# Copyright Notice

# Neural Network Training

TensorFlow implementation

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

①

②

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

③

epochs: number of steps
        in gradient descent

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

$\vec{x}$

1 unit

15 units

25 units

Given set of (x,y) examples

How to build and train this in code?

Andrew Ng

# Model Training Steps

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

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

② specify loss and cost

$$L\left(f_{\overrightarrow{\mathbf{w}},b}(\overrightarrow{\mathbf{x}}), y\right) \quad \text{1 example}$$

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

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

## 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

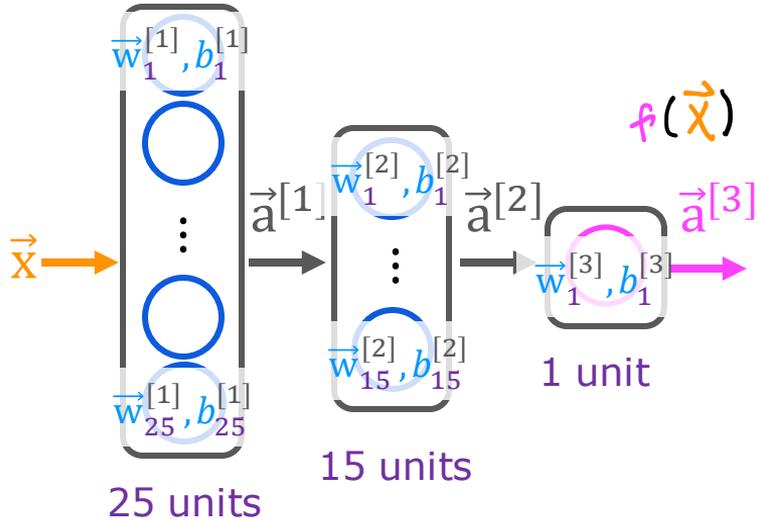```
model.compile(
loss=BinaryCrossentropy())



