

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



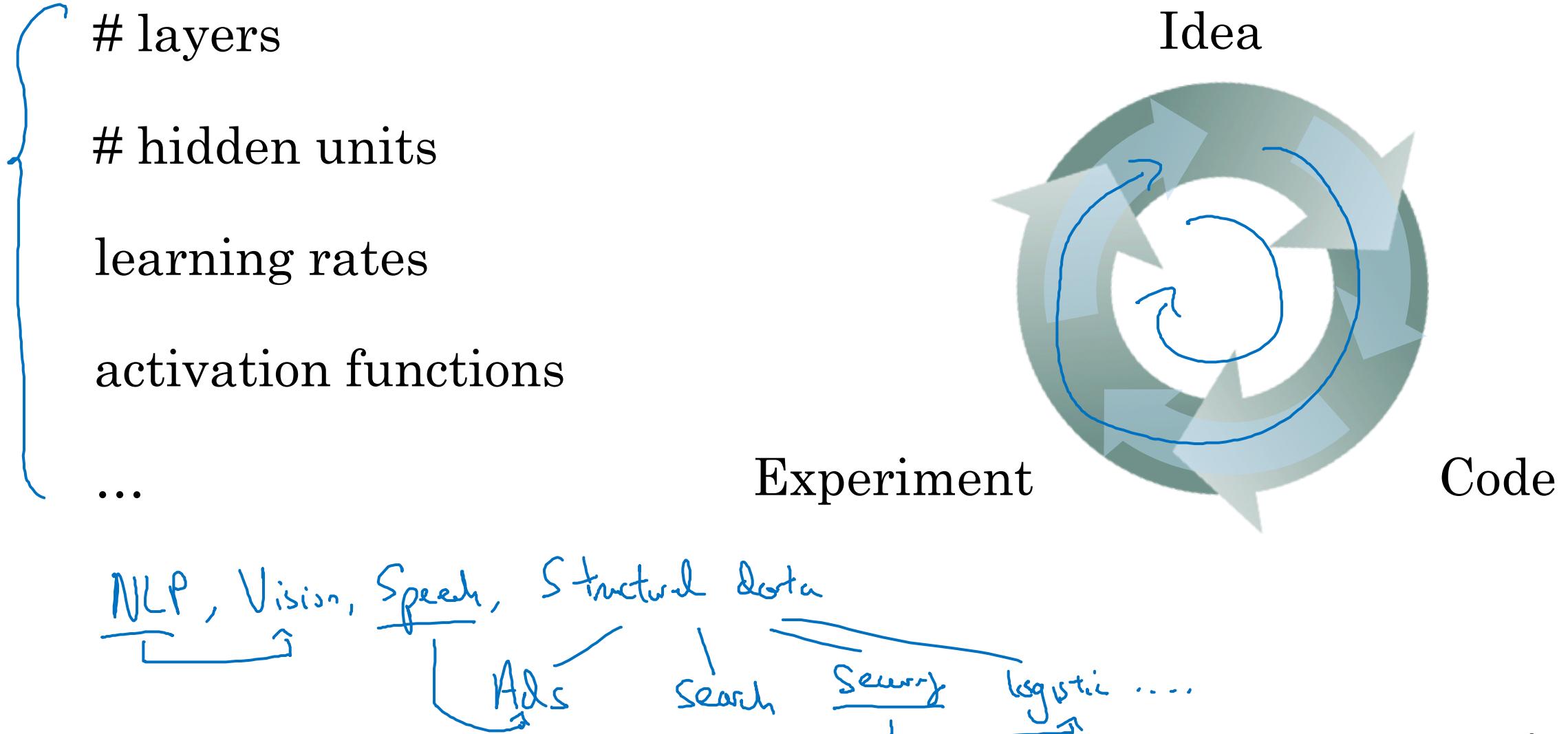
deeplearning.ai

Setting up your  
ML application

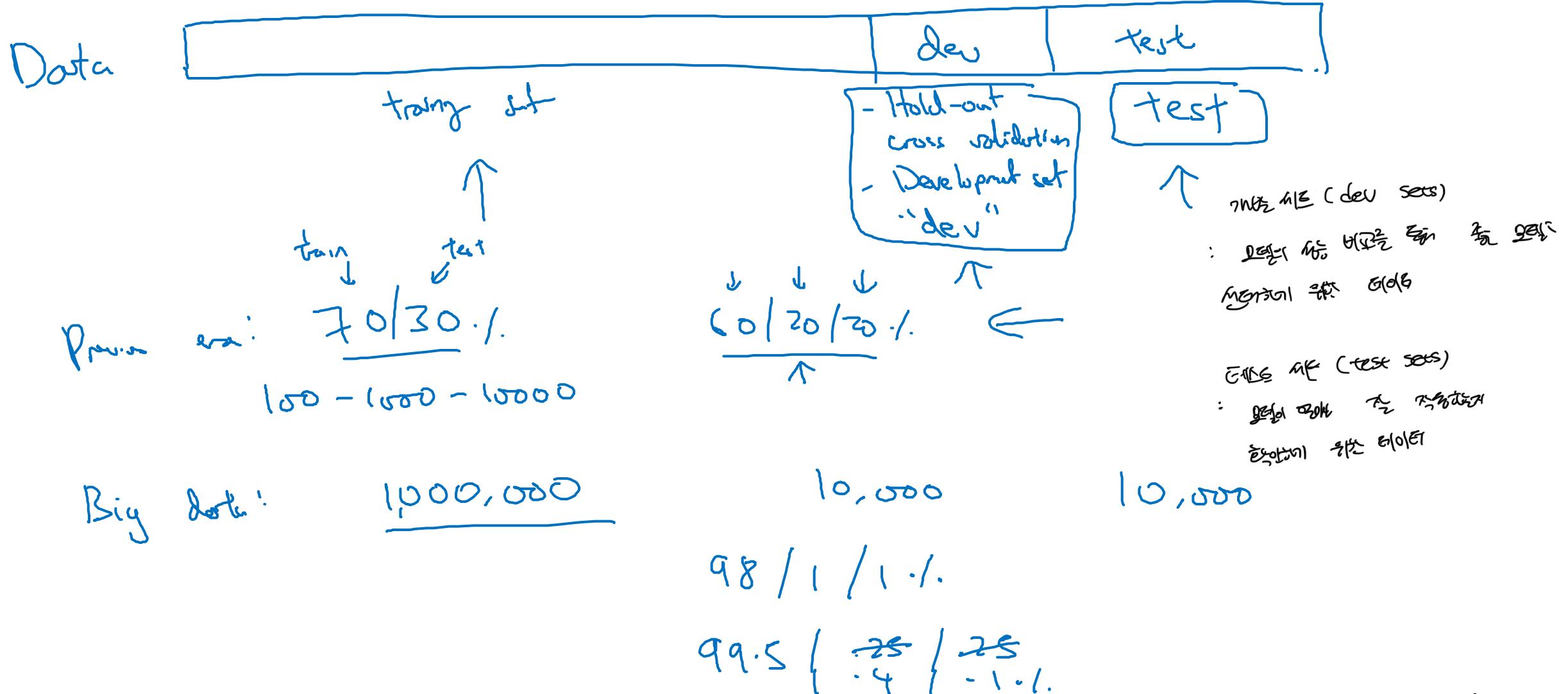
---

Train/dev/test  
sets

# Applied ML is a highly iterative process



# Train/dev/test sets



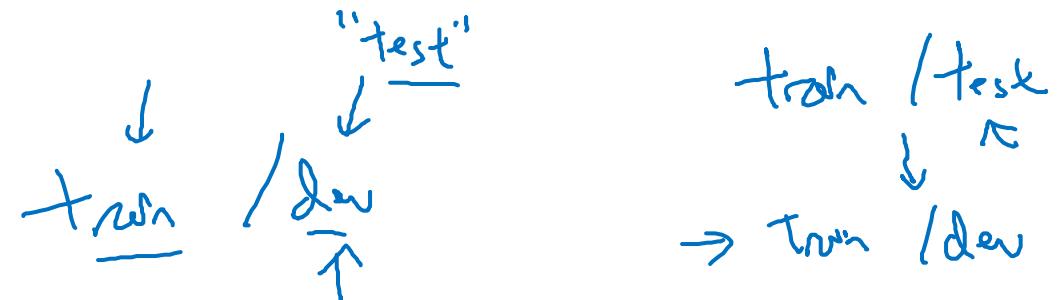
# Mismatched train/test distribution

Conts

Training set:  
Cat pictures from }  
webpages

Dev/test sets:  
Cat pictures from }  
users using your app

→ Make sure dev and test come from same distribution.



Not having a test set might be okay. (Only dev set.)



