# 第一周 - Andrew Ng.

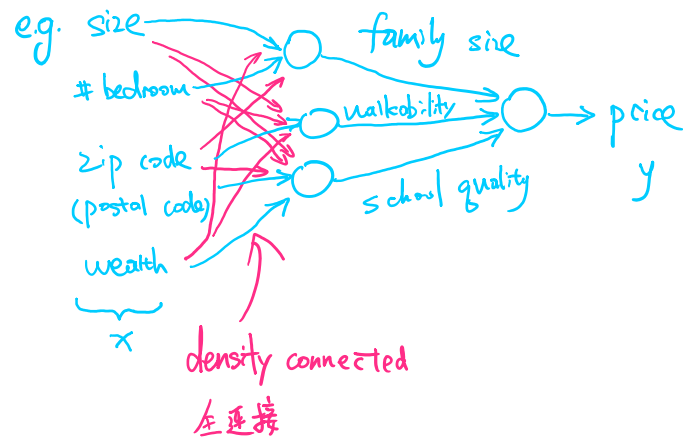==Welcome to the deep learning specialization==

1. Neural networks and deep learning
2. Improve deep neural networks: hyperparameter tuning, regularization, and optimization
3. Structure your machine learning project
4. Convolutional neural networks CNN 卷积神经网络
5. Natural language processing: building sequence models (Recurrent neural network RNN Long short term memory model LSTM) ⇒ NLP

rectified linear regression (ReLu)
线性整流函数     max(0, y)

size → O → price
x    "neuron"  y

==Supervised learning with neural networks==

| input (x) | output (y) | application | |
|---|---|---|---|
| home features | price | real estate | } standard NN |
| ad, user info | click on ad? (0/1) | online ads | |
| (DL) image | object (1,....,1000) | photo tagging | CNN 卷积神经网络 |
| audio | text transcript | Speech recognition | } RNN (sequence data) |
| English | Chinese | Machine translation | |
| image, radar | position of car | autonomous driving | customer hybrid |

e.g. size
# bedroom
zip code (postal code)
wealth
→ x

family size
walkability
school quality
→ price y

density connected
全连接

Structured data and unstructured data
size. #bedrooms. price        audio. image

Scale drives deep learning progress

performance
Large NN
medium NN
small NN
traditional learning algo (SVM, logistic regression)
amount of data

Why is deep learning taking off?

region
sigmoid function        ReLu
gradient descent ≈ 0.
algo with be slow

# 第二周

## Logistic regression as a neural network

### Binary classification

RGB. if a image is $64 \times 64$ pixel, $\rightarrow 3 \uparrow 64 \times 64$ matrix

unroll all pixel intensity: $X = \begin{cases} \text{:} \} \text{ red} \\ \text{:} \} \text{ green} \\ \text{:} \} \text{ blue} \end{cases}$ $\sim 64 \times 64 \times 3$ rows

$x \in \mathbb{R}^{n_x}$, $y \in \{0,1\}$    $n_x = 12288$, $x \rightarrow y$

$m \rightarrow$ tr examples $\{(x^{(1)}, y^{(1)}), \cdots (x^{(m)}, y^{(m)})\}$

$m = m_{\text{train}}$. $m_{\text{test}} = \#$ test examples

$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(1)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix} \updownarrow$ n rows (features)

$\underbrace{\phantom{XXXXX}}_{\text{m}}$
columns

**GOOD!**
in this way $n \times m$
easier for Neural Network

$X \in \mathbb{R}^{n_x \times m}$. $X.\text{shape} = (n_x, m)$

$Y = [y^{(1)} \cdots y^{(m)}]$    $Y \in \mathbb{R}^{1 \times m}$, $Y.\text{shape} = (1, m)$

### Logistic regression

Given $x$, want $\hat{y}$ (y两估计值) $\hat{y} = P(y=1|x)$, $0 \leq \hat{y} \leq 1$

$x \in \mathbb{R}^{n_x}$

parameter of logistic regression: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

Output $\hat{y} = \sigma(w^T x + b)$    $\sigma(z) = \frac{1}{1 + e^{-z}}$  if $z$ large $\sigma(z) \approx 1$
                                                            small $\sigma(z) = 0$

