# Train a Neural Network in TensorFlow



$\vec{x} \rightarrow$

$\vec{a}^{[1]}$  $\vec{a}^{[2]}$  $\vec{a}^{[3]}$

1 unit

15 units
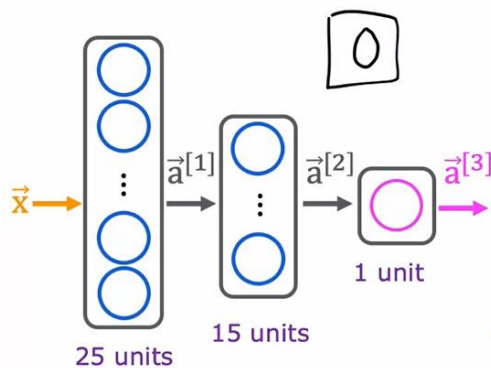
25 units

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
    model = Sequential([
        Dense(units=25, activation='sigmoid'),
        Dense(units=15, activation='sigmoid'),
        Dense(units=1, activation='sigmoid'),
        ])
from tensorflow.keras.losses import
BinaryCrossentropy
    model.compile(loss=BinaryCrossentropy())    ②

    model.fit(X,Y,epochs=100)    ③
```

①

epochs: number of steps
in gradient descent

Given set of $(x,y)$ examples

How to build and train this in code? gradient descent you
may want to run.

DeepLearning.AI   Stanford ONLINE                    Andrew Ng

❚❚ ◀) 2:31 / 3:35

---

# Model Training Steps    Tensor Flow

① specify how to
compute output
given input x and
parameters w,b
(define model)

$f_{\vec{w},b}(\vec{x}) = ?$

② specify loss and cost

$L\left(f_{\vec{w},b}(\vec{x}), y\right)$  1 example

$J(\vec{w}, b) = \dfrac{1}{m} \sum_{i=1}^{m} L\left(f_{\vec{w},b}(\vec{x}^{(i)}), y^{(i)}\right)$

③ Train on data to
minimize $J(\vec{w}, b)$

### logistic regression

```
z = np.dot(w,x)+ b



f_x = 1/(1+np.exp(-z))
```

logistic loss

```
loss = -y * np.log(f_x)
      -(1-y) * np.log(1-f_x)




w = w - alpha * dj_dw
b = b - alpha * dj_db
```

### neural network

```
model = Sequential([
    Dense(...)
    Dense(...)
    Dense(...)
    ])
```

binary cross entropy

```
model.compile(
loss=BinaryCrossentropy())


model.fit(X,y,epochs=100)
```
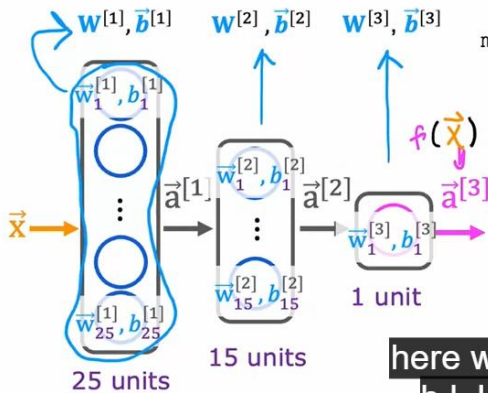
Let's look in greater detail in

DeepLearning.AI   Stanford ONLINE                    Andrew Ng

▶ ◀) 4:49 / 13:27

# 1. Create the model

define the model

$f(\vec{x}) = ?$

$w^{[1]}, \vec{b}^{[1]}$    $w^{[2]}, \vec{b}^{[2]}$    $w^{[3]}, \vec{b}^{[3]}$

$\vec{w}_1^{[1]}, b_1^{[1]}$

$\vec{a}^{[1]}$    $\vec{w}_1^{[2]}, b_1^{[2]}$   $\vec{a}^{[2]}$

$\vec{x}$

$f(\vec{x})$

$\vec{w}_1^{[3]}, b_1^{[3]}$   $\vec{a}^{[3]}$

$\vec{w}_{25}^{[1]}, b_{25}^{[1]}$    $\vec{w}_{15}^{[2]}, b_{15}^{[2]}$   1 unit

25 units      15 units

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense

model = Sequential([
    Dense(units=25, activation='sigmoid'),
    Dense(units=15, activation='sigmoid'),
    Dense(units=1, activation='sigmoid'),
])
```

here we have written w I and b I. Let's go on to step 2.

