

## <WEEK2> NN Implemental techniques 신경망 구조와 학습 기법.

### OVERVIEW

- ex) ①  $\rightarrow$  in training set.  $\rightarrow$  loop  
 $\Rightarrow$  explicit  $\rightarrow$  not for loop!  $\rightarrow$  for loop 쓰지 않고 training
- ② forward pass, forward propagation computation in NN.  
 $\text{State } g_2 = \text{Input}$
- why?  
 backward pass, backward propagation 계산방법. why 이런 계산?  
 역방향 정리 역전파

logistic regression.  $\rightarrow$  2차원적 회귀.

$\hookrightarrow$  이진 분류 (Binary Classification)을 위한 알고리즘

고양이 예제 Input  $x$   
 output  $y$ .

0|1|2| - RGB.  $\rightarrow$   $\underbrace{64 \times 64 \times 3}_{\text{이미지 } 4 \times 1^3} \rightarrow \underbrace{5 \times 4}_{\text{이미지 } 4 \times 1^3}$ .

intensity  $\rightarrow$  feature vector.  $\rightarrow$  펼쳐내 all pixel into a input feature vector  $x$   
 $\hookrightarrow$  to unroll  $\rightarrow$  find feature vector  $x$  리스트에 일렬로 나열.

:  $64 \times 64$   
 $\downarrow$  RGB

$\hookrightarrow 64 \times 64 \times 3 = 12288$ .

$N = N_x = 12288$ .

$\hookrightarrow$  dimension of input feature  $x$   
 $\downarrow$  차원

42488  $\div 128$

$N/48$ .

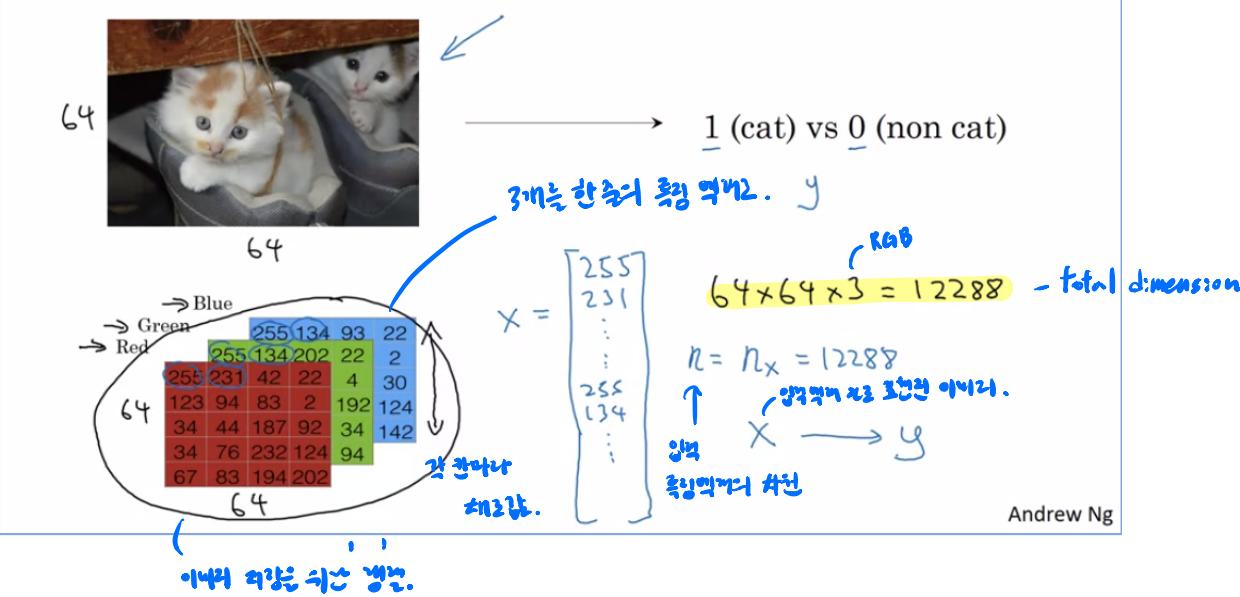
이진분류

입력벡터  $x \rightarrow$  출력으로 주어 분류기 학습  
 classifier

output  $y \rightarrow 1$  or  $0$ ?

cat noncat.

## Binary Classification



$\Delta z \rightarrow 0$  ဆိုရင်  $0 \leq z \leq 1$

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

여기 x가 리는 책을.

$x \rightarrow p: L \Rightarrow$  dimensional vector       $x \in \mathbb{R}^{n_x}$

$\hat{y} \rightarrow$  고양이일 확률.

parameters:  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$ .

Real num.

Input  $x_2$ , para w.b.  $\Rightarrow \hat{y}$  ?

$$\text{output } \hat{y} = w^T x + b$$

↳ 0 ≤ x ≤ 1 단위 단위에 있어 가능하게 풀리 x.

회귀분석을 가르기는데 쓰는 초기 양식  
신경망에 대한 초기 양식

< Notation  $\rightarrow$  표기법 >

single train example  $\rightarrow$   $(x, y)$ .  
 $x \in \mathbb{R}^{n_x}$ ,  $y \in \{0, 1\}$   
 $x$  차수를 가진 투명벡터 | Table.

$m$  training examples :  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$

M train.      M test  
 ↳ 훈련용 데이터  $m_s$ .      테스트용 데이터  $m_t$ .

implement - 구현하기

Matrix      input  $x$  - stack in columns

$X = \begin{bmatrix} | & | & | \\ x_1 & x_2 & \dots & x_m \\ | & | & | \end{bmatrix} \quad \leftarrow m \text{개의 행}\rightleftharpoons n_x \text{개의 열}$

$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$   
 $Y \in \mathbb{R}^{1 \times m} \quad Y.\text{shape} = (1, m)$ .

$X \in \mathbb{R}^{n_x \times m}$        $X.\text{shape} (n_x, m)$   
 훈련용 행렬을 일렬로 쌓은 행렬.  
 ↳ 각 행은.

각 행은 예제와 관련된 데이터.      ↳ 같은 방식으로 전부.

## Notation

$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$

$m$  training examples :  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$M = M_{\text{train}}$        $M_{\text{test}} = M_{\text{test}}$        $M_{\text{train}} + M_{\text{test}} = \# \text{train examples}$   
 훈련용 훈련셋      테스트용 훈련셋

$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \quad \leftarrow m \text{개의 행} \quad \leftarrow n_x \text{개의 열}$

$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$   
 $Y \in \mathbb{R}^{1 \times m}$   
 $Y.\text{shape} = (1, m)$

$X \in \mathbb{R}^{n_x \times m}$   
 ↳ 각 행은 예제와 관련된 행렬.  
 $X.\text{shape} = (n_x, m)$

Andrew Ng

수학 2/3).

↳  $y \rightarrow$  이진 분류의 예제 0이나 1. 이진 분류.