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

# 1. Create the model

define the model

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

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



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

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

$f(\vec{x})$

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

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

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

$\vec{x}$

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

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

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

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')
                    ])
```

# 2. Loss and cost functions

Mnist digit classification problem

binary classification

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

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

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

$f_{W,B}(\vec{x})$

Compare prediction vs. target

logistic loss
also known as binary cross entropy

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

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

regression
(predicting numbers and not categories)

Mean squared error

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

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

# 3. Gradient descent



$J(w)$

minimum

$w$

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

repeat {

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

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

}

Compute derivatives
for gradient descent
using "back propagation"

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

# Neural network libraries
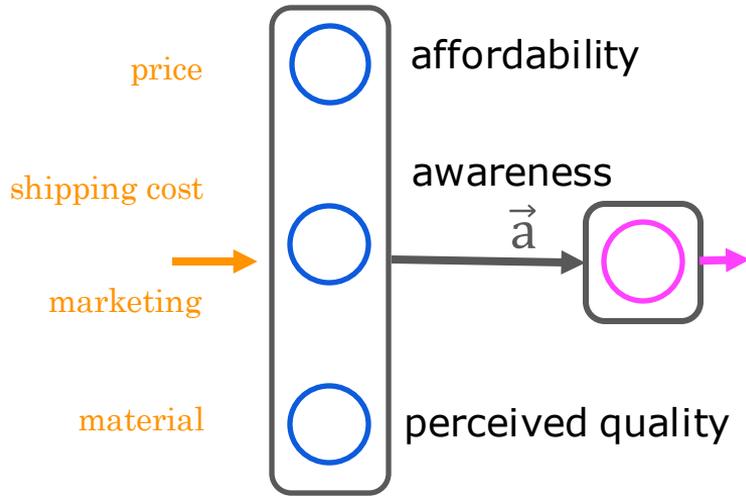
