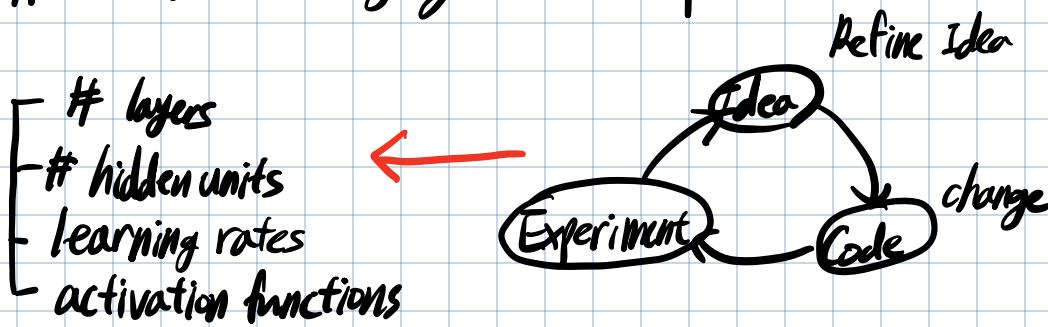


WEEK 1. Setting up your Machine Learning Application

1. Setting up your Machine Learning Applications

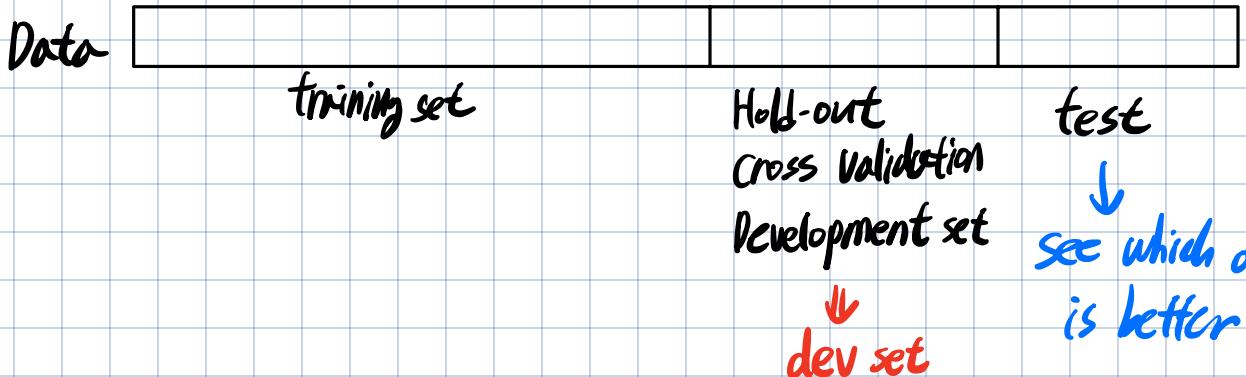
- Train / Dev / Test sets

- Applied ML is a highly iterative process



→ go round this cycle repeatedly & efficiently to find best hyperparameters

- Train / dev / test sets for efficient cycle



- Previous era: split data to 70/30 %
(Machine learning)
60/20/20%

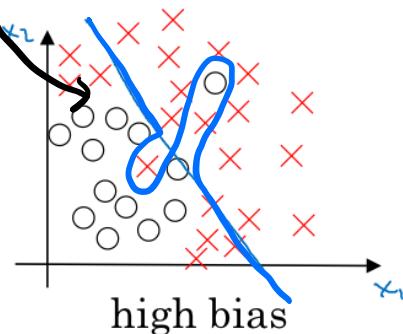
- Big data era: dev & test set → to quickly decide which algorithm is (1,000,000) better → 98/1/1 % or 99/0.5/0.5 %
e.g. 1,000,000 → 10,000 is enough for dev set
10,000 " test set