$B = \text{inter-spectrum}$

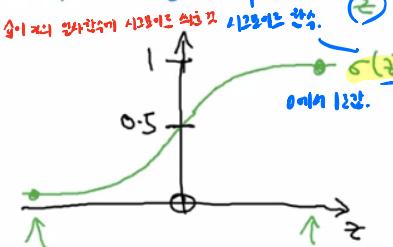
## Logistic Regression

Given feature vector  $x$ , want  $\hat{y} = P(y=1|x)$

Parameters:  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$ .

단위  $n_x$ -dimensional vector  $w$   
parameter: vector  $n_x$ -dimensional vector  $b$

$$\hat{y} = \sigma(w^T x + b)$$



$0 \leq \hat{y} \leq 1$ 이기 때문에 sigmoid 함수를 취한다.

( $x$ 가  $y$ 에 대한 확률)

$y$ 가  $x$ 가 되는 확률

이산화  
고양이인 확률?

$0 \leq \hat{y} \leq 1$

<--- 확률로는 합의  $x$

나타내는 확률

나타내는 확률

일반화

binary classification  
1과 0의 확률

승은 부가 1이 되는

why?

$$x_0 = 1, \quad x \in \mathbb{R}^{n_x+1}$$

$$\hat{y} = \sigma(w^T x)$$

$$\Theta = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \\ \vdots \\ \Theta_{n_x} \end{bmatrix} \quad \left. \begin{array}{l} \Theta_0 \text{ } \leftarrow \text{bias} \\ \Theta_1, \dots, \Theta_{n_x} \text{ } \leftarrow \text{계수} \end{array} \right\} w \leftarrow \text{계수} \quad \left. \begin{array}{l} \Theta_0 \text{ } \leftarrow \text{bias} \\ \Theta_1, \dots, \Theta_{n_x} \text{ } \leftarrow \text{계수} \end{array} \right\} w \leftarrow \text{계수}$$

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

If  $z$  large  $\sigma(z) \approx \frac{1}{1+0} = 1$

If  $z$  large negative number

$$\sigma(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+\text{Bignum}} \approx 0$$

Andrew Ng

근가 작으면  $\sigma(z)$ 는 0에 가까운

즉  $y=1$  확률은 축소되는지.

to think param  $w$  &  $b \rightarrow$  model



define  
cost func.

sigmoid?

model training set에서 param  $w$  와  $b$  찾음.

prediction  $\hat{y}^{(i)}$

$$\hat{y}^{(i)} \approx y^{(i)}$$

sigmoid func

super

식 설명: 광고 단위 있는 것의 각각의 인덱스에 해당.  
or 한바탕으로 각각의 training set을 나눠놓은 것

superscript = subscript

한바탕, 여러 글자.

$$\delta(\quad)$$

↳ 미분 가능한지?

광고 허락 - 예상 인덱스.

least squares error.  
↳ 각각의 풀  $w$ .



$$\langle \text{Loss func} \rangle \mathcal{L}(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2.$$

↳  $\frac{1}{2}$  미분.

$$(X).$$

non-convex.  $\rightarrow$  비별곡.

gradient descent 잘 안나온.

(비슷한게 있어서)

복잡한 최적화 문제.

$\Rightarrow$  optimization with multiple local optima.

복수 최적화 optimization.

↳ gradient descent  $\rightarrow$  may not find global optima

전역 최적값을 못 찾을 수 있음.

유감을 가지고 있을 때

얼마나 정복이 좋을

선수에게 나에게는 좋음.

gradient descent  
Loss func ① 2차 소거법은 찾을 수 있다.  
Loss func ② 4차.

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

↳ 이거 학습률 좋음.

$$\text{if } y=1 : \mathcal{L}(\hat{y}, y) = -\frac{\log \hat{y}}{\log 1!} \text{ 라이어스니까.}$$

•  $\hat{y}$  까아

아이가 놀고 춤추는데

주제 아름다워!

$\dots$

↳ 428이면 func의 결과  
124 출석 X.

$\Rightarrow$   $\hat{y}$ 는 1에 최적은 가깝지.

$$\text{if } y=0 : \mathcal{L}(\hat{y}, y) = -\log(1-\hat{y}).$$

▷ 해야!

$\therefore \hat{y}$ 가 가령 0.1이면

$\Rightarrow 0$ 에 가까워야.

↳  $y$ 가 1이면  $\hat{y}$ 는 가령 크지.

$y$ 가 0이면  $\hat{y}$ 는 가령 작지.

en 이를 나누는가?  $\rightarrow$  속도.

따라서 single example를 끝이.

cost func :  $J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$

(entire set)  
cost of parameters

$\uparrow$   
prediction output.

$$= -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right]$$

Loss func  $\rightarrow$  미분 가능.

$\Rightarrow$  Loss func는 single training example에 적용.

cost func는 각각의 예의 비중을 나타낸다.

Logistic regression 모델을 훈련

$\rightarrow$  weight  $w$ 와 bias  $b$ 를 찾고, cost func는 비중  $J$ 를 증가나가는 것.

$\Rightarrow$  cost func는 loss func의 평균

이걸 통해  $w$ 와  $b$ 를 찾아내기

$\hat{y} \rightarrow$  어떤  $x$ 는 선형 합성에 sigmoid 취한 것.

## Logistic Regression cost function

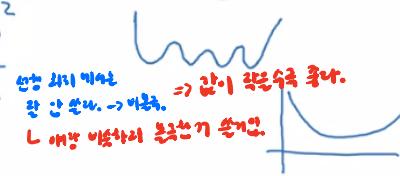
$$\rightarrow \hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \text{ where } \sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}} \quad z^{(i)} = w^T x^{(i)} + b$$

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

( $i$ -th example) ; non training example

Loss (error) function:

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$



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

loss function

If  $y=1$ :  $L(\hat{y}, y) = -\log \hat{y} \leftarrow$  want  $\log \hat{y}$  large, want  $\hat{y}$  large.

$y=1$  일 때  $\hat{y}$ 가 작아야 한다.

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

$y=0$  일 때  $\hat{y}$ 가 작아야 한다.

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

Andrew Ng

Loss func single training example로 나눠보자.

Ratidah's effect

cost: loss  $\hat{y}$ 와  $y$ .

cost of parameter

find  $w$  and  $b$  that minimize cost func  $J$ 's value

## Clarification about Upcoming Gradient Descent Video

Please note that in the next video at the second slide, there is a missing parenthesis.

The negative sign should apply to the entire cost function (both terms in the summation).

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

Also Andrew mentions: "To make this easier to draw, I'm going to ignore  $b$  for now, just to make this a one-dimensional plot instead of a high-dimensional plot." This plot is actually two-dimensional; it shows the loss  $J(w)$  for values of  $w$ .

