

Week 2



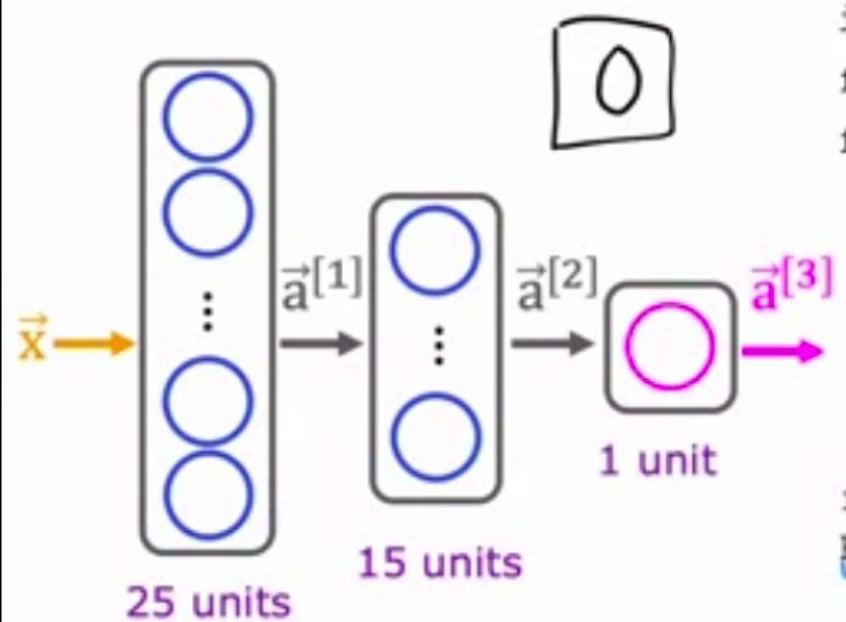
Press Esc to exit full screen

# Neural Network Training

---

## TensorFlow implementation

# Train a Neural Network in TensorFlow



Given set of  $(x, y)$  examples

How to build and train this in code?

```
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())
model.fit(X, Y, epochs=100) ③
    epochs: number of steps
    in gradient descent
```

The code shows how to build and train a neural network in TensorFlow. It imports the necessary modules and defines a sequential model with three layers, each having 25, 15, and 1 unit respectively, using sigmoid activation. It then compiles the model with binary crossentropy loss and fits it to training data X and Y for 100 epochs. A green circle highlights the 'epochs' parameter.



# Neural Network Training

---

## Training Details

# Model Training Steps

Tensor Flow

①

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

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

②

specify loss and cost

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

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

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

logistic regression

$$\begin{aligned} z &= np.dot(w, x) + b \\ f_x &= 1 / (1 + np.exp(-z)) \end{aligned}$$

logistic loss

$$\begin{aligned} \text{loss} &= -y * np.log(f_x) \\ &- (1-y) * np.log(1-f_x) \end{aligned}$$

$$\begin{aligned} w &= w - \alpha * dj_dw \\ b &= b - \alpha * dj_db \end{aligned}$$

neural network

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

binary cross entropy

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

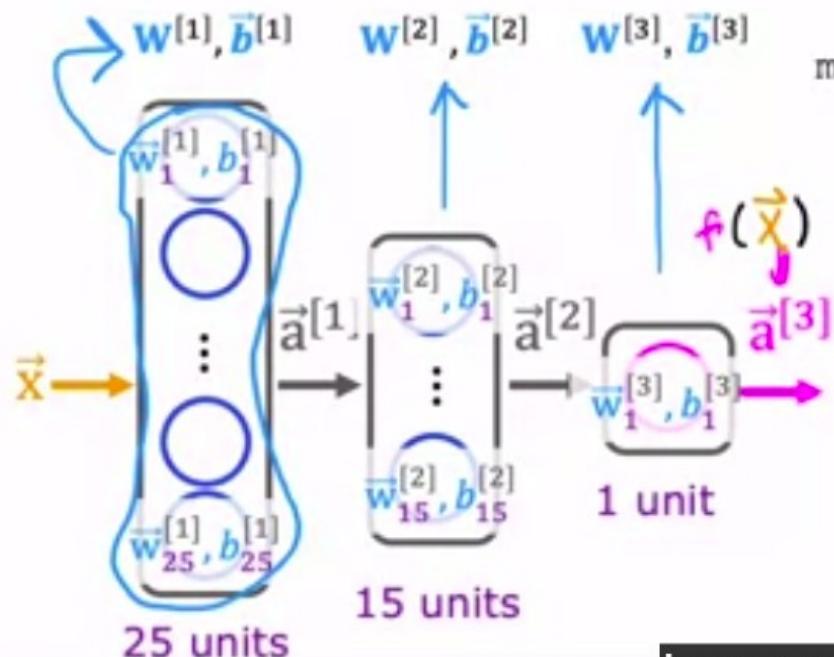
```
model.fit(X, y, epochs=100)
```

the context of training  
a neural network.

# 1. Create the model

define the model

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



```
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  $w^l$  and  $b^l$ . Let's go on to step 2.

## 2. Loss and cost functions

handwritten digit classification problem      ↗ binary classification

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

compare prediction vs. target

logistic loss

also known as binary cross entropy

model.compile(loss= BinaryCrossentropy())

regression

(predicting numbers and not categories)

model.compile(loss= MeanSquaredError())

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

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

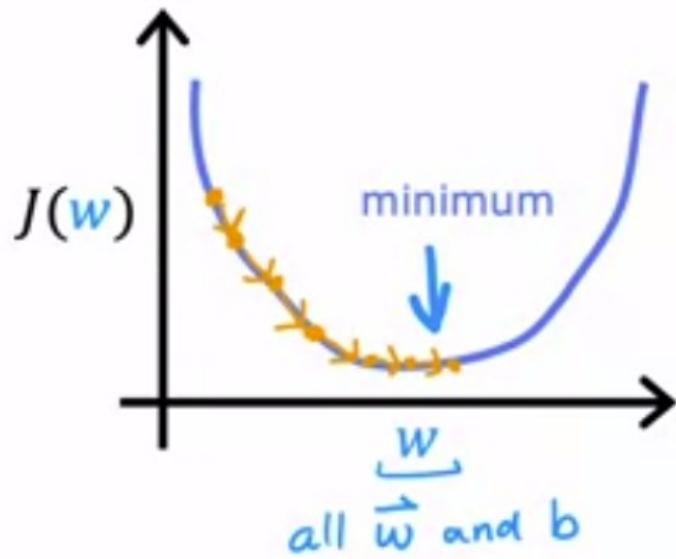
from tensorflow.keras.losses import  
BinaryCrossentropy



from tensorflow.keras.losses import  
MeanSquaredError

That's the loss function  
and the cost function.

### 3. Gradient descent



repeat {

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

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

} Compute derivatives  
for gradient descent  
using "back propagation"

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

enough that most developers  
will use these libraries,

# Neural network libraries

Use code libraries instead of coding "from scratch"



Good to understand the implementation  
(for tuning and debugging).

which still does with  
today's libraries,

[Back](#)

## Practice quiz: Neural Network Training

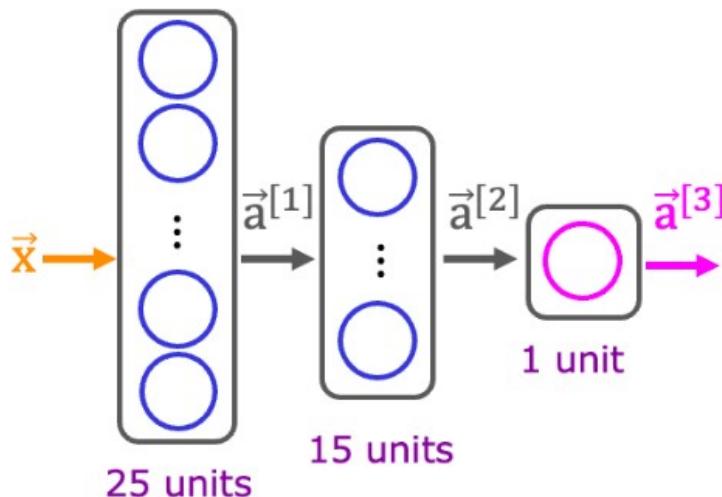
Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

1.

1 / 1 point

# Train a Neural Network in TensorFlow



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

A blue back arrow icon.

## Practice quiz: Neural Network Training

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

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

Here is some code that you saw in the lecture:

...

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

...

For which type of task would you use the binary cross entropy loss function?

- regression tasks (tasks that predict a number)
- BinaryCrossentropy() should not be used for any task.
- binary classification (classification with exactly 2 classes)
- A classification task that has 3 or more classes (categories)

A green circular icon with a white checkmark inside.

Correct

Yes! Binary cross entropy, which we've also referred to as logistic loss, is used for classifying between two classes (two categories).

[Back](#)

## Practice quiz: Neural Network Training

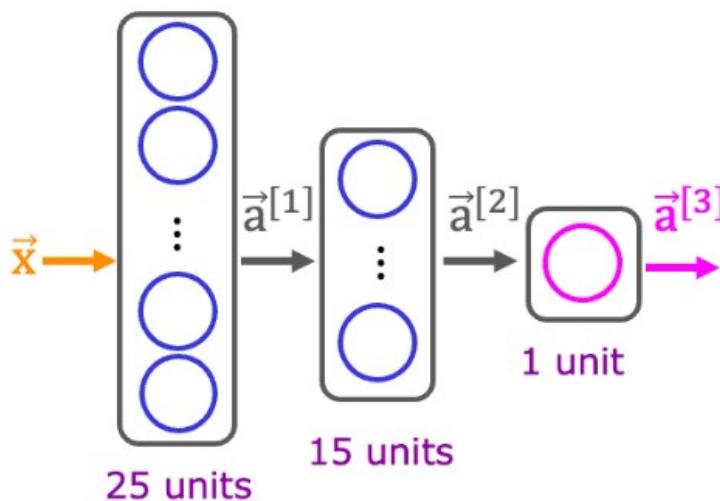
Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

2.

1 / 1 point

# Train a Neural Network in TensorFlow



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

[Back](#)

## Practice quiz: Neural Network Training

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

Here is code that you saw in the lecture:

...

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

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

...

Which line of code updates the network parameters in order to reduce the cost?

- model = Sequential([...])

[Back](#)

## Practice quiz: Neural Network Training

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

])

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

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