DeepLearning.AI   Stanford | ONLINE     Andrew Ng

5:45 / 13:27

---

# 2. Loss and cost functions

handwritten digit classification problem    binary classification

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^{m} L\left(f(\vec{x}^{(i)}), y^{(i)}\right)$$

$L(f(\vec{x}), y) = -y\log(f(\vec{x})) - (1-y)\log(1 - f(\vec{x}))$

$\mathbf{W}^{[1]}, \mathbf{W}^{[2]}, \mathbf{W}^{[3]}$    $\vec{b}^{[1]}, \vec{b}^{[2]}, \vec{b}^{[3]}$    $f_{W,B}(\vec{x})$

Compare prediction vs. target

logistic loss
also known as binary cross entropy

```
model.compile(loss= BinaryCrossentropy())
```

regression
(predicting numbers and not categories)

mean squared error

```
model.compile(loss= MeanSquaredError())
```

```
from tensorflow.keras.losses import
    BinaryCrossentropy
```
K Keras

```
from tensorflow.keras.losses import
    MeanSquaredError
```

on all the parameters in

DeepLearning.AI   Stanford | ONLINE     Andrew Ng

9:34 / 13:27

# 3. Gradient descent



$J(w)$

minimum

all $\vec{w}$ and $b$

repeat {

$$w_j^{[l]} = w_j^{[l]} - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b_j^{[l]} = b_j^{[l]} - \alpha \frac{\partial}{\partial bj} J(\vec{w}, b)$$

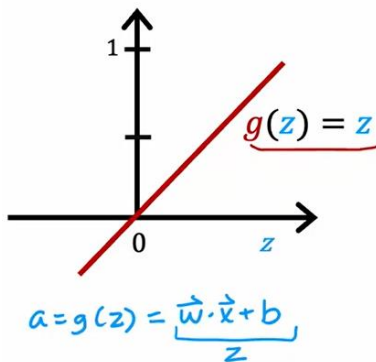} Compute derivatives for gradient descent using "back propagation"

```
model.fit(X,y,epochs=100)
```

In fact, what you see later is that TensorFlow can use
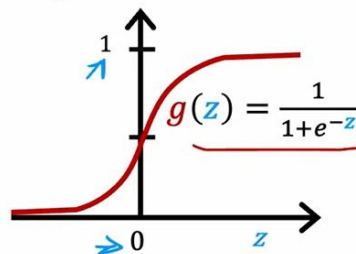
# Examples of Activation Functions

"No activation function"

Linear activation function

$$a_2^{[1]} = g(\vec{w}_2^{[1]} \cdot \vec{x} + b_2^{[1]})$$

Sigmoid

Later: softmax activation

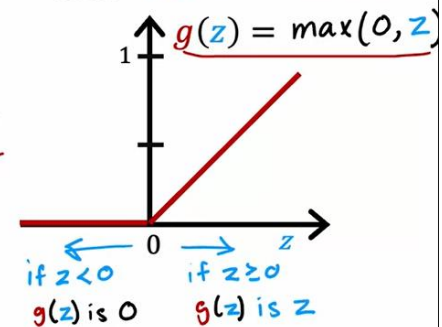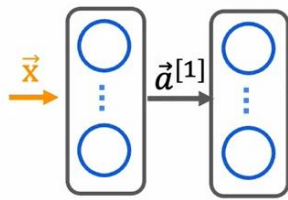ReLU  Rectified Linear Unit

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

$g(z) = z$

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

$$a = g(z) = \vec{w} \cdot \vec{x} + b$$
z

$0 < g(z) < 1$

if z<0    if z≥0
g(z) is 0   g(z) is z

a rich variety of powerful neural networks.

5:10 / 5:29

# Output Layer



$\vec{x}$ → $\vec{a}^{[1]}$ → $\vec{a}^{[2]}$ → $\vec{a}^{[3]} = f(\vec{x})$

$$f(\vec{x}) = a_1^{[3]} = g(z_1^{[3]})$$

Choosing $g(z)$ for output layer?

**Binary classification**
Sigmoid
y = 0/1

**Regression**
Linear activation function
y = +/-

**Regression**
ReLU
Y = 0 or +

as well for the output
layer of a neural network.