※: Loss - single set에서의 Loss

cost - 여러 set에서의 평균으로서의 Loss

cost는 min set에서 Loss의 평균

전체  $\frac{1}{m} \sum_{i=1}^m$  gradient descent에서 전부의 최적값을 찾을 수 있게  
log 쓰는 Loss func 좋다.

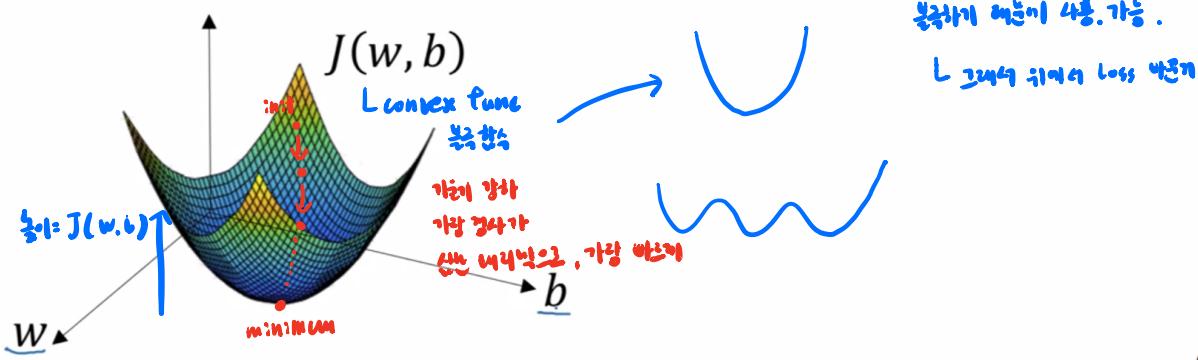
## Gradient Descent

Gradient Descent

Recap:  $\hat{y} = \sigma(w^T x + b)$ ,  $\sigma(z) = \frac{1}{1+e^{-z}}$  ← 0과 1 사이의 값이 출력하기 위한 sigmoid 함수.

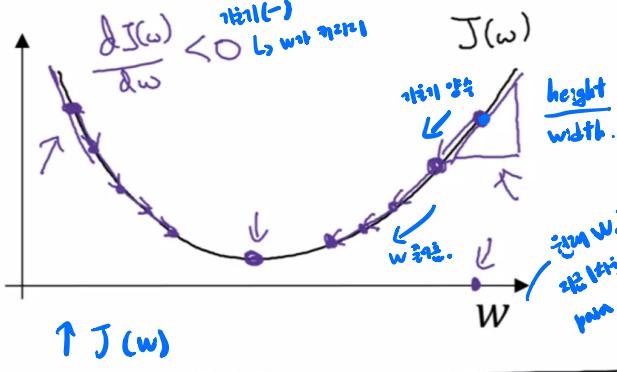
cost func  $J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$  average of loss func

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



Andrew Ng

## Gradient Descent



$$J(\omega, b)$$

$$\omega := \omega - \alpha \frac{\partial J(\omega, b)}{\partial \omega}$$

$$b := b - \alpha$$

$$\frac{\partial J(\omega, b)}{\partial \omega}$$

$$\frac{\partial J(\omega, b)}{\partial \omega_k}$$

$$\boxed{\frac{2J(\omega, b)}{2b}}$$

A diagram illustrating the relationship between  $d\omega$  and  $db$ . It features two blue circles labeled  $\omega$  and  $b$ , each with a small arrow indicating rotation. A red arrow points from the bottom right towards the top left, passing through a yellow oval containing the text  $d\omega$ . Another red arrow points from the bottom right towards the top left, passing through a yellow oval containing the text  $db$ .

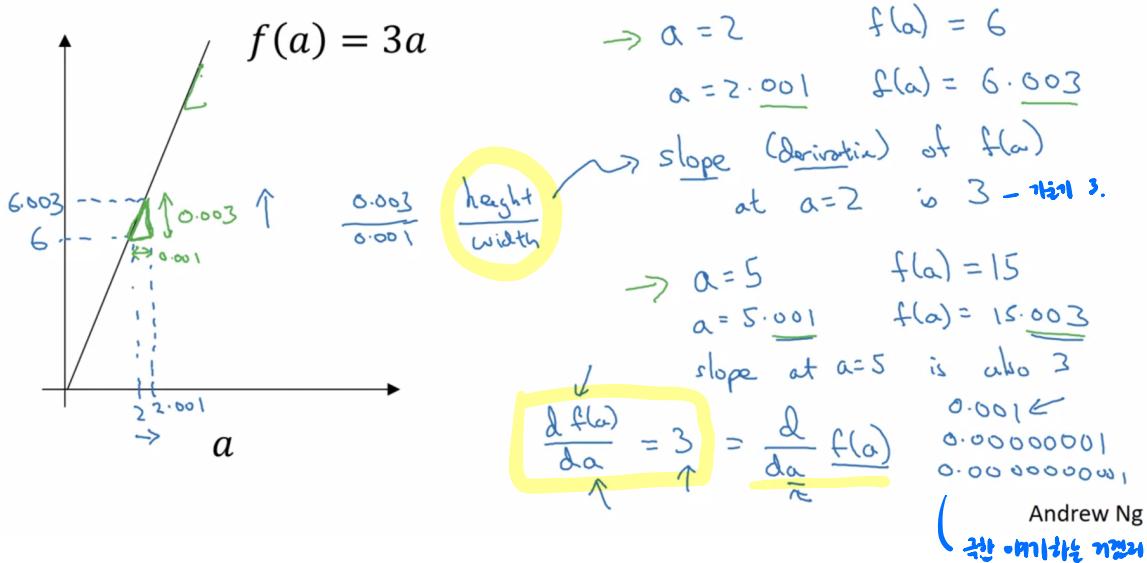
Andrew Ng

순간 변화 .

