Repeat these steps:

    1. Forward propagate an input

    2. Compute the cost function

    3. Compute the gradients of the cost with respect to parameters using backpropagation

    4. Update each parameter using the gradients, according to the optimization algorithm

# Summary

- Universal function approximator
- Gradient descent: a method to find weights and biases that minimize J
- Calculate $\nabla J$
  - One example
  - A full dataset
- Weight initialization
- Forward and backpropagation (a clever way to get gradients of J)

# Neural Networks for Pattern Recognition

## Christopher M. Bishop

**Chapter 4
The Multi-layer Perceptron**