# Choosing Activation Summary



$\vec{x}$ → $\vec{a}^{[1]}$ → $\vec{a}^{[2]}$ → $\vec{a}^{[3]} = f(\vec{x})$

ReLU hidden layers

binary classification
activation='sigmoid'

regression  y negative/positive
activation='linear'

regression  y ≥ 0
activation='relu'

```
from tf.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),     layer1
    Dense(units=15, activation='relu'),     layer2
    Dense(units=1,  activation='sigmoid')   layer3
])
```

or 'linear'
or 'relu'

that shows the syntax for it.

# Example



$g(z) = z$

$$\vec{a}^{[4]} = \vec{w}_1^{[4]} \cdot \vec{a}^{[3]} + b_1^{[4]}$$

all linear (including output)
↳ equivalent to linear regression

$$\vec{a}^{[4]} = \frac{1}{1+e^{-(\vec{w}_1^{[4]} \cdot \vec{a}^{[3]} + b_1^{[4]})}}$$

output activation is sigmoid
(hidden layers still linear)
↳ equivalent to logistic regression

Don't use linear activations in hidden layers

in the hidden layers
of the neural network.

⏸ 🔊  4:33 / 5:30

---

Logistic regression
(2 possible output values)
$z = \vec{w} \cdot \vec{x} + b$

0.71

✗ $a_1 = g(z) = \frac{1}{1+e^{-z}} = P(y = 1|\vec{x})$

○ $a_2 = 1 - a_1 = P(y = 0 | \vec{x})$

0.29

Softmax regression
(N possible outputs) $y = 1, 2, 3, ..., N$

$z_j = \vec{w}_j \cdot \vec{x} + b_j$  $j = 1, ..., N$

parameters $w_1, w_2, ..., w_N$
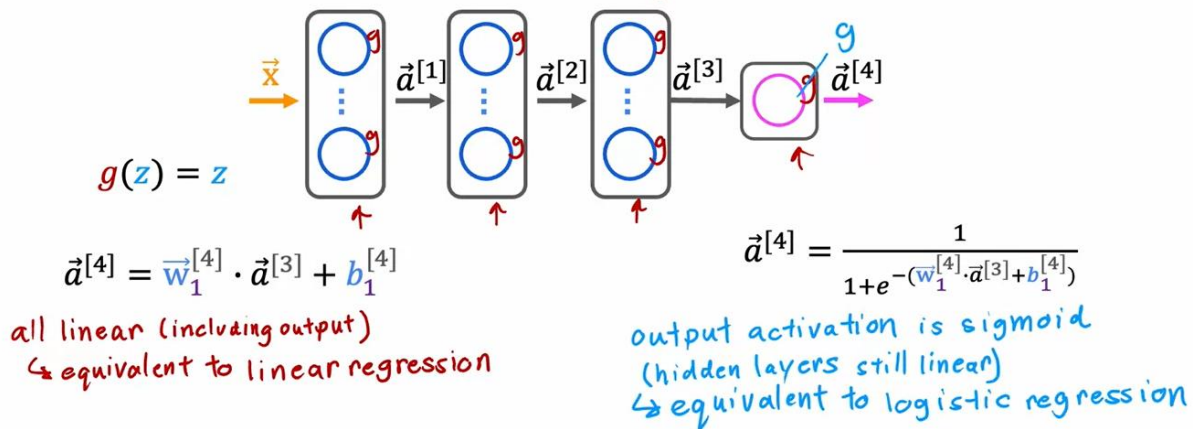$b_1, b_2, ..., b_N$

$a_j = \frac{e^{z_j}}{\sum_{k=1}^{N} e^{z_k}} = P(y = j|\vec{x})$

note: $a_1 + a_2 + ... + a_N = 1$

Softmax regression (4 possible outputs) $y = 1, 2, 3, 4$

✗ $z_1 = \vec{w}_1 \cdot \vec{x} + b_1$    $a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$

✗ ○ □ △

$= P(y = 1|\vec{x})$ 0.30

○ $z_2 = \vec{w}_2 \cdot \vec{x} + b_2$    $a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$

$= P(y = 2|\vec{x})$ 0.20

□ $z_3 = \vec{w}_3 \cdot \vec{x} + b_3$    $a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$

$= P(y = 3|\vec{x})$ 0.15

△ $z_4 = \vec{w}_4 \cdot \vec{x} + b_4$    $a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$

$= P(y = 4|\vec{x})$ 0.35

thing as logistic regression.