Use code libraries instead of coding "from scratch"



Good to understand the implementation

# Activation Functions

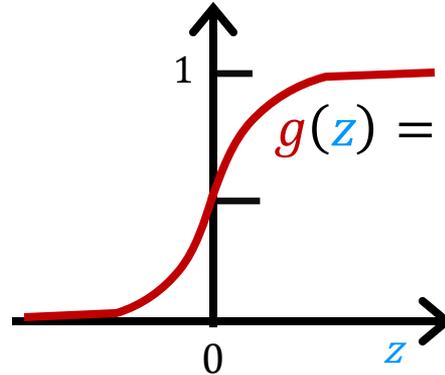Alternatives to the sigmoid activation

# Demand Prediction Example



price

shipping cost

marketing

material

affordability

awareness
$\vec{a}$

perceived quality

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

Sigmoid

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

$0 < g(z) < 1$
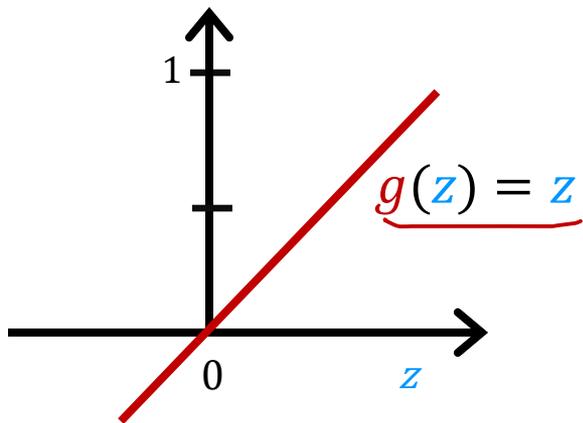
ReLU   Rectified Linear Unit

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

if z < 0
g(z) is 0

if z ≥ 0
g(z) is z

Andrew Ng

# Examples of Activation Functions

"No activation function"

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

### Linear activation function



$$g(z) = z$$

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

### Sigmoid



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

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

### ReLU  Rectified Linear Unit



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

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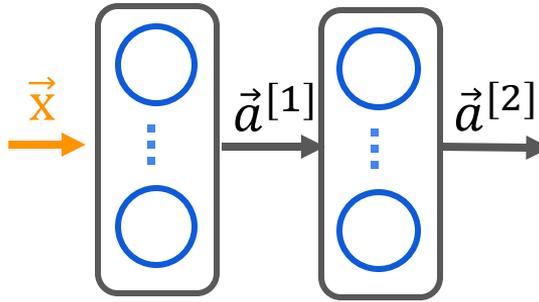g(z) is 0   g(z) is z

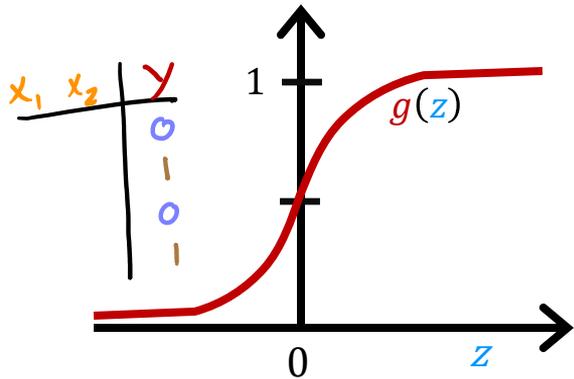Andrew Ng

# Output Layer



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

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

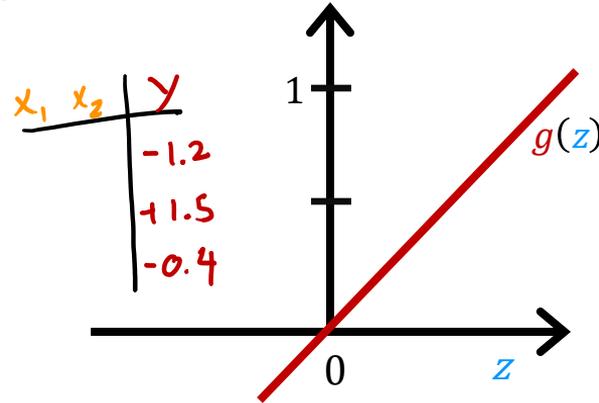Choosing $g(z)$ for output layer?

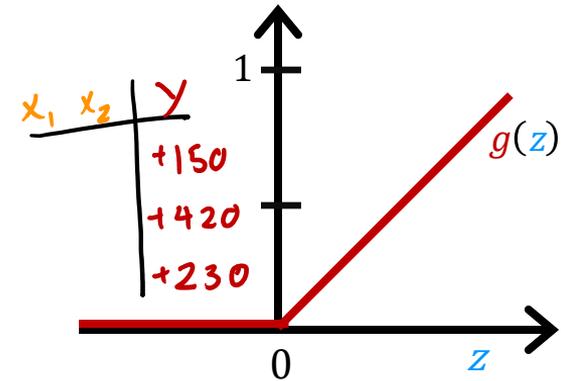## Binary classification
Sigmoid
y=0/1
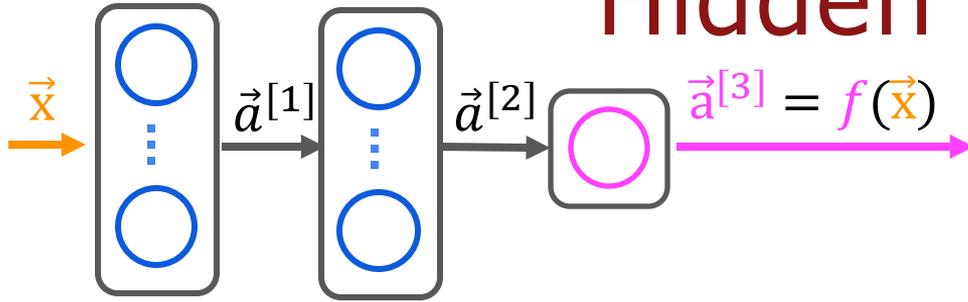


## Regression
Linear activation function
y = +/-



## Regression
ReLU
Y = 0 or +



Andrew Ng

# Hidden Layer



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

Choosing $g(z)$ for hidden layer

**Sigmoid**

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

