



Setting up your optimization problem

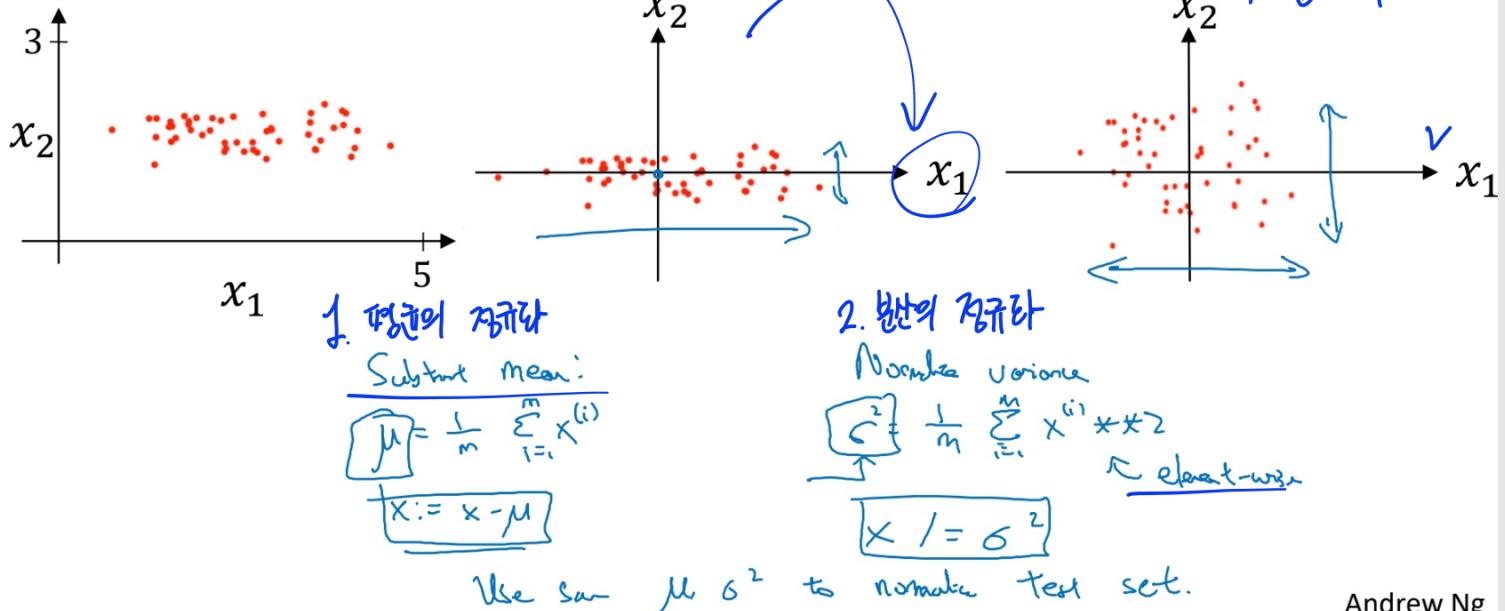
Normalizing inputs

deeplearning.ai

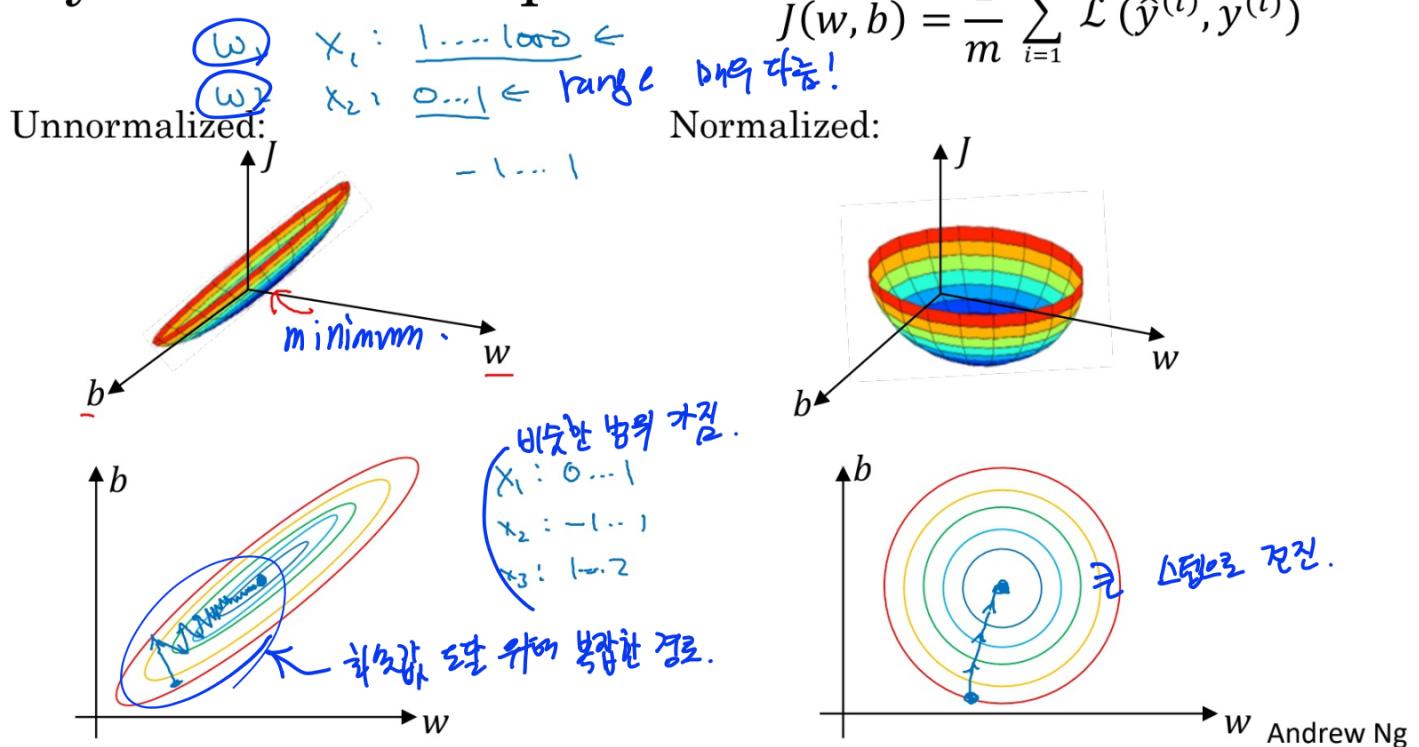
입력의 정규화.

Normalizing training sets

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$



Why normalize inputs?





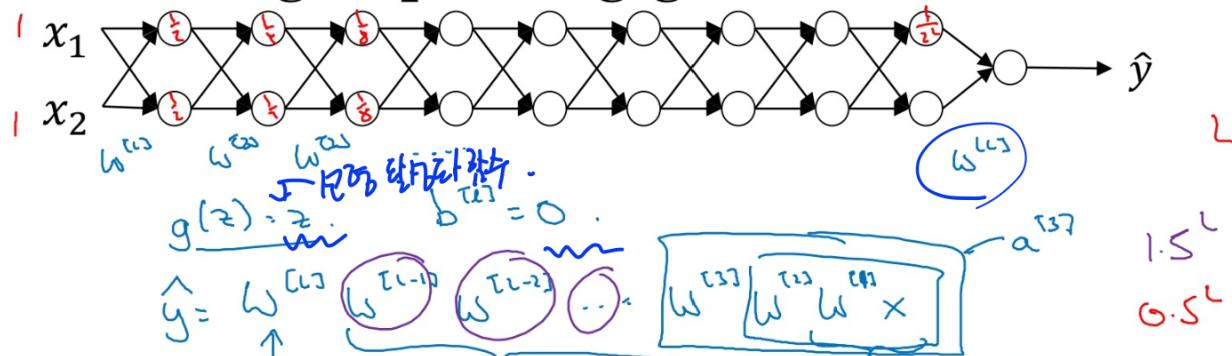
Setting up your optimization problem

Vanishing/exploding gradients

대부분의 경우 신경망을 훈련할 때
미분값은 매우 작아서 커질 수 있다. → 무한대 가중치 초기화

$$L = 150 \dots 832.$$

Vanishing/exploding gradients



$$w^{(l)} > 1$$

$$w^{(l)} < 1 \quad [0.9 \quad 0.9]$$

$$w^{(l)} = \begin{bmatrix} 0.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

$$\hat{y} = w^{(l)} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^{-1} x$$