Overfitting to training set
(low generalization)
⇒ high variance

It didn't even fit training set
(Underfitting)
⇒ high bias

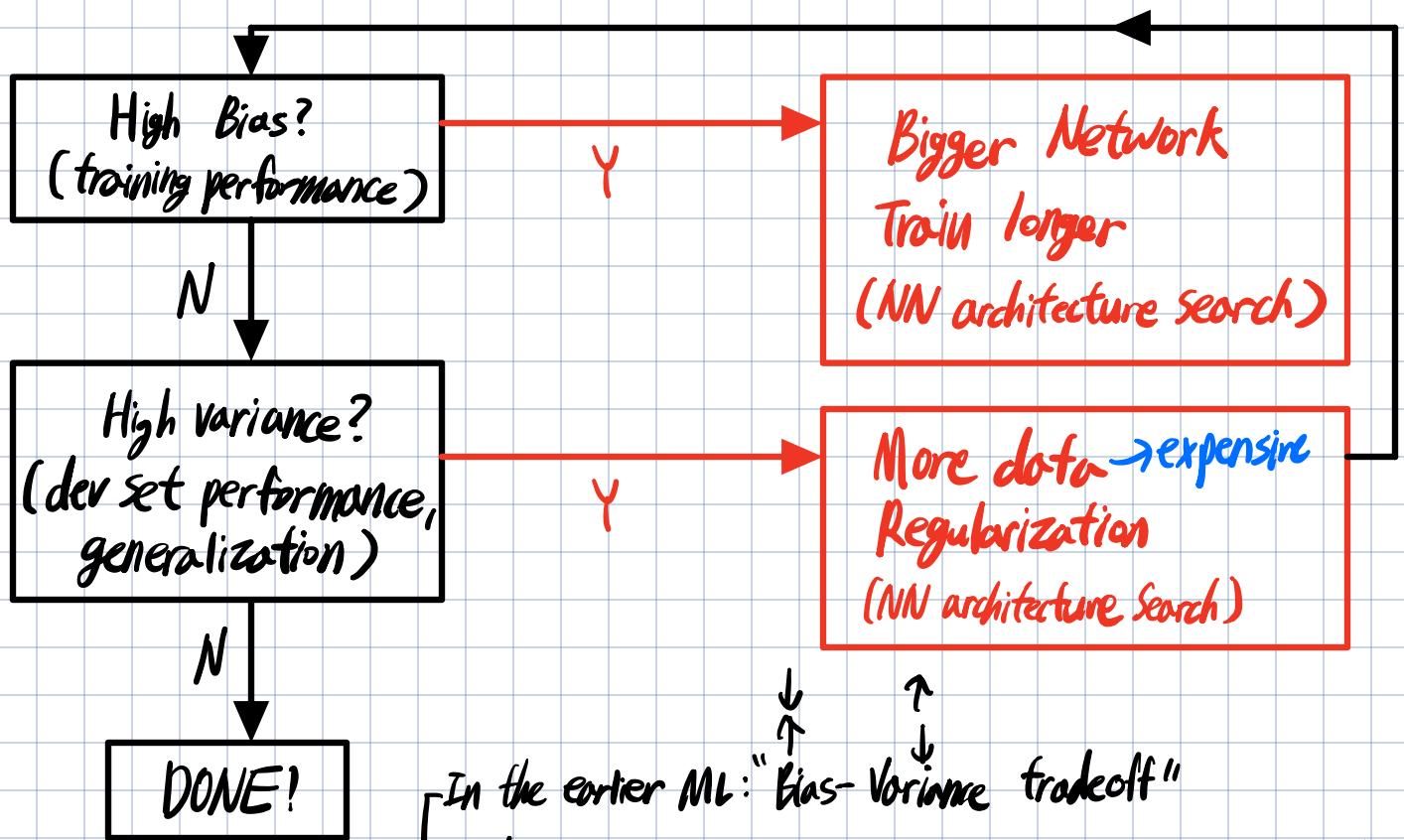
high bias low variance
& high variance

↓
overfitting to a part of data



• Basic "recipe" for machine learning

- After training initial model, check bias & variance by using training/dev set
- ⇒ Depending on its bias and variance, solution can be different



2. Regularizing your neural network.

• Regularization

- Overfitting due to High variance \Rightarrow **Regularization**

- Regularization of Logistic Regression

$$J(w, b) = \underbrace{\frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y)}_{\text{cross-entropy loss}} + \frac{\lambda}{2m} \|w\|_2^2 + \cancel{\frac{b^2}{2m}}$$

mbd

hyperparameter
used in cross-validation
and dev set

omit
: almost all parameters
are w to fit well
(high dimensional)

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \quad (\text{Euclidean norm of } w)$$

↳ L2 norm of vector w

↳ L1 Regularization $\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1 \Rightarrow w$ often be **sparse**
= w vector has many zeros

↳ L2 Regularization $\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$

• Regularization of Neural Network

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

- **Frobenius norm** : norm of matrix w ($n^{[L-1]}, n^{[1]}$)

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

↳ square of each element in matrix w

↳ weight update \neq all
 $w \leftarrow \text{gradient}$ 만드는

- if λ is too large \rightarrow oversmooth (model with high bias)

- Backpropagation

$$dw^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]}$$

$$w^{[l]} := w^{[l]} - \alpha dw^{[l]}$$

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

extra term $\frac{\lambda}{m} w^{[l]}$
penalty $\frac{\lambda}{m}$ dict.

$$= w^{[l]} - \alpha (\text{from backprop}) - \alpha \frac{\lambda}{m} w^{[l]} \Rightarrow \text{weight ends up smaller}$$

$$= w^{[l]} \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha (\text{from backprop})$$

\Rightarrow make $w^{[l]}$ a little bit smaller by multiplying something < 1

\therefore "L2 Regularization = Weight decay"

• Why regularization reduces overfitting

- How does regularization prevent overfitting?