```
...
```

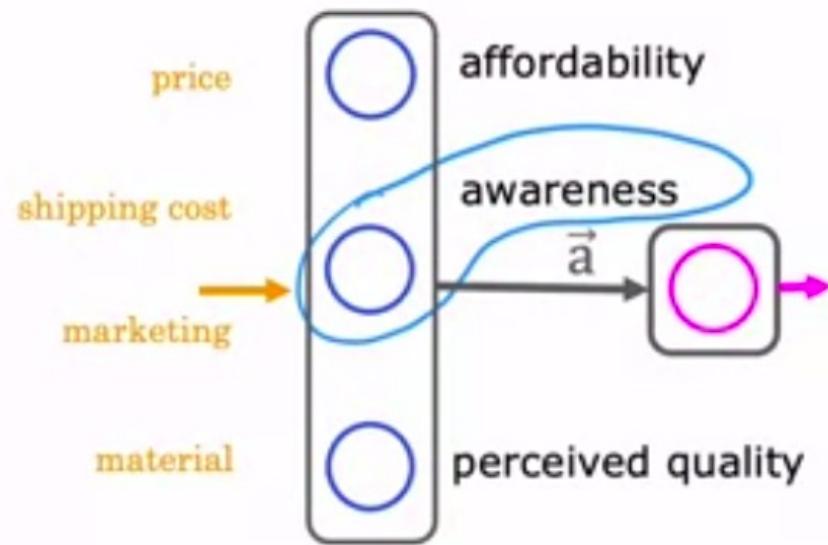
Which line of code updates the network parameters in order to reduce the cost?

- model = Sequential([...])
- None of the above -- this code does not update the network parameters.
- model.fit(X,y,epochs=100)
- model.compile(loss=BinaryCrossentropy())

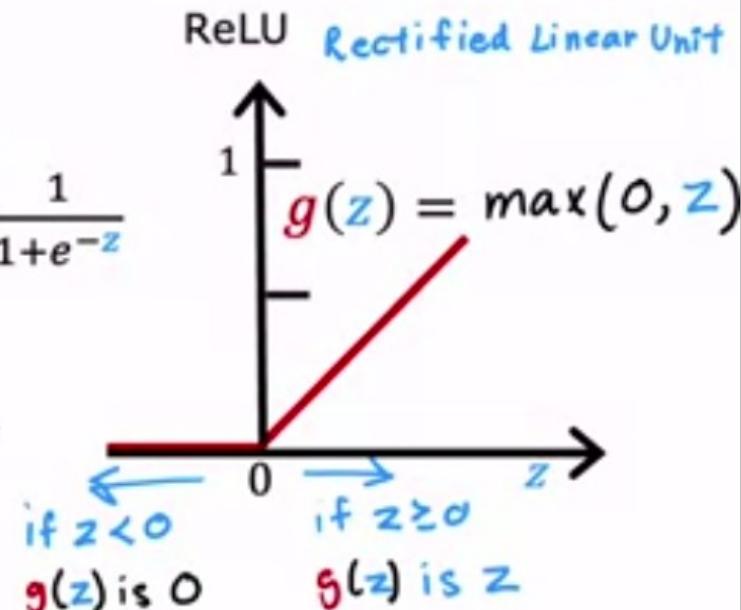
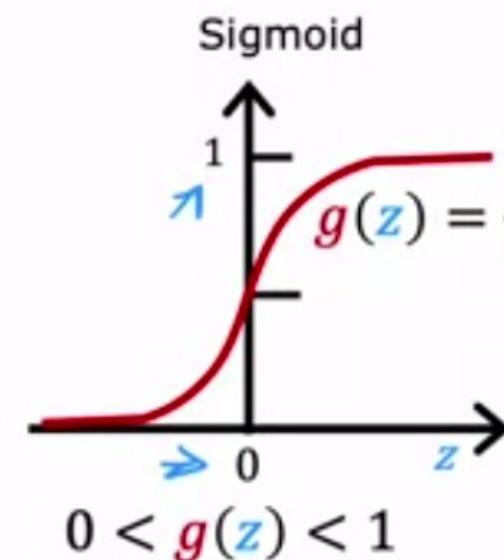
Correct

Yes! The third step of model training is to train the model on data in order to minimize the loss (and the cost)

# Demand Prediction Example



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

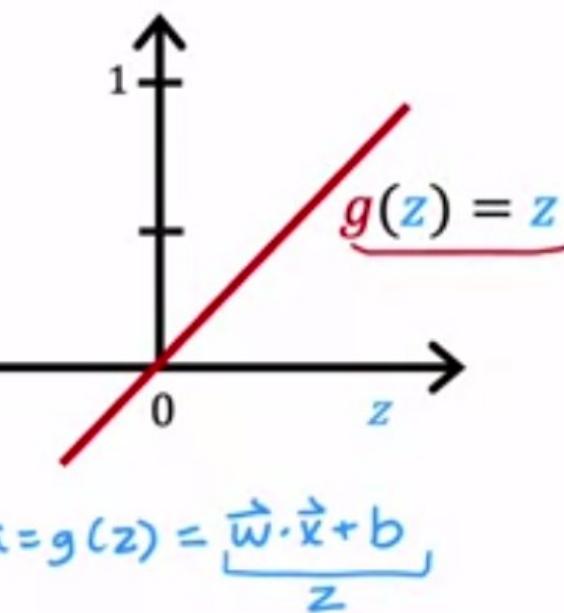


More generally you have a choice  
of what to use for  $g(z)$  and

# Examples of Activation Functions

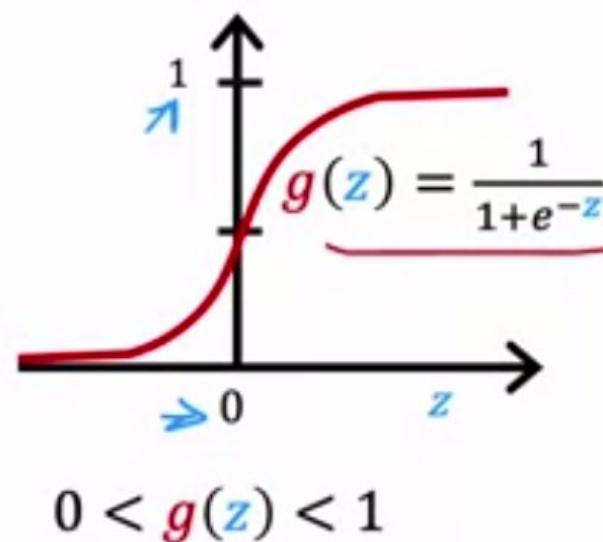
"No activation function"

Linear activation function



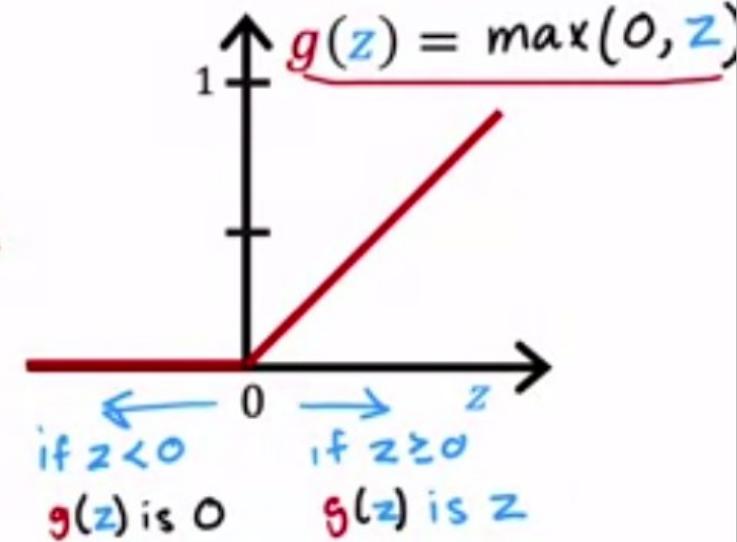
$$a_2^{[1]} = g(\overbrace{\vec{w}_2^{[1]} \cdot \vec{x}}^z + b_2^{[1]})$$

Sigmoid



ReLU

Rectified Linear Unit



# Examples of Activation Functions

"No activation function"

Linear activation function

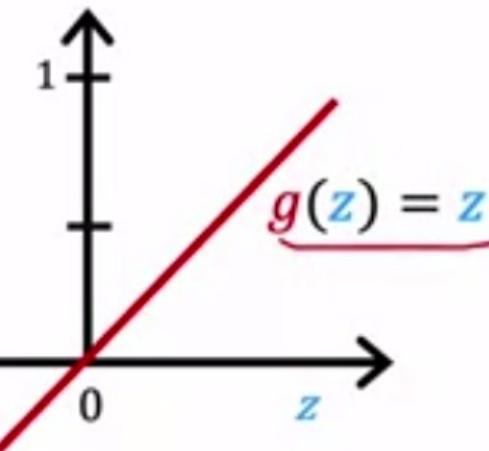
$$a_2^{[1]} = g(\overbrace{\vec{w}_2^{[1]} \cdot \vec{x}}^z + b_2^{[1]})$$

Sigmoid

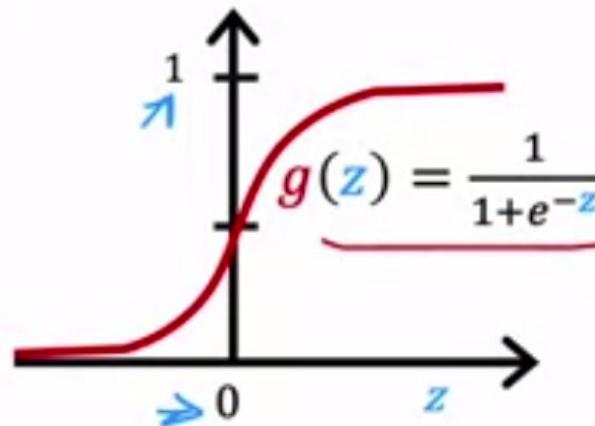
Later: softmax activation

ReLU

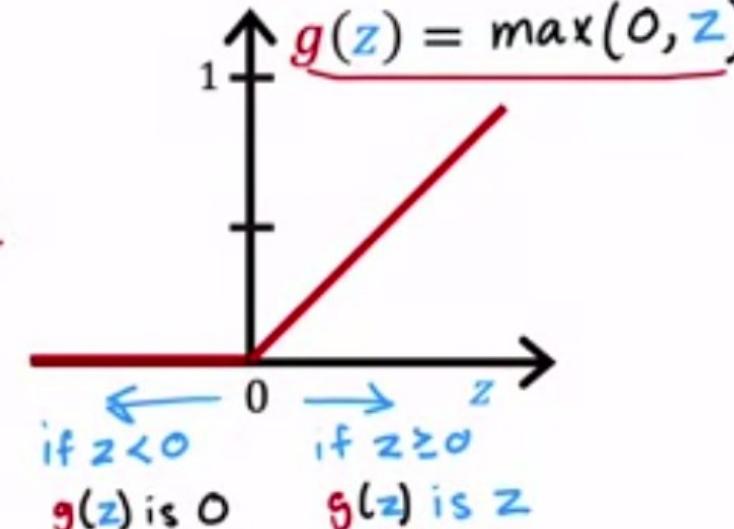
Rectified Linear Unit



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



$$0 < g(z) < 1$$



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

if  $z \geq 0$   
 $g(z)$  is  $z$

But with these activation  
functions you'll be able to build



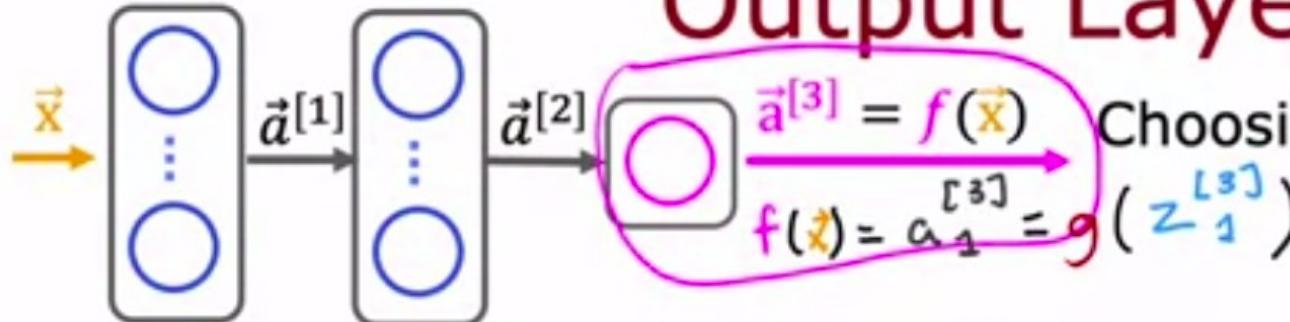
# Activation Functions

---

## Choosing activation functions

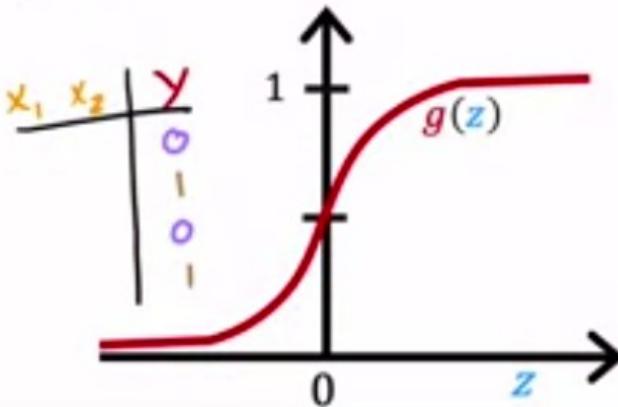
Let's take a look at  
how you can choose

# Output Layer



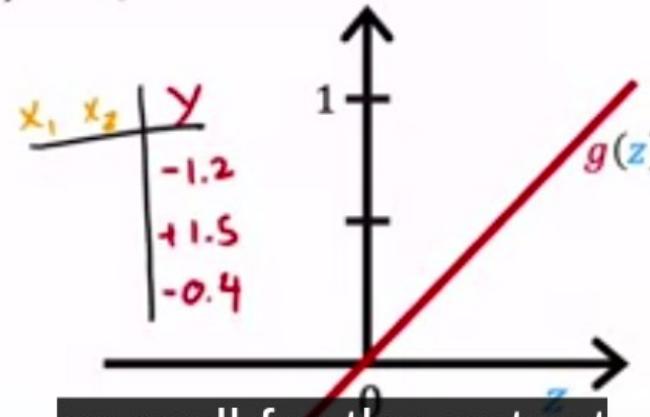
Binary classification

Sigmoid  
 $y=0/1$



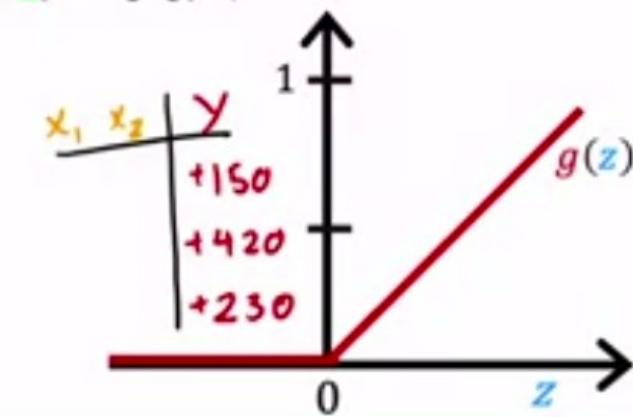
Regression

Linear activation function  
 $y = +/-$



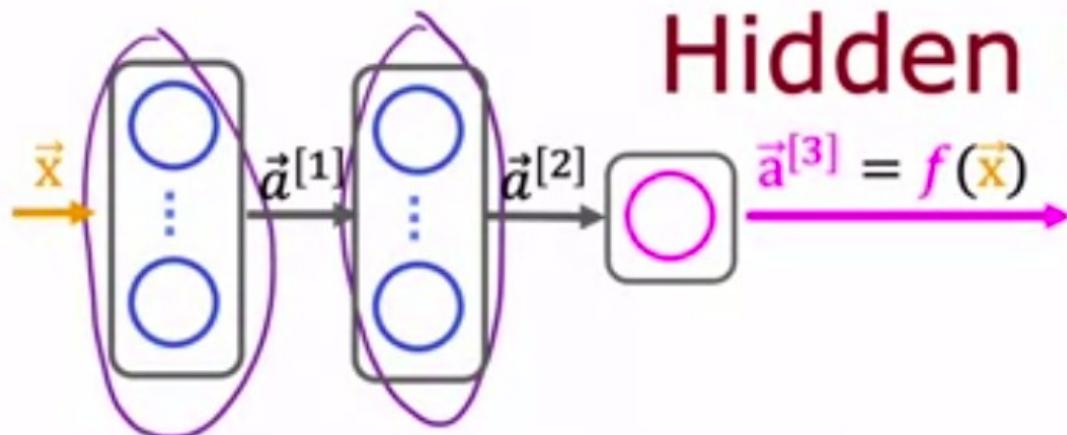
Regression

ReLU  
 $y = 0 \text{ or } +$

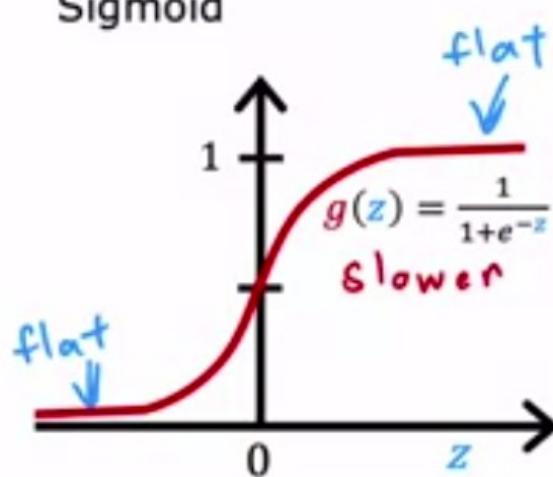


as well for the output  
layer of a neural network.

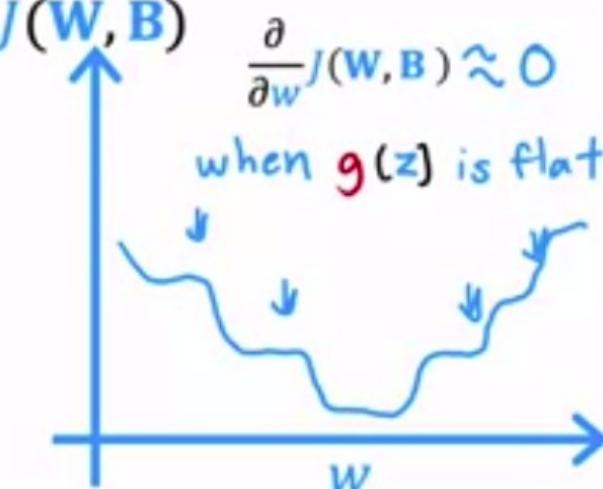
# Hidden Layer



Sigmoid



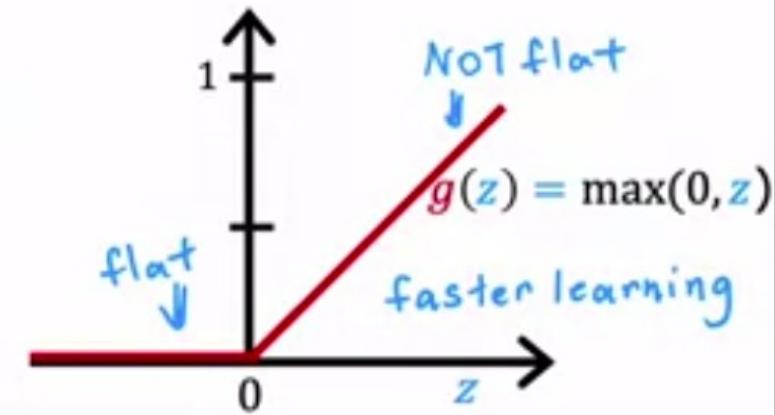
$J(W, B)$



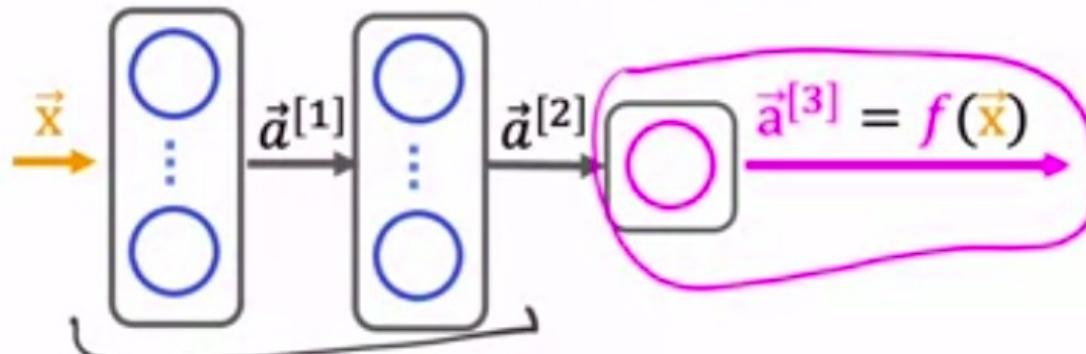
most common choice

ReLU

faster



# Choosing Activation Summary



binary classification

activation='sigmoid'

regression  $y$  negative/  
positive  
activation='linear'

regression  $y \geq 0$   
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'
```