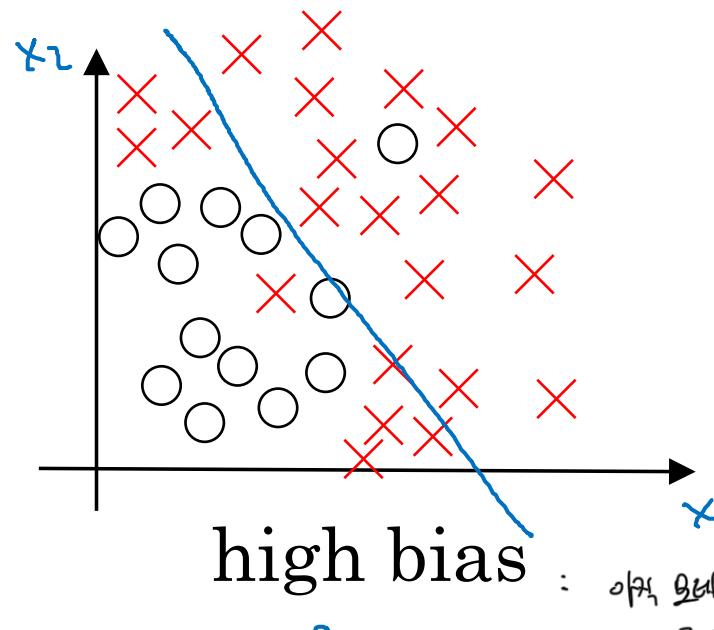
deeplearning.ai

Setting up your  
ML application

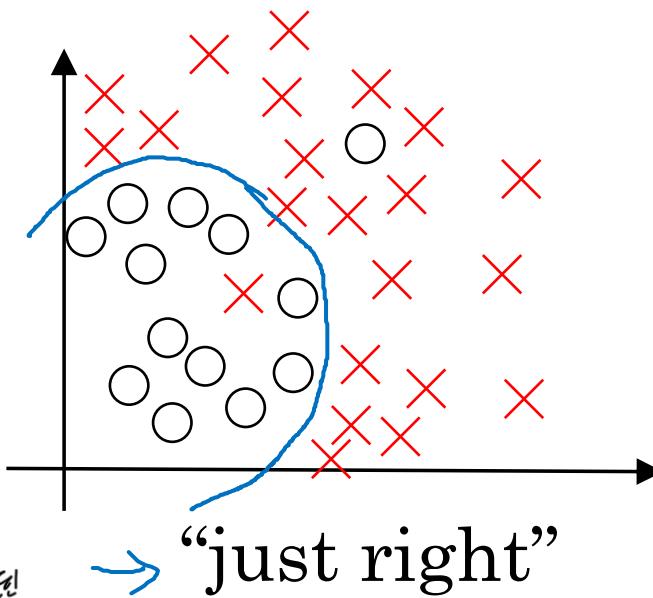
---

Bias/Variance

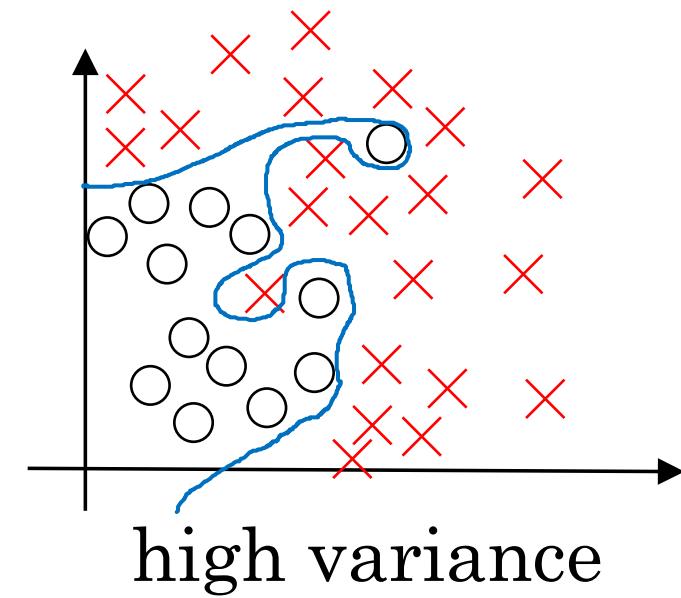
# Bias and Variance



Bias:  $\text{편향}$   $\text{직선의 거리}$   
Variance:  $\text{예측값의 범위}$   $\text{거리}$



→ “just right”



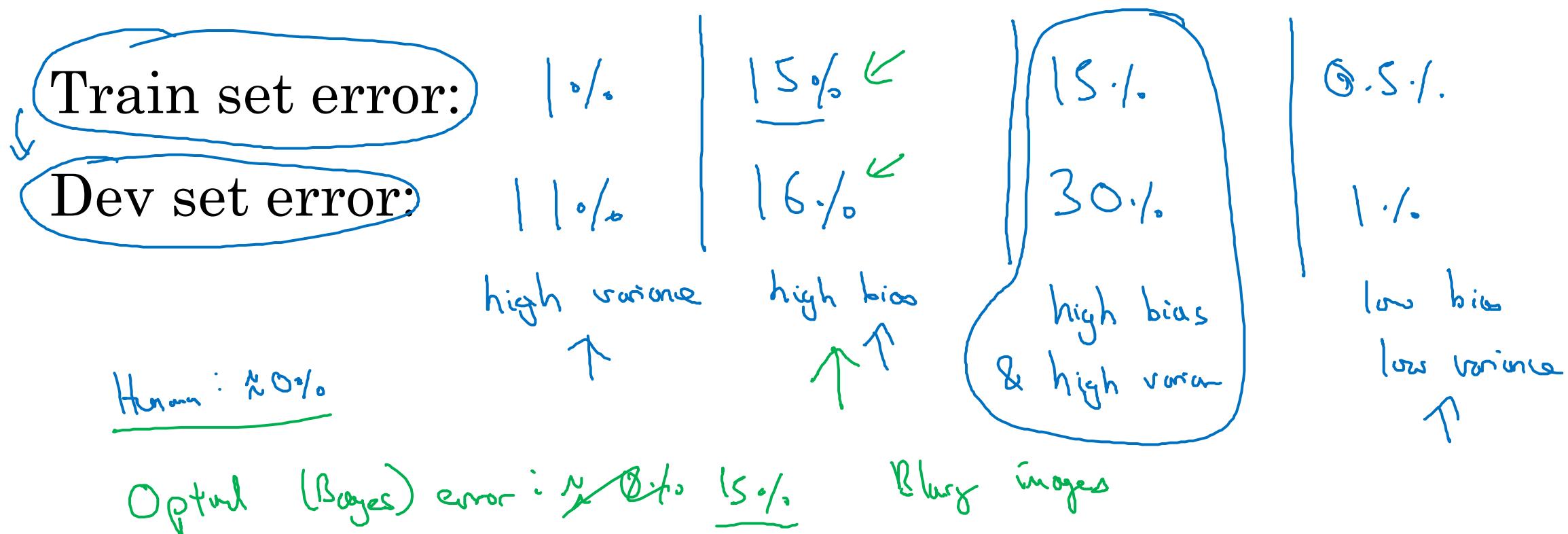
high variance

Overfitting

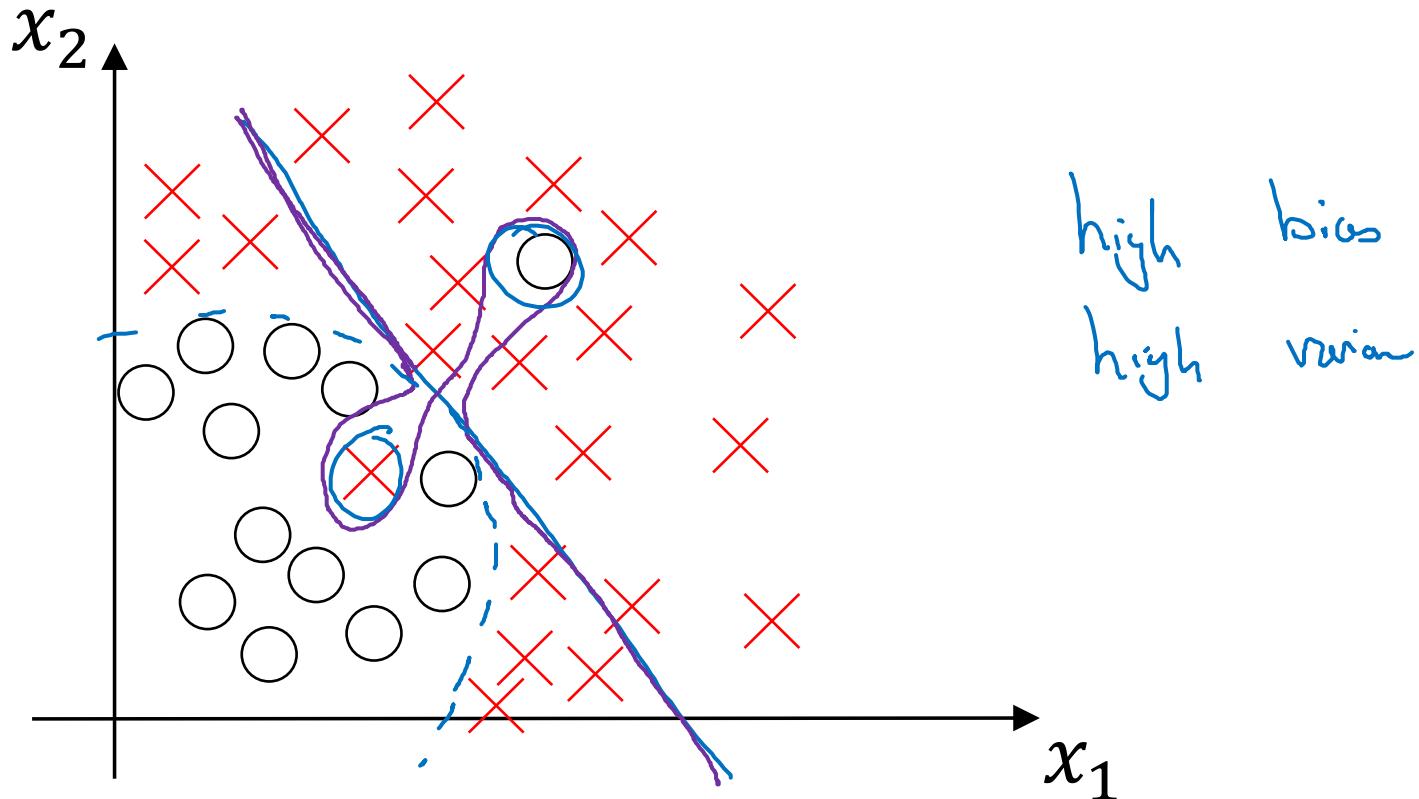
=  $\text{과도한 편향}$   $\text{직선의 거리}$   
 $\text{점마다 거리}$