$x_0 = 1$. $x \in \mathbb{R}^{n_x+1}$. $\hat{y} = \sigma(\theta^T x)$    $\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \begin{matrix} \leftarrow b \\ \} w \end{matrix}$

#### Logistic regression cost function

$\hat{y} = \sigma(w^T x + b)$, where $\sigma(z) = \frac{1}{1 + e^{-z}}$. $z^{(i)} = w^T x^{(i)} + b$

Given $\{(x^{(1)}, y^{(1)}) \cdots (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$

#### Loss (error) function

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

(得了对此, 不面足 machine learning 学到的).
一样有意义, 以是凸起出试不同



Cost function
Logistic regression:

Neural network:

if $y=1$. $\mathcal{L}(\hat{y}, y) = -\log(\hat{y})$. goal: want $\hat{y}$ large
if $y=0$. $\mathcal{L}(\hat{y}, y) = -\log(1-\hat{y})$. goal: want $\hat{y}$ small

Cost function $J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^{m} (y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)}))$

The loss function computes the error for a single training examples, and the cost function is the average of the loss function of the

## Python and vectorization

### Vectorization

$z = w^T x + b$    $w = \begin{bmatrix} \vdots \end{bmatrix}$ $x = \begin{bmatrix} \vdots \end{bmatrix}$ $w \in \mathbb{R}^{n_x}$, $x \in \mathbb{R}^{n_x}$

vectorized:

$z = \underbrace{np.\text{dot}(w,x)}_{w^T x} + b$

### More vectorization examples

Vectors and matrix valued functions

```
Import numpy as np
U = np.exp(v)
U = np.log(v)
U = np.abs(v)
U = np.maximum(v,0)
```

$J = 0$; $dw_1 = 0$; $dw_2 = 0$; $db = 0$     $dw = np.\text{zeros}(n_x, 1)$

For $i = 1$ to $m$:  the first loop
   $z^{(i)} = w^T x^{(i)} + b$
   $a^{(i)} = \sigma(z^{(i)})$
   $J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$
   $dz^{(i)} = a^{(i)} - y^{(i)}$              $dw += x^{(i)} \cdot dz^{(i)}$

for $j=1 \cdots n_x$ second loop    $dw_1^{(i)} += x_1^{(i)} dz^{(i)}$    $n = 2 (2 \text{ feature})$    $w_1 := w_1 - \alpha dw_1$
$dw_j += x_j^{(i)} dz^{(i)}$       $dw_2^{(i)} += x_2^{(i)} dz^{(i)}$                                $w_2 := w_2 - \alpha dw_2$
                                  $db += dz^{(i)}$                                                  $b := b - \alpha db$

$J /= m$
$dw_1 /= m$; $dw_2 /= m$    $db /= m$

### Vectorizing logistic regression

$z^{(1)} = w^T x^{(1)} + b$       $z^{(m)} = w^T x^{(m)} + b$
$a^{(1)} = \sigma(z^{(1)})$  $\cdots$  $a^{(m)} = \sigma(z^{(m)})$

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

$z = [z^{(1)}; \cdots z^{(m)}] = w^T X + [b, \cdots, b]$  $(1, m)$

in python: $np.\text{dot}(w.T, X) + b$
   $\underset{\text{a } 1 \times 1 \text{ real number}}{}$
expand $b$ to $1 \times m$.
call broadcasting

$A = [a^{(1)}, \cdots, a^{(m)}] = \sigma(z)$

### Vectorizing logistic regression's gradient output

$dz^{(1)} = a^{(1)} - y^{(1)}$
$dz = [dz^{(1)} \cdots dz^{(m)}]$  $(1 \times m)$
$A = [a^{(1)} \cdots a^{(m)}]$ $Y = [y^{(1)} \cdots y^{(m)}]$
$dz = A - Y$
$dw += x^{(i)} dz^{(i)}$      $db += dz^{(i)}$
$dw /= m$                     $db /= m$

$db = \frac{1}{m} \sum_{i=1}^{m} dz^{(i)}$
   $= \frac{1}{m} np.\text{sum}(dz)$
$dw = \frac{1}{m} X dz^T$
   $= \frac{1}{m} [x^{(1)} dz^{(1)} + \cdots + x^{(m)} dz^{(m)}]$
                    $n \times 1$

The loss function computes the error for a single training examples, and the cost function is the (average) of the loss function of the entire training set.