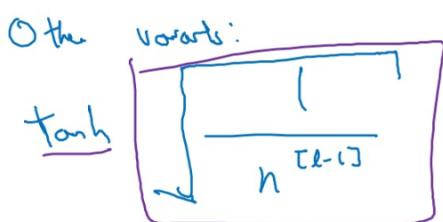
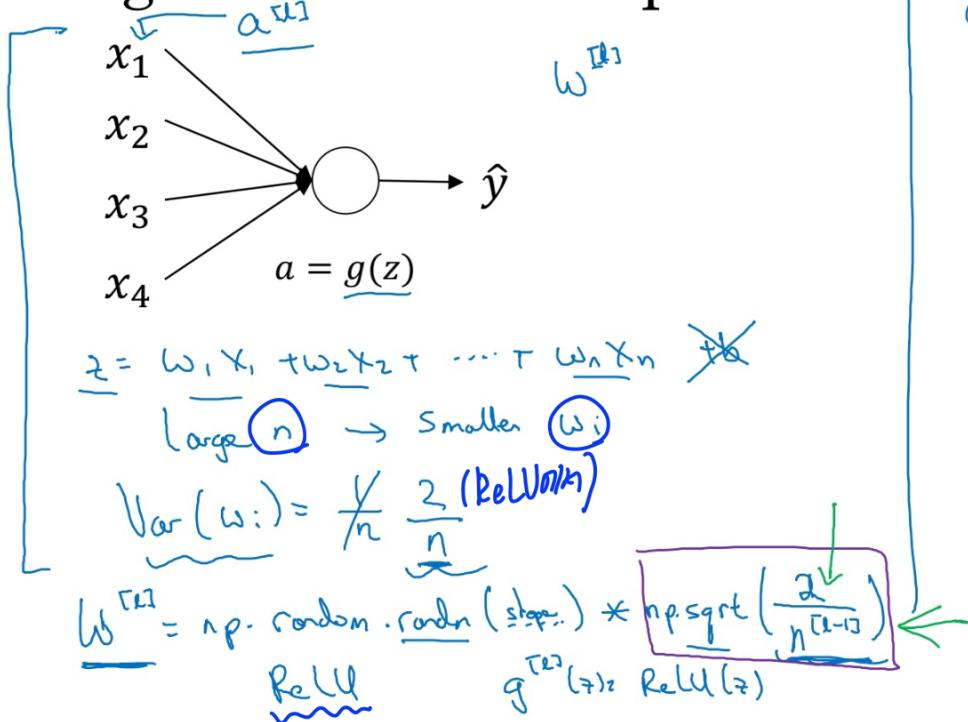
$$\begin{aligned} z^{(l)} &= w^{(l)} x \\ a^{(l)} &= g(z^{(l)}) = z^{(l)} \\ a^{(l)} &= g(z^{(l)}) = g(w^{(l)} a^{(l-1)}) \end{aligned}$$

$$1.5 \times \leftarrow \text{폭발} \quad 0.5 \times \leftarrow \text{소멸.}$$

Andrew Ng

Weight initialization for Deep network

Single neuron example



Xavier initialization

$$\frac{2}{n^{(l-1)} + n^{(l)}}$$

↑

Andrew Ng



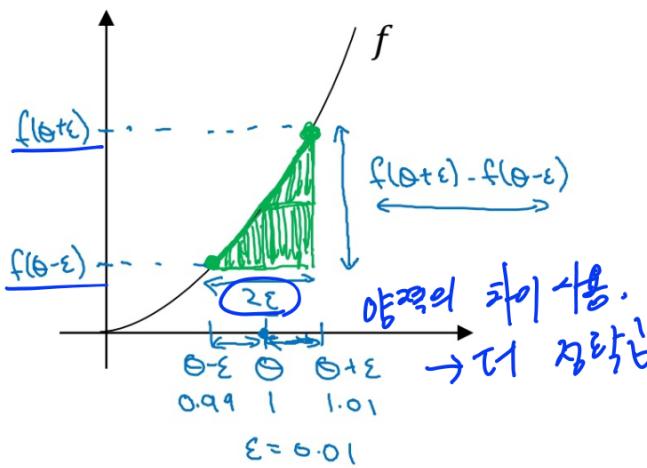
deeplearning.ai

Setting up your optimization problem

Numerical approximation of gradients

Checking your derivative computation

$$f(\theta) = \theta^3$$



$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta)$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = \underline{\underline{3.0001}} \approx 3$$