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

\Rightarrow reduce L2/Frobenius norm \Rightarrow less overfitting

- Intuitions

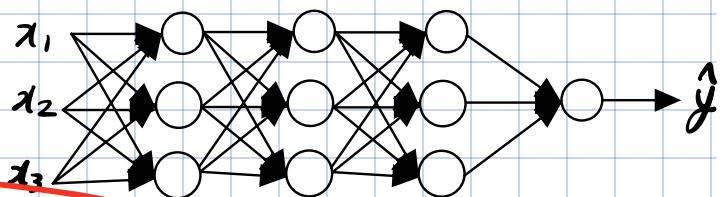
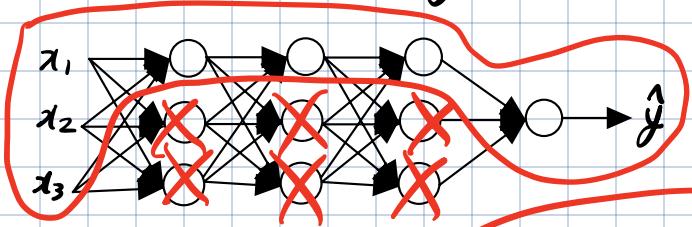
① λ big: make $w^{[l]} \approx 0$ \Rightarrow set many hidden units 0

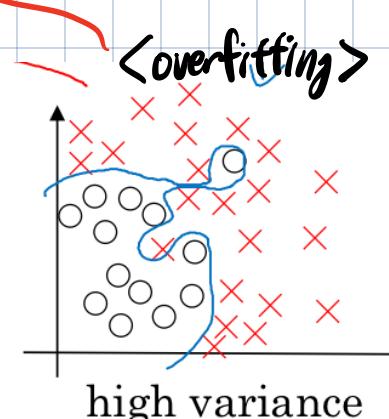
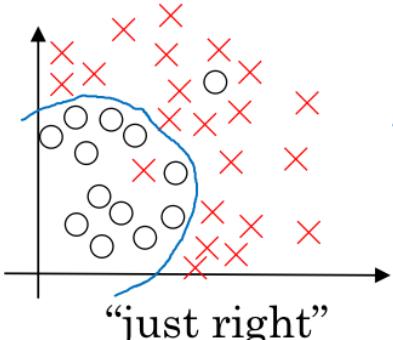
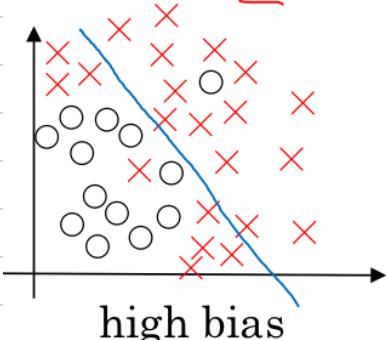
\Rightarrow reduces impacts of hidden units

\Rightarrow close to LR unit (simpler network)

\Rightarrow reduces overfitting

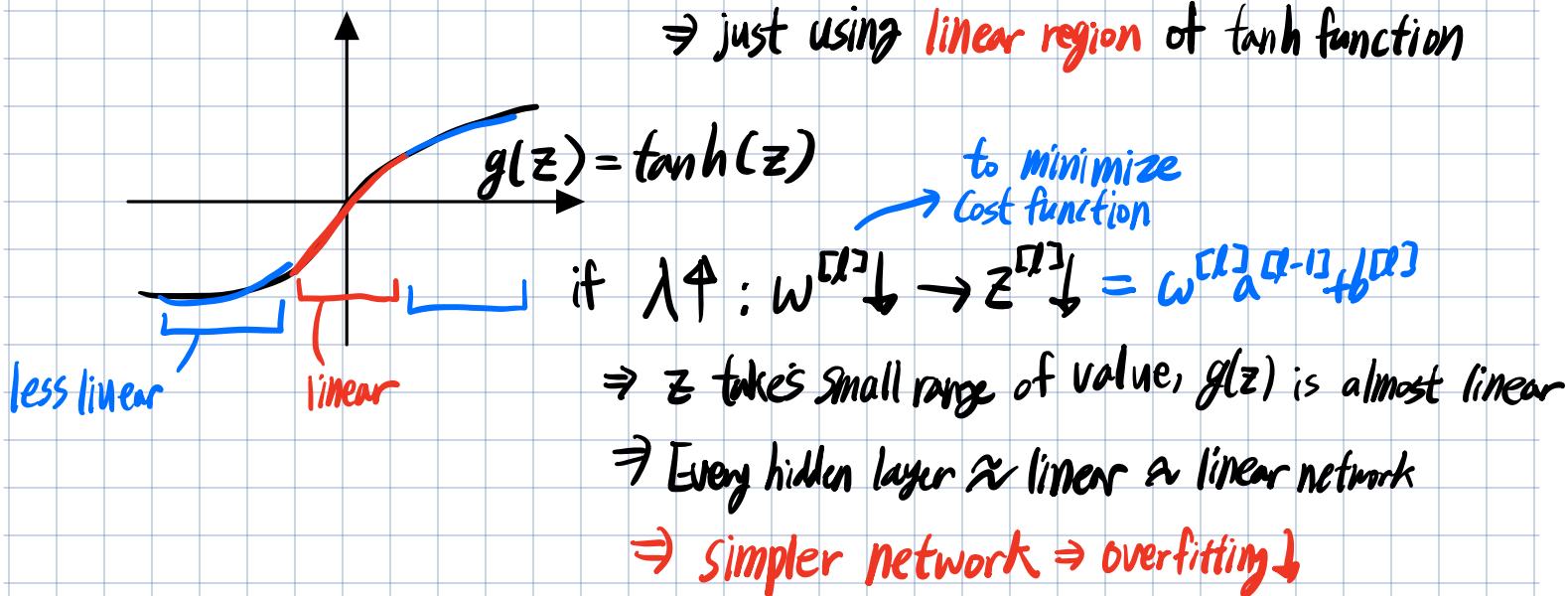
- It can hurt training set performance but ultimately it gives better test accuracy