## Gradient descent

$$\hat{y} = \sigma(w^T x + b) \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^{m} (y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)}))$$

$\cup$ convex function

to find $w, b$ that minimize $J(w,b)$
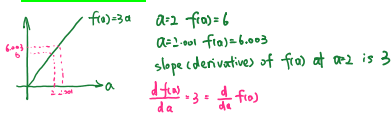


repeat {
$$w := w - \alpha \frac{dJ(w)}{dw}$$
} $w := w - \alpha\, dw$

$J(w,b)$ ① $w := w - \alpha \frac{dJ(w,b)}{dw}$   derivative

② $b := b - \alpha \frac{dJ(w,b)}{db}$

$\partial$: partial derivative
$d$: $J$ (cost function)

① update $dw$  ② update $db$

## Derivatives



$f(a) = 3a$   $a=2$  $f(a)=6$
$a=2.001$  $f(a)=6.003$
slope (derivative) of $f(a)$ at $a=2$ is 3

$$\frac{d f(a)}{da} = 3 = \frac{d}{da} f(a)$$

## More derivative examples



$f(a) = a^2$   $a=2$, $f(a)=4$
$a=2.001$, $f(a)=4.004$
slope (derivative) of $f(a)$ at $a=2$ is 4
while $a=5$  $f(a)=25$
$a=5.001$  $f(a)=25.010$
$\frac{d}{da} f(a) = 10$ when $a=5$

$$(\ln(a))' = \frac{1}{a}$$

## Computation graph

$$J(a,b,c) = 3(a + bc)$$

$u = bc$
$v = a+u$
$J = 3v$



## Derivative with a computation graph



$J'$ called $\frac{d\ \text{Final Output Var}}{d\ var}$ ⇒ "$dvar$" means the derivative of a final output variables with respect to various intermediate quantities

$$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da} = 3 \cdot 1 = 3$$

$\frac{dJ}{da}$ (abbr.) → $da$   $\frac{dJ}{dv}$ (abbr.) → $dv$



## Logistic regression gradient descent

$z = w^T x + b$
$\hat{y} = a = \sigma(z)$
$\mathcal{L}(a,y) = -(y \log(a) + (1-y) \log(1-a))$



$dw_1 = x_1 \cdot dz$
$dw_2 = x_2 \cdot dz$
$db = dz$

$dz = a - y$

$\frac{da}{} = -\frac{y}{a} + \frac{1-y}{1-a}$

$\frac{d\mathcal{L}}{da} \cdot \frac{da}{dz}$ denote by $dz$

since $(\sigma(z))' = \sigma(z) \cdot (1-\sigma(z))$, so $dz = (-\frac{y}{a} + \frac{1-y}{1-a}) a(1-a)$
$= -y + ya + a - ay$
$= a - y$

---

$dw += x^{(i)} dz^{(i)}$    $db += dz^{(i)}$
$dw\ /= m$    $db\ /= m$

## Implementing logistic regression

$Z = w^T x + b$
$\quad = np.dot(w.T, x) + b$
$A = \sigma(Z)$
$dZ = A - Y$
$dw = \frac{1}{m} X dZ^T$
$db = \frac{1}{m} np.sum(dZ)$
$w := w - \alpha\, dw$
$b := b - \alpha\, db$

## Broadcasting in python

```
import numpy as np
A = np.array([[56.0,0.0,4.4,68.0],
              [1.2,104.0,52.0,8.0],
              [1.8,135.0,99.0,0.9]])
cal = A.sum(axis=0)
percentage = 100*A/cal.reshape(1,4)
```

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix} \quad \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

### General principle

$(m,n)$ matrix   $\begin{matrix} + \\ - \\ * \\ / \end{matrix}$   $(1,n) \rightsquigarrow (m,n)$
  $(m,1) \rightsquigarrow (m,n)$    matlab: similar function: bsxfun