# Bias and Variance

## Cat classification



# High bias and high variance





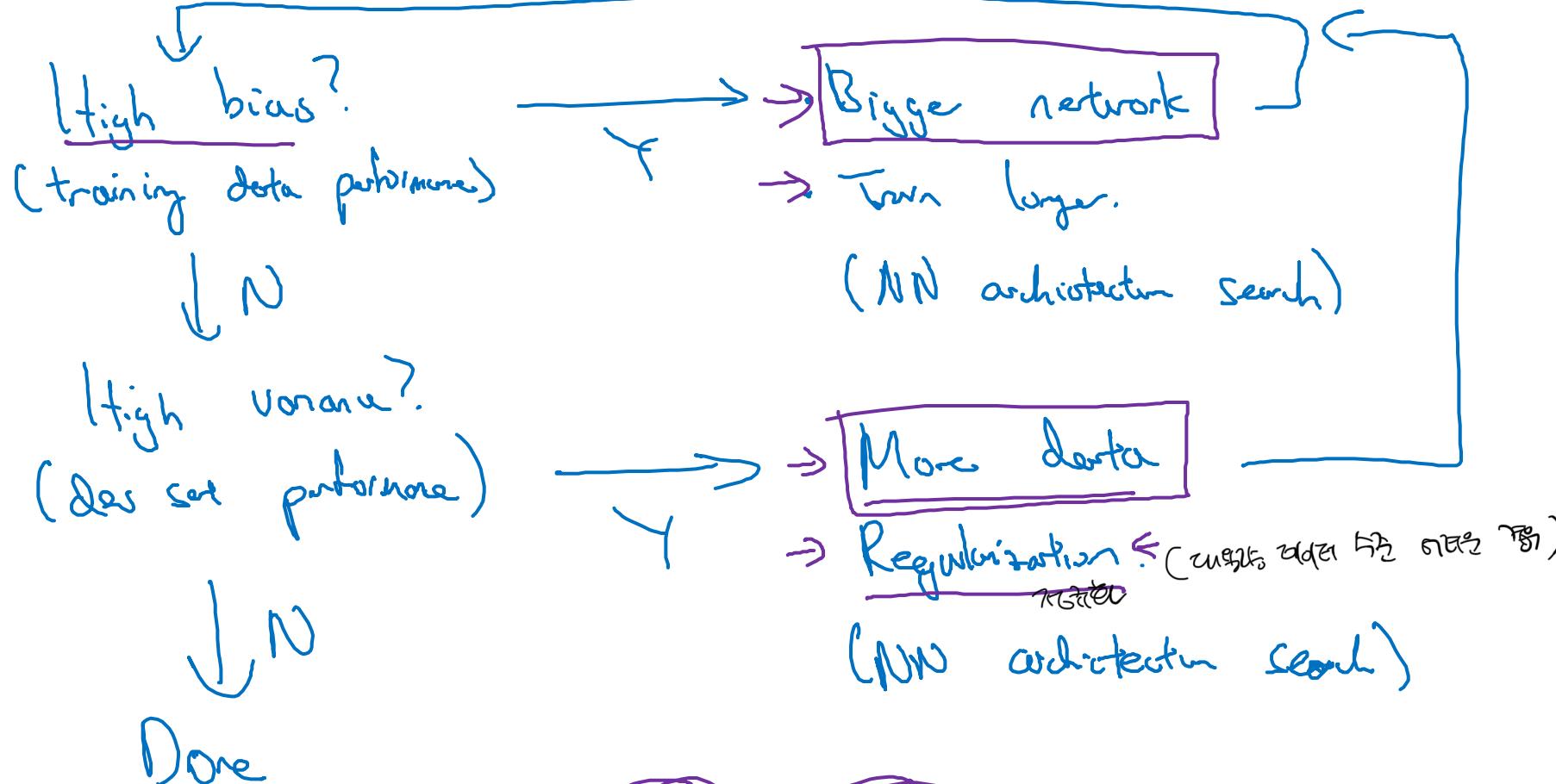
deeplearning.ai

# Setting up your ML application

---

## Basic “recipe” for machine learning

# Basic recipe for machine learning





deeplearning.ai

Regularizing your  
neural network

---

Regularization

# Logistic regression

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

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

$\lambda$  = regularization parameter  
 lambda  
 lambd

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

L<sub>2</sub> regularization

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

L<sub>1</sub> regularization

$$\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$$

~~$$+ \frac{\lambda}{2m} b^2$$~~

:  $w$  ( $L_2$ )를 계산할 때, 정규화를 사용하여 얻는 이유는  
 $w$ 는 항상  $0$ 이 아니고,  $b$ 는 항상  $0$ 이 아니기 때문

$w$  will be sparse

무한远处 0과 1의 결합으로 sparse  
 예제로는 0과 1의 결합으로 sparse.  
 가능합니다. 그 외에 다른 것은 없습니다.  
 유용한 경우는 그 이상입니다.  
 예제로는 0과 1의 결합으로 sparse.  
 특히 weight.

# Neural network

$$\rightarrow J(w^{(0)}, b^{(0)}, \dots, w^{(L)}, b^{(L)}) = \underbrace{\frac{1}{m} \sum_{i=1}^m f(\hat{y}^{(i)}, y^{(i)})}_{n^{(0)} \times n^{(L-1)}} + \underbrace{\frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2}_{\text{regularization}}$$

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^{n^{(l)}} \sum_{j=1}^{n^{(l-1)}} (w_{ij}^{(l)})^2$$

$w^{(l)} \in (n^{(l)}, n^{(l-1)})$

$w^{(l)}: (n^{(l)}, n^{(l-1)})$   
 $\uparrow \quad \uparrow$

"Frobenius norm"

$$\|\cdot\|_2^2 \quad \|\cdot\|_F^2$$

$$dW^{(l)} = \boxed{(\text{from backprop}) + \frac{\lambda}{m} w^{(l)}}$$

$$\rightarrow w^{(l)} := w^{(l)} - \alpha dW^{(l)}$$

$$\frac{\partial J}{\partial w^{(l)}} = dw^{(l)}$$

$$w^{(l)} := w^{(l)} - \alpha \left[ (\text{from backprop}) + \frac{\lambda}{m} w^{(l)} \right]$$

L2: "Weight decay" at 50% off

$\therefore \text{weight of } (1 - \frac{\alpha \lambda}{m}) \text{ of } ((1 - \frac{\alpha \lambda}{m}) \text{ of } \dots)$

$$\begin{aligned} w^{(l)} &= w^{(l)} - \frac{\alpha \lambda}{m} w^{(l)} - \alpha (\text{from backprop}) \\ &= \underbrace{\left(1 - \frac{\alpha \lambda}{m}\right)}_{<1} \underbrace{w^{(l)}}_{>1} - \alpha (\text{from backprop}) \end{aligned}$$