## < Derivative (미분학) >

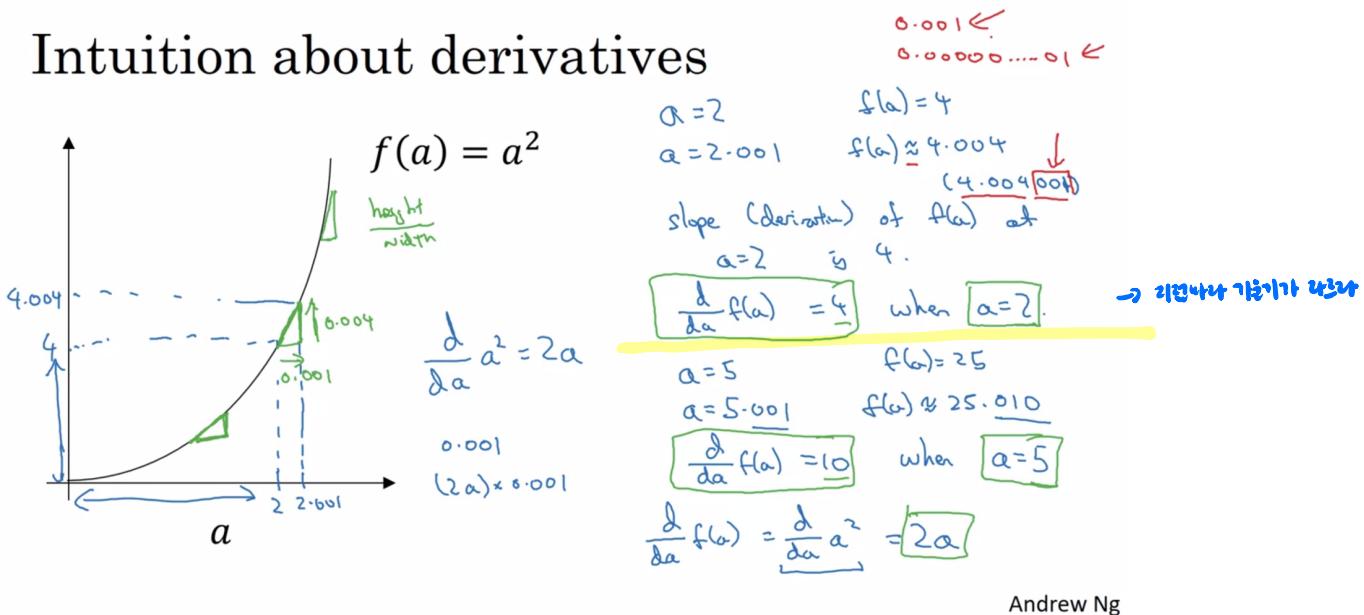
22224번 과제 할 수 있는 21m. + forward backward 4m. on 4m.

# Intuition about derivatives



기울기 변화는?

# Intuition about derivatives



# More derivative examples

$$f(a) = a^2 \quad \frac{d}{da} f(a) = \underline{\underline{2a}} \quad a=2 \quad f(a)=4$$

$$a=2.001 \quad f(a) \approx 4.004$$

$$f(a) = a^3 \quad \frac{d}{da} f(a) = \underline{\underline{3a^2}} = 12 \quad a=2 \quad f(a)=8$$

$$a=2.001 \quad f(a) \approx 8.012$$

수학적 계산,  
기울기 변화

$$f(a) = \frac{\log_e(a)}{\ln(a)} \quad \frac{d}{da} f(a) = \frac{1}{a} \quad a=2 \quad f(a) \approx 0.4315$$

$$a=2.001 \quad f(a) \approx 0.69365$$

$$\frac{d}{da} f(a) = \frac{1}{a^2} \quad 0.0005 \quad 0.0005$$

$\ln(a)$   
 $0.0005$

Andrew Ng



deeplearning.ai

## Basics of Neural Network Programming

### Computation Graph

computations of NN  $\rightarrow$  forward pass / propagation step

↑ 전방향

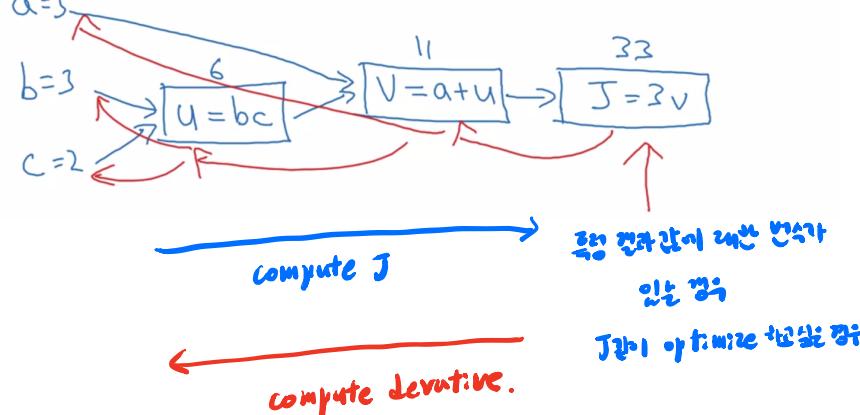
전방향 계산을 backward pass 통해 거친 계산

↓ why? example.

### Computation Graph $\rightarrow$ compute J

$$J(a, b, c) = 3(a + bc) = 3(\underbrace{a}_{\text{cost func}} + \underbrace{bc}_{v}) = 33$$