## A note on python / numpy vectors

$a = np.random.randn(5)$
$a.shape = (5,)$  "rank 1 array" 被为1的秩阵阵

## Explanation of logistic regression cost function

$\hat{y} = \sigma(w^T x + b)$  where  $\sigma(z) = \frac{1}{1+e^{-z}}$

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

if  $y=1$:  $P(y|x) = \hat{y}$ ;  if $y=0$  $P(y|x) = 1-\hat{y}$

$P(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$  两边取对数  $\log P(y|x) = y\log\hat{y} + (1-y)\log(1-\hat{y})$

if $y=1$, $P(y|x) = \hat{y}$    因为 $\log P(y|x)$ 是 ↑, 所以 $-\log P(y|x)$ 是 ↓

$y=0$, $P(y|x) = 1-\hat{y}$

### Cost on m examples

$P(\text{labels in tr set}) = \log \prod_{i=1}^{m} P(y^{(i)}|x^{(i)})$

$\log P(\cdots) = \sum_{i=1}^{m} \log P(y^{(i)}|x^{(i)})$    maximum likelihood estimation
  $-\mathcal{L}(\hat{y}^{(i)}, y^{(i)})$    最大似然估计

$= -\sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(a^{(i)}, y) \qquad (x^{(i)}, y^{(i)})$$

$$\rightarrow a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b) \qquad dw_1^{(i)} \cdots dw_2^{(i)} \cdots db^{(i)}$$

$$\frac{\partial}{\partial w_1} J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \underbrace{\frac{\partial}{\partial w_1} \cdot \mathcal{L}(a^{(i)}, y^{(i)})}_{dw_1^{(i)} \rightarrow (x^{(i)}, y^{(i)})}$$

$$J=0; \; dw_1 = 0; \; dw_2 = 0; \; db = 0 \qquad \text{假设只有两个变量 } w_1, w_2$$

For i = 1 to m : the first loop

$$z^{(i)} = w^T x^{(i)} + b$$

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

$$J \mathrel{+}= -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$$

second loop
$dw_1$
$\vdots$
$dw_n$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1^{(i)} \mathrel{+}= x_1^{(i)} \cdot dz^{(i)} \qquad n=2 \; (2 \text{ feature})$$

$$dw_2^{(i)} \mathrel{+}= x_2^{(i)} \cdot dz^{(i)}$$

$$db \mathrel{+}= dz^{(i)}$$

$$w_1 := w_1 - \alpha \, dw_1$$
$$w_2 := w_2 - \alpha \, dw_2$$
$$b := b - \alpha \, db$$

$$J \mathrel{/}= m$$

$$dw_1 \mathrel{/}= m \qquad dw_2 \mathrel{/}= m$$

# 第三周

## Shallow neural network

### Neural networks overview

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

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

### Neural network representation

input layer　hidden layer　output layer

$$eg: a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}$$

"2 layer NN"

$$a^{[0]} = X \qquad a^{[1]}$$

### Computing a neural network's output

$$z = W^T X + b$$
$$a = \sigma(z)$$

$$\boxed{W^T X + b \mid \sigma(z)} \quad a = \hat{y}$$

$$z_1^{[1]} = W_1^{[1]T} X + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$
$$\vdots$$
$$z_4^{[1]} = W_4^{[1]T} X + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

$$z^{[1]} = \begin{bmatrix} - W_1^{[1]T} - \\ - W_2^{[1]T} - \\ - W_3^{[1]T} - \\ - W_4^{[1]T} - \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} W_1^{[1]T}x + b_1^{[1]} \\ W_2^{[1]T}x + b_2^{[1]} \\ W_3^{[1]T}x + b_3^{[1]} \\ W_4^{[1]T}x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

Given input $x$:

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

### Vectoring across multiple examples

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

$$X = \begin{bmatrix} x^{(1)} & \cdots & x^{(m)} \end{bmatrix}$$
$$X.shape = (n_x, m)$$

$$Z^{[1]} = \begin{bmatrix} z^{[1](1)} & \cdots & z^{[1](m)} \end{bmatrix} \qquad A^{[1]} = \begin{bmatrix} a^{[1](1)} & \cdots & a^{[1](m)} \end{bmatrix} \quad \text{hidden units number on the 1st layer}$$