# Cost

## Logistic regression

$$z = \vec{w} \cdot \vec{x} + b$$

$$a_1 = g(z) = \frac{1}{1 + e^{-z}} \quad = P(y = 1|\vec{x})$$

$$a_2 = 1 - a_1 \quad = P(y = 0|\vec{x})$$

$a_2$

$$loss = -y \log a_1 - (1 - y) \log(1 - a_1)$$

if y = 1     if y = 0

$$J(\vec{w}, b) = \text{average loss}$$

## Softmax regression

$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \cdots + e^{z_N}} \quad = P(y = 1|\vec{x})$$

$$\vdots$$

$$a_N = \frac{e^{z_N}}{e^{z_1} + e^{z_2} + \cdots + e^{z_N}} \quad = P(y = N|\vec{x})$$

*Crossentropy loss*

$$loss(a_1, \ldots, a_N, y) = \begin{cases} -\log a_1 & \text{if } y = 1 \\ -\log a_2 & \text{if } y = 2 \\ \quad \vdots \\ -\log a_N & \text{if } y = N \end{cases}$$

$L$

$loss = -\log a_j \text{ if } y = j$

$a_j \downarrow L \uparrow$

For example, if y was equal to 2,

0.5        1        $a_j$

# Neural Network with Softmax output

$$z_1^{[3]} = \vec{w}_1^{[3]} \cdot \vec{a}^{[2]} + b_1^{[3]} \qquad a_1^{[3]} = \frac{e^{z_1^{[3]}}}{e^{z_1^{[3]}} + \cdots + e^{z_{10}^{[3]}}}$$

$$= P(y = 1|\vec{x})$$

$$\vdots$$

$$z_{10}^{[3]} = \vec{w}_{10}^{[3]} \cdot \vec{a}^{[2]} + b_{10}^{[3]} \qquad a_{10}^{[3]} = \frac{e^{z_{10}^{[3]}}}{e^{z_1^{[3]}} + \cdots + e^{z_{10}^{[3]}}}$$

$$= P(y = 10|\vec{x})$$

$\vec{x}$    $\vec{a}^{[1]}$    $\vec{a}^{[2]}$    $\vec{a}^{[3]}$

relu    softmax

25 units    15 units    10 units

10 classes

logistic regression

$$a_1^{[3]} = g\left(z_1^{[3]}\right) \qquad a_2^{[3]} = g\left(z_2^{[3]}\right)$$

softmax

$$\vec{a}^{[3]} = \left(a_1^{[3]}, \ldots a_{10}^{[3]}\right) = g\left(z_1^{[3]}, \ldots, z_{10}^{[3]}\right)$$

# MNIST with softmax

① specify the model

$f_{\vec{w},b}(\vec{x}) = ?$

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
from tensorflow.keras.losses import
    SparseCategoricalCrossentropy
model.compile(loss= SparseCategoricalCrossentropy() )
model.fit(X,Y,epochs=100)
```

② specify loss and cost

$L(f_{\vec{w},b}(\vec{x}), y)$

③ Train on data to minimize $J(\vec{w}, b)$

you can train a neural network on a multi class classification problem.

DeepLearning.AI    Stanford | ONLINE                                        Andrew Ng

6:21 / 7:24

---

# Numerical Roundoff Errors

More numerically accurate implementation of logistic loss:

$1 + \dfrac{1}{10,000}$     $1 - \dfrac{1}{10,000}$

Logistic regression:

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

Original loss

$loss = -y \log(\textcircled{a}) - (1-y)\log(1 - \textcircled{a})$

More accurate loss (in code)

$loss = -y \log\left(\dfrac{1}{1+e^{-z}}\right) - (1-y)\log(1 - \dfrac{1}{1+e^{-z}})$

logit: z

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),  'linear'
    Dense(units=1, activation='sigmoid')
])
model.compile(loss=BinaryCrossEntropy() )
model.compile(loss=BinaryCrossEntropy(from_logits=True) )
```

worse when it comes to softmax.