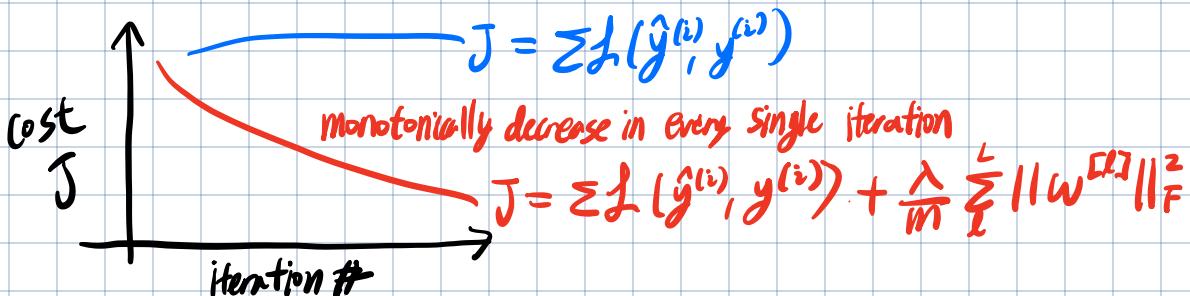


② If z takes on only small ranges of parameters

\Rightarrow just using **linear region** of tanh function

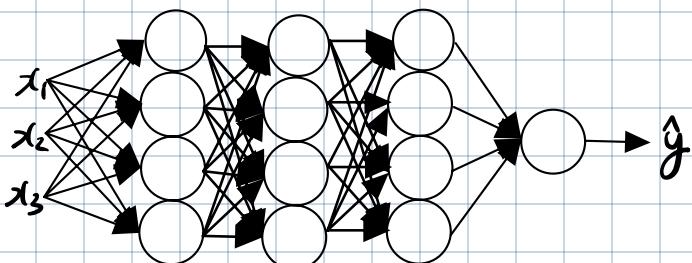


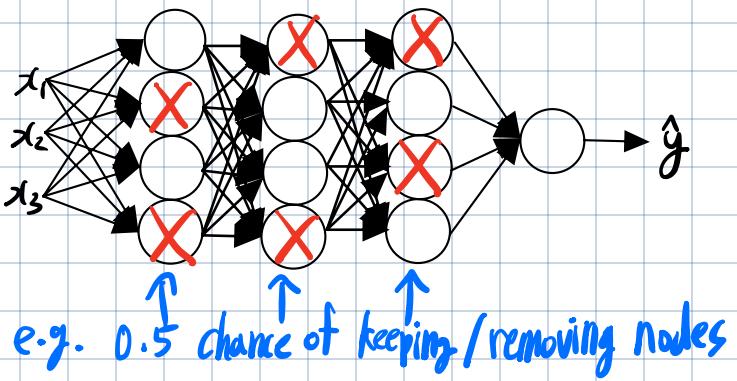
- Tips for implementation of gradient descent



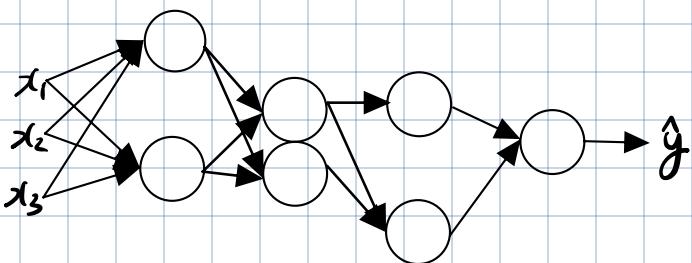
• Dropout Regularization

- Let's train this overfitting network





① go through each of the layers in the network and set some probability of eliminating the node of neural network



② remove every in/outgoing links from removed nodes and perform backpropagation in the smaller Network

③ repeat for every training samples

- Implementing dropout ("Inverted dropout")

e.g. $l=3$ keep-prob = 0.8

↳ layer probability of hidden unit to be kept for layer 3

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

↳ dropout vector d_3 는 각각의 노드가 각 유닛에 대해

for layer 3

zero out hidden unit

$\begin{bmatrix} 0.8\text{의 확률로 } 1 \\ 0.2\text{의 확률로 } 0 \end{bmatrix}$ 이 된다.