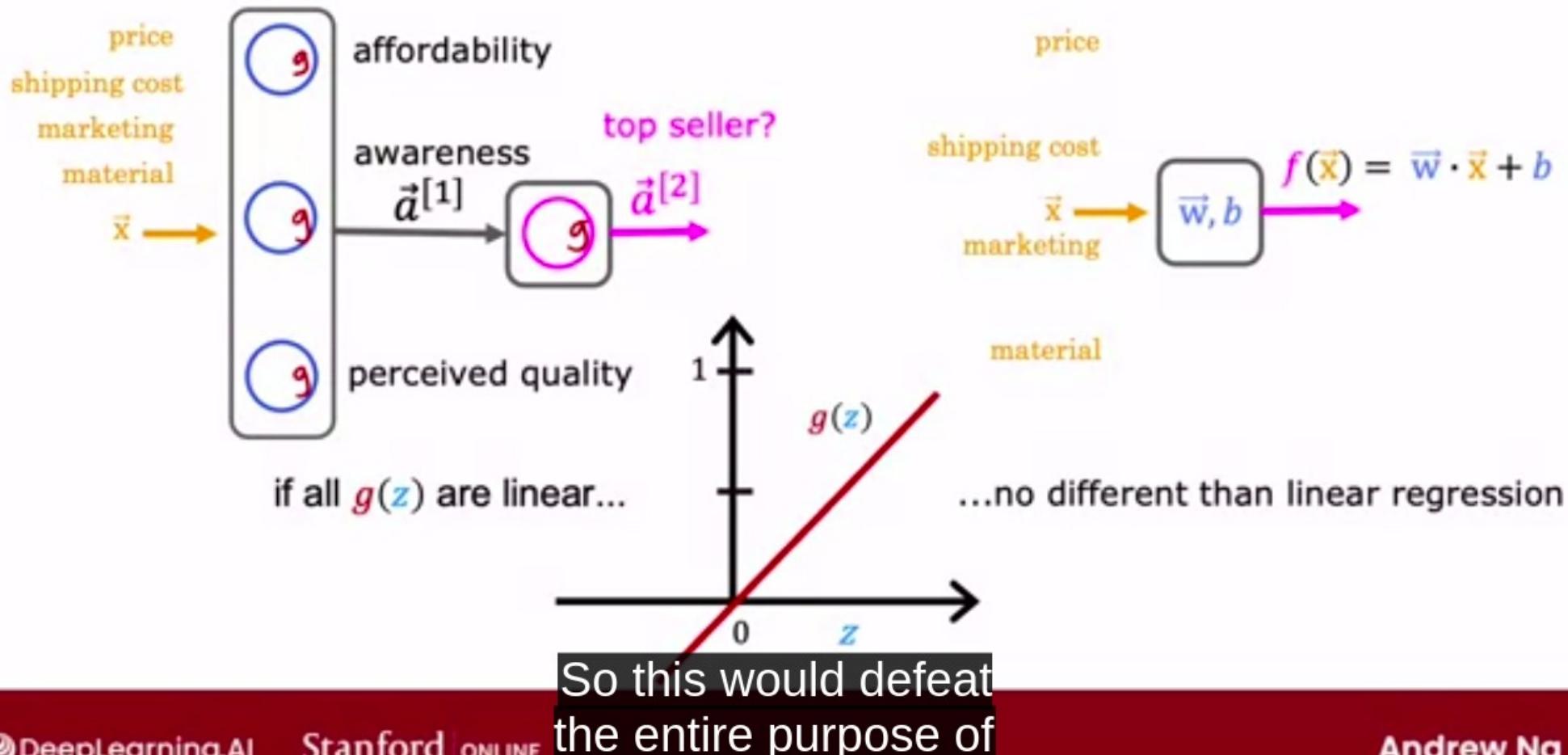
Press **Esc** to exit full screen

# Activation Functions

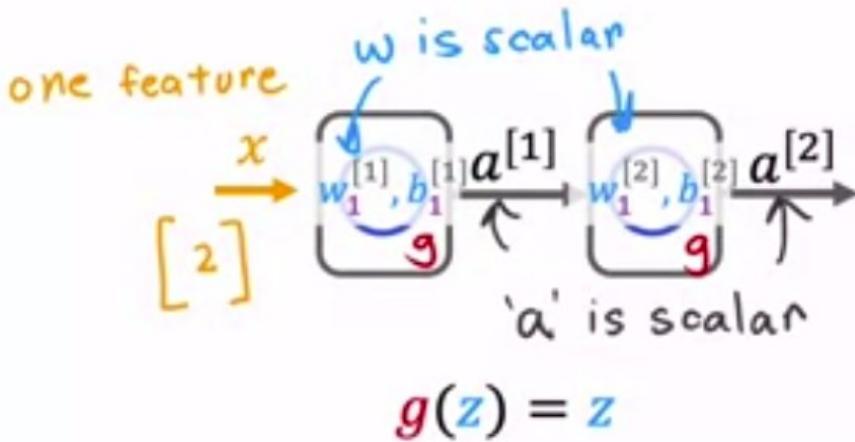
---

Why do we need  
activation functions?  
Let's take a look at why

# Why do we need activation functions?



# Linear Example



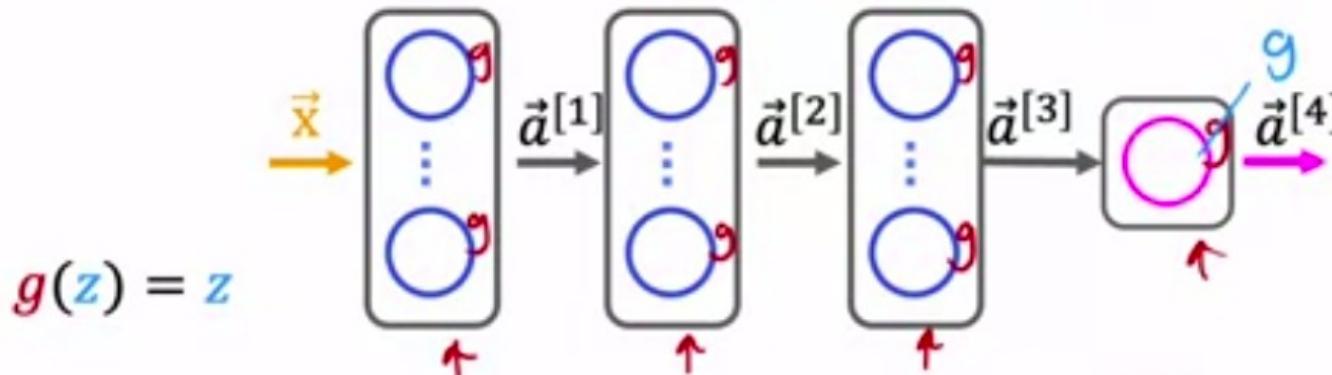
$$\begin{aligned}a^{[1]} &= \underbrace{w_1^{[1]} x + b_1^{[1]}}_{\downarrow} \\a^{[2]} &= w_1^{[2]} a^{[1]} + b_1^{[2]} \\&= w_1^{[2]} (w_1^{[1]} x + b_1^{[1]}) + b_1^{[2]} \\\vec{a}^{[2]} &= (\underbrace{\vec{w}_1^{[2]} \vec{w}_1^{[1]}}_{\omega}) x + \underbrace{w_1^{[2]} b_1^{[1]} + b_1^{[2]}}_{b}\end{aligned}$$

$$\vec{a}^{[2]} = w x + b$$

$$f(x) = w x + b \text{ linear regression}$$

anything more complex than  
just a linear function.

# Example



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

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 (use ReLU)

the ReLU activation function  
should do just fine.

[Back](#)

## Practice quiz: Activation Functions

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

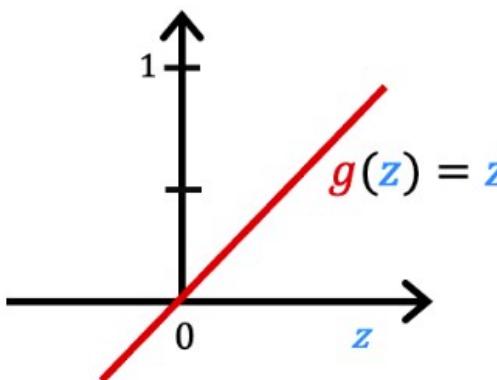
1.

# Examples of Activation Functions

1 / 1 point

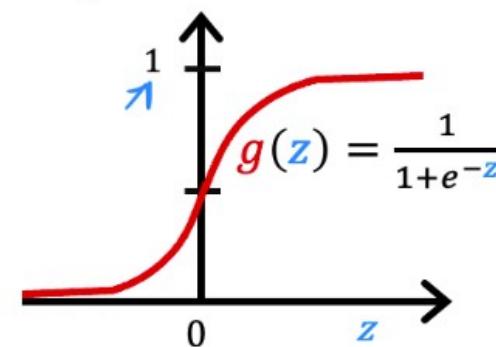
"No activation function"

Linear activation function



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

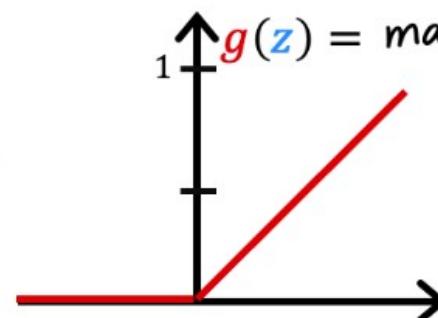
Sigmoid



ReLU

Rectified Linear Unit

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



Which of the following activation functions is the most common choice for the hidden layers of a neural network?

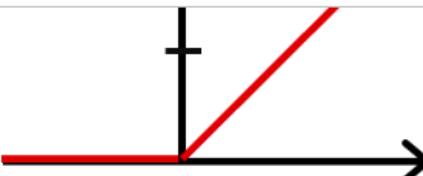
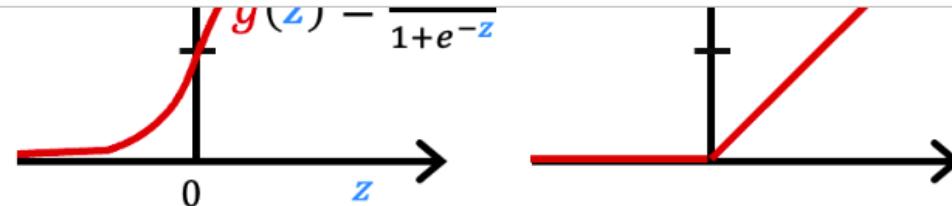
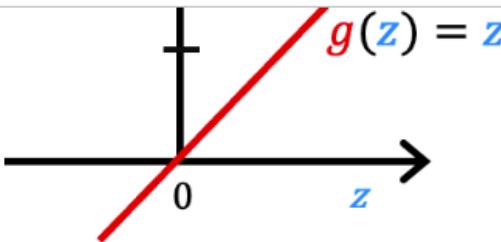
- ReLU (rectified linear unit)

[Back](#)

## Practice quiz: Activation Functions

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST



Which of the following activation functions is the most common choice for the hidden layers of a neural network?

- ReLU (rectified linear unit)
- Sigmoid
- Linear
- Most hidden layers do not use any activation function

**Correct**

Yes! A ReLU is most often used because it is faster to train compared to the sigmoid. This is because the ReLU is only flat on one side (the left side) whereas the sigmoid goes flat (horizontal, slope approaching zero) on both sides of the curve.

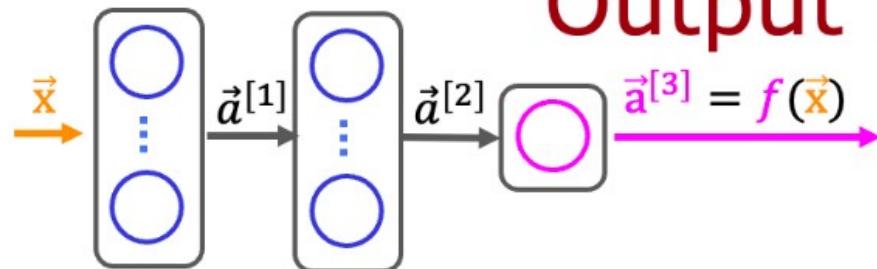
[Back](#)

## Practice quiz: Activation Functions

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

2.