DeepLearning.AI    Stanford | ONLINE                                        Andrew Ng

5:14 / 9:11

# More numerically accurate implementation of softmax

**Softmax regression**

$$(a_1, \ldots, a_{10}) = g(z_1, \ldots, z_{10})$$

$$\text{Loss} = L(\vec{a}, y) = \begin{cases} -\log a_1 & \text{if } y = 1 \\ \vdots \\ -\log a_{10} & \text{if } y = 10 \end{cases}$$

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
    ])
```
'linear'

```
model.compile(loss=SparseCategoricalCrossEntropy() )
```

**More Accurate**

$$L(\vec{a}, y) = \begin{cases} -\log \dfrac{e^{z_1}}{e^{z_1} + \cdots + e^{z_{10}}} & \text{if } y = 1 \\ \vdots \\ -\log \dfrac{e^{z_{10}}}{e^{z_1} + \cdots + e^{z_{10}}} & \text{if } y = 10 \end{cases}$$

```
model.compile(loss=SparseCategoricalCrossEntropy(from_logits=True))
```

and this whole computation of

# MNIST (more numerically accurate)

**model**
```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='linear') ])
```

**loss**
```
from tensorflow.keras.losses import
    SparseCategoricalCrossentropy
model.compile(...,loss=SparseCategoricalCrossentropy(from_logits=True) )
```

**fit**
```
model.fit(X,Y,epochs=100)
```

**predict**
```
logits = model(X)
```
not $a_1 \ldots a_{10}$
is $z_1 \ldots z_{10}$
```
f_x = tf.nn.softma
```
It is instead of putting $z_1$ through $z_{10}$.

8:27 / 9:11

# logistic regression
# (more numerically accurate)

model
```
model = Sequential([
    Dense(units=25, activation='sigmoid'),
    Dense(units=15, activation='sigmoid'),
    Dense(units=1, activation='linear')
    ])
from tensorflow.keras.losses import
    BinaryCrossentropy
```

loss
```
model.compile(..., BinaryCrossentropy(from_logits=True))

model.fit(X,Y,epochs=100)
```

fit
```
logit = model(X)
```

predict
```
f_x = tf.nn.sigmoid(logit)
```

to actually get the probability.

8:47 / 9:11

# Multi-label Classification



Is there a car?          yes          $y = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$      no          $y = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$      yes          $y = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$
Is there a bus?          no                          no                          yes
Is there a pedestrian    yes                         yes                         no

The second one to
detect buses and

2:00 / 4:19

# Multi-label Classification



Alternatively, train one neural network with three outputs



$$\vec{a}^{[3]} = \begin{bmatrix} a_1^{[3]} \\ a_2^{[3]} \\ a_3^{[3]} \end{bmatrix} \begin{array}{l} \text{car} \\ \text{bus} \\ \text{pedestrian} \end{array}$$

sigmoid activations

and no pedestrians in the image.

3:11 / 4:19

---

# Gradient Descent

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

learning rate



minimum

$J(\vec{w}, b)$

$w_2$

$w_1$    start

"Adam" algorithm

$J(\vec{w}, b)$

$w_2$

$w_1$

Go faster – increase $\alpha$          Go slower – decrease $\alpha$

a smaller learning rate Alpha.

2:41 / 6:25

# MNIST Adam

### model

```
model = Sequential([
        tf.keras.layers.Dense(units=25, activation='sigmoid'),
        tf.keras.layers.Dense(units=15, activation='sigmoid'),
        tf.keras.layers.Dense(units=10, activation='linear')
])
```

### compile

$$\alpha = 10^{-3} = 0.001$$

```
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=1e-3),
   loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True))
```

### fit