$a_3 = \text{np. multiply}(a_3, d_3)$ or $a_3 *= d_3$

$a_3 /= \text{keep-prob}$ scaling step

inverted-dropout technique

⇒ keep-prob을 얻기 위해 keep-prob을 다시 나누어야 a_3 의 기대값을 같게 유지

즉, 이를 cost function에 적용할 때마다 동일한 같은 유지하고자 마지막에 나눠준다.

e.g. 50 units in 3rd layer $\xrightarrow{\text{keep-prob}=0.8}$ 10 units shut off

$$z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}$$

↑ 20% of elements are zeroed out

⇒ dropout 전 A_3 의 기대값은 A_3 과 같다, but dropout $\Rightarrow A_3$ 의 기대값은 80% off. $\Rightarrow \text{keep-prob} \leq 4/5$ 까닭은 여기

$\Rightarrow d$ 는 결과가 정답 한번 봄마다 0이 되는 원인 유발들이 달라진다. (보통에 훈련 세트의 $d \neq$ 정답이)

- Making predictions at test time

$$a^{[0]} = X \quad (=X_{\text{test}})$$

test에서는 no drop out \Rightarrow random 몇몇 hidden unit을 zero-out하는 경우

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]} \quad (\text{output } \hat{y} \text{ of random 몇몇 } X, \text{ dropout을 적용하면 test의 noise를 추가한다.})$$

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

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \dots$$

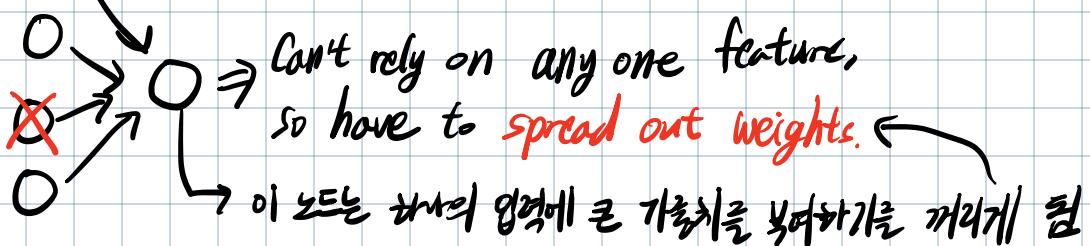
\downarrow
 \hat{y}

test에서 drop-out을 고려하지 않아도 크기는 변하지 않기 때문에
keep-prob scaling을 할 필요 X

• Understanding dropout

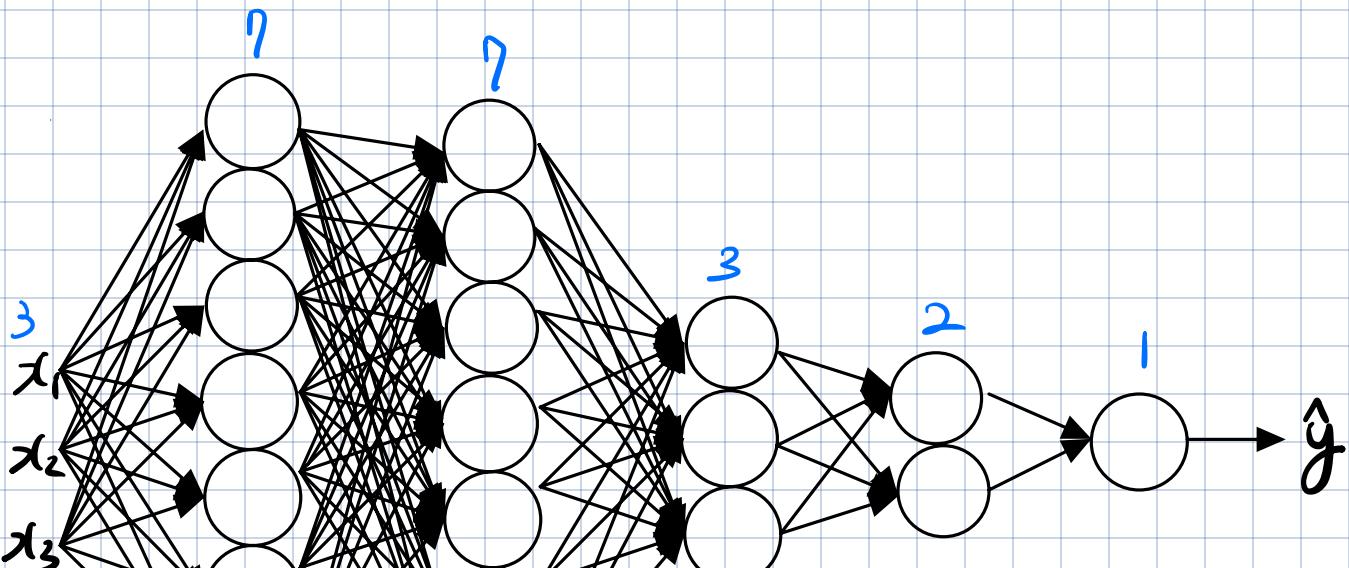
- Why does drop-out Work?