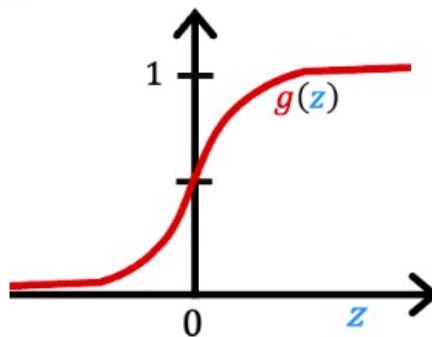


## Output Layer

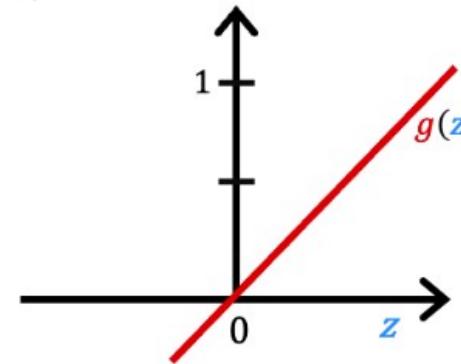
1 / 1 point

Choosing  $g(z)$  for output layer?

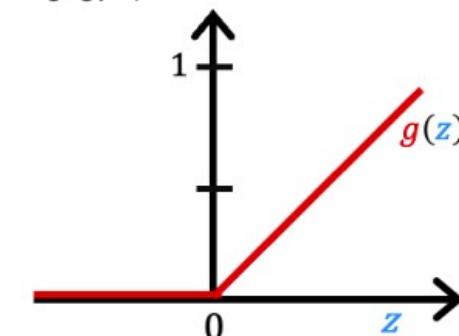
Binary classification

Sigmoid  
 $y=0/1$ 

Regression

Linear activation function  
 $y = +/-$ 

Regression

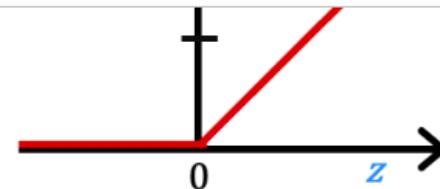
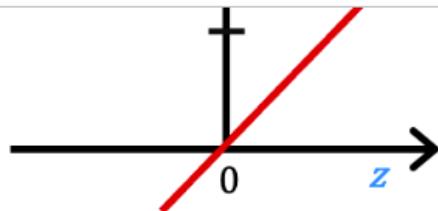
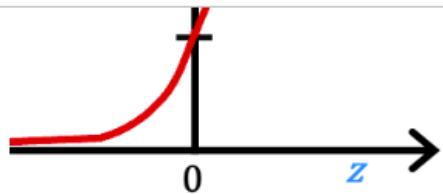
ReLU  
 $Y = 0 \text{ or } +$ 

[Back](#)

## Practice quiz: Activation Functions

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST



For the task of predicting housing prices, which activation functions could you choose for the output layer? Choose the 2 options that apply.

 Sigmoid linear **Correct**

Yes! A linear activation function can be used for a regression task where the output can be both negative and positive, but it's also possible to use it for a task where the output is 0 or greater (like with house prices).

 ReLU **Correct**

Yes! ReLU outputs values 0 or greater, and housing prices are positive values.

[Back](#)

## Practice quiz: Activation Functions

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

 linear **Correct**

Yes! A linear activation function can be used for a regression task where the output can be both negative and positive, but it's also possible to use it for a task where the output is 0 or greater (like with house prices).

 ReLU **Correct**

Yes! ReLU outputs values 0 or greater, and housing prices are positive values.

3. True/False? A neural network with many layers but no activation function (in the hidden layers) is not effective; that's why we should instead use the linear activation function in every hidden layer.

1 / 1 point

 False True **Correct**

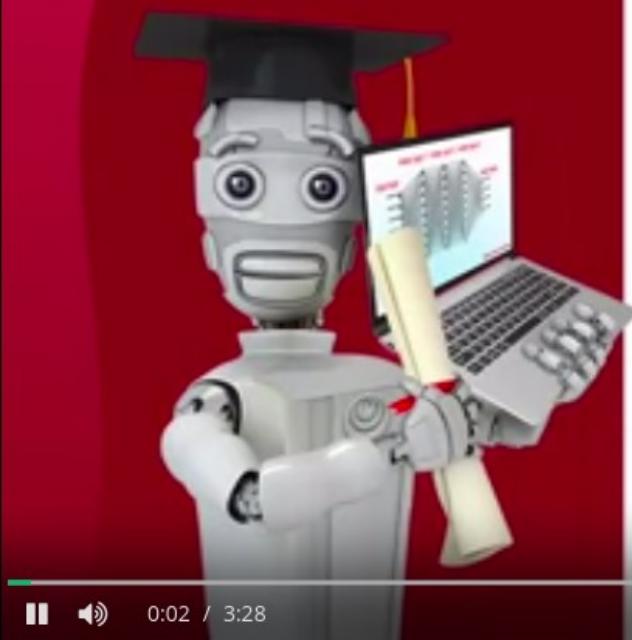
Yes! A neural network with many layers but no activation function is not effective. A linear activation is the same as "no activation function".

# Multiclass Classification

---

## Multiclass

Multiclass classification refers  
to classification problems where



# MNIST example

0 1 2 3 4 5 6 7 8 9  
 $y = 0$       1      2      3      4      5      6      7      8      9

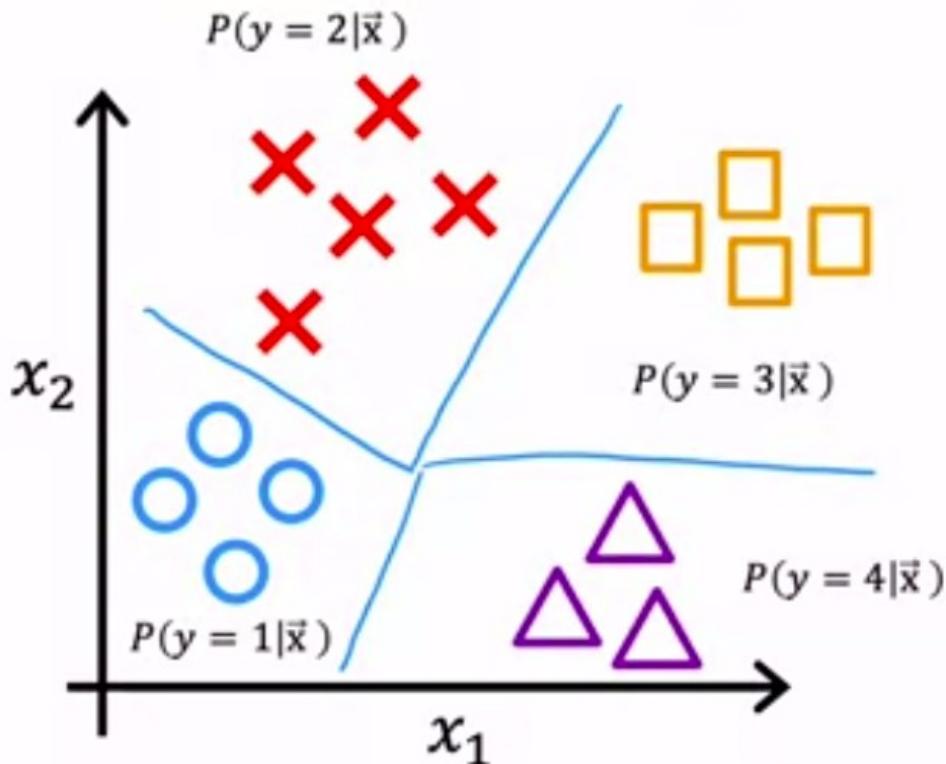
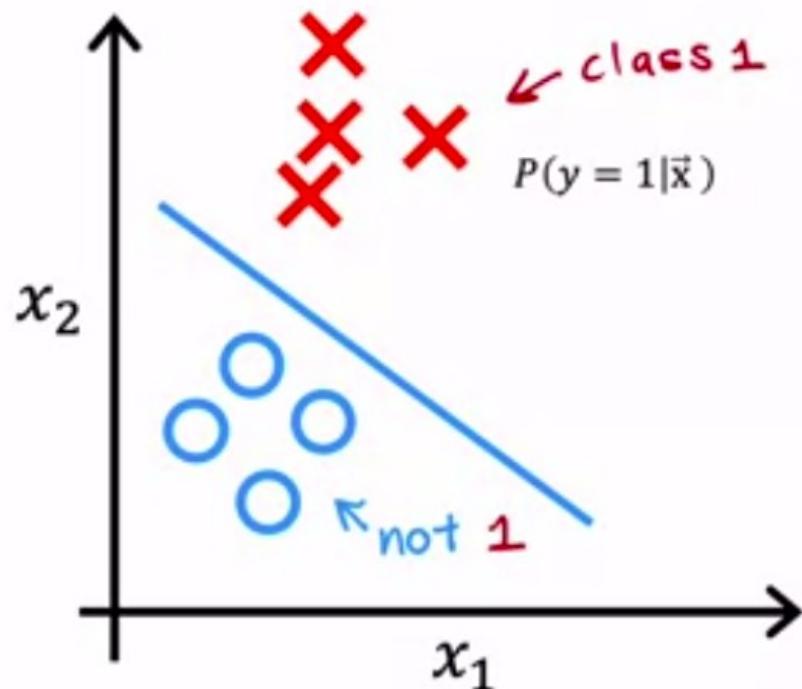
$x \xrightarrow{\text{?}} 7$        $y = 7$

multiclass classification problem:

target  $y$  can take on more than two possible values

but now  $y$  can take on more  
than just two possible values.

# Multiclass classification example



exploded next to into four categories  
rather than just two categories.



# Multiclass Classification

---

## Softmax

a generalization of  
logistic regression,

## Question

What do you think  $a_4$  is equal to?

$a_1 = g(z_1)$

$a_2 = \dots$

$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$



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

$$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

$$a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

Skip

Continue

## Question

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

$$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

$$a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

$$a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

[Skip](#)[Continue](#)

## Question

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

$$a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

- 0.35
- 0.40
- 0.40

✓ **Correct**

This is correct. Please continue the video to see why!

Skip

Continue

## Logistic regression (2 possible output values)

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

0.11

$\times \alpha_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1|\vec{x})$

$\circlearrowleft \alpha_2 = 1 - \alpha_1 = P(y=0|\vec{x})$   
0.29

## Softmax regression (N possible outputs) $y=1, 2, 3, \dots, N$

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

parameters  $w_1, w_2, \dots, w_N$   
 $b_1, b_2, \dots, b_N$

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

Note:  $\alpha_1 + \alpha_2 + \dots + \alpha_N = 1$

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

$\times z_1 = \vec{w}_1 \cdot \vec{x} + b_1$

$$\alpha_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$\times \circlearrowleft \square \triangle$

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

$\circlearrowleft z_2 = \vec{w}_2 \cdot \vec{x} + b_2$

$$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

$\square z_3 = \vec{w}_3 \cdot \vec{x} + b_3$

$$a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

$\triangle z_4 = \vec{w}_4 \cdot \vec{x} + b_4$

$$a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

# Cost

## Logistic regression

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

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

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

$$\text{loss} = -y \underbrace{\log a_1}_{\text{if } y=1} - (1-y) \underbrace{\log(1-a_1)}_{\text{if } y=0}$$

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

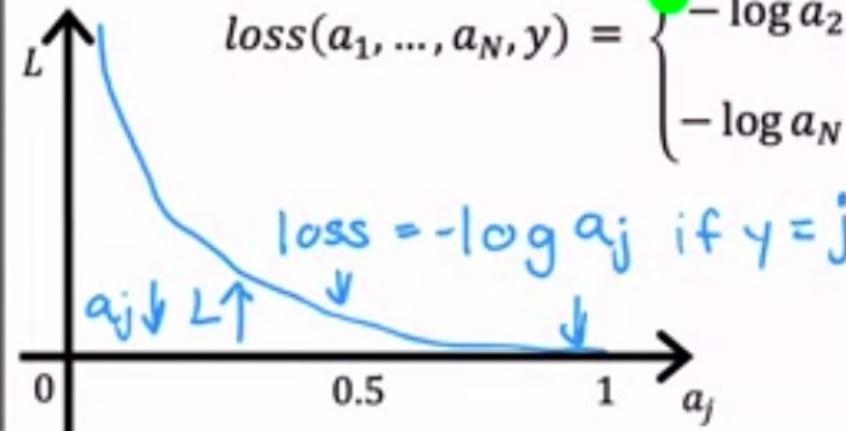
## Softmax regression

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

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

Crossentropy loss

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



# Cost

## Logistic regression

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

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

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

$$\text{loss} = -y \underbrace{\log a_1}_{\text{if } y=1} - (1-y) \underbrace{\log(1-a_1)}_{\text{if } y=0}$$