```
model.fit(X,Y,epochs=100)
```

see what gives you the
fastest learning performance.

# Convolutional Layer



Each neuron only looks at
part of the previous layer's outputs.

Why?
- Faster computation
- Need less training data
  (less prone to
  overfitting)

of how to get convolutional layers
to work and popularized their use.

# Convolutional Neural Network

EKG

$90°$

height

$x_1\ x_2\ x_3\ \cdots\cdots\ x_{100}$

2D $\uparrow\!\!\to$
1D $\longleftrightarrow$

convolutional layer

presence or absence of heart disease

$x_1$
$x_2$
$x_3$
$\vdots$
$x_{100}$

$x_1 - x_{20}$
$x_{11} - x_{30}$
$x_{21} - x_{40}$
$\vdots$
$x_{81} - x_{100}$

$\vec{x}$

$\vec{a}^{[1]}$

$a_1^{[1]} - a_5^{[1]}$
$a_3^{[1]} - a_7^{[1]}$
$a_5^{[1]} - a_9^{[1]}$

9 units

$\vec{a}^{[2]}$

$\vec{a}^{[3]}$

sigmoid

3 units

how many neurons should layer have.

7:24 / 8:55

---

# More Derivative Examples

$w = 3$ $\qquad J(w) = w^2 = 9$ $\qquad w \uparrow 0.001$ $\qquad J(w) = J(3.001) = 9.006001$ $\qquad \frac{\partial}{\partial w}J(w) = 6$

$J(w) \uparrow 6 \times 0.001$

$w = 2$ $\qquad J(w) = w^2 = 4$ $\qquad w \uparrow 0.001$ $\qquad J(w) = J(2.001) = 4.004001$ $\qquad \frac{\partial}{\partial w}J(w) = 4$

$J(w) = w^2$

$J(w) \uparrow 4 \times 0.001$

$\downarrow 0.006$

$w = -3$ $\qquad J(w) = w^2 = 9$ $\qquad w \uparrow 0.001$ $\qquad J(w) = J(-2.999) = 8.994001$ $\qquad \frac{\partial}{\partial w}J(w) = -6$

$J(w) \downarrow 6 \times 0.001$

$-3 + 0.001$

$J(w) \uparrow -6 \times 0.001$

$J(w)$

$-6$ $\quad -9$ $\quad 6$ height

$4$ $\quad 4$ width

$-3$ $\quad 2$ $\quad 3$ $\quad w$

| Calculus | $w$ | $\partial J(w)/\partial w$ |
|---|---|---|
| $\frac{\partial}{\partial w}J(w) = 2w$ | 3 | $2 \times 3 = 6$ |
| | 2 | $2 \times 2 = 4$ |
| | -3 | $2 \times -3 = -6$ |

2 turns out to give you the derivative.

12:48 / 22:56

File   Edit   View   Insert   Cell   Kernel   Widgets   Help                    Trusted   ✏  Python 3 ○

⊞ Run   Code ▾   Validate