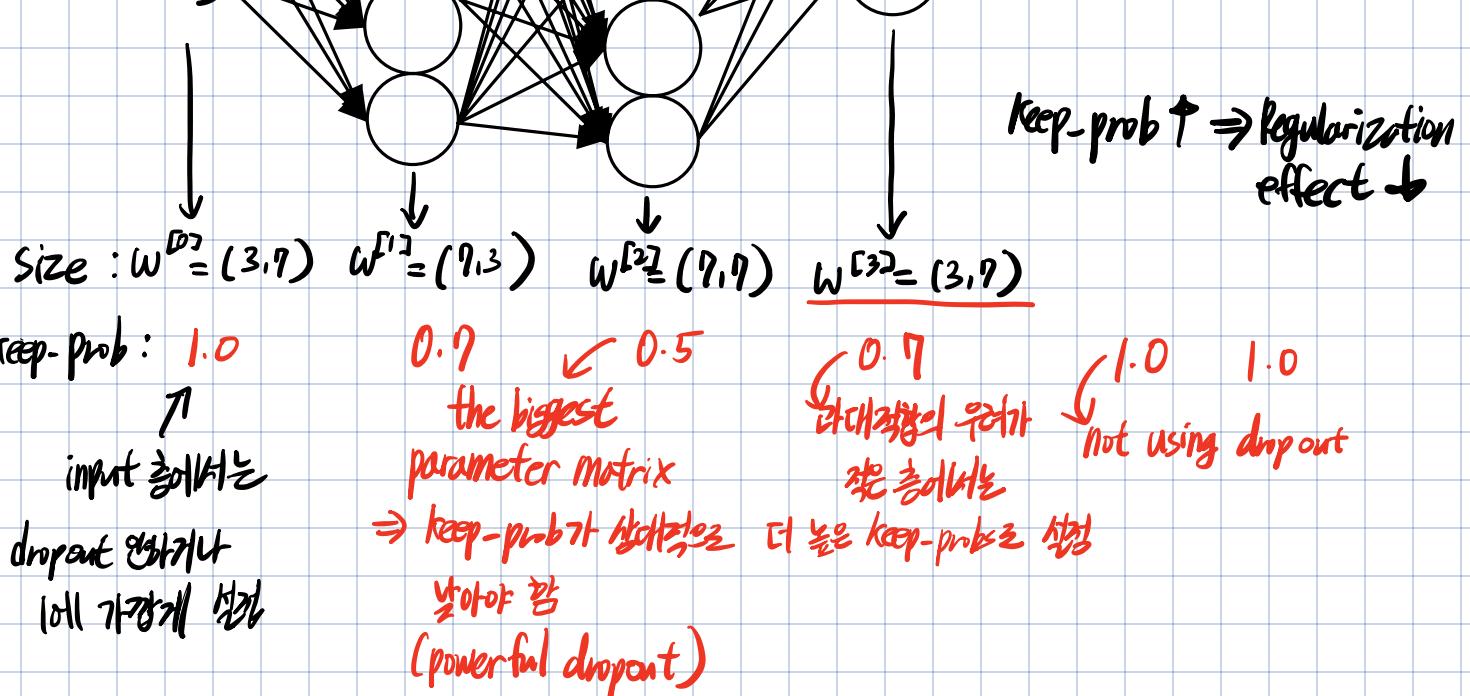
~~X~~ \rightarrow randomly drop out



\Rightarrow Weight를 spread-out함으로써 $\|W^{[l]}\|^2$ 이 줄어든다. \leftarrow L2 정규화처럼
가중치를 줄여 overfitting ~~X~~

- keep-prob을 각 층에 대해 설정할 수 있다. (but, more hyperparameter for CV)





- Computer vision에서 drop-out 구현을怎么做?