$$g(\theta) = 3\theta^2 = \underline{\underline{3}}$$

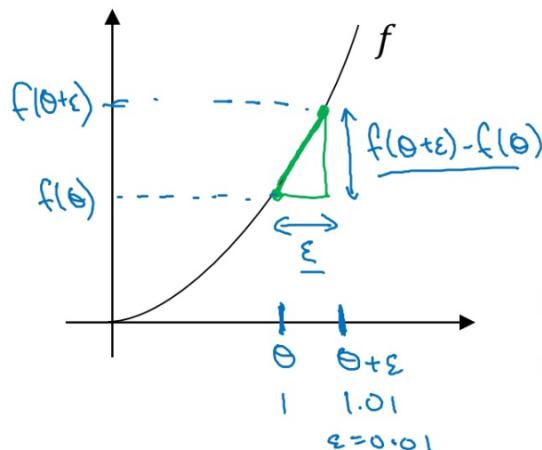
approx error: 0.0001
(prev slide: 3.0301. error: 0.03)

$$\left\{ f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \right. \quad \begin{matrix} O(\epsilon^2) \\ 0.01 \\ 0.0001 \end{matrix} \quad \left| \quad \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \right. \quad \text{error: } O(\frac{1}{\epsilon}) \quad \begin{matrix} 0.01 \end{matrix}$$

Andrew Ng

Checking your derivative computation

$$\boxed{f(\theta) = \theta^3} \quad \theta \in \mathbb{R}$$



$$g(\theta) = \frac{d}{d\theta} f(\theta) = f'(\theta)$$

$$\boxed{g(\theta) = 3\theta^2.}$$

$$g(\theta) = 3 \cdot (1)^2 = 3 \quad \text{when } \theta = 1$$

$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \approx g(\theta)$$

$$\frac{(1.01)^3 - 1^3}{0.01} = \underline{\underline{3.0301}} \approx \underline{\underline{3}}$$

$\theta = 1$
 $\theta + \epsilon = 1.01$

0.0301
3.1
3.2

Andrew Ng



deeplearning.ai

Setting up your optimization problem

Gradient Checking implementation notes

Gradient check for a neural network

큰 차이의 미지 변수 →

Take $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$ and reshape into a big vector θ .

concatenate

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$$

Take $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$.

concatenate

Is $d\theta$ the gradient of $J(\theta)$?

비교로 정답과 주면

Gradient checking (Grad check)

for each i :

$$\rightarrow \underline{d\theta_{\text{approx}}[i]} = \frac{J(\theta_0, \theta_1, \dots, \theta_i + \epsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\approx \underline{d\theta[i]} = \frac{\partial J}{\partial \theta_i}$$

$d\theta_{\text{approx}}$ ≈ $d\theta$
θ의 각도 차이.

Check
두 벡터의
가까운지 여부

$$\rightarrow \frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$$

$\epsilon = 10^{-7}$

$$\times \left[\frac{10^{-7}}{10^{-5}} - \text{great!} \right] \leftarrow$$

$\rightarrow 10^{-3}$ - worry. ←
너그의 가능성의 상한은 많.

Andrew Ng



deeplearning.ai

Setting up your optimization problem

Gradient Checking

Gradient checking implementation notes

- Don't use in training – only to debug

$$\frac{\partial \theta_{approx}[i]}{\uparrow} \longleftrightarrow \frac{\partial \theta[i]}{\uparrow}$$

↙ 매우 느리다!

- If algorithm fails grad check, look at components to try to identify bug.

$$\frac{\partial b^{[l]}}{\uparrow} \quad \frac{\partial w^{[l]}}{\uparrow}$$

$$J(\theta) = \frac{1}{m} \sum_i \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_l \|w^{(l)}\|_F^2$$

$\delta\theta = \text{gradt of } J \text{ wrt. } \theta$

- Remember regularization.

- Doesn't work with dropout.

$$J \quad \underline{\text{keep-prob} = 1.0}$$

- Run at random initialization; perhaps again after some training.

$$\underline{w, b \approx 0}$$

Andrew Ng