flat

flat

slower

$J(\mathbf{W}, \mathbf{B})$

$$\frac{\partial}{\partial w} J(\mathbf{W}, \mathbf{B}) \approx 0$$

when $g(z)$ is flat

$w$

most common choice

**ReLU**

faster

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

NOT flat

flat

faster learning

Andrew Ng

# Choosing Activation Summary



ReLU hidden layers

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'

Andrew Ng

# Why do we need activation functions?

price
shipping cost
marketing
material

$\vec{x}$

affordability

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

perceived quality

top seller?
$\vec{a}^{[2]}$

if all $g(z)$ are linear…

price

shipping cost

$\vec{x}$

marketing

material

$\vec{w}, b$

$f(\vec{x}) = \vec{w} \cdot \vec{x} + b$

$g(z)$

1

0

$z$

…no different than linear regression

# Linear Example

one feature

$x$

$[2]$

$w$ is scalar



'a' is scalar

$g(z) = z$

$$a^{[1]} = w_1^{[1]} x + b_1^{[1]}$$

$$a^{[2]} = w_1^{[2]} a^{[1]} + b_1^{[2]}$$

$$= w_1^{[2]} \left( w_1^{[1]} x + b_1^{[1]} \right) + b_1^{[2]}$$

$$\vec{a}^{[2]} = \left( \vec{w}_1^{[2]} \vec{w}_1^{[1]} \right) x + \underbrace{w_1^{[2]} b_1^{[1]} + b_1^{[2]}}_{b}$$

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

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

Andrew Ng

# 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

$$\vec{a}^{[4]} = \frac{1}{1+e^{-(\vec{w}_1^{[4]} \cdot \vec{a}^{[3]} + b_1^{[4]})}}$$
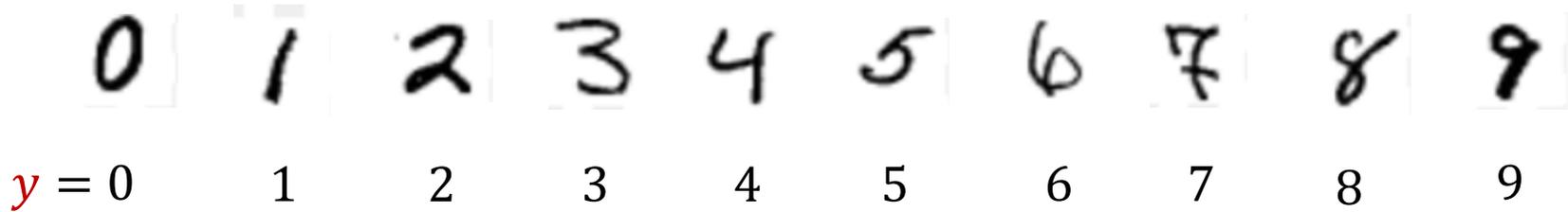
output activation is sigmoid
(hidden layers still linear)
↳ equivalent to logistic regression

Don't use linear activations in hidden layers (use ReLU)

# Multiclass Classification

---

Multiclass

# MNIST example



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

$x \to 7$      $y = 7$

multiclass classification problem:
target $y$ can take on more than two possible values

Andrew Ng

# Multiclass classification example



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

class 1

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

not 1

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

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

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

Andrew Ng

# Multiclass Classification

Softmax

## 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) $y = 1, 2, 3, ..., N$

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

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

Andrew Ng

# Cost

## Logistic regression

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

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

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

$a_2$

$$loss = -y \underbrace{\log a_1} - (1 - y) \underbrace{\log(1 - a_1)}$$

if $y=1$      if $y=0$

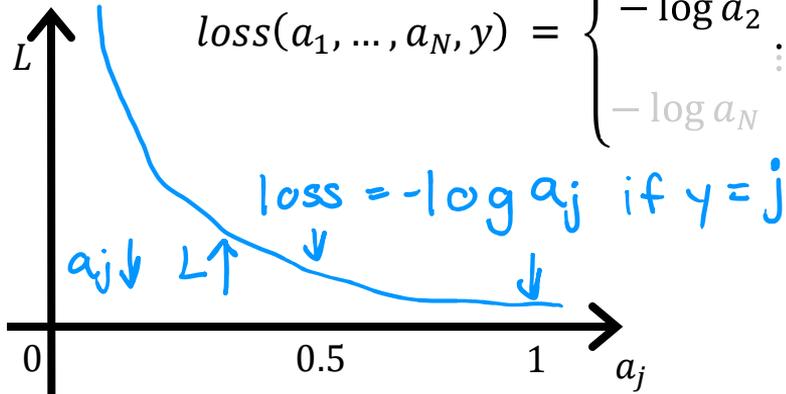$$J(\vec{w}, b) = \text{average loss}$$

## Softmax regression

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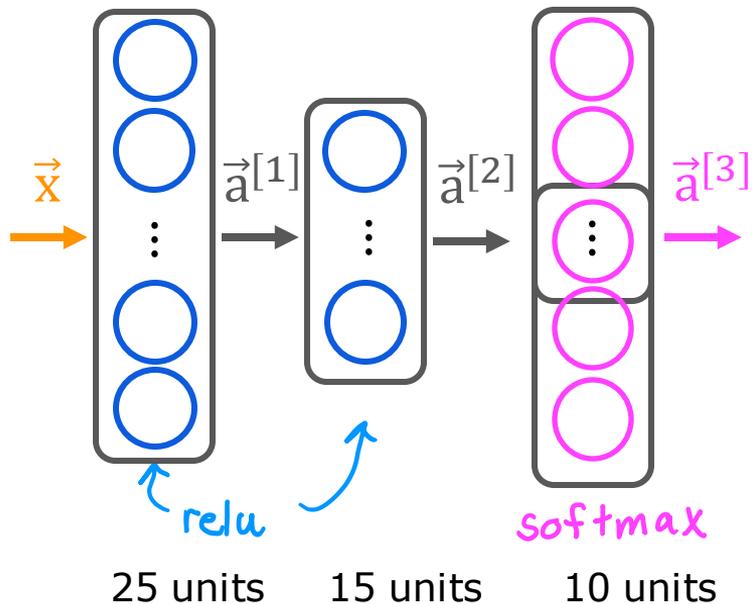
$$\vdots$$

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

*Crossentropy loss*

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

$L$

loss $= -\log a_j$ if $y = j$

$a_j \downarrow$  $L \uparrow$

0      0.5      1      $a_j$

Andrew Ng

# Neural Network with Softmax output



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

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

$$\vdots$$

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

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

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

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

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