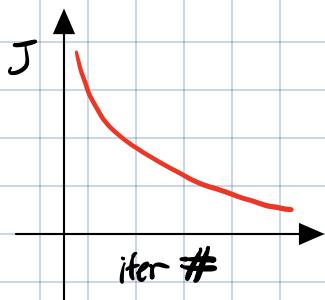
∴ 아직 많은 pixel 사용 → no enough data

⇒ 하지만 이 용도에서는 over-fitting 때문에 drop-out 사용X

- Downside of drop-out

⇒ 반복마다 무작위로 노드를 제거하기 때문에 cost function J의 차이를 확인하기 어렵음
(따라서 디버깅 어려움)

(drop out을 끄고 학습하는 drop out을 끄고 J를 확인 후 다시 turn-on)



• Other regularization Methods

① Data augmentation

- 데이터를 다양하게 바꾸 / 편집하여 사용을 데이터로 얻기

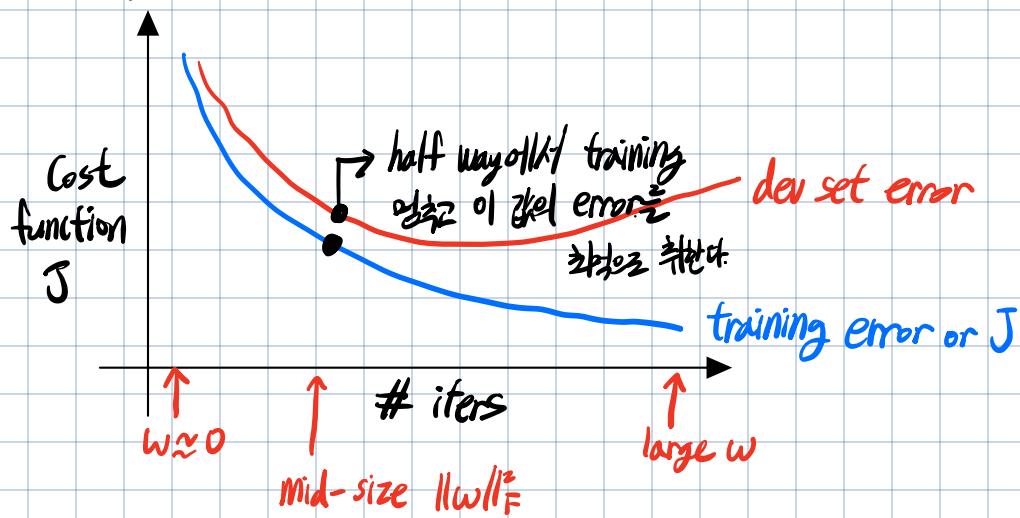
e.g. 이미지의 확장, 학대, 편집

⇒ 증명된 사실이 많아서 좋지 않을 수 있으며, 동양적인 metric이 이미지를 구하는 것보다는 bad.

② Early Stopping

- Early stopping은 training 중 최소의 loss가 되면 stop

→ L2 정제화 비슷하게 W norm이 작은 선형망을 선택함으로써 cost function을 optimize 한다.



- L2 regularization vs early stopping

	L2	Early Stopping
Good	모델을 잘 훈련할 수 있어 J 최적화 가능	한번에 경사하강법으로 small w, mid w, big w를 얻을 수 있어 많은 것을 시도해볼 필요 있다.
Bad	계산 비용이 커 시간 필요 ⇒ computationally expensive	J 최적화를 어렵게 함 (orthogonalization X) ⇒ ①, ②를 복잡정수로 수행 불가

Orthogonalization [① optimize J : Gradient descent, momentum ... ⇒ optimize w & b
 ② Avoid overfitting : augmentation, getting more data, ... ⇒ reduce variance]
 ↓ 이 두를 한 번에

think one thing at a time ⇒ but early stopping mix these two things up

이 두는 양자나
다른 걸

→ 컴퓨터 강당 가능하면 L2 Regularization 이용하는 걸 추천

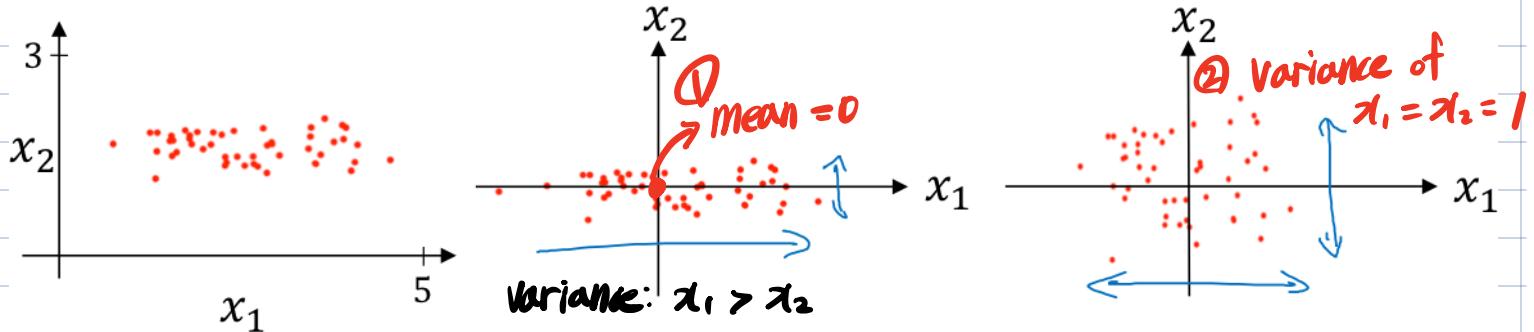
3. Setting up your optimization problem

- Normalizing inputs

- How to speed up training? \Rightarrow Normalize the inputs!