```
In [2]: J, w = sympy.symbols('J,w')
```

```
In [3]: J = w**2
        J
```

Out[3]: $w^2$

```
In [5]: dJ_dw = sympy.diff(J,w)
        dJ_dw
```

Out[5]: $2w$

```
In [6]: dJ_dw.subs([(w,2)])
```

Out[6]: 4

```
In [ ]:
```

Let's look at some other examples.

15:27 / 22:56

---

# Even More Derivative Examples

$w = 2$

$J(w) = w^2 = 4$   $\frac{\partial}{\partial w} J(w) = 2w = 4$   $w \uparrow 0.001 \; \varepsilon$   $J(w) = 4.004001$ $J(w) \uparrow 4 \times \varepsilon$

$J(w) = w^3 = 8$   $\frac{\partial}{\partial w} J(w) = 3w^2 = 12$   $w \uparrow \varepsilon$   $J(w) = 8.012006$ $J(w) \uparrow 12 \times \varepsilon$

$J(w) = w = 2$   $\frac{\partial}{\partial w} J(w) = 1$   $w \uparrow \varepsilon$   $J(w) = 2.001$ $J(w) \uparrow 1 \times \varepsilon$

$J(w) = \frac{1}{w} = \frac{1}{2} = 0.5$   $\frac{\partial}{\partial w} J(w) = \frac{-1}{w^2} = \frac{-1}{4}$   $w \uparrow \varepsilon$ $w = \frac{1}{2.001}$   $-0.25 \times 0.001$ $0.5 - 0.00025$ $J(w) = 0.49975$ $J(w) \uparrow -\frac{1}{4} \times \varepsilon$

$\frac{\partial}{\partial w} J(w)$   $w \uparrow \varepsilon$   $J(w) \uparrow k \times \varepsilon$

DeepLearning.AI   Stanford ONLINE                 Andrew Ng

19:53 / 22:56

# Computing the Derivatives

Forward prop →
← Back prop

$w = 2 \quad b = 8 \quad x = -2 \quad y = 2 \quad a = wx + b \quad J = \frac{1}{2}(a - y)^2$

w ②  $c = wx$  ④  $a = c + b$  ④  $d = a - y$  ②  $J = \frac{1}{2}d^2$  ②  → J
(-4)        (2)        (2)        (2)
b (2) (8)

$\frac{\partial J}{\partial w} = -4 \qquad \frac{\partial J}{\partial c} = 2 \qquad \frac{\partial J}{\partial b} = 2 \qquad \frac{\partial J}{\partial a} = 2 \qquad \frac{\partial J}{\partial d} = 2$

w=2.001  C=-4.002
w ↑ 0.001   c ↓ 2×0.001
          c ↑ −2×0.001
J ↑ −4×0.001

$\frac{\partial J}{\partial w} = \frac{\partial c}{\partial w} \times \frac{\partial J}{\partial c}$
−2    2

c ↑ 0.001  a ↑ 0.001  J ↑ 0.002

$\frac{\partial J}{\partial c} = \frac{\partial a}{\partial c} \times \frac{\partial J}{\partial a}$
1    2

b ↑ 0.001  a ↑ 0.001  J ↑ 0.002

$\frac{\partial J}{\partial b} = \frac{\partial a}{\partial b} \times \frac{\partial J}{\partial a}$
1    2

a ↑ 0.001   d ↑ 0.001
a = 4.001  d = 2.001

J ↑ 0.002

$\frac{\partial J}{\partial a} = \frac{\partial d}{\partial a} \times \frac{\partial J}{\partial d}$
1    2

d ↑ 0.001  J ↑ 0.002
         ε         2ε

chain rule
(calculus)

DeepLearning.AI    Stanford ONLINE               Andrew Ng

13:26 / 19:19

---

# Computing the Derivatives

$w = 2 \quad b = 8 \quad x = -2 \quad y = 2 \quad a = wx + b \quad J = \frac{1}{2}(a - y)^2$

w  2  $c = wx$  -4  $a = c + b$  4  $d = a - y$  2  $J = \frac{1}{2}d^2$  2  → J
(-4)        2              2              2
8  2
b

$\frac{\partial J}{\partial w} = -4 \qquad \frac{\partial J}{\partial c} = 2 \qquad \frac{\partial J}{\partial b} = 2 \qquad \frac{\partial J}{\partial a} = 2 \qquad \frac{\partial J}{\partial d} = 2$

$J = \frac{1}{2}\big((wx + b) - y\big)^2 = \frac{1}{2}\big((2 \times -2 + 8) - 2\big)^2 = 2$

w ↑ 0.001   $J = \frac{1}{2}\big((2.001 \times -2 + 8) - 2\big)^2 = 1.996002$

J ↓ 4 × 0.001
J ↑ −4 × 0.001

$\frac{\partial J}{\partial w} = -4$

predicted by this
derivative calculation.

DeepLearning.AI    Stanford ONLINE               Andrew Ng