계산  
↓  
 $u = bc$   
 $v = au$   
 $J = 3v$

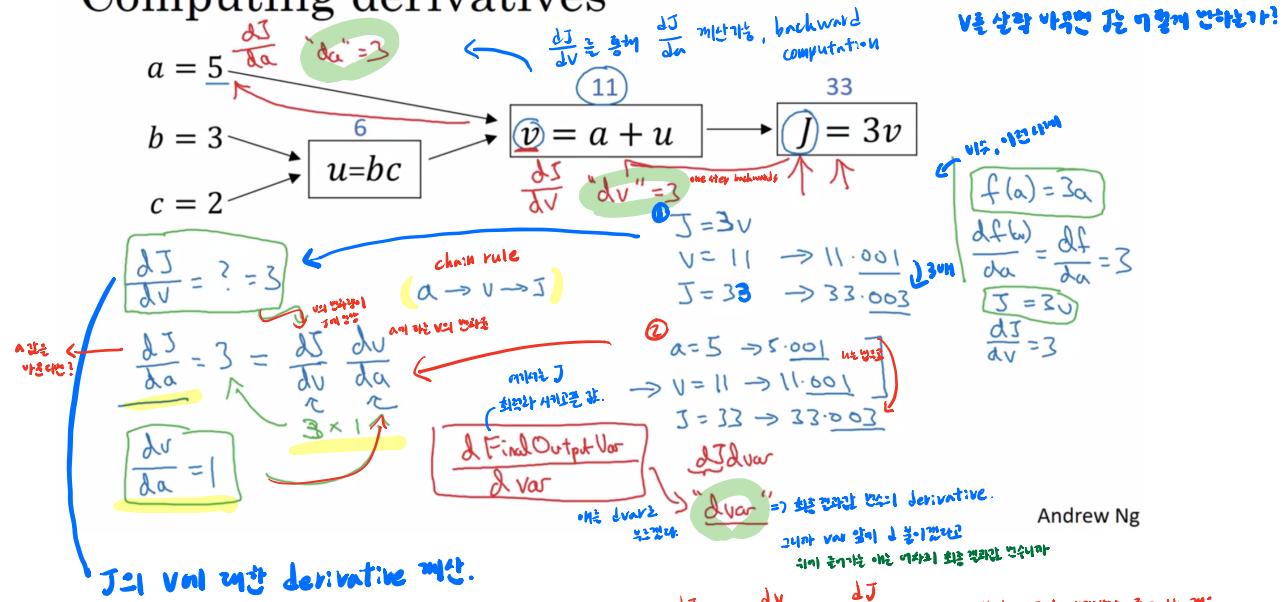


deeplearning.ai

## Basics of Neural Network Programming

### Derivatives with a Computation Graph

## Computing derivatives

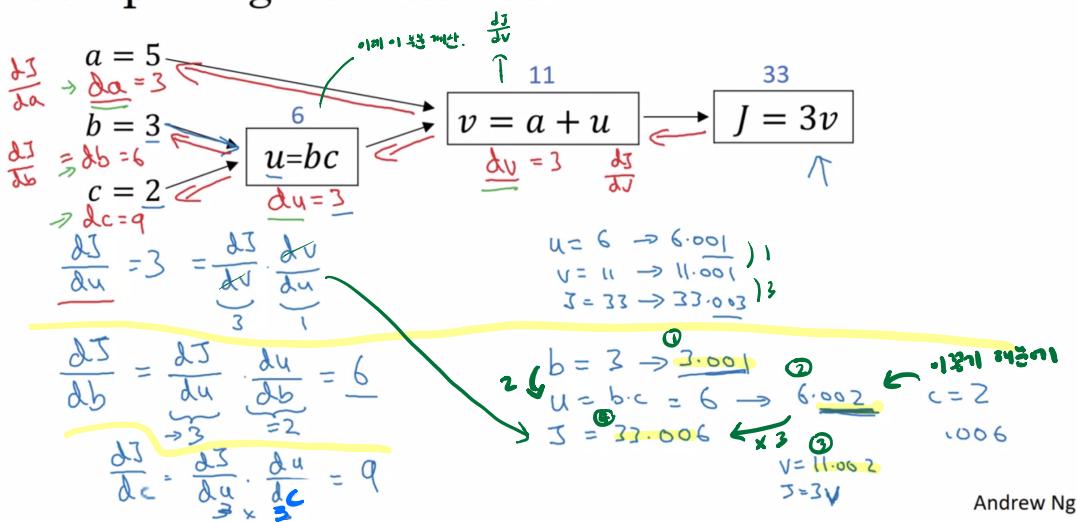


② change  $a \rightarrow$  change  $v \rightarrow$  change  $J$        $\frac{dJ}{da} = \frac{dv}{da} \times \frac{dJ}{dv}$       비례분이니 같은 배수변수를かける정리

Andrew Ng

$$\text{review} \quad \frac{dJ}{Jv} = 3 \quad \frac{dJ}{dn} = \frac{dJ}{Jv} \times \frac{dv}{dn} \quad \text{a변수에 따른 v의 변화.}$$

# Computing derivatives



chain rule 쉽게 생각하면 되는 걸 암기해라.

$\Rightarrow$  chain rule 用来反向传播梯度

$\Delta V$ 를 찾으면  $\Delta h_f$   $\Delta h_i$  찾을 수 있다.  $\Delta v \rightarrow \Delta h$ ,  $\Delta c$  compute.

$J$ 를 계산하는 방식은 forward, derivative를 계산하는 방식은 backward.

→ Logistic regression model의 derivative 미분.