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

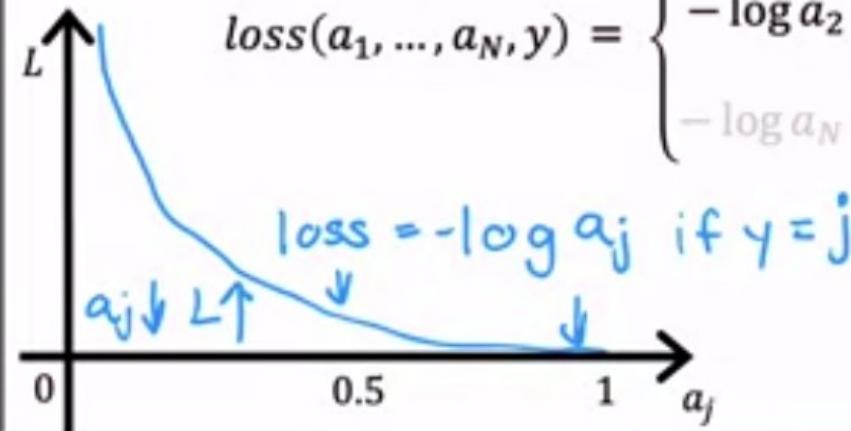
## Softmax regression

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

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

### Crossentropy loss

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



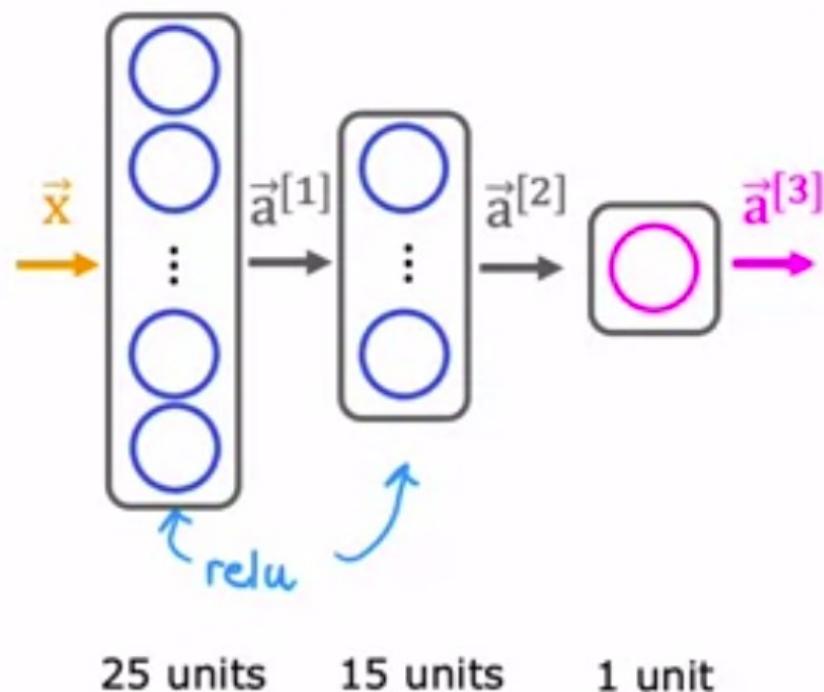


# Multiclass Classification

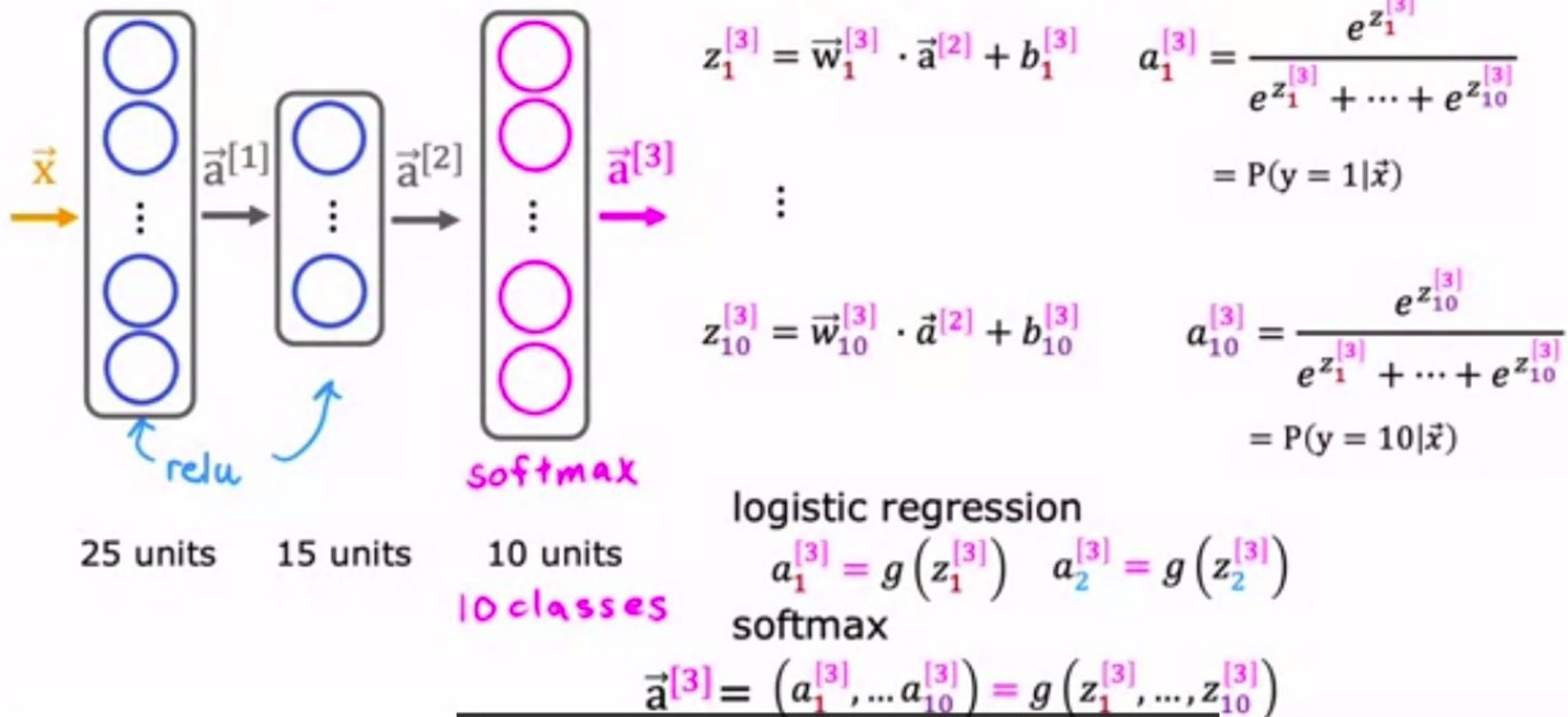
---

Neural Network with  
Softmax output

# Neural Network with Softmax output



# Neural Network with Softmax output



If you want to implement the neural network that I've shown here on

# 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')
])
```

②

specify loss and cost

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

```
from tensorflow.keras.losses import
SparseCategoricalCrossentropy
model.compile(loss= SparseCategoricalCrossentropy() )
model.fit(X, Y, epochs=100)
```

③

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

Note: better (recommended) version later.

Don't use the version shown here!



Press **Esc** to exit full screen

# Multiclass Classification

---

Improved implementation  
of softmax

# Multiclass Classification

---

Improved implementation  
of softmax



# Numerical Roundoff Errors

option 1

$$x = \frac{2}{10,000}$$

option 2

$$x = \underbrace{\left(1 + \frac{1}{10,000}\right)} - \underbrace{\left(1 - \frac{1}{10,000}\right)} =$$



File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

Cell Kernel Help Trusted Python 3

```
In [1]: x1 = 2.0 / 10000
print(f"{x1:.18f}") # print 18 digits to the right of decimal point
```

0.00020000000000000000

```
In [2]: x2 = 1 + (1/10000) - (1 - 1/10000)
print(f"{x2: .18f}")
```

0.000199999999999978

```
In [ ]:
```

# Numerical Roundoff Errors

More numerically accurate implementation of logistic loss:

Logistic regression:

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

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='sigmoid')
])
model.compile(loss=BinaryCrossEntropy())
```

Original loss

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

# Numerical Roundoff Errors

More numerically accurate implementation of logistic loss:

Logistic regression:

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

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

Original loss

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

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

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

More accurate loss (in code)

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

logit: z

## More numerically accurate implementation of softmax

Softmax regression

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

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

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
model.compile(loss=SparseCategoricalCrossEntropy())
```

# More numerically accurate implementation of softmax

Softmax regression

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

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

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

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

More Accurate

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

~~model.compile(loss=SparseCategoricalCrossEntropy(from\_logits=True))~~



# 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  $a_1 \dots a_{10}$ 
         f_x = tf.nn.softmax(logits)           is  $z_1 \dots z_{10}$ 
```

# 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)      Z ↗
predict   f_x = tf.nn.sigmoid(logit)
```



# Multi-label Classification

---

Classification with  
multiple outputs  
(Optional)

# Multi-label Classification



Is there a car?

$$\begin{matrix} \text{yes} \\ \text{no} \end{matrix} \quad y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

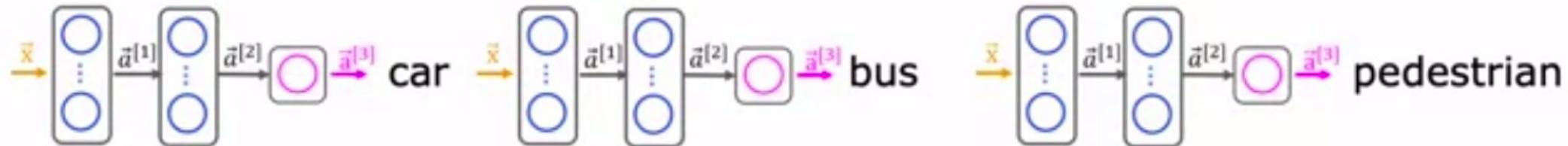
Is there a bus?

$$\begin{matrix} \text{no} \\ \text{no} \\ \text{yes} \end{matrix} \quad y = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

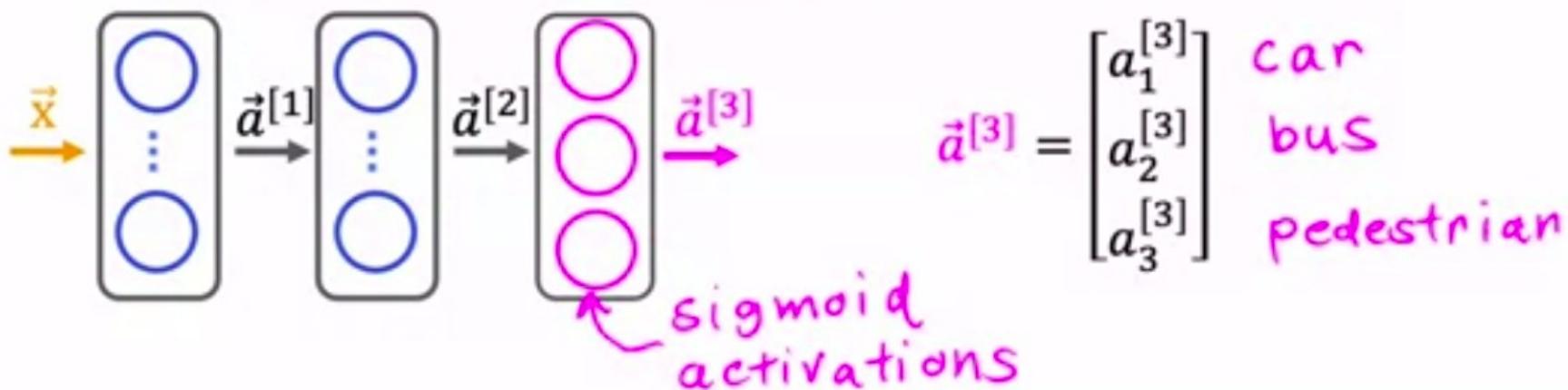
Is there a pedestrian

$$\begin{matrix} \text{yes} \\ \text{yes} \\ \text{no} \end{matrix} \quad y = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$$

# Multi-label Classification



Alternatively, train one neural network with three outputs



[Back](#)

## Practice quiz: Multiclass Classification

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

1.

1 / 1 point

## Softmax regression (4 possible outputs)

$$z_1 = \vec{w}_1 \cdot \vec{x} + b_1$$

$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$
$$= P(y = 1 | \vec{x}) \quad 0.30$$

$$z_2 = \vec{w}_2 \cdot \vec{x} + b_2$$

$$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$
$$= P(y = 2 | \vec{x}) \quad 0.20$$

$$z_3 = \vec{w}_3 \cdot \vec{x} + b_3$$

$$a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$
$$= P(y = 3 | \vec{x}) \quad 0.15$$

[Back](#)

## Practice quiz: Multiclass Classification

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

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

$$z_3 = \vec{w}_3 \cdot \vec{x} + b_3$$

$$a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

$$z_4 = \vec{w}_4 \cdot \vec{x} + b_4$$

$$a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

For a multiclass classification task that has 4 possible outputs, the sum of all the activations adds up to 1. For a multiclass classification task that has 3 possible outputs, the sum of all the activations should add up to ....

 Less than 1

[Back](#)

## Practice quiz: Multiclass Classification

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

$$\Delta \quad z_4 = \vec{w}_4 \cdot \vec{x} + b_4$$

$$a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$
$$= P(y = 4 | \vec{x}) \text{ o.35}$$

For a multiclass classification task that has 4 possible outputs, the sum of all the activations adds up to 1. For a multiclass classification task that has 3 possible outputs, the sum of all the activations should add up to ....

- Less than 1
- More than 1
- 1
- It will vary, depending on the input x.

 Correct

Yes! The sum of all the softmax activations should add up to 1. One way to see this is that if  $e^{z_1} = 10, e^{z_2} = 20, e^{z_3} = 30$ , then the sum of  $a_1 + a_2 + a_3$  is equal to  $\frac{e^{z_1}+e^{z_2}+e^{z_3}}{e^{z_1}+e^{z_2}+e^{z_3}}$  which is 1.

[Back](#)

## Practice quiz: Multiclass Classification

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

2.

1 / 1 point

# Cost

## Logistic regression

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

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

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

$$\text{loss} = -y \underbrace{\log a_1}_{\text{if } y=1} - (1-y) \underbrace{\log(1-a_1)}_{\text{if } y=0}$$

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

## Softmax regression

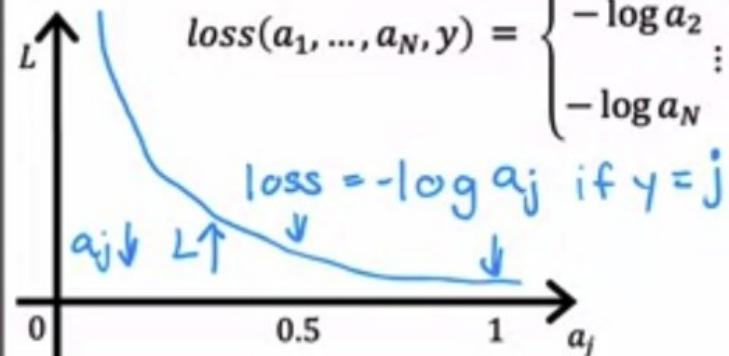
$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}} = P(y = 1 | \vec{x})$$

$$\vdots$$

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

### Crossentropy loss

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



[Back](#)

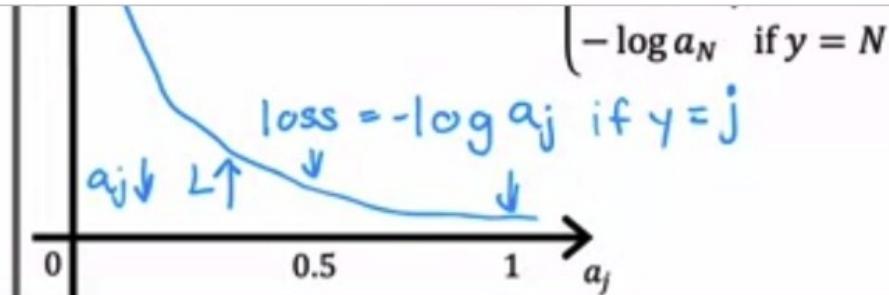
## Practice quiz: Multiclass Classification

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

$$\text{loss} = \underbrace{-y \log a_1}_{\text{if } y=1} - \underbrace{(1-y) \log(1-a_1)}_{\text{if } y=0}$$

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



For multiclass classification, the cross entropy loss is used for training the model. If there are 4 possible classes for the output, and for a particular training example, the true class of the example is class 3 ( $y=3$ ), then what does the cross entropy loss simplify to? [Hint: This loss should get smaller when  $a_3$  gets larger.]

- z\_3
- $\frac{-\log(a_1) + -\log(a_2) + -\log(a_3) + -\log(a_4)}{4}$
- $z_3/(z_1+z_2+z_3+z_4)$
- $-\log(a_3)$

✓ **Correct**

Correct. When the true label is 3, then the cross entropy loss for that training example is just the negative of the log of the activation for the third neuron of the softmax. All other terms of the cross entropy loss equation ( $-\log(a_1)$ ,  $-\log(a_2)$ , and  $-\log(a_4)$ ) are ignored

A blue back arrow icon.

## Practice quiz: Multiclass Classification

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

3.

1 / 1 point

# 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)
         f_x = tf.nn.softmax(logits)
```

A blue back arrow icon.