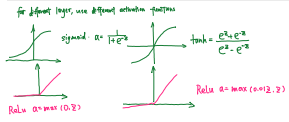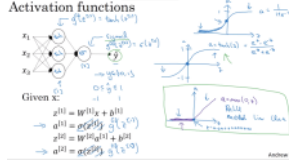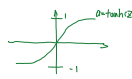### Explanation for vectorised implementation

Justification for vectorized implementation

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

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

## Activation functions

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

Activation functions

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

for different layers, use different activation functions

$$\text{sigmoid} \quad a = \frac{1}{1+e^{-z}} \qquad \tanh = \frac{e^z + e^{-z}}{e^z - e^{-z}}$$

Relu $a = \max(0, z)$ 　　Relu $a = \max(0.01z, z)$

### Why do you need non-linear activation functions?

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

### Derivatives of activation functions

→ Sigmoid
$$g(z) = \frac{1}{1+e^{-z}}$$
$$g'(z) = \frac{d}{dz} g(z) = g(z)(1 - g(z))$$
$z = 10, \ g(z) \approx 1 \quad z = -10, \ g(z) \approx 0$
$g(z) \approx 0 \qquad g'(z) \approx 0$
$g(z) = \frac{1}{2}, \ g'(z) = \frac{1}{4}$

→ Tanh
$$g(z) = \tanh(z)$$
$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$
$$g'(z) = \frac{d}{dz} g(z) = 1 - (\tanh(z))^2$$
$z = 10, \ \tanh(z) \approx 1 \quad z = -10, \ \tanh(z) \approx -1$
$g(z) \approx 0 \qquad g'(z) \approx 0$
$z = 0 \ \tanh = 0$
$g'(z) = 1$

→ Relu
$$g(z) = \max(0, z)$$
$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \end{cases}$$

→ Leaky Relu
$$g(z) = \max(0.01z, z)$$
$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \end{cases}$$

## Gradient descent for neural networks

parameter: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$ 　$n_x = n^{[0]}$ input units
　$(n^{[1]}, n^{[0]}), (n^{[1]}, 1)$
　$(n^{[2]}, n^{[1]}), (n^{[2]}, 1)$ 　$n^{[1]}$ hidden units

cost function: $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]})$ 　$n^{[2]}$ output units
$$= \frac{1}{m} \sum_{i=1}^{m} L(\hat{y}, y)$$

Gradient descent:
Repeat { compute predicts $(\hat{y}^{(i)}, i=1, \ldots, m)$
$$dW^{[1]} = \frac{dJ}{dW^{[1]}}, \quad db^{[1]} = \frac{dJ}{db^{[1]}}, \ldots$$
$$W^{[1]} = W^{[1]} - \alpha \, dW^{[1]}$$
$$b^{[1]} = b^{[1]} - \alpha \, db^{[1]}$$
$\vdots$ }

### Formulas for computing derivatives

Forward propagation:
$$z^{[1]} = W^{[1]} X + b^{[1]}$$
$$A^{[1]} = g^{[1]}(z^{[1]})$$
$$z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$
$$A^{[2]} = g^{[2]}(z^{[2]}) = \sigma(z^{[2]})$$

Backpropagation:
1. $dz^{[2]} = A^{[2]} - Y \quad Y = [Y^{(1)} \cdots Y^{(m)}]$
2. $dW^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$
3. $db^{[2]} = \frac{1}{m} np.sum(dz^{[2]}, axis=1, keepdims=True)$
4. $dz^{[1]} = \underbrace{W^{[2]T} dz^{[2]}}_{(n^{[1]}, m)} * \underbrace{g'^{[1]}(z^{[1]})}_{(n^{[1]}, m)}$
5. $dW^{[1]} = \frac{1}{m} dz^{[1]} x^T$
6. $db^{[1]} = \frac{1}{m} np.sum(dz^{[1]}, axis=1, keepdims=True)$
　$(n^{[1]}, 1)$

### Backpropagation intuition

logistic regression
$$x, w, b \rightarrow \boxed{z = Wx+b} \rightarrow \boxed{a = \sigma(z)} \rightarrow \boxed{L(a,y)}$$
$$L = -y \log a - (1-y) \log(1-a)$$
$$dz = a - y \qquad da = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{dL}{dz} = \frac{dL}{da} \cdot \frac{da}{dz}$$
$$= W^{[2]T} dz^{[2]} \cdot \frac{da^{[1]}}{dz^{[1]}}$$
$$= W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$\frac{dL}{da} \cdot \frac{da}{dz} = \frac{dz^{[1]}}{da^{[1]}} \cdot \frac{da^{[1]}}{dz^{[1]}}$$