deeplearning.ai

# Basics of Neural Network Programming

---

## Logistic Regression Gradient descent

What you need to implement

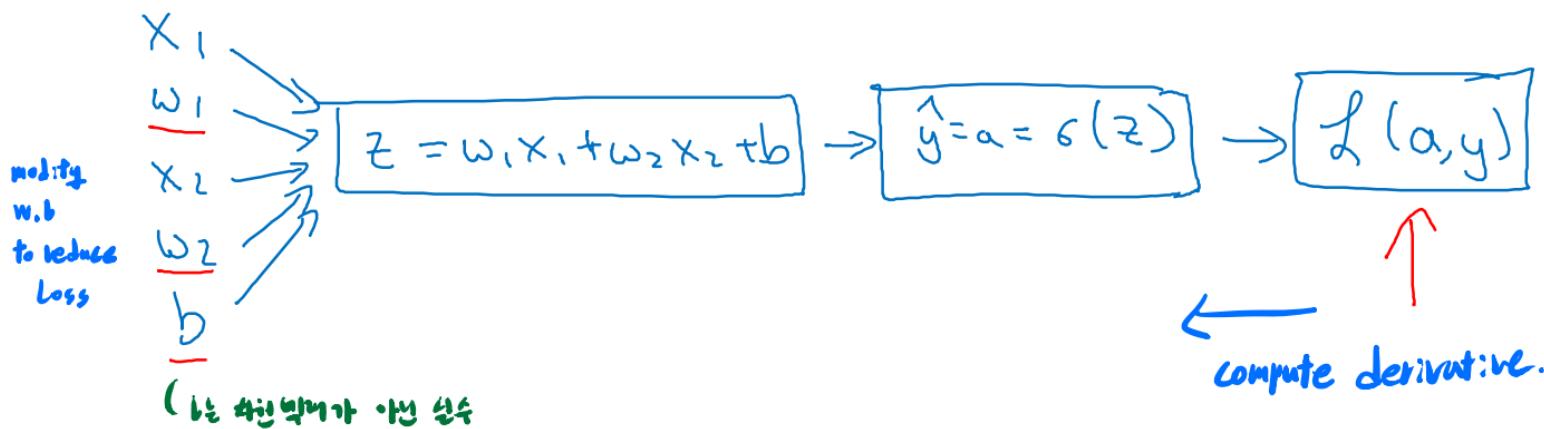
그래프학습을 위한 기본 기능을 조합해 사용하는 주제 공식

# Logistic regression recap

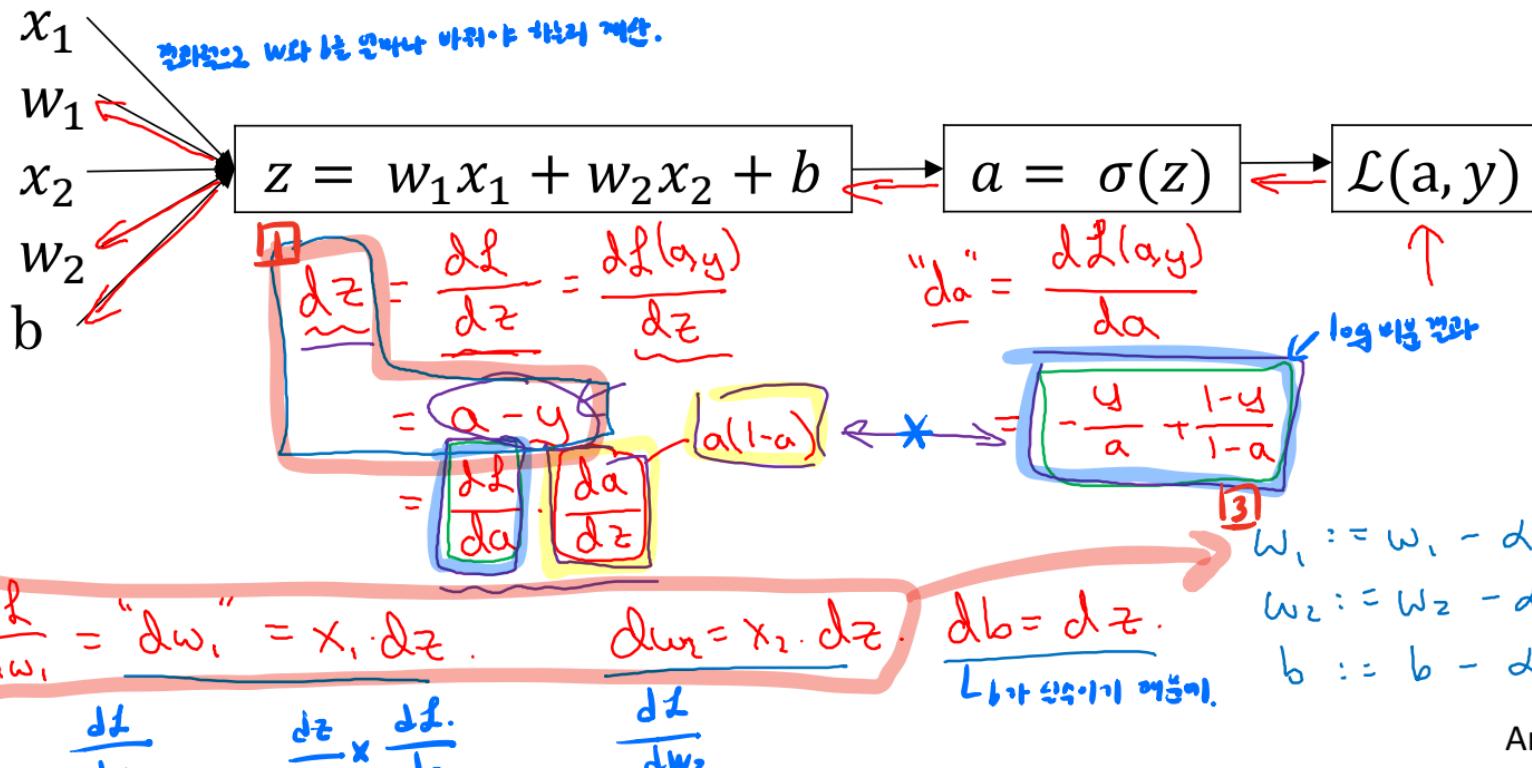
$$\rightarrow z = w^T x + b \quad y: \text{ground truth table}$$

prediction  $\rightarrow \hat{y} = a = \sigma(z)$

loss  $\rightarrow \mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$



# Logistic regression derivatives





deeplearning.ai

$dW_1 \quad dz$   
Lagrangian w.r.t.  $w_1$  and  $z$

$$\frac{dz}{d\lambda_1}$$

1회 example of  $m=1$  gradient descent  
를 위한 예  $\square \rightarrow \square \rightarrow \square$   
but 예의 example 있는.

# Basics of Neural Network Programming

---

## Gradient descent on $m$ examples

# Logistic regression on $m$ examples

cost func

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m l(a^{(i)}, y^{(i)})$$
$$\rightarrow \hat{a}^{(i)} = \hat{y}^{(i)} = g(z^{(i)}) = g(\omega^\top x^{(i)} + b)$$

$(x^{(i)}, y^{(i)})$

$\underline{d\omega_1^{(i)}}, \underline{d\omega_2^{(i)}}, \underline{db^{(i)}}$

1 training example

$$\frac{\partial}{\partial \omega_1} J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial \omega_1} l(a^{(i)}, y^{(i)})}_{d\omega_1^{(i)}} - (x^{(i)}, y^{(i)})$$

$d\omega_1^{(i)}$ 는  $d\omega_1$ 의 정수값이 결과로 나온다.

# Logistic regression on $m$ examples

$$J=0; \underline{\Delta w_1}=0; \underline{\Delta w_2}=0; \underline{\Delta b}=0$$

$\rightarrow$  For  $i = 1$  to  $m$

$$z^{(i)} = \omega^\top x^{(i)} + b$$

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

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

$$\underline{\Delta z^{(i)}} = a^{(i)} - y^{(i)}$$

$$\begin{aligned} \textcircled{1} \quad \Delta w_1 &+= x_1^{(i)} \Delta z^{(i)} \\ \Delta w_2 &+= x_2^{(i)} \Delta z^{(i)} \\ \Delta b &+= \Delta z^{(i)} \end{aligned}$$

$\uparrow \downarrow \quad \uparrow \downarrow$  기여 확장판 합계를  
가정하여

내부의 정한값.

$$J_t/m \leftarrow$$

$$\Delta w_1/m; \Delta w_2/m; \Delta b/m \leftarrow$$

$$\Delta w_1 = \frac{\partial J}{\partial w_1}$$

$\Delta w_2$  &  $\Delta b$  를 정의

gradient descent의 2단

$$w_1 := w_1 - \alpha \frac{\partial J}{\partial w_1}$$

$$w_2 := w_2 - \alpha \frac{\partial J}{\partial w_2}$$

$$b := b - \alpha \frac{\partial J}{\partial b}$$

기준마다 | 가지 단계의 가중치를  
증가하는 것.

Vectorization

조합한 것.



이것이 여러 번 반복  
multiple steps of  
gradient descent

Andrew Ng

계산에의 2가지 방법.

방법 for loop 사용 ① m training example  $\rightarrow$  to compute average.  
② all the features  $n$ 개의 features

$\Rightarrow$  알고리즘 복잡성을 줄여. 여기서 세트 처리는데.



Vectorization  $\rightarrow$  for loop 사용.

↑      ↗  
speedy.      매우 빠르게 처리할 수 있는 방법.



deeplearning.ai

# Basics of Neural Network Programming

---

## Vectorization

코드에서 for loop 사용  
제거해 놓고 예제에 빠는 코드 쓰고

# What is vectorization?

$$z = \underbrace{w^T x}_{} + b$$

Non-vectorized:

$z = 0$       0이 더해지는 경우  
for i in range(n - x):  
 $z += w[i] * x[i]$

$$z += b$$

slow!

column vector.

$$w = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad x = \begin{bmatrix} : \\ : \\ : \end{bmatrix}$$

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

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

Vectorized

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

$\Rightarrow$  GPU      } SIMD - single instruction  
 $\Rightarrow$  CPU      } multiple data.

parallelization instruction이 있음

↳ numpy가 이를 처리하는 경우, 즉, GPU가 그 작업을 할 때.

Andrew Ng



deeplearning.ai

# Basics of Neural Network Programming

---

## More vectorization examples

# Neural network programming guideline

Whenever possible, ~~avoid explicit for-loops.~~

Using built-in func.  
for speed.