## Practice quiz: Multiclass Classification

Graded Quiz • 30 min

Due Nov 27, 11:59 PM IST

```
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)
          f_x = tf.nn.softmax(logits)
```

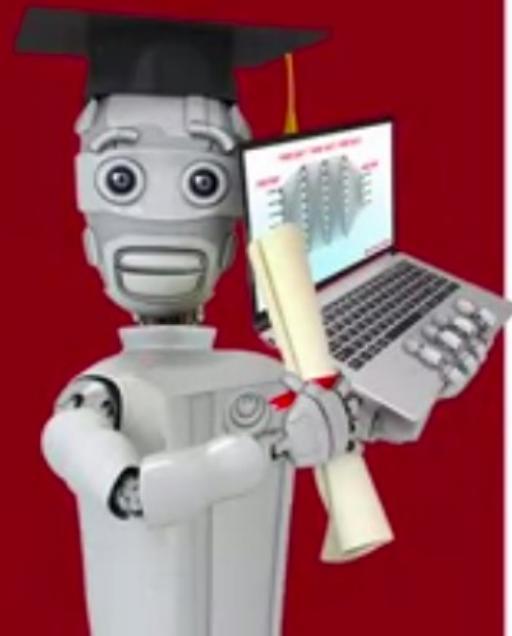
For multiclass classification, the recommended way to implement softmax regression is to set `from_logits=True` in the loss function, and also to define the model's output layer with...

- a 'linear' activation
- a 'softmax' activation

A green checkmark icon inside a circle.

Correct

Yes! Set the output as linear, because the loss function handles the calculation of the softmax with a more numerically stable method.



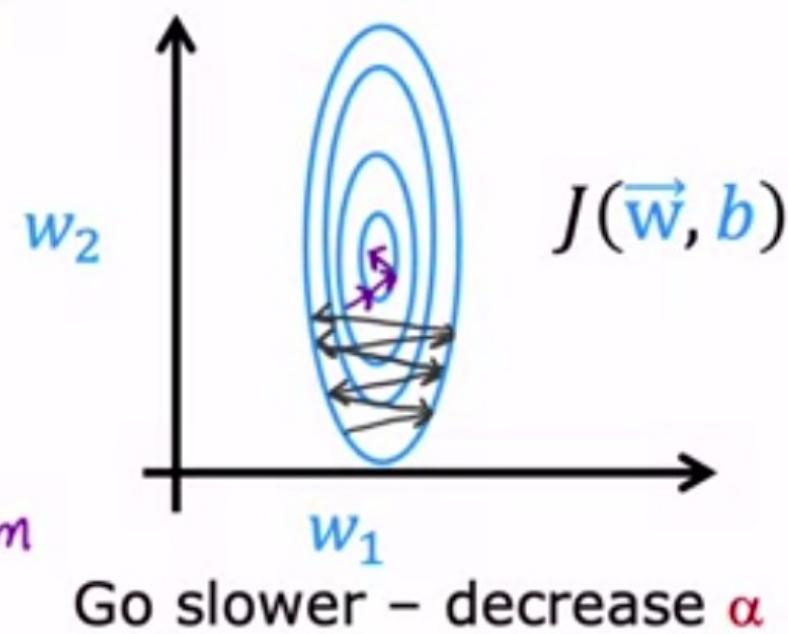
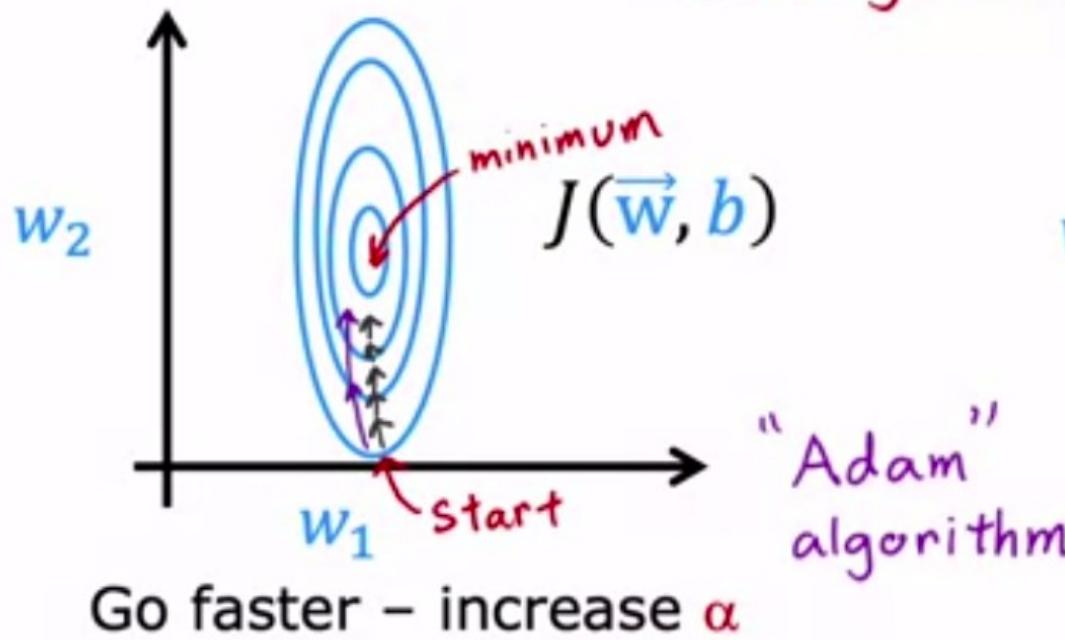
## Additional Neural Network Concepts

### Advanced Optimization

# Gradient Descent

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

↑  
learning rate



a smaller learning rate Alpha.

# Adam Algorithm Intuition

Adam: Adaptive Moment estimation      not just one  $\alpha$

$$w_1 = w_1 - \underbrace{\alpha_1}_{\text{not just one } \alpha} \frac{\partial}{\partial w_1} J(\vec{w}, b)$$

⋮

$$w_{10} = w_{10} - \underbrace{\alpha_{10}}_{\text{not just one } \alpha} \frac{\partial}{\partial w_{10}} J(\vec{w}, b)$$

$$b = b - \underbrace{\alpha_{11}}_{\text{not just one } \alpha} \frac{\partial}{\partial b} J(\vec{w}, b)$$

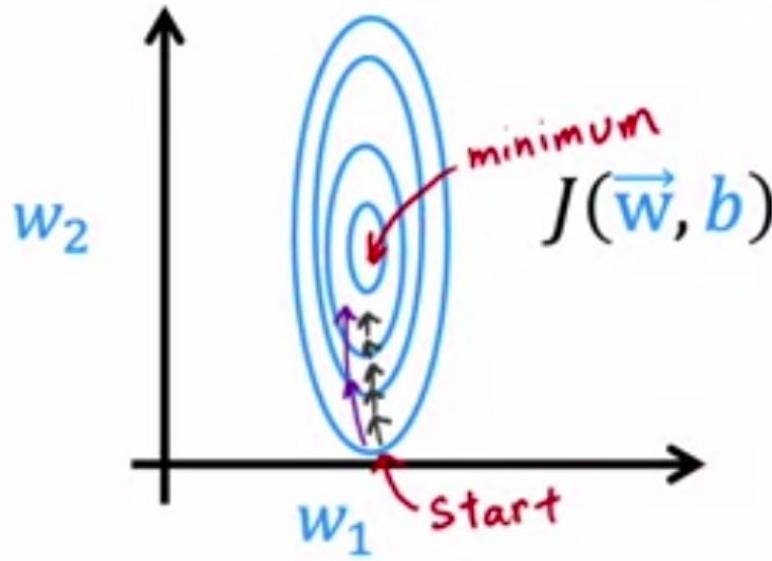
as well as I'll call it

Alpha\_11 for the parameter b.

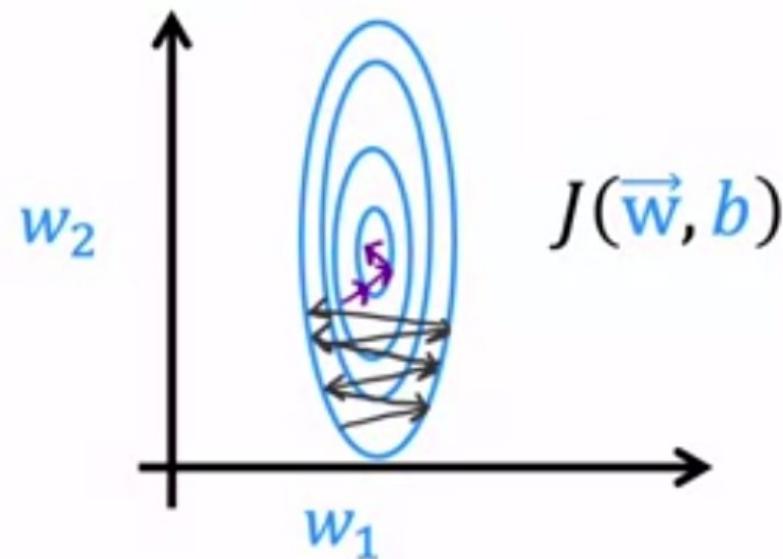
# MNIST Adam

but in codes this is  
how you implement it.

# Adam Algorithm Intuition



If  $w_j$  (or  $b$ ) keeps moving  
in same direction,  
increase  $\alpha_j$ .



If  $w_j$  (or  $b$ ) keeps oscillating,  
reduce  $\alpha_j$ .

you may learn more  
about the details

# 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

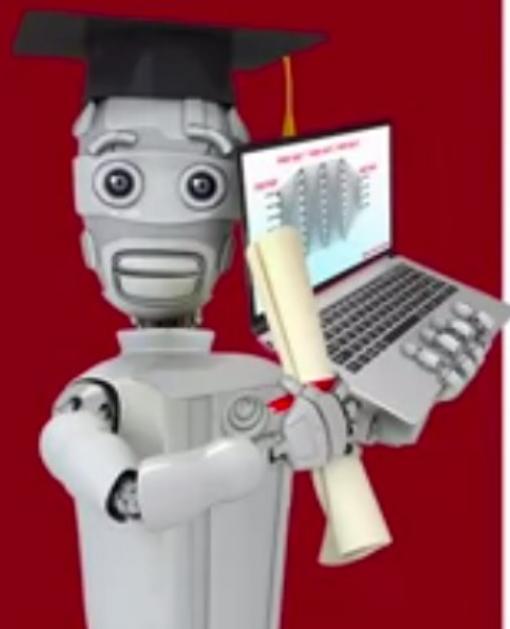
$$\alpha = 10^{-3} = 0.001$$

```
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)
```

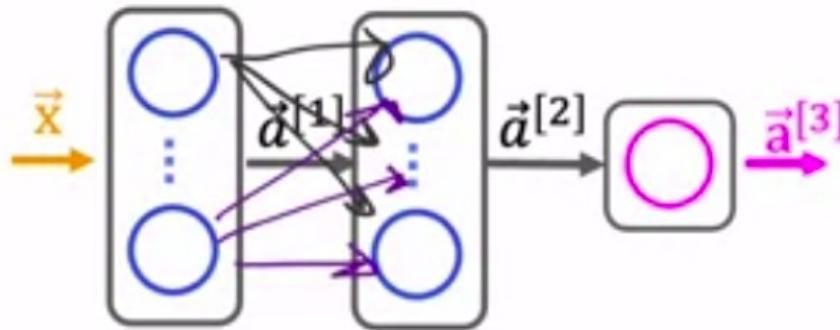
if you can get somewhat  
faster learning.



## Additional Neural Network Concepts

### Additional Layer Types

# Dense Layer

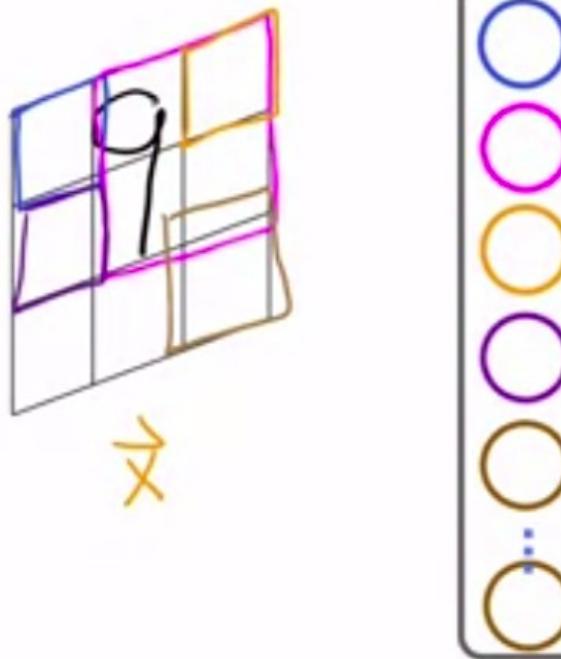


Each neuron output is a function of  
all the activation outputs of the previous layer.

$$\vec{a}_1^{[2]} = g \left( \vec{w}_1^{[2]} \cdot \vec{a}^{[1]} + b_1^{[2]} \right)$$

someone designing a neural network may  
choose to use a different type of layer.