# Backprop is an efficient way to compute derivatives

$$w \xrightarrow{\;2\;} \boxed{c = wx} \xrightarrow{\;-4\;} \boxed{a = c + b} \xrightarrow{\;4\;} \boxed{d = a - y} \xrightarrow{\;2\;} \boxed{J = {}^1/_2\, d^2} \xrightarrow{\;2\;} J$$

$$\frac{\partial J}{\partial w} \qquad \frac{\partial J}{\partial c} \quad \overset{8}{\underset{b}{\phantom{x}}} \quad \frac{\partial J}{\partial b} \qquad \frac{\partial J}{\partial a} \qquad \frac{\partial J}{\partial d}$$

Compute $\frac{\partial J}{\partial a}$ once and use it to compute both $\frac{\partial J}{\partial w}$ and $\frac{\partial J}{\partial b}$.

If N nodes and P parameters, compute derivatives in roughly $N + P$ steps rather than $N \times P$ steps.
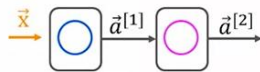
| N | P | N + P | N × P |
|---|---|---|---|
| 10,000 | 100,000 | $1.1 \times 10^5$ | $10^9$ |

---

# Neural Network Example

$x = 1 \quad y = 5$ 　　　　 $w^{[1]} = 2, b^{[1]} = 0$ 　　 ReLU activation

$$\vec{x} \rightarrow \boxed{\bigcirc} \xrightarrow{\vec{a}^{[1]}} \boxed{\bigcirc} \xrightarrow{\vec{a}^{[2]}}$$
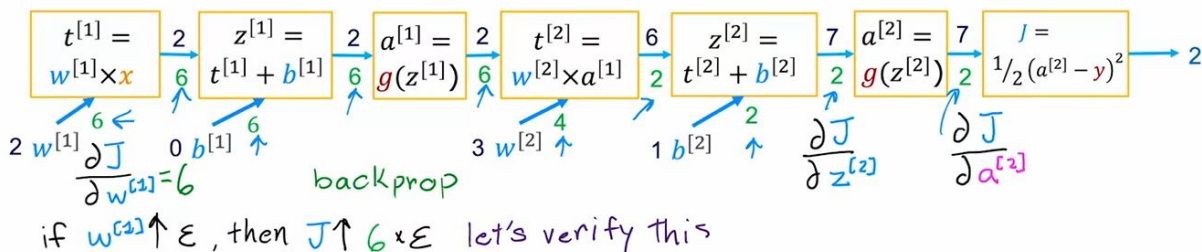
$w^{[2]} = 3, b^{[2]} = 1$ 　　 $g(z) = \max(0, z)$

$a^{[1]} = g(w^{[1]} x + b^{[1]}) = \overbrace{w^{[1]} x + b^{[1]}}^{z^{[1]}} = 2 \times 1 + 0 = 2$

$a^{[2]} = g(w^{[2]} a^{[1]} + b^{[2]}) = \overbrace{w^{[2]} a^{[1]} + b^{[2]}}^{z^{[2]}} = 3 \times 2 + 1 = 7$

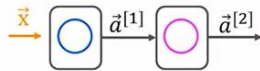$J(w, b) = \frac{1}{2}(a^{[2]} - y)^2 = \frac{1}{2}(7 - 5)^2 = 2$

$$\boxed{\begin{array}{c} t^{[1]} = \\ w^{[1]} \times x \end{array}} \xrightarrow[6]{2} \boxed{\begin{array}{c} z^{[1]} = \\ t^{[1]} + b^{[1]} \end{array}} \xrightarrow[6]{2} \boxed{\begin{array}{c} a^{[1]} = \\ g(z^{[1]}) \end{array}} \xrightarrow[6]{2} \boxed{\begin{array}{c} t^{[2]} = \\ w^{[2]} \times a^{[1]} \end{array}} \xrightarrow[2]{6} \boxed{\begin{array}{c} z^{[2]} = \\ t^{[2]} + b^{[2]} \end{array}} \xrightarrow[2]{7} \boxed{\begin{array}{c} a^{[2]} = \\ g(z^{[2]}) \end{array}} \xrightarrow[2]{7} \boxed{\begin{array}{c} J = \\ {}^1/_2 (a^{[2]} - y)^2 \end{array}} \xrightarrow{2}$$

$2\, w^{[1]} \quad 0\, b^{[1]} \qquad\qquad\qquad 3\, w^{[2]} \qquad 1\, b^{[2]} \qquad \frac{\partial J}{\partial z^{[2]}} \qquad \frac{\partial J}{\partial a^{[2]}}$

$\frac{\partial J}{\partial w^{[2]}} = 6$

backprop

if $w^{[2]} \uparrow \varepsilon$, then $J \uparrow 6 \times \varepsilon$   let's verify this

# Neural Network Example

$x = 1$  $y = 5$

$w^{[1]} = 2, b^{[1]} = 0$

$w^{[2]} = 3, b^{[2]} = 1$

ReLU activation

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



$a^{[1]} = g(w^{[1]} x + b^{[1]}) = w^{[1]} x + b^{[1]} = 2 \times 1 + 0 = 2$   2.001   2.001

$a^{[2]} = g(w^{[2]} a^{[1]} + b^{[2]}) = w^{[2]} a^{[1]} + b^{[2]} = 3 \times 2 + 1 = 7$   2.001   7.003

$J(w, b) = \frac{1}{2}(a^{[2]} - y)^2 = \frac{1}{2}(7 - 5)^2 = 2$   1   7.003   2.006   005

$\frac{(2.003)^2}{2}$

$w^{[2]} \uparrow 0.001$

$J \uparrow 6 \times 0.001$   $\frac{\partial J}{\partial w^{[2]}} = 6$

$\frac{\partial J}{\partial w^{[1]}}$   $\frac{\partial J}{\partial b^{[1]}}$

$\frac{\partial J}{\partial w^{[2]}}$   $\frac{\partial J}{\partial b^{[2]}}$

N nodes ☐→☐→☐

P parameters
$w_1, b_1, w_2, b_2 \cdots$

inefficient way

$N \times P$

efficient way (backprop)

$N + P$