# Neural network programming guideline

Whenever possible, avoid explicit for-loops.

$$u = Av$$

A matrix의 곱셈.

2m\*1 for 2m\*1.

$$\text{정답이 } u_i = \sum_i \sum_j A_{ij} * v_j$$

$$u = np.zeros((n,))$$

for i ...      ↪

  for j ...      ↪

$$u[i] += A[i][j] * v[j]$$

$$u = np.dot(A, v)$$

# Vectors and matrix valued functions

Say you need to apply the exponential operation on every element of a matrix/vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

vector 0:2 출력  
 $\rightarrow u = np.zeros((n, 1))$   
for i in range(n):  $\leftarrow$   
 $\rightarrow u[i] = \text{math.exp}(v[i])$

import numpy as np  
u = np.exp(v)  $\leftarrow$   
↑ ↗ **행렬의 선형연산: 행과 열에 따른 원소별 연산.**  
np.log(v)  $\rightarrow$  element-wise log.  
np.abs(v)  
np.maximum(v, 0)  $\downarrow$  **element-wise maximum**  
v\*\*2  $\downarrow$  **v의 모든 원소에 적용**  
1/v  $\downarrow$  **회전값은 0:2**  
↳ v 원소의 원소에 적용  
element-wise의 연산.

# Logistic regression derivatives

$J = 0, \boxed{dw_1 = 0, dw_2 = 0}, db = 0$

$\rightarrow \text{for } i = 1 \text{ to } n:$

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

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

$\text{주의 } dz^{(i)} = \frac{da}{dz}$

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

$dz^{(i)} = a^{(i)}(1 - a^{(i)})$

$\downarrow$

$\text{for } j=1 \dots n_x$

$dw_1 += x_1^{(i)} dz^{(i)}$

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

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

$J = J/m, \boxed{dw_1 = dw_1/m, dw_2 = dw_2/m}, db = db/m$

$dw /= m$

$dw = np.zeros((n_x, 1))$

$n_x \in 1+인 수면2.$

$\uparrow$   
이거 충분히 많아.



deeplearning.ai

# Basics of Neural Network Programming

---

## Vectorizing Logistic Regression

to compute forward propagation  $\Rightarrow \alpha \in \mathbb{R}^{n+1}$

부정법률.

# Vectorizing Logistic Regression

$$\begin{array}{l} \rightarrow z^{(1)} = w^T x^{(1)} + b \\ \rightarrow a^{(1)} = \sigma(z^{(1)}) \end{array}$$

$$\begin{array}{l} \rightarrow z^{(2)} = w^T x^{(2)} + b \\ \rightarrow a^{(2)} = \sigma(z^{(2)}) \end{array}$$

$$\begin{array}{l} \rightarrow z^{(3)} = w^T x^{(3)} + b \\ \rightarrow a^{(3)} = \sigma(z^{(3)}) \end{array}$$

$$\overline{X} = \left[ \begin{array}{c} x^{(1)} \\ | \\ x^{(2)} \\ | \\ \dots \\ | \\ x^{(m)} \end{array} \right] \quad \begin{matrix} (n_x, m) \\ \overline{\mathbb{R}^{n_x \times m}} \end{matrix}$$

$$\overline{w}^\top \left[ \begin{array}{c} x^{(1)} \\ | \\ x^{(2)} \\ | \\ \dots \\ | \\ x^{(m)} \end{array} \right]$$

$$\overline{z} = \left[ \begin{array}{c} z^{(1)} \\ | \\ z^{(2)} \\ | \\ \dots \\ | \\ z^{(m)} \end{array} \right] = \underbrace{w^\top \overline{X}}_{1 \times m} + \underbrace{\left[ \begin{array}{c} b \\ | \\ b \\ | \\ \dots \\ | \\ b \end{array} \right]}_{1 \times m} = \left[ \begin{array}{c} w^\top x^{(1)} + b \\ | \\ w^\top x^{(2)} + b \\ | \\ \dots \\ | \\ w^\top x^{(m)} + b \end{array} \right] \quad \text{"Broadcasting"}$$

$$A = \left[ \begin{array}{c} a^{(1)} \\ | \\ a^{(2)} \\ | \\ \dots \\ | \\ a^{(m)} \end{array} \right] = \sigma(\overline{z})$$

programming을 vector valued sigmoid로 한다.

각각의 원소를 대상으로 A 선호.

=> Vectorization을 사용하여 효율적으로 모든 activation을 신호하는 방법.



deeplearning.ai

# Basics of Neural Network Programming

---

## Vectorizing Logistic Regression's Gradient Computation

# Vectorizing Logistic Regression

$$dz^{(1)} = a^{(1)} - y^{(1)}$$

$$dz^{(2)} = a^{(2)} - y^{(2)}$$

...

$$d\bar{z} = \begin{bmatrix} dz^{(1)} & dz^{(2)} & \dots & dz^{(m)} \end{bmatrix} \quad (1 \times m)$$

$$A = [a^{(1)} \dots a^{(m)}], \quad Y = [y^{(1)} \dots y^{(m)}]$$

$$\rightarrow d\bar{z} = A - Y = [a^{(1)} - y^{(1)} \quad a^{(2)} - y^{(2)} \quad \dots]$$

$$\begin{cases} \cancel{dw_1} \\ \cancel{dw_2} \\ \vdots \\ dw_l = m \end{cases}$$

$$\begin{cases} dw_1 + = \frac{x^{(1)} dz^{(1)}}{m} \\ dw_2 + = \frac{x^{(2)} dz^{(2)}}{m} \\ \vdots \\ dw_l + = \frac{x^{(l)} dz^{(l)}}{m} \end{cases}$$

$$\begin{cases} db = 0 \\ db + = dz^{(1)} \\ db + = dz^{(2)} \\ \vdots \\ db + = dz^{(m)} \\ db/m = m \end{cases}$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

$$= \frac{1}{m} \underbrace{\text{np. sum}(d\bar{z})}_{\text{계산}}$$

$$dw = \frac{1}{m} \times d\bar{z}^T$$

$$= \frac{1}{m} \left[ \underbrace{\begin{matrix} x^{(1)} & \dots & x^{(m)} \\ 1 & \dots & 1 \end{matrix}}_{n \times 1} \right] \left[ \begin{matrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{matrix} \right]$$

$$= \frac{1}{m} \left[ \underbrace{x^{(1)} dz^{(1)}}_n + \dots + \underbrace{x^{(m)} dz^{(m)}}_n \right]$$

 $n \times 1$ 

dw - parameter.

# Implementing Logistic Regression

$$J = 0, dw_1 = 0, dw_2 = 0, db = 0$$

for  $i = 1$  to  $m$ :

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

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

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

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

$$\left. \begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \end{array} \right\} dw = X^{(i)} * dz^{(i)}$$

vector  
value.