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

$z^{[1]}, dz^{[1]} (n^{[1]}, 1)$ 　$W^{[1]} (n^{[1]}, n^{[0]})$
$z^{[2]}, dz^{[2]} (n^{[2]}, 1)$

$$dz^{[1]} = \underbrace{W^{[2]T} dz^{[2]}}_{(n^{[1]}, n^{[2]})(n^{[2]}, 1)} * \underbrace{g'^{[1]}(z^{[1]})}_{(n^{[1]}, 1)}$$

$$dz^{[2]} = da^{[2]} - Y$$
$$dW^{[2]} = dz^{[2]} a^{[1]T}$$
$$db^{[2]} = dz^{[2]}$$

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

### Summary of gradient descent

1. $dz^{[2]} = a^{[2]} - y$
2. $dw^{[2]} = dz^{[2]} \cdot a^{[1]T}$
3. $db^{[2]} = dz^{[2]}$
4. $dz^{[1]} = W^{[2]T} dz^{[2]} \cdot g^{[1]'}(z^{[1]})$
5. $dw^{[1]} = dz^{[1]} x^T$
6. $db^{[1]} = dz^{[1]}$

Vectorization
$$dz^{[2]} = A^{[2]} - Y$$
$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$
$$db^{[2]} = \frac{1}{m} np.sum(dz^{[2]}, axis=1, keepdims=True)$$
$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]}) \quad (n^{[1]}, m)$$
$$dw^{[1]} = \frac{1}{m} dz^{[1]} x^T$$
$$db^{[1]} = \frac{1}{m} np.sum(dz^{[1]}, axis=1, keepdims=True)$$

### Random initialization

$$W^{[1]} = np.random.randn((2,2)) * 0.01$$
$$b^{[1]} = np.zeros((2,1))$$
$$W^{[2]} = np.random.randn((1,2)) * 0.01$$
$$b^{[2]} = np.zeros((2,1))$$

2017年10月26日　　10:18

## Deep-layer neural network



What is a deep neural network?

logistic regression — 1 hidden layer
2 hidden layers — 5 hidden layers

Andrew Ng



$L = 4$ (# layers)
$n^{[l]}$ = # units in layer $l$, $n^{[1]} = 5$, $n^{[2]} = 5$, $n^{[3]} = 3$, $n^{[4]} = 1$, $n^{[0]} = n_x = 3$
$a^{[l]}$ = activations in layer $l$
$a^{[l]} = g^{[l]}(z^{[l]})$, $W^{[l]}$ = weight for $z^{[l]}$
$X = a^{[0]}$

## Forward propagation in a deep network



$z^{[1]} = W^{[1]} X + b^{[1]}$　　$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$
$a^{[1]} = g^{[1]}(z^{[1]})$　　$a^{[1]} = g^{[1]}(z^{[1]})$
$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$
$a^{[2]} = g^{[2]}(z^{[2]})$
$\vdots$
$z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}$
$a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$

Vectorized:
$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$　⎫ for loop
$A^{[l]} = g^{[l]}(Z^{[l]})$　　　　　⎬ $l = 1, \ldots, 4$
$\hat{y} = g^{[4]}(Z^{[4]}) = A^{[4]}$　⎭

## Getting your matrix dimensions right

Parameter $W^{[l]}$ and $b^{[l]}$



$dW^{[l]}, W^{[l]} : (n^{[l]}, n^{[l-1]})$
$db^{[l]}, b^{[l]} : (n^{[l]}, 1)$

$(n^{[1]}, n^{[0]})$
$z^{[1]} = W^{[1]} \cdot X + b^{[1]}$　the vector of activations
$(n^{[1]}, 1)$　　$(n^{[0]}, 1)$