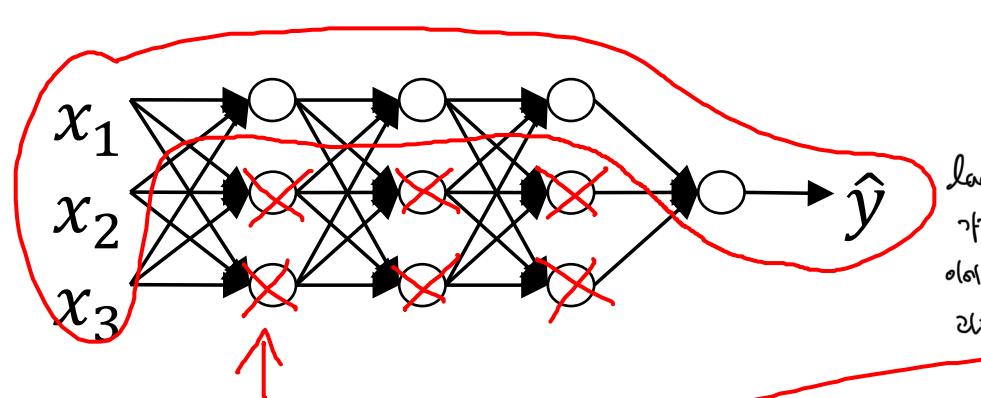
deeplearning.ai

# Regularizing your neural network

---

## Why regularization reduces overfitting

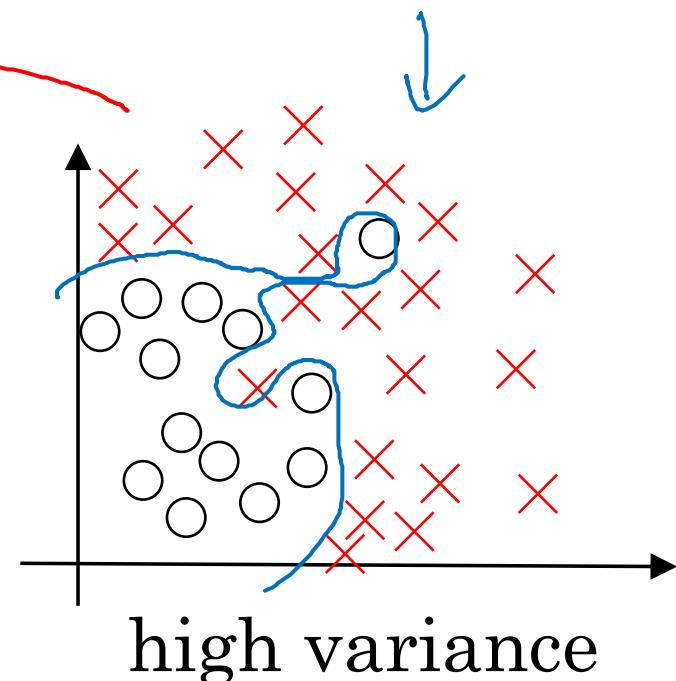
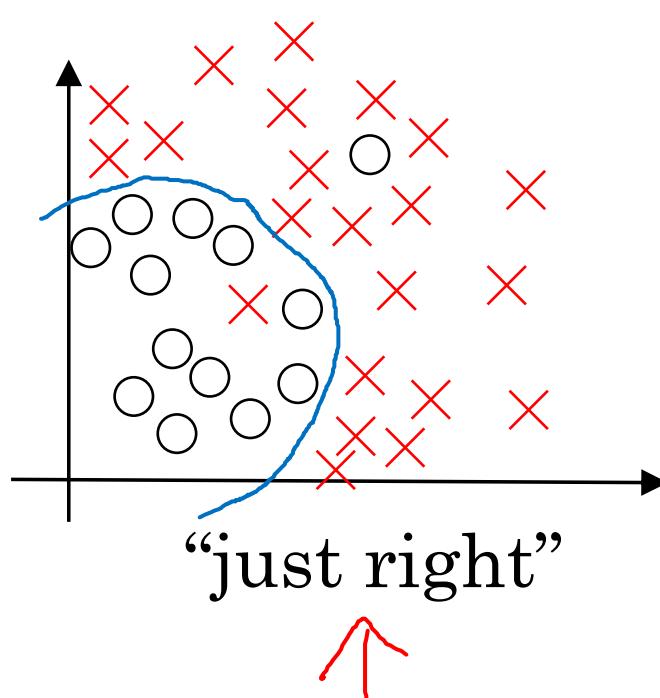
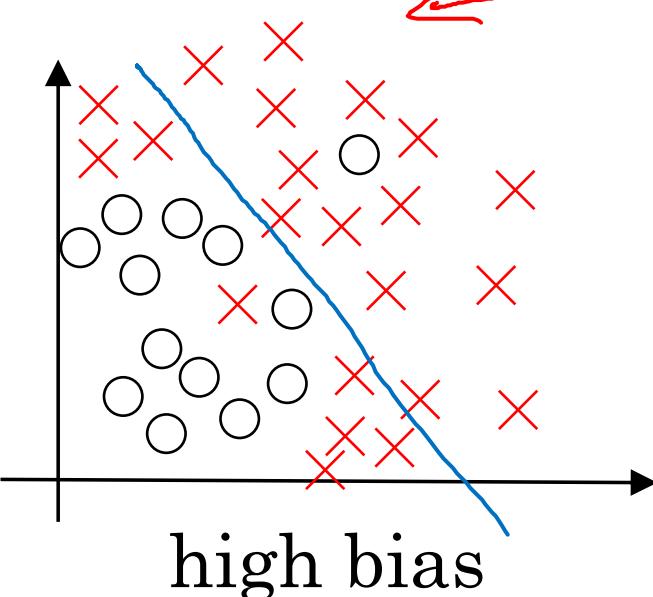
# How does regularization prevent overfitting?



$$J(\boldsymbol{w}^{(1)}, \boldsymbol{b}^{(1)}) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\boldsymbol{w}^{(l)}\|_F^2$$

Handwritten notes in Korean:

- lambda 는 정제화에 기여하는  $\boldsymbol{w}^{(1)}$ 에 대한 오류에 대한 허용 범위를 줄여주는 것 같다.
- hidden unit 레이어에 대해서는 오류를 줄여주는 역할을 한다.
- 이제  $\boldsymbol{w}^{(1)} \approx 0$ 이다.



# How does regularization prevent overfitting?

tanh 흡수를 사용하면 (lambda 값이 커지면) 과적합 방지 가능  
가중치  $w$ 는 제한된 범위,  $z = w \cdot x + b$  이 약간 조정되는 듯.

결과로 일정한 평균을  $\lambda \rightarrow 0$ 로

$w$ 의 범위는 제한된 범위로 제한된다.

자신 부위를 제한하는 듯, tanh 흡수는

제한된 범위를 확장한다.

$\Rightarrow$  Activation function  $\approx$  linear when decision of separation

$$\lambda \uparrow$$

$$w^{[l]} \downarrow$$

overfitting  $\approx$  A 문제 같다.

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

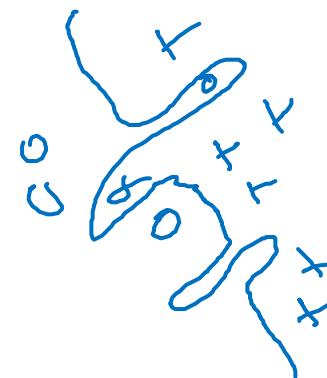
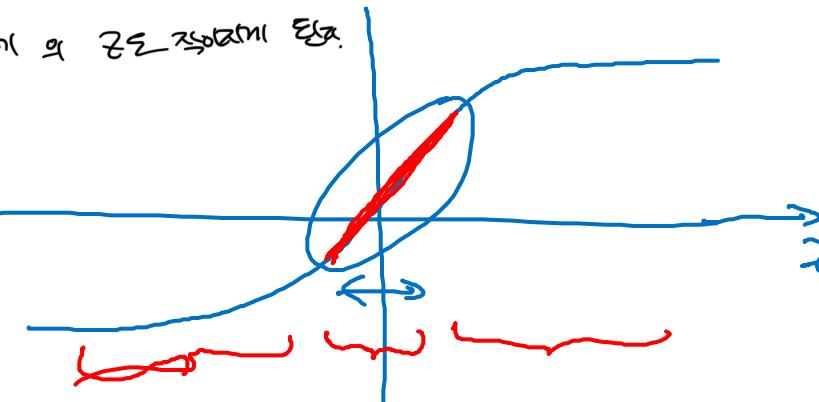
Every layer  $\approx$  linear.

$$J(\dots) = \left[ \sum_i L(\hat{y}^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2m} \sum_l \|w^{[l]}\|_F^2$$



tanh

$g(z) = \tanh(z)$





deeplearning.ai

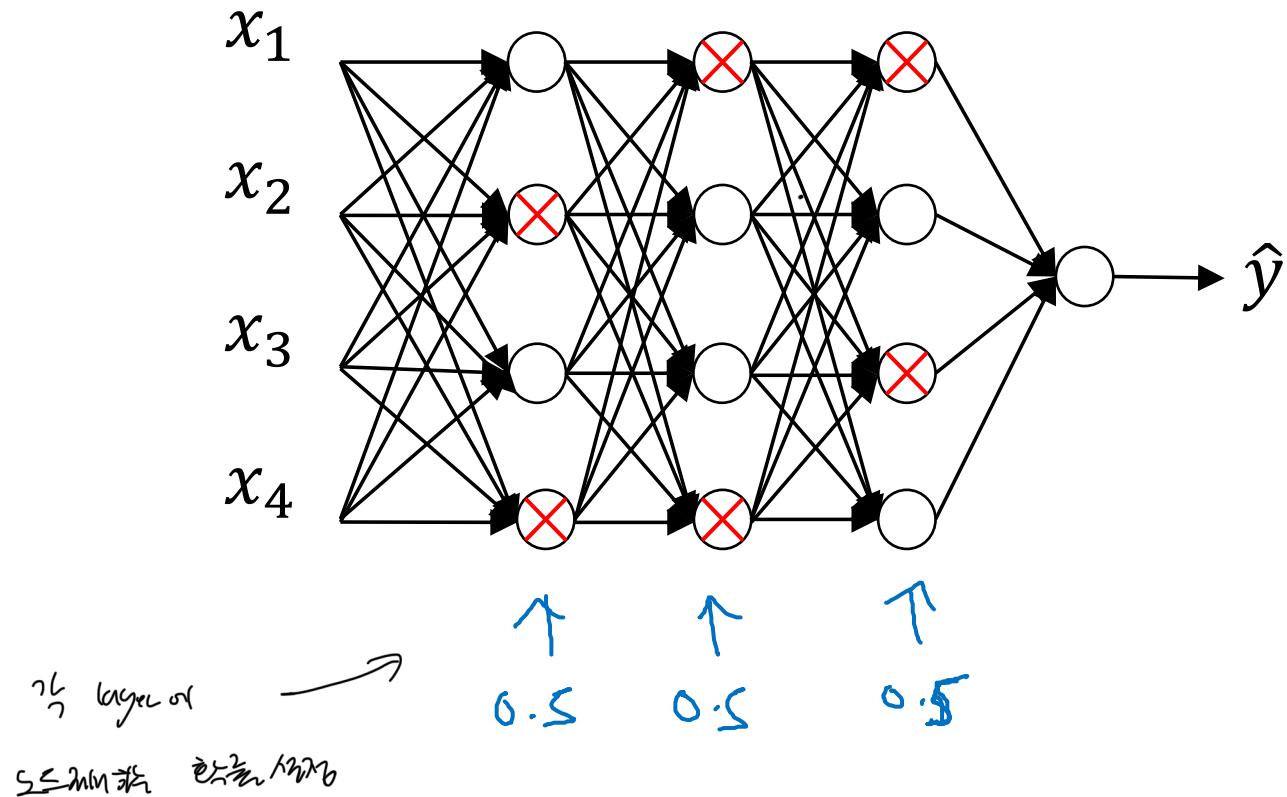
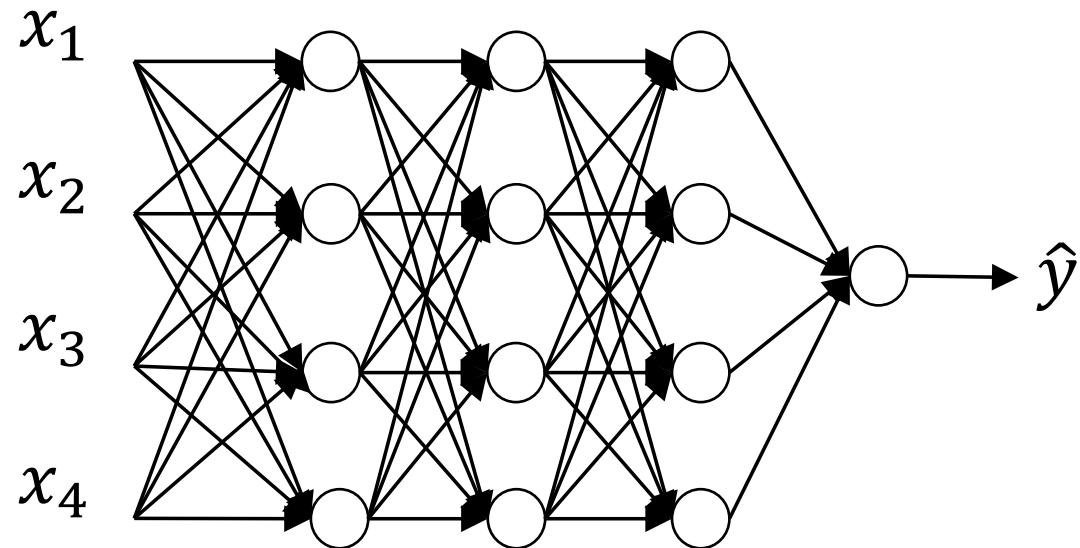
Regularizing your  
neural network