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

$$J = J/m, dw_1 = dw_1/m, dw_2 = dw_2/m$$

$$db = db/m$$

for iter in range(1000):  
 $Z = w^T X + b$   
 $= np.dot(w.T, X) + b$   
 $A = \sigma(Z)$   
 $dZ = A - Y$   
 $dw = \frac{1}{m} X^T dZ$   
 $db = \frac{1}{m} np.sum(dZ)$

$$w := w - \alpha dw$$

$$b := b - \alpha db$$

Logistic regression in single iteration.

Logistic for  $X$

Andrew Ng



deeplearning.ai

# Basics of Neural Network Programming

---

## Broadcasting in Python

# Broadcasting example

Calories from Carbs, Proteins, Fats in 100g of different foods:

$$\begin{array}{c} \text{Apples} \quad \text{Beef} \quad \text{Eggs} \quad \text{Potatoes} \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ \begin{matrix} \text{Carb} & [56.0 & 0.0 & 4.4 & 68.0] \\ \text{Protein} & [1.2 & 104.0 & 52.0 & 8.0] \\ \text{Fat} & [1.8 & 135.0 & 99.0 & 0.9] \end{matrix} = A_{(3,4)} \\ \downarrow^0 \quad \downarrow^1 \end{array}$$

59 cal  $\frac{56}{59} \approx 94.9\%$

Calculate % of calories from Carb, Protein, Fat. / Can you do this without explicit for-loop?  
→ 59 cal ≈ 94.9%, 104.0 cal.

cal = A.sum(axis = 0) ↗ 1x1 1x4 ↗ 모양 맞았으면 그냥 써자.  
percentage = 100\*A/(cal.sum(axis = 0)) ↗ 1x1 1x4 ↗ 비율 맞아 압니다.

$\uparrow (3,4) / (1,4)$

# Broadcasting example

L과 같은 차원이면 연산하는 행렬이 맞춰

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \xrightarrow{\text{각 줄은 행렬로 확장 가능}} \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = 100$$

그냥 +가 +  
끼고 각 +  
각 행의 모든 요소에  
적용

$3 \times 4$  matrix  $\times \frac{1}{2}$   $1 \times 4$  벡터로 43.

L A      L cal

$$= \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

각 줄은 4x1 벡터로 확장.  
L broadcasting

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m,n)} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}_{(1,n) \rightsquigarrow (m,n)} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

각 줄은 행렬로  
적용

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

각 줄은 행렬로  
적용

# General Principle

신경망을 조합한 때 사용하는 주요 broadcasting 유형.

행을 m번 복사

$$\begin{array}{c} (m, n) \\ \text{matrix} \\ \hline + \\ \times \\ / \end{array} \quad \begin{array}{c} (1, n) \\ \rightsquigarrow (m, n) \\ \text{열을 n번 복사} \\ (m, 1) \\ \rightsquigarrow (m, n) \end{array}$$

$$\begin{array}{ccc} (m, 1) & + & \mathbb{R} \\ \left[ \begin{smallmatrix} 1 \\ 2 \\ 3 \end{smallmatrix} \right] & + & 100 \\ \left[ \begin{smallmatrix} 1 & 2 & 3 \end{smallmatrix} \right] & + & 100 \end{array} = \begin{array}{c} (m, 1) \text{ matrix가 나올 때까지 } m\text{번 복사} \\ \left[ \begin{smallmatrix} 101 \\ 102 \\ 103 \end{smallmatrix} \right] \\ = \left[ \begin{smallmatrix} 101 & 102 & 103 \end{smallmatrix} \right] \end{array}$$

Matlab/Octave: bsxfun

similar as broadcasting.



deeplearning.ai

# Basics of Neural Network Programming

---

A note on python/  
numpy vectors

↑ 유동 but 불변 가능.

# Python Demo

# Python / numpy vectors

```
import numpy as np
```

```
a = np.random.randn(5)      a.shape(5, )      } Don't use.  
                            "rank 1 array"
```

```
a = np.random.randn((5,1))  a.shape = (5,1).  
                           col vector
```

```
a = np.random.randn((1,5))  a.shape = (1,5)  
                           row vector.
```

```
assert(a.shape = (5,1))
```

L 벡터의 dimension을 모르는 경우. (5,1) 벡터가 2D인.

rank 1 array인 경우. a=a.reshape ((5,1))

# Explanation of Logistic regression cost function

< 왜 이걸 써-트는지.

## Logistic regression cost function

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

Interpret  $\hat{y} = p(y=1|x)$   $\hat{y}$ 가  $x$ 에 대해  $y=1$ 인 확률.

$$\text{If } y=1 : p(y|x) = \hat{y}$$

$$\text{If } y=0 : p(y|x) = 1 - \hat{y} \quad y가 0인 확률.$$

## Logistic regression cost function

$y$ 는 반드시 0 혹은 1이어야 한다.

$\hat{y}$ 는  
반면에

$$\begin{aligned} &\rightarrow \boxed{\text{If } y=1: p(y|x) = \hat{y}} \\ &\rightarrow \boxed{\text{If } y=0: p(y|x) = 1 - \hat{y}} \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} p(y|x)$$

$$p(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)} \quad \leftarrow$$

$$\text{If } y=1: p(y|x) = \hat{y}^1 (1-\hat{y})^{(1-1)} = 1 \cdot (1-\hat{y}) = \hat{y}$$

$$\text{If } y=0: p(y|x) = \hat{y}^0 (1-\hat{y})^{(1-0)} = 1 \cdot (1-\hat{y}) = 1 - \hat{y}$$

$$\uparrow \log p(y|x) = \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y})$$

$$= -f(\hat{y}, y) \quad \downarrow \text{Loss 함수의 바이너스 값과 동일.} \rightarrow \text{Loss는 최소화}$$

$\log$  함수는 감소함수이고 증가하는 함수 (부등식으로)  
 $\log(p(y|x))$ 를 최대화하면 optimizing  $p(y|x)$ 와 비슷한 결과.)

Andrew Ng  
제로 일정 것.

최종의  $\log$ 을 최대화  
하는 것과 동일.

최종 속성의 최적화 → 그동안 만들어주는 parameter를 찾아야.

# Cost on $m$ examples

예측 가능한 곳?

$$\frac{\log p(\text{labels in training set})}{\log p(\dots)} = \log \prod_{i=1}^m p(y^{(i)} | x^{(i)}) \leftarrow$$

$$\begin{aligned} \log p(\dots) &= \sum_{i=1}^m \underbrace{\log p(y^{(i)} | x^{(i)})}_{-L(\hat{y}^{(i)}, y^{(i)})} \\ &= -\sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \end{aligned}$$

$$\text{Cost: } J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$\Rightarrow J(w, b)$ 를 최소화하여 고리스틱 최적화법을 사용해 최대 확률 추정법 적용.

Andrew Ng