# Convolutional Layer



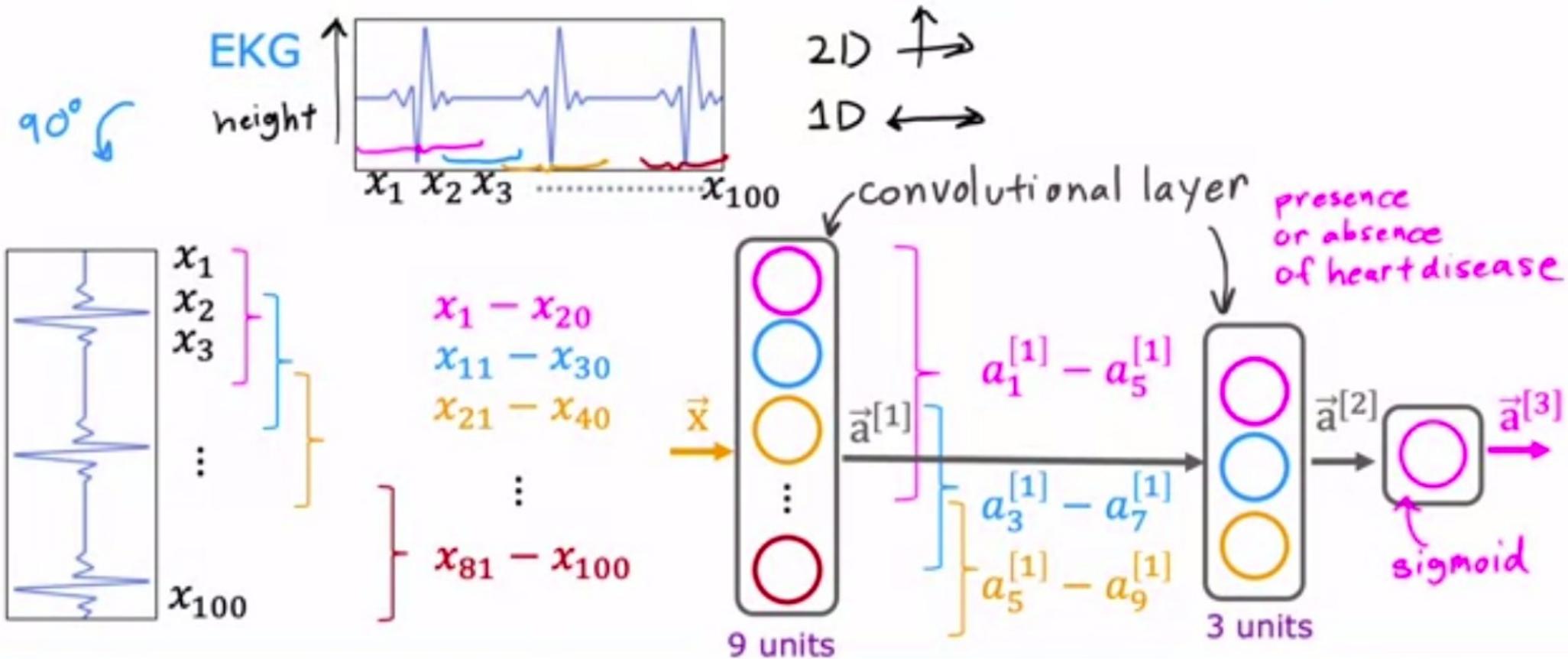
Each Neuron only looks at part of the previous layer's inputs.

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



[Back](#)

## Practice quiz: Additional Neural Network Concepts

Due Mar 26, 11:59 PM IST

Graded Quiz • 30 min

1.

1 / 1 point

# 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)
```

[Back](#)

## Practice quiz: Additional Neural Network Concepts

Graded Quiz • 30 min

Due Mar 26, 11:59 PM IST

```
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)
```

The Adam optimizer is the recommended optimizer for finding the optimal parameters of the model. How do you use the Adam optimizer in TensorFlow?

- The call to `model.compile()` will automatically pick the best optimizer, whether it is gradient descent, Adam or something else. So there's no need to pick an optimizer manually.
- When calling `model.compile`, set `optimizer=tf.keras.optimizers.Adam(learning_rate=1e-3)`.
- The call to `model.compile()` uses the Adam optimizer by default
- The Adam optimizer works only with Softmax outputs. So if a neural network has a Softmax output layer, TensorFlow will automatically pick the Adam optimizer.

A green checkmark icon indicating the answer is correct.

Correct. Set the optimizer to Adam.

[Back](#)

## Practice quiz: Additional Neural Network Concepts

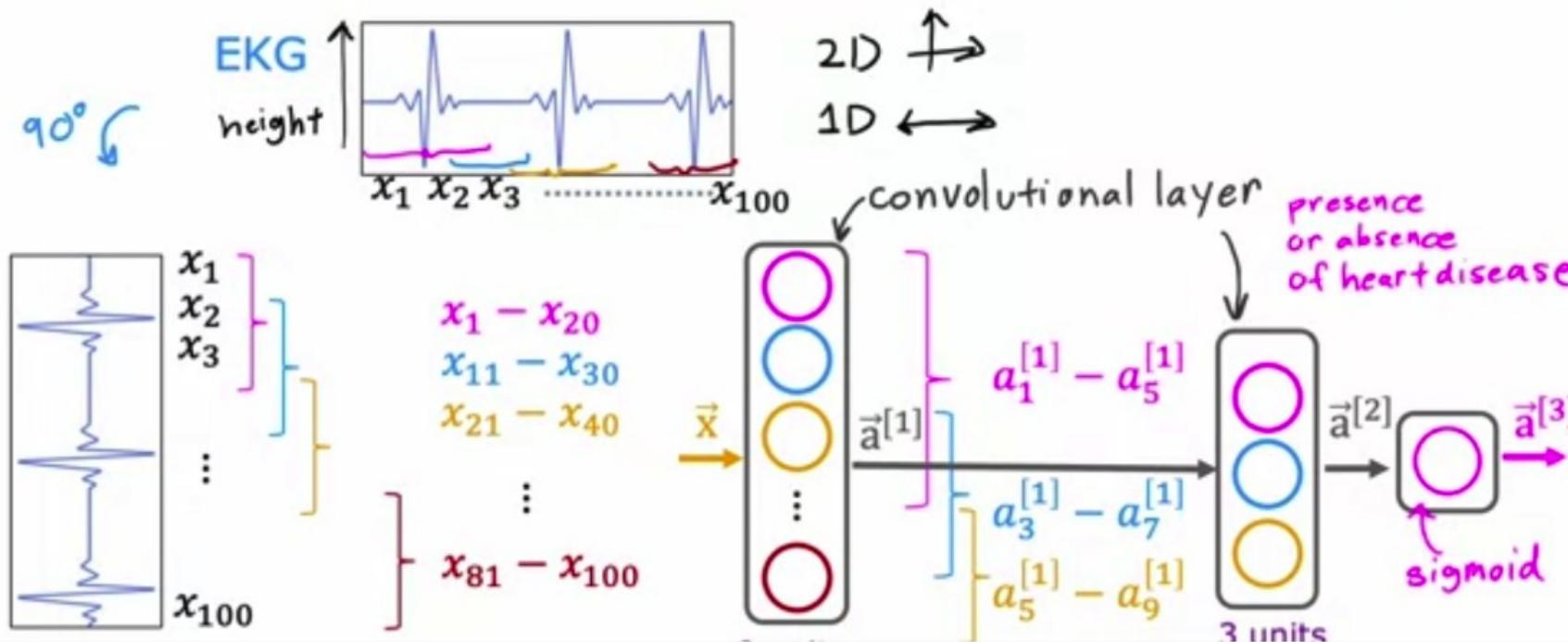
Graded Quiz • 30 min

Due Mar 26, 11:59 PM IST

2.

1 / 1 point

## Convolutional Neural Network

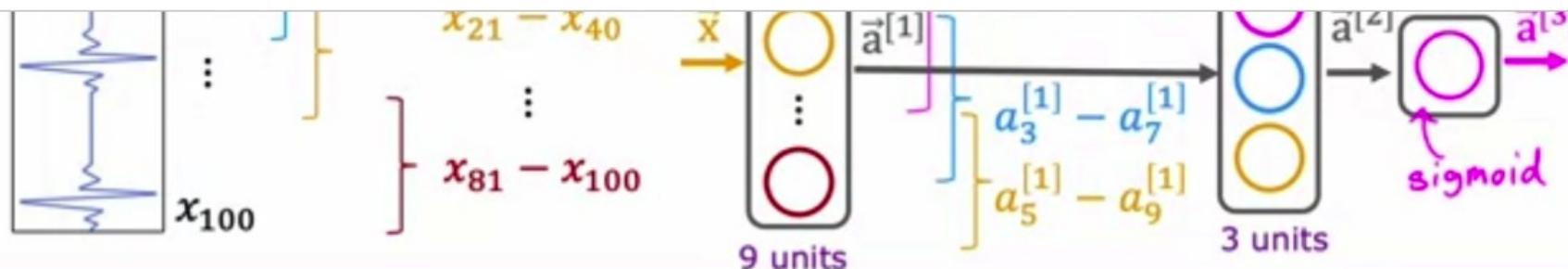


[Back](#)

## Practice quiz: Additional Neural Network Concepts

Graded Quiz • 30 min

Due Mar 26, 11:59 PM IST



The lecture covered a different layer type where each single neuron of the layer does not look at all the values of the input vector that is fed into that layer. What is this name of the layer type discussed in lecture?

- convolutional layer
- Image layer
- A fully connected layer
- 1D layer or 2D layer (depending on the input dimension)

A green checkmark icon indicating the answer is correct.

Correct. For a convolutional layer, each neuron takes as input a subset of the vector that is fed into that layer.



## Backprop Intuition (Optional)

**What is a derivative?**

# Derivative Example

Cost function

$$J(w) = w^2$$

Say  $w = 3$

$$J(w) = 3^2 = 9$$

$$\epsilon = 0.002$$

If we increase  $w$  by a tiny amount  $\epsilon = 0.001$  how does  $J(w)$  change?

$$w = 3 + \cancel{0.001} \overset{0.002}{\cancel{0.001}}$$

$$J(w) = w^2 = \cancel{9.006001}$$

$$\underbrace{\cancel{9.012} \overset{0.004}{\cancel{0.001}}}_{9.012}$$

$$\text{If } \underline{w} \uparrow \overset{0.002}{\cancel{0.001}} \quad \epsilon \leftarrow$$

$$J(w) \uparrow 6 \times \overset{0.002}{\cancel{0.001}} \quad 6 \times \epsilon \quad 6 \times 0.002 = 0.012$$

$$\frac{\partial}{\partial w} J(w) = 6 \quad \uparrow$$

# Informal Definition of Derivative

If  $w \uparrow \varepsilon$  causes  $J(w) \uparrow k \times \varepsilon$  then

$$\frac{\partial}{\partial w} J(w) = k$$

Gradient descent

repeat {

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

}

*learning rate*

If derivative is small, then this update step will make a small update to  $w_j$

If the derivative is large, then this update step will make a large update to  $w_j$

# More Derivative Examples

$w = 3$

$J(w) = w^2 = 9$

$w \uparrow 0.001$

$J(w) = J(3.001) = 9.006001$

$\frac{\partial}{\partial w} J(w) = 6$

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

$w = 2$

$J(w) = w^2 = 4$

$w \uparrow 0.001$

$J(w) = J(2.001) = 4.004001$

$\frac{\partial}{\partial w} J(w) = 4$

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

$J(w) = w^2$

↓ 0.006

$w = -3$

$J(w) = w^2 = 9$

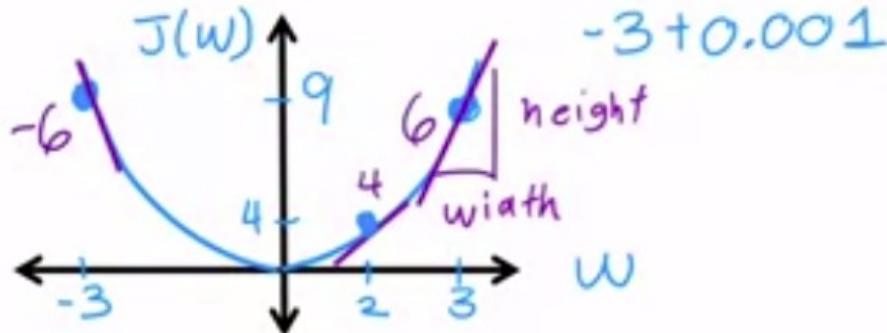
$w \uparrow 0.001$

$J(w) = J(-2.999) = 8.994001$

$\frac{\partial}{\partial w} J(w) = -6$

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

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



Calculus

$\frac{\partial}{\partial w} J(w) = 2w$

$w$

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

$3$

$2 \times 3 = 6$

$2$

$2 \times 2 = 4$

$-3$

$2 \times -3 = -6$

# Even More Derivative Examples

$$w = 2 \left[ \begin{array}{l} J(w) = w^2 = 4 \\ \frac{\partial}{\partial w} J(w) = 2w = 4 \\ w \uparrow \underbrace{0.001}_{\varepsilon} \\ J(w) = 4.004001 \\ J(w) \uparrow 4 \times \varepsilon \\ \\ J(w) = w^3 = 8 \\ \\ J(w) = w = 2 \\ \\ J(w) = \frac{1}{w} = \frac{1}{2} = 0.5 \end{array} \right]$$

# jupyter Using code to get derivatives (unsaved changes)



File Edit View Insert Cell Kernel Widgets Help Trusted | Python 3



Run



Code



Validate

In [1]: `import sympy`

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

In [27]: `J = w**3`|  
J

Out[27]:  $w^3$

In [28]: `dJ_dw = sympy.diff(J,w)`  
`dJ_dw`

Out[28]:  $3w^2$

In [29]: `dJ_dw.subs([(w,2)])`

Out[29]: 12

In [ ]:

# jupyter Using code to get derivatives (autosaved)



File Edit View Insert Cell Kernel Widgets Help Trusted Python 3



```
In [4]: J, w = sympy.symbols('J, w')
```

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

```
Out[3]: w2
```

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

```
Out[5]: 2w
```

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

```
Out[6]: 4
```

```
In [ ]:
```

# jupyter Using code to get derivatives (unsaved changes)



File Edit View Insert Cell Kernel Widgets Help Trusted Python 3



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

```
In [7]: J=w**3 #J = w**2  
J
```

```
Out[7]: w3
```

```
In [8]: dJ_dw = sympy.diff(J,w)  
dJ_dw
```

```
Out[8]: 3w2
```

```
In [9]: dJ_dw.subs([(w,2)])
```

```
Out[9]: 12
```

```
In [ ]:
```

# jupyter Using code to get derivatives (unsaved changes)



File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

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

```
In [10]: J=w #J=w**3 #J = w**2  
J
```

```
Out[10]: w
```

```
In [11]: dJ_dw = sympy.diff(J,w)  
dJ_dw
```

```
Out[11]: 1
```

```
In [12]: dJ_dw.subs([(w,2)])
```

```
Out[12]: 1
```

```
In [ ]:
```

# jupyter Using code to get derivatives (unsaved changes)



File Edit View Insert Cell Kernel Widgets Help Trusted Python 3



Code



Validate

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