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

- ① 모드 축의 평균에 대해 **mean** \bar{x}_i (zero-out mean)
- ② 모드 축의 편차에 대해 **Variance** σ_i^2



Use same μ & σ^2 to normalize your test set

$$\text{① } \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\text{② } \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

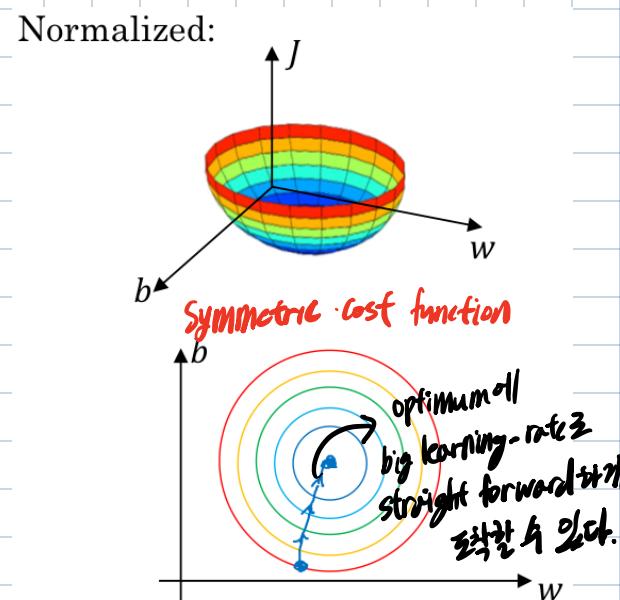
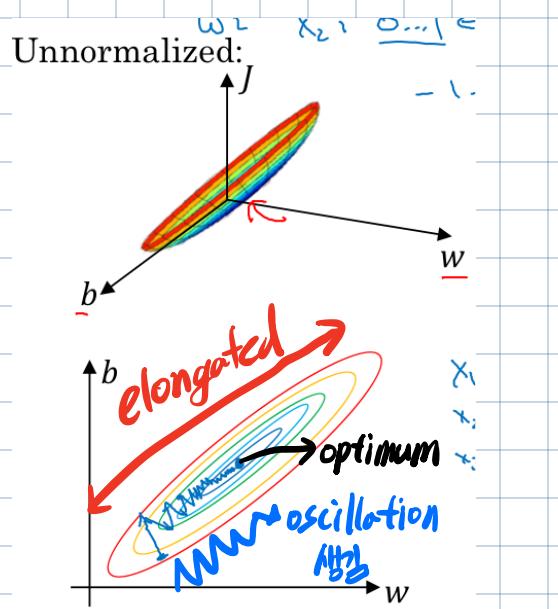
(element-wise)

$$x := x - \mu$$

$$x / \sigma^2$$

- Why do we normalize input features?

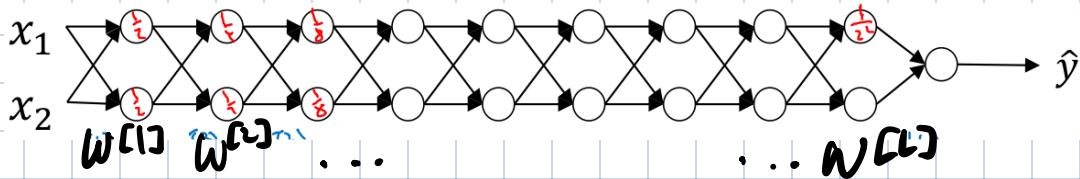
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)})$$



\Rightarrow feature 간 scale이 너무 다르면 w_1 & w_2 의 비가 너무
달라져 **elongated cost function**을 갖게 되기 때문에
이유는 learning-rate η gradient descent 속도에
(비교적 더 빠르거나 느리다)

• Vanishing/exploding gradients

— 깊은 neural network 훈련 시 derivative가 점점 감소하거나 증가한다.



$$g(z) = z \text{ (linear activation)}, b^{[l]} = 0$$

$$\hat{y} = w^{[L]} \cdot w^{[L-1]} \dots w^{[2]} w^{[1]} x$$

$\Rightarrow \hat{y}$ 는 모든 w 행렬의 곱이 됨

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

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

if $w^{[l]} > 1 \Rightarrow$ Deep NN makes activations explode

$$\left[\begin{array}{l} \text{e.g. } w^{[1]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \Rightarrow \hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} z \\ = 1.5^{(L-1)} z \end{array} \right]$$

if $w^{[l]} < 1 \Rightarrow$ Deep NN makes activations very small

$$\left[\begin{array}{l} \text{e.g. } w^{[1]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \Rightarrow \hat{y} = 0.5^{L-1} z \end{array} \right]$$

\Rightarrow You can apply this intuition to derivative exploded/vanish

• Weight initialization for deep networks

- Partial solution for gradient vanishing/exploding