---

Dropout  
regularization

# Dropout regularization

: 신경망 구조의 종류에 따라 노드를 선택하는 확률 ( $(1 - \text{keep\_prob})$ ) 을 적용하는 것  
이다. 선택된 노드는 학습과 예측 과정으로 신경망에서 선택된 노드의 풍자로  
제거와 함께 그 위치를 모두 재설정한다.  
이에 따라 가중치 모델이 다른 차고 간접화된 신경망과 같은  
아트워크를 훈련을 진행하는데 유용



# Implementing dropout (“Inverted dropout”)

Illustrate with layer  $l=3$ .  $\text{keep-prob} = \frac{0.8}{x}$   $\underline{\underline{0.2}}$

$$\rightarrow d_3 = \underbrace{\text{np.random.rand}(a_3.shape[0], a_3.shape[1]) < \text{keep-prob}}_{d_3}$$

$$\underbrace{a_3}_{a_3} = \text{np.multiply}(a_3, d_3) \quad \# a_3 * d_3.$$

$$\rightarrow \underbrace{a_3 /=\cancel{\text{keep-prob}}}_{\downarrow} \leftarrow$$

50 units.  $\rightsquigarrow$  10 units shut off

여기서 Inverted dropout은  
드롭아웃과 같은 역할을 한다.

단계 3에서 했던 것과 같은 역할을 한다.

단계 4(keep-prob)는 훈련과 테스트

단계 5(정규화)와 같은 역할을 한다.

단계 6(정규화)와 같은 역할을 한다.

$$z^{(4)} = w^{(4)} \cdot \frac{a^{(3)}}{\cancel{c}} + b^{(4)}$$

$\cancel{c}$  reduced by 20%.

Test

$$1 = \underline{\underline{0.8}}$$

# Making predictions at test time

$$a^{(0)} = X$$

No drop out.

$$\uparrow z^{(1)} = W^{(1)} \underline{a^{(0)}} + b^{(1)}$$

$$a^{(1)} = g^{(1)}(z^{(1)})$$

$$z^{(2)} = W^{(2)} \underline{a^{(1)}} + b^{(2)}$$

$$a^{(2)} = \dots$$

$$\downarrow \hat{y}$$

$\lambda$  = keep-prob



deeplearning.ai

Regularizing your  
neural network

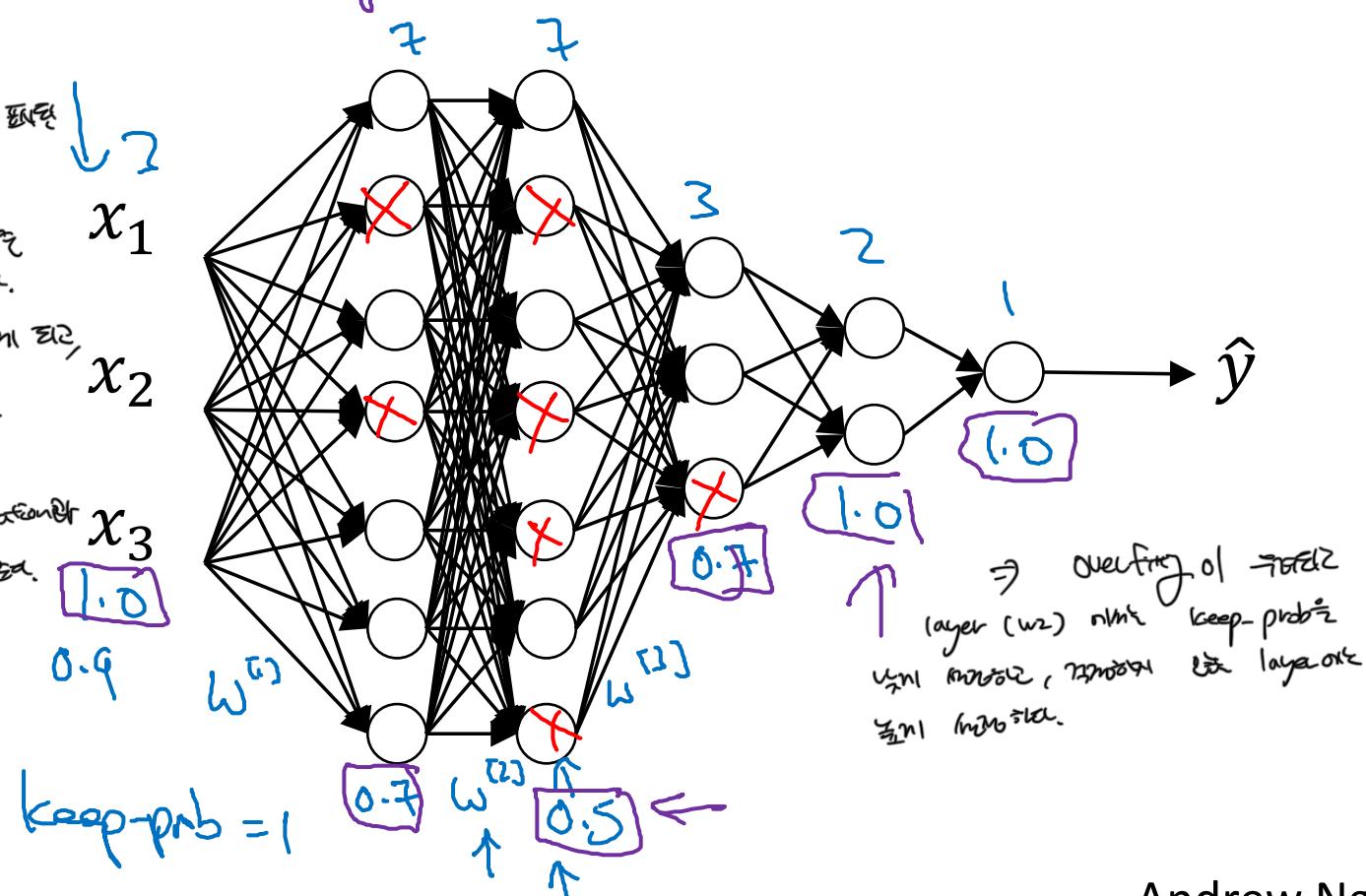
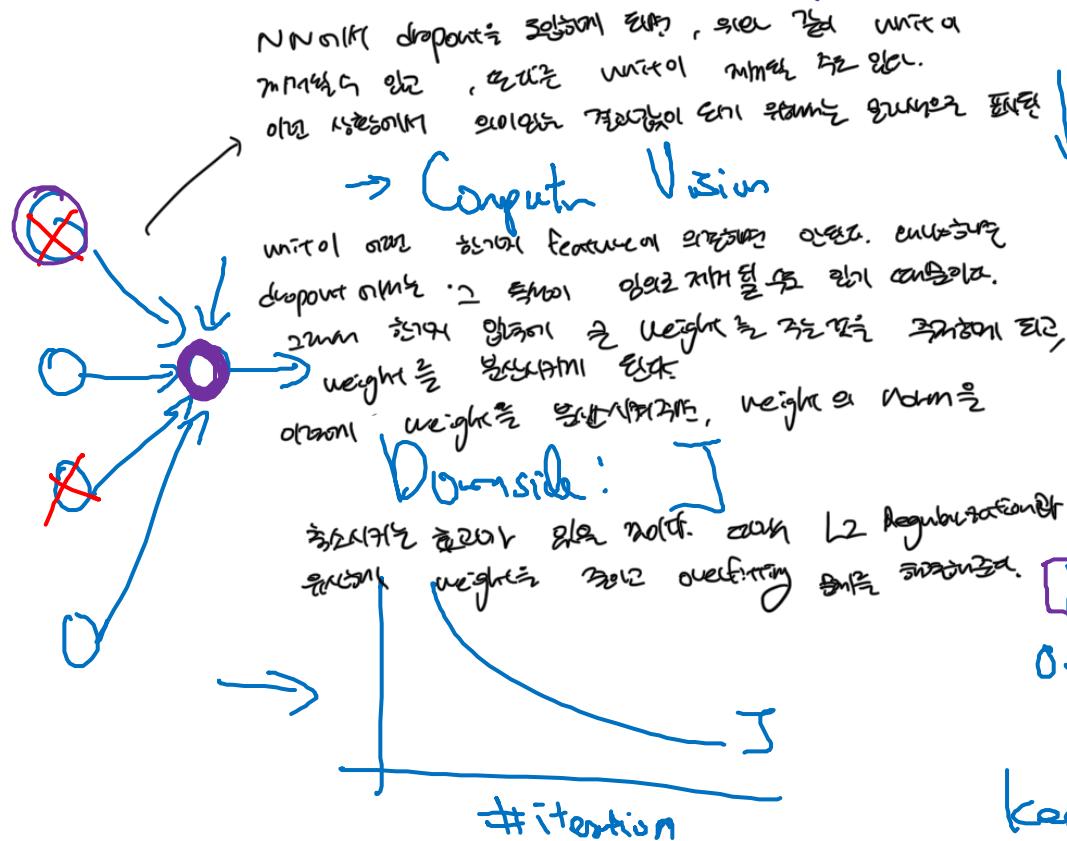
---

Understanding  
dropout

# Why does drop-out work?

: 학습률을 낮춰 차게 만들면, 학습의 흐름이 일정하게 유지된다  
가중치를 다른 끝에 봉인시키는 효과가 있다.