```
In [13]: J=1/w #J=w #J=w**3 #J = w**2  
J
```

```
Out[13]:  $\frac{1}{w}$ 
```

```
In [14]: dJ_dw = sympy.diff(J,w)  
dJ_dw
```

```
Out[14]:  $-\frac{1}{w^2}$ 
```

```
In [15]: dJ_dw.subs([(w,2)])
```

```
Out[15]:  $-\frac{1}{4}$ 
```

```
In [ ]:
```

# jupyter Using code to get derivatives (unsaved changes)



File Edit View Insert Cell Kernel Widgets Help

Trusted Python 3



Code



Validate

In [1]: `import sympy`

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

In [13]: `J=1/w #J=w #J=w**3 #J = w**2`  
J

Out[13]:  $\frac{1}{w}$

In [14]: `dJ_dw = sympy.diff(J,w)`  
`dJ_dw`

Out[14]:  $-\frac{1}{w^2}$

In [15]: `dJ_dw.subs([(w,2)])`

Out[15]:  $-\frac{1}{4}$

In [ ]:

# Even More Derivative Examples

$w = 2$	$J(w) = w^2 = 4$	$\frac{\partial}{\partial w} J(w) = 2w = 4$	$w \uparrow \underbrace{0.001}_{\varepsilon}$	$J(w) = 4.004001$ $J(w) \uparrow 4 \times \varepsilon$
	$J(w) = w^3 = 8$	$\frac{\partial}{\partial w} J(w) = 3w^2 = 12$	$w \uparrow \varepsilon$	$J(w) = 8.012006$ $J(w) \uparrow 12 \times \varepsilon$
	$J(w) = w = 2$	$\frac{\partial}{\partial w} J(w) = 1$	$w \uparrow \varepsilon$	$J(w) = 2.001$ $J(w) \uparrow 1 \times \varepsilon$
	$J(w) = \frac{1}{w} = \frac{1}{2} = 0.5$	$\frac{\partial}{\partial w} J(w) = -\frac{1}{w^2} = -\frac{1}{4}$	$w \uparrow \varepsilon$ $w = \frac{1}{2.001}$	$-0.25 \times 0.001$ $0.5 - \underline{0.00025}$ $J(w) = 0.49975$ $J(w) \uparrow -\frac{1}{4} \times \varepsilon$
	$\frac{\partial}{\partial w} J(w)$	$w \uparrow \varepsilon$	$J(w) \uparrow k \times \varepsilon$	

# A note on derivative notation

If  $J(w)$  is a function of one variable ( $w$ ),

$$d \quad \frac{d}{dw} J(w)$$

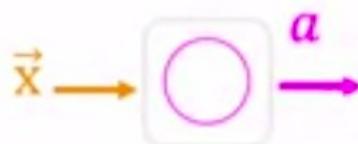
If  $J(w_1, w_2, \dots, w_n)$  is a function of more than one variable,

$$\partial \quad \bullet \frac{\partial}{\partial w_i} J(w_1, w_2, \dots, w_n) \quad \frac{\partial J}{\partial w_i} \quad \text{or} \quad \frac{\partial}{\partial w_i} J$$

"partial derivative"

notation used  
in these courses

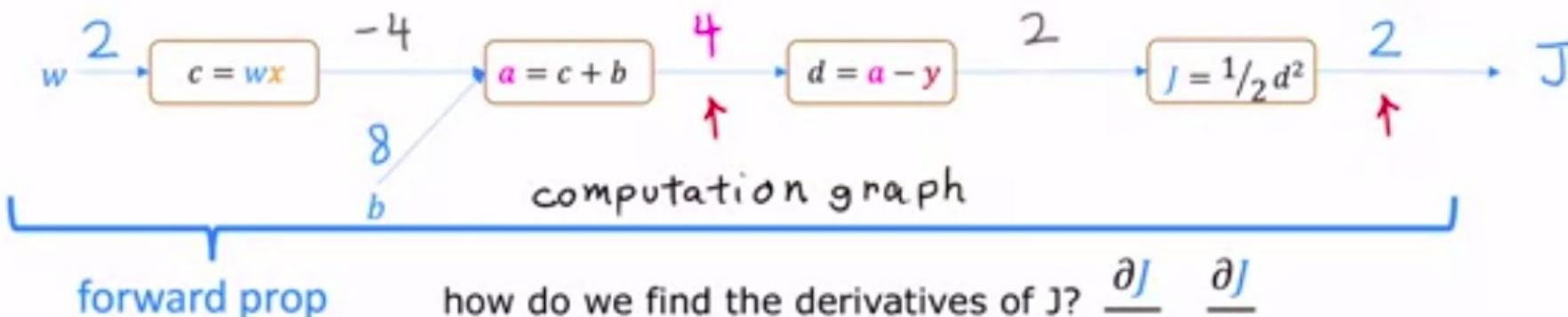
# Small Neural Network Example



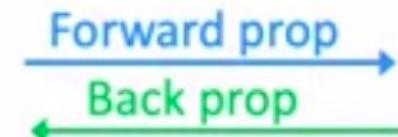
$$w=2 \quad b=8 \quad x=-2 \quad y=2$$

$$a = wx + b \quad \text{linear activation} \quad a = g(z) = z$$

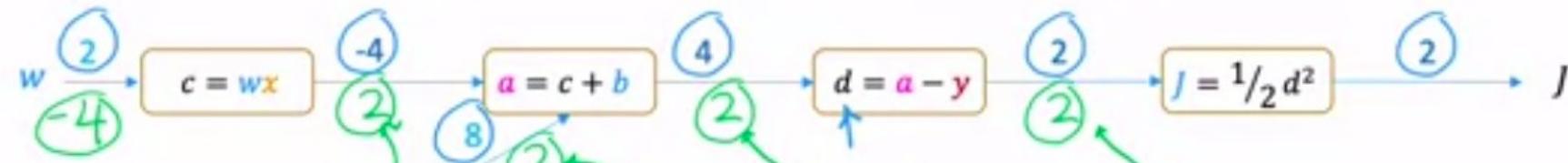
$$J(w,b) = \frac{1}{2} (a - y)^2$$



# Computing the Derivatives



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



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

$$\frac{\partial J}{\partial c} = 2$$

$$\frac{\partial J}{\partial b} = 2$$

$$\frac{\partial J}{\partial a} = 2$$

$$\frac{\partial J}{\partial d} = 2$$

$$w=2.001 \quad c=-4.002$$

$$w \uparrow 0.001 \quad c \downarrow 2 \times 0.001 \\ c \uparrow -2 \times 0.001$$

$$J \uparrow -4 \times 0.001$$

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

$$c \uparrow 0.001 \quad a \uparrow 0.001 \quad J \uparrow 0.002 \\ \frac{\partial J}{\partial c} = \frac{\partial a}{\partial c} \times \frac{\partial J}{\partial a} \\ 1 \quad 2$$

$$b \uparrow 0.001 \quad a \uparrow 0.001 \quad J \uparrow 0.002 \\ \frac{\partial J}{\partial b} = \frac{\partial a}{\partial b} \times \frac{\partial J}{\partial a} \\ 1 \quad 2$$

$$a \uparrow 0.001 \quad d \uparrow 0.001 \\ a=4.001 \quad d=2.001$$

$$J \uparrow 0.002$$

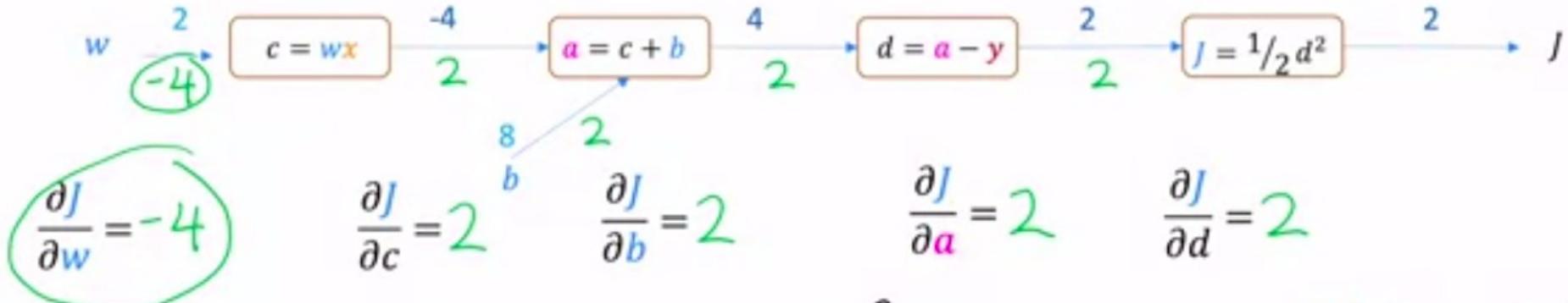
$$\frac{\partial J}{\partial a} = \frac{\partial d}{\partial a} \times \frac{\partial J}{\partial d} \\ 1 \quad 2$$

$$d \uparrow \underbrace{0.001}_{\epsilon} \quad J \uparrow \underbrace{0.002}_{2\epsilon}$$

chain rule  
(calculus)

# Computing the Derivatives

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



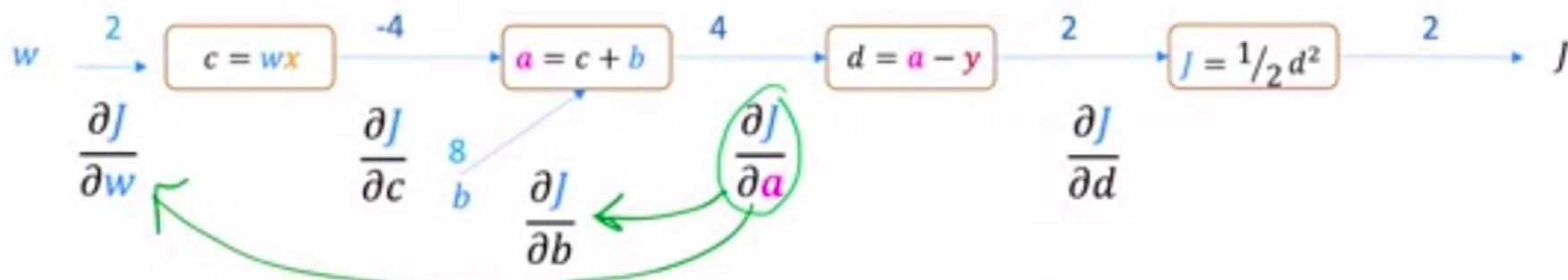
$$J = \frac{1}{2}((wx + b) - y)^2 = \frac{1}{2}((2 \times -2 + 8) - 2)^2 = 2$$

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

$J \downarrow 4 \times 0.001$   
 $J \uparrow -4 \times 0.001$

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

# Backprop is an efficient way to compute derivatives



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.

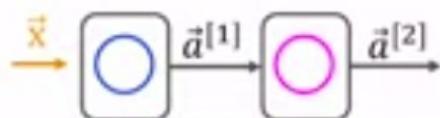
$N$	$P$	$N+P$	$N \times P$
10,000	100,000	$1.1 \times 10^5$	$10^9$

# Neural Network Example

$$x = 1 \quad y = 5$$

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

ReLU activation



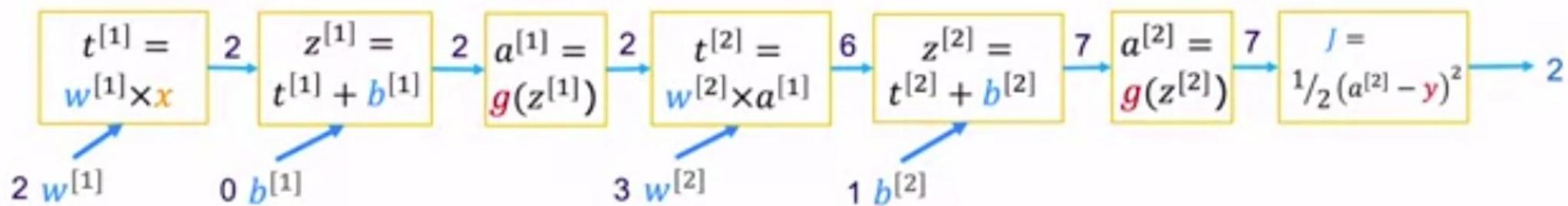
$$w^{[2]} = 3, b^{[2]} = 1$$

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

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

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

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

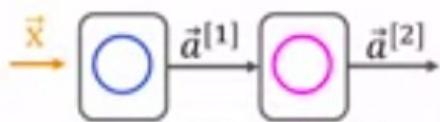


# Neural Network Example

$$x = 1 \quad y = 5$$

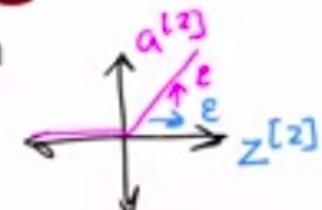
$$w^{[1]} = 2, b^{[1]} = 0$$

ReLU activation



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

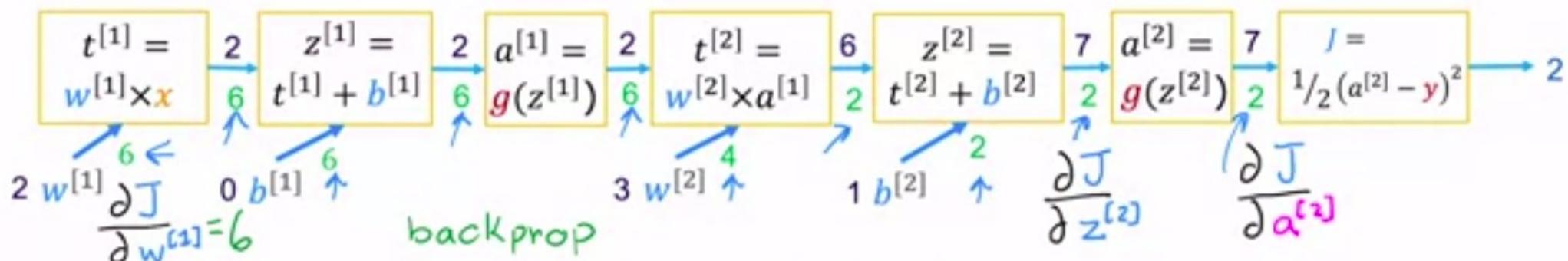
$$g(z) = \max(0, z)$$



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

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

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



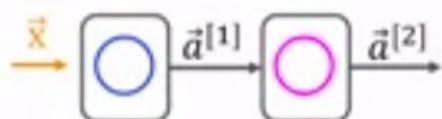
if  $w^{[2]} \uparrow \varepsilon$ , then  $J \uparrow 6 \times \varepsilon$  let's verify this!

# Neural Network Example

$$x = 1 \quad y = 5$$

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

ReLU activation



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

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

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

2.001

2.001

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

2.001

7.003

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

$$\frac{(2.003)^2}{2} = 2.006005$$

$w^{[1]} \uparrow 0.001$

$J \uparrow 6 \times 0.001$

$$\frac{\partial J}{\partial w^{[1]}} = 6$$

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

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

N nodes  $\square \rightarrow \square \rightarrow \square$

inefficient way

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

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

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

$N \times P$

efficient way (backprop)

$N + P$