Vectorized implementation
$Z^{[1]} = W^{[1]} X + b^{[1]}$　$(n^{[1]}, m)$　$m$ is the number of dataset
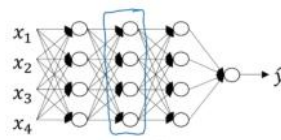　　by broadcasting. $b^{[1]} (n^{[1]}, 1) \Longrightarrow (n^{[1]}, m)$
$Z^{[l]}, A^{[l]} : (n^{[l]}, m)$

## Why deep representation?

Circuit theory and deep learning

Informally, there are functions you can compute with a small L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

## Building blocks of deep neural networks



Layer L　　eg

layer $l$ : $W^{[l]}, b^{[l]}$
Forward : input $a^{[l-1]}$, output $a^{[l]}$
Backward : input $da^{[l]}$, output $da^{[l-1]}$, $dW^{[l]}, db^{[l]}$

$z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$
$a^{[l]} = g^{[l]}(z^{[l]})$

## Forward and backward propagation

### Forward propagation for layer l

input $a^{[l-1]}$
output $a^{[l]}$, cache $(z^{[l]})$

$z^{[l]} = W^{[l]} \cdot a^{[l-1]} + b^{[l]}$
$a^{[l]} = g^{[l]}(z^{[l]})$

Vectorized:
$Z^{[l]} = W^{[l]} \cdot A^{[l-1]} + b^{[l]}$
$A^{[l]} = g^{[l]}(Z^{[l]})$

### Backward propagation for layer l

input $da^{[l]}$
output $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$dz^{[l]} = da^{[l]} * g^{[l]'}(z^{[l]})$
$dW^{[l]} = dz^{[l]} \cdot a^{[l-1]}$
$db^{[l]} = dz^{[l]}$
$da^{[l-1]} = W^{[l]T} \cdot dz^{[l]}$
$dz^{[l]} = W^{[l+1]T} dz^{[l+1]} * g^{[l]'}(z^{[l]})$

Vectorized:
$dZ^{[l]} = dA^{[l]} * g^{[l]'}(Z^{[l]})$
$dW^{[l]} = \frac{1}{m} dZ^{[l]} \cdot A^{[l-1]T}$
$db^{[l]} = \frac{1}{m} np.sum(dZ^{[l]}, axis=1, keepdims=True)$
$dA^{[l-1]} = W^{[l]T} dZ^{[l]}$

### Summary



$X \rightarrow \boxed{ReLu} \rightarrow \boxed{ReLu} \rightarrow \boxed{sigmoid} \rightarrow \hat{y} \rightarrow \mathcal{L}(\hat{y}, y)$

$da^{[l]} = -\frac{y}{a} + \frac{1-y}{1-a}$　$L'(a, y)$
$\mathcal{L}(a, y) = -(y\log a + (1-y)\log(1-a))$

## Parameters vs hyperparameters

parameters : $W^{[l]}, b^{[l]} \ldots$

hyperparameters: learning rate $\alpha$
　　　# iterations
　　　# hidden layer $L$
　　　# hidden units $n^{[1]}, n^{[2]} \ldots$
　　　choice of activation function
other: momentum, mini-batch size, regularization

Applied deep learning is a very empirical process
Application: vision, speech recognition, NLP, online advertising, recommendation

## What does this have to do with human brain?

forward propagation

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$
$$A^{[1]} = g^{[1]}(Z^{[1]})$$
$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$
$$A^{[2]} = g^{[2]}(Z^{[2]})$$
$$\vdots$$
$$A^{[L]} = g^{[L]}(Z^{[L]}) = \hat{Y}$$

backward propagation

$$dZ^{[L]} = A^{[L]} - Y$$
$$dW^{[L]} = \frac{1}{m} dZ^{[L]} A^{[L]T}$$
$$db^{[L]} = \frac{1}{m} np.\,sum(dZ^{[L]}, axis = 1, keepdims = True)$$
$$dZ^{[L-1]} = dW^{[L]T} dZ^{[L]} g'^{[L]}(Z^{[L-1]})$$
$$\vdots$$
$$dZ^{[1]} = dW^{[L]T} dZ^{[2]} g'^{[1]}(Z^{[1]})$$
$$dW^{[1]} = \frac{1}{m} dZ^{[1]} A^{[1]T}$$
$$db^{[1]} = \frac{1}{m} np.\,sum(dZ^{[1]}, axis = 1, keepdims = True)$$