Intuition: Can't rely on any one feature, so have to spread out weights.  $\rightsquigarrow$  Shrink weights.  $b_2$



< Dropout의 역사>

1. Cross Validation이 AI 학습 성과를 향상시키는 데에 기여한 점이다.
2. Cost Function J가 대체로 미분 가능한 점이다.
3. 디버깅 코드를 풀면서 gradient descent, Gradient Descent 알고리즘을 개선하기 어렵다.



deeplearning.ai

Dropout은 Computer Vision에서 성공했다.

Computer Vision에서 성공한 것은 바로 커널,  
필터 편집을 이용하는데 그 이유는 복잡한 계산 때문이다.

# Regularizing your neural network

---

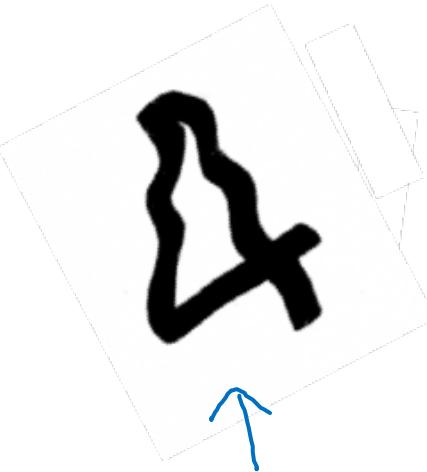
## Other regularization methods

# Data augmentation

↑ 사진 데이터를 확장하는 기법을  
는 자동화된 **augmentation**을  
하는 방법이다.

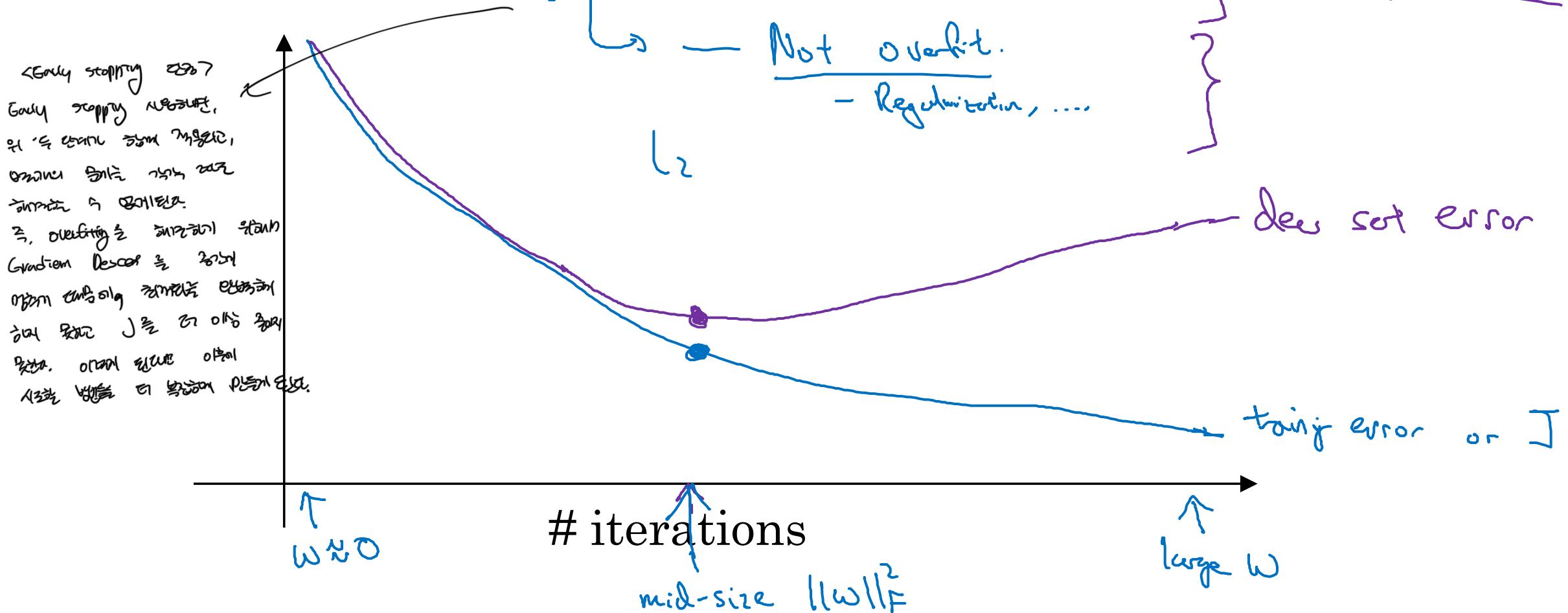


4



# Early stopping

• 학생들은 어떤 학생인가? 그는 어떤 학생인가? 그는 어떤 학생인가?  
그는 어떤 학생인가? 그는 어떤 학생인가? 그는 어떤 학생인가?  
그는 어떤 학생인가? 그는 어떤 학생인가? 그는 어떤 학생인가?





deeplearning.ai

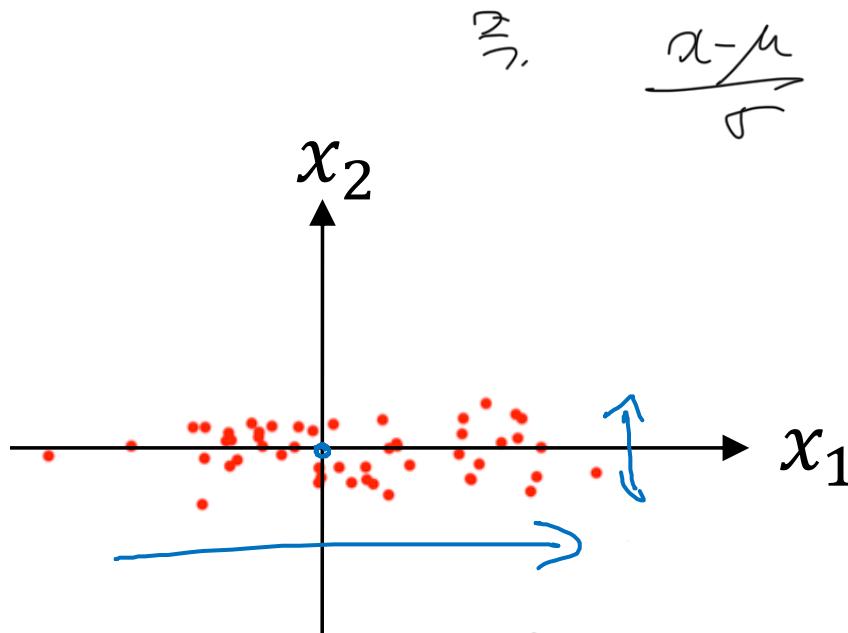
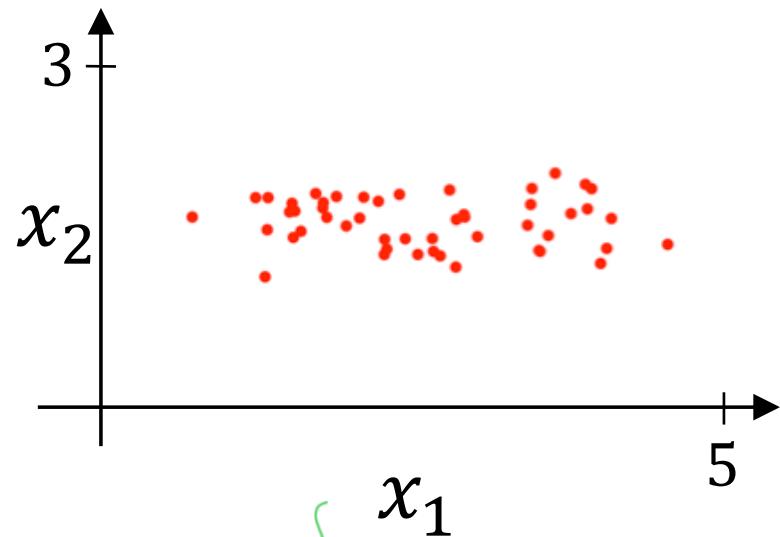
Setting up your  
optimization problem