```
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

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

```
from tensorflow.keras.losses import
    SparseCategoricalCrossentropy

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

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

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

Note: better (recommended) version later.

Andrew Ng

# Multi-label Classification



$\vec{x}$

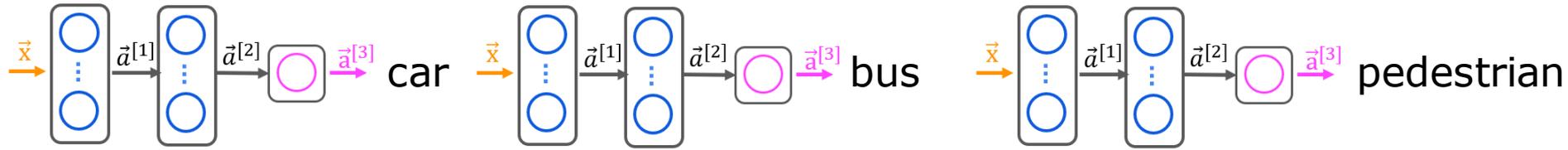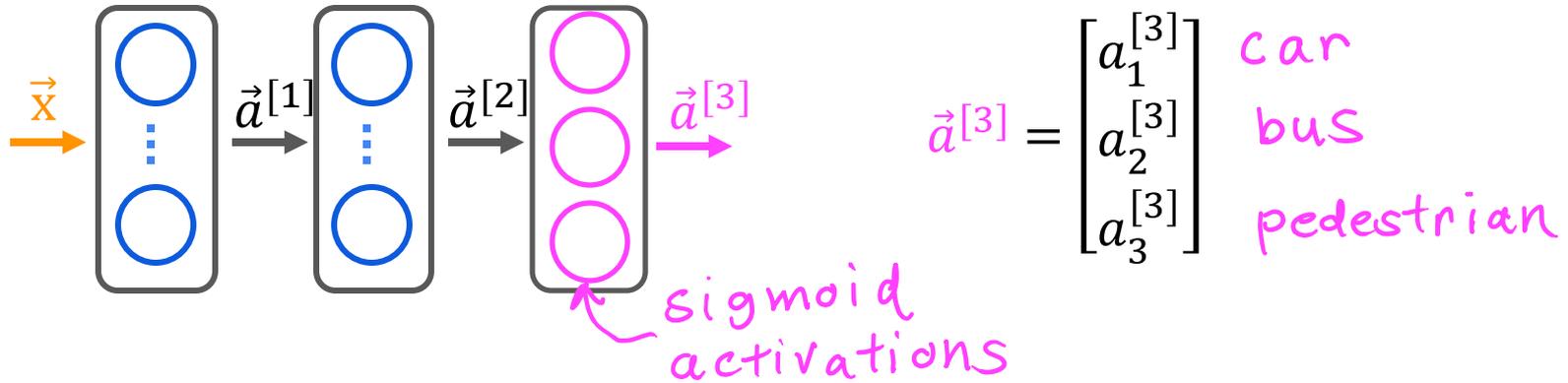Is there a car?      yes     $y = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$     no     $y = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$     yes     $y = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$
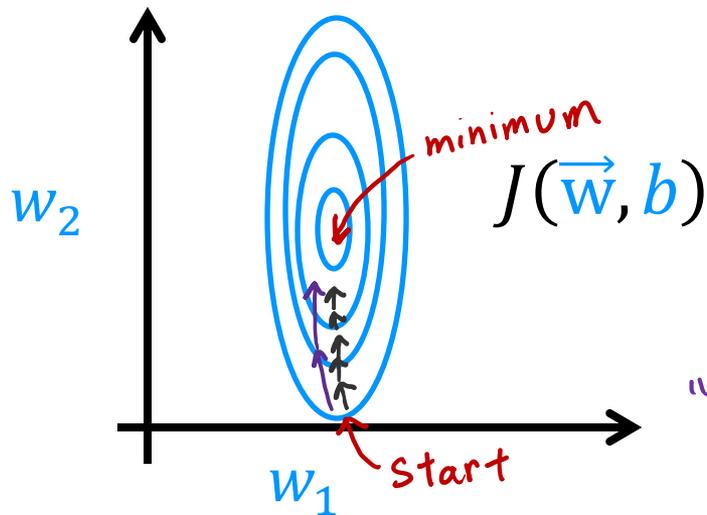
Is there a bus?     no     no     yes

Is there a pedestrian     yes     yes     no

Andrew Ng

# Multiple classes



Alternatively, train one neural network with three outputs

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

sigmoid activations

Andrew Ng

# Gradient Descent

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

learning rate



$w_2$

minimum

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

start

$w_1$

"Adam" algorithm

Go faster – increase $\alpha$

$w_2$

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

$w_1$

Go slower – decrease $\alpha$

Andrew Ng

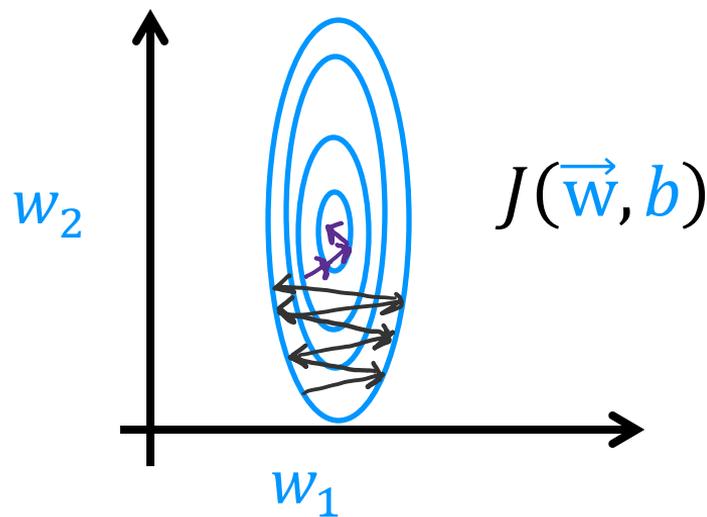# Adam Algorithm Intuition

Adam: <u>A</u>daptive <u>M</u>oment estimation     <span style="color:red">not just one $\alpha$</span>

$$w_1 = w_1 - \alpha_1 \frac{\partial}{\partial w_1} J(\vec{w}, b)$$

$$\vdots$$

$$w_{10} = w_{10} - \alpha_{10} \frac{\partial}{\partial w_{10}} J(\vec{w}, b)$$

$$b = b - \alpha_{11} \frac{\partial}{\partial b} J(\vec{w}, b)$$

# Adam Algorithm Intuition



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

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

Andrew Ng

# 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)
```

Andrew Ng