---

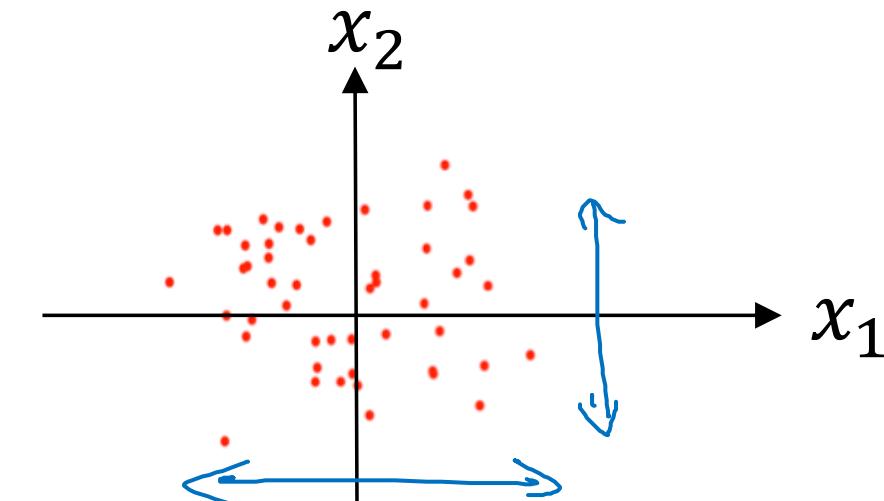
Normalizing inputs

# Normalizing training sets

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



∴  $\frac{x - \mu}{\sigma}$  *(Normalized)*

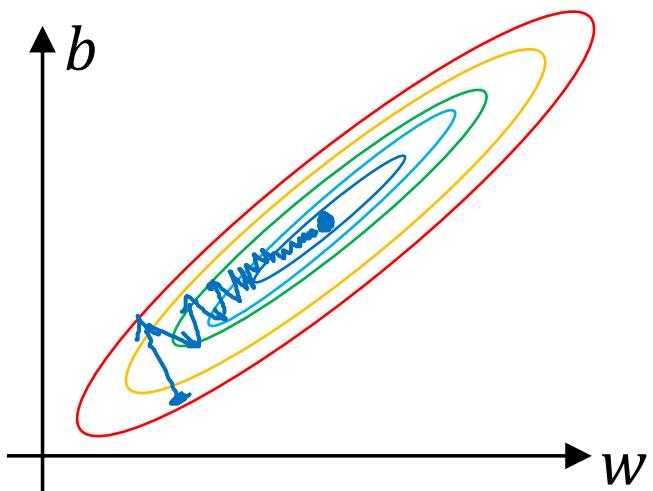
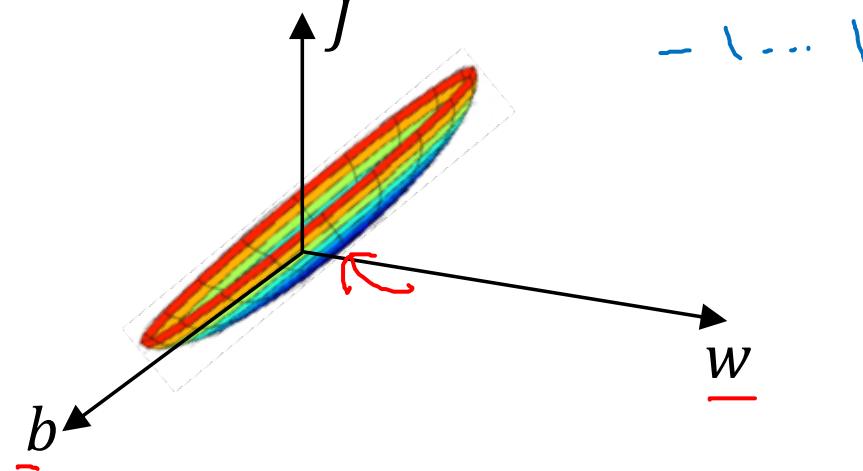


Use same  $\mu, \sigma^2$  to normalize test set.

# Why normalize inputs?

$\omega_1 \quad x_1: \frac{1 \dots 1000}{0 \dots 1} \leftarrow$   
 $\omega_2 \quad x_2: \frac{0 \dots 1}{-1 \dots 1} \leftarrow$

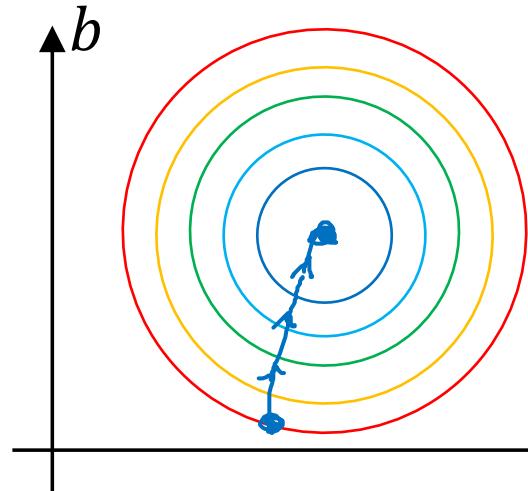
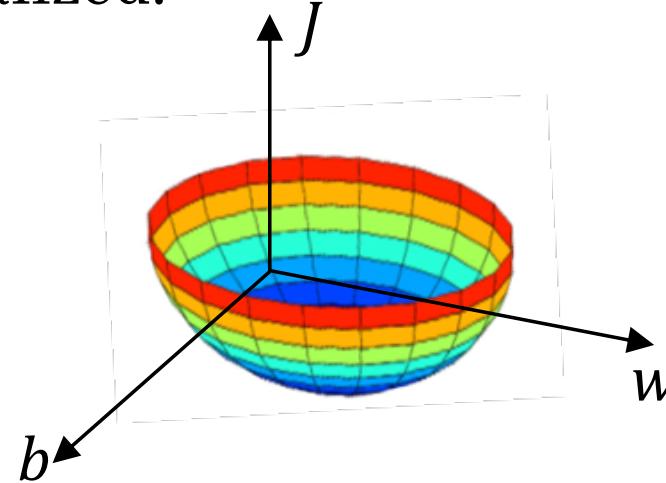
Unnormalized:



$x_1: 0 \dots 1$   
 $x_2: -1 \dots 1$   
 $x_3: 1 \dots 2$

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

Normalized:



$w$  Andrew Ng



deeplearning.ai

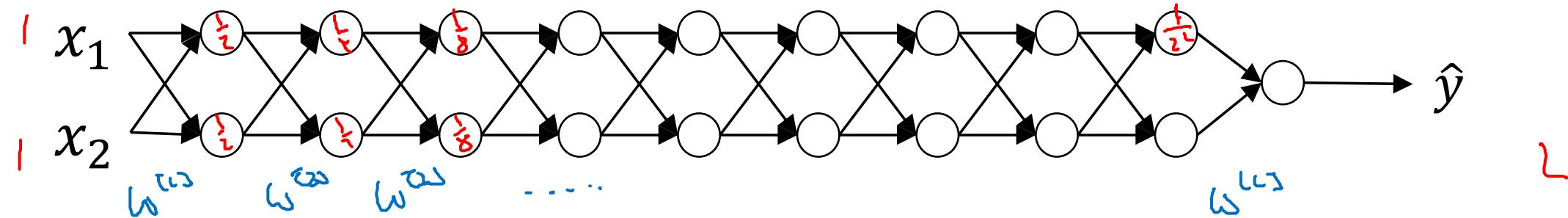
Setting up your  
optimization problem

---

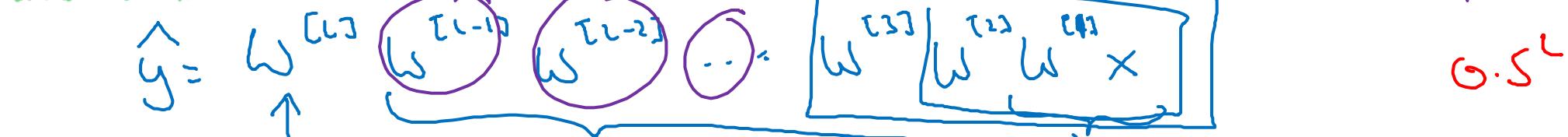
Vanishing/exploding  
gradients

# Vanishing/exploding gradients

$L = 150$



activation function  $g(z) = z^2$ . Layer  $b^{[L]} = 0$ .



$$w^{[1]} > 1$$

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

$$w^{[L]} = \begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix} \quad \begin{bmatrix} 0 & 0 \\ 0 & 1.5 \end{bmatrix} \quad \begin{bmatrix} 0 & 0 \\ 0 & 6.5 \end{bmatrix}$$

$z^{[1]} = w^{[1]} x$

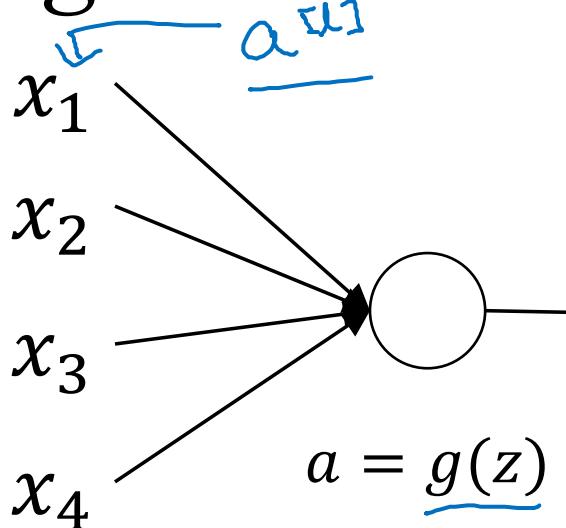
$a^{[1]} = g(z^{[1]}) = z^{[1]}$

$a^{[2]} = g(z^{[2]}) = g(w^{[2]} a^{[1]})$

$1.5^{L-1} \times$  : 1.5^L for NN, 6.5^L for deep NN  
 $6.5^{L-1} \times$  : 6.5^L for NN, 1.5^L for deep NN

Andrew Ng

# Single neuron example



$w^{[l]}$  < 가중치 페인트>  
 : 가중치는 높은 번호를 갖는 뉴런  
 적용되는 경우 학습률이 빠르게 줄어들 수 있다.  
 적용되는 경우 학습률이 빠르게 줄어들 수 있다.

$$z = \underline{w_1 x_1 + w_2 x_2 + \dots + w_n x_n} \quad \cancel{\text{X}}$$

Large  $n \rightarrow$  Smaller  $w_i$

$$\text{Var}(w_i) = \frac{1}{n} \frac{2}{n}$$

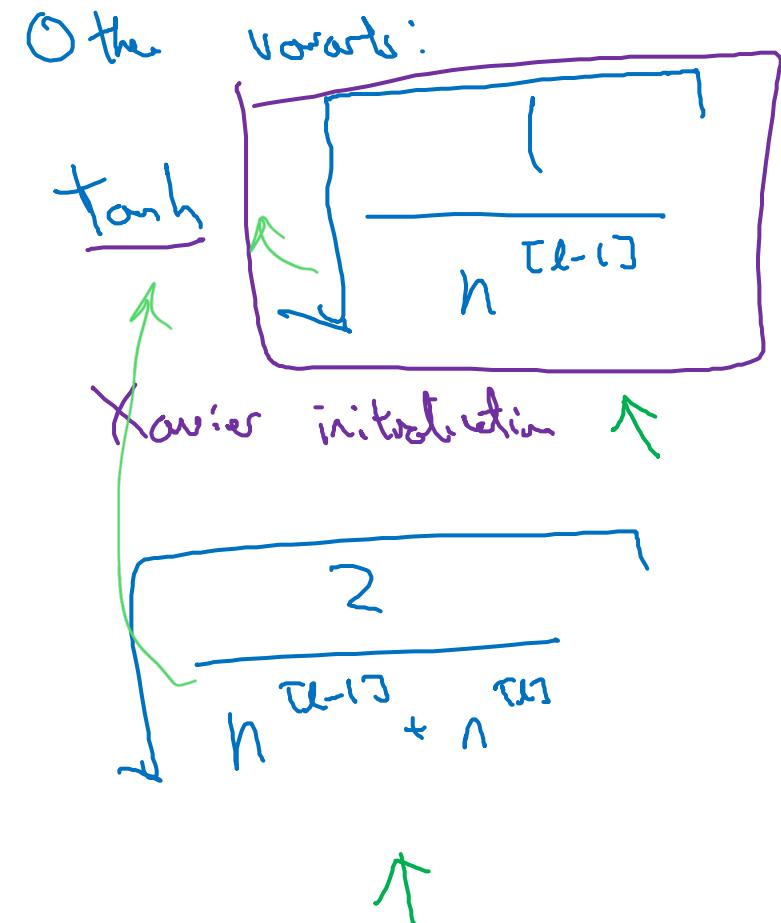
( $w_i$  평균  $\frac{1}{n} \approx 0$ )  
 $n$ : 뉴런 개수

$$\underline{w^{[l]}} = \text{np.random.randn}(\text{shape}) * \text{np.sqrt}\left(\frac{2}{n^{[l-1]}}\right)$$

ReLU

$$g^{[l]}(z) = \text{ReLU}(z)$$

ReLU 활성화 함수  
 $w_{avg} = \frac{2}{n^{[l-1]}}$





deeplearning.ai

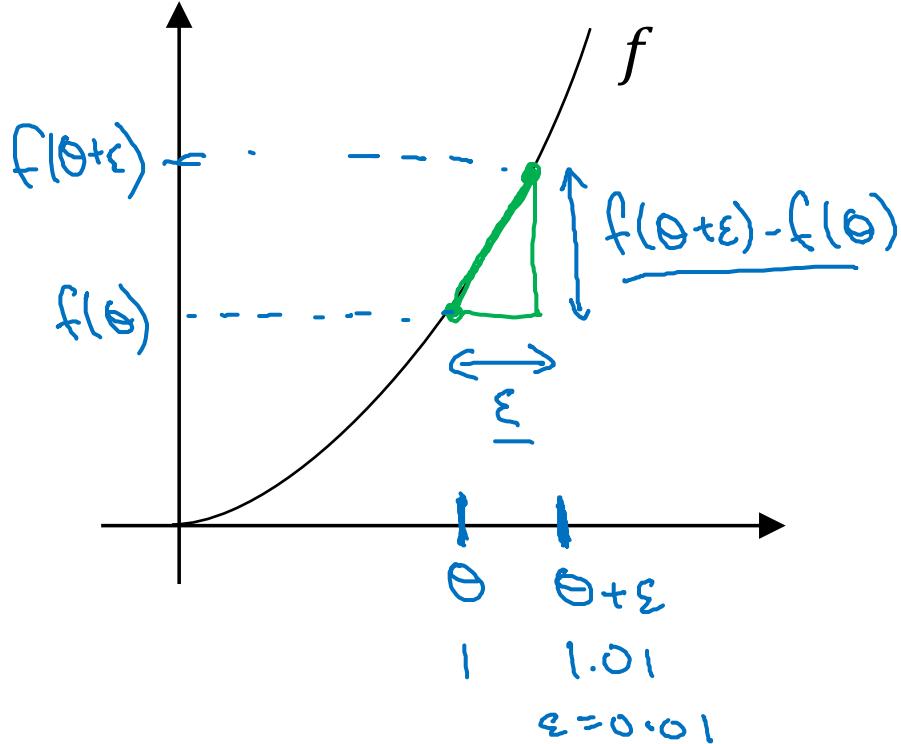
# Setting up your optimization problem

---

## Numerical approximation of gradients

# Checking your derivative computation

$$\begin{aligned} f(\theta) &= \underline{\theta^3} \\ \theta &\in \mathbb{R}. \\ \text{I} \end{aligned}$$



$$\begin{aligned} g(\theta) &= \frac{d}{d\theta} f(\theta) = f'(\theta) \\ g(\theta) &= 3\theta^2. \\ g(1) &= 3 \cdot (1)^2 = 3 \\ \text{when } \theta &= 1 \\ \frac{dw}{db} \end{aligned}$$

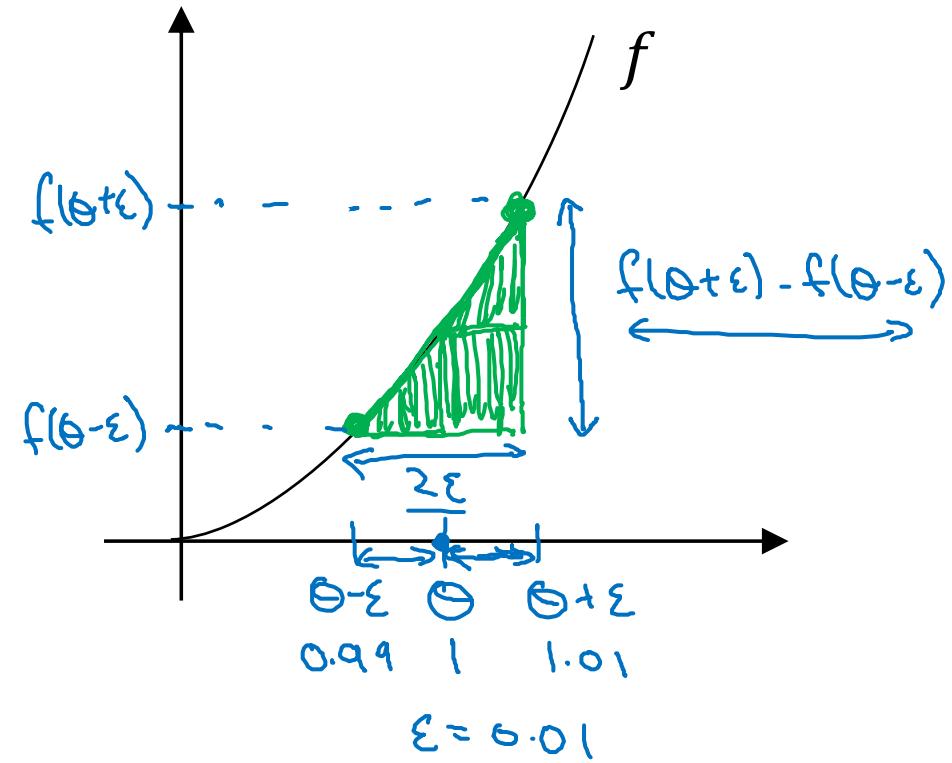
$$\begin{aligned} \frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} &\approx g(\theta) \\ \frac{(1.01)^3 - 1^3}{0.01} &= 3.0301 \\ \frac{3.0301}{0.0301} &\approx 3 \end{aligned}$$

$$\begin{aligned} \theta &= 1 \\ \theta + \epsilon &= 1.01 \end{aligned}$$

# Checking your derivative computation

$$\underline{f(\theta) = \theta^3}$$

Gradient checking in multiple dimensions



$$\left[ \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \right] \approx g(\theta)$$

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

$$g(\theta) = 3\theta^2 = 3$$

approx error: 0.0001

(prev slide: 3.0301. error: 0.03)

|  |   |  |  |
|--|---|--|--|
| $f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$ | $\frac{\mathcal{O}(\epsilon^2)}{\epsilon}$<br>$\frac{0.01}{0.0001}$ | $\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}$<br>$\uparrow \quad \uparrow$ | error: $\mathcal{O}(\epsilon)$<br>$0.01$ |
|--|---|--|--|



deeplearning.ai

Setting up your  
optimization problem

---

Gradient Checking

# Gradient check for a neural network

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

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

1. 모든 계층의 모든 가중치와 편향을 big vector에 포함시킨다.

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

$$J(\underline{dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}})$$

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

2.  $d\theta$ 는  $J(\theta)$ 의 기울기이다.

# Gradient checking (Grad check)

$$J(\theta) = J(\theta_0, \theta_1, \theta_2, \dots)$$

for each  $i$ :

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

가능한 경우에 따른

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

|  $d\theta_{\text{approx}} \approx d\theta$

각 단계에서  $d\theta$ 의 차이

Check

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

$$\varepsilon = 10^{-7}$$

$$\approx [10^{-7} - \text{great!}] \leftarrow$$

$10^{-5}$

$$\rightarrow 10^{-3} - \text{worry.} \leftarrow$$



deeplearning.ai

Setting up your  
optimization problem

---

Gradient Checking  
implementation notes

# Gradient checking implementation notes

- Don't use in training – only to debug

또한 노인은 가족과 구별해 구하는 것은 예전 시기의 풍습이었지만 최근에는 노인들이 대체로 혼자 살고, 자녀들이,

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

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

$\frac{\partial b}{\partial \gamma}$   $\frac{\partial w}{\partial \gamma}$  만약 앞에서 Gradient check이 통과되었을 때, 같이 초기  $d\gamma$  를 주어지면 layer와 바운더리 층을 순회  
 → 흐름을 따라가면서 차를 계산할 수 있음.

- Remember regularization.

Regulation을 하기 위해서는 그동안 이루어온

$$J(\theta) = \frac{1}{m} \sum_i f(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_i \|\omega^{(i)}\|_F^2$$

- Doesn't work with dropout.

Dropout rate는 훈련 단계에서 흔히 사용되는 Gradient checking은 학습률을 조정하는 데 사용된다.

J keep-prob = 1.0

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

$w, b \sim \mathcal{O}$  Back Propagation 은 웨이트 및 바이어스를 업데이트하는 과정입니다. 이를 통해 예측값과 실제값 간의 차이를 줄여나가는 과정입니다. 웨이트와 바이어스는 학습 과정에서 최적화되는 값입니다. 예전에는 흔히 Gradient Descent 라는 random initialization 이후에, 이를 통해 최적화를 수행합니다. 대신 Gradient checking 을 수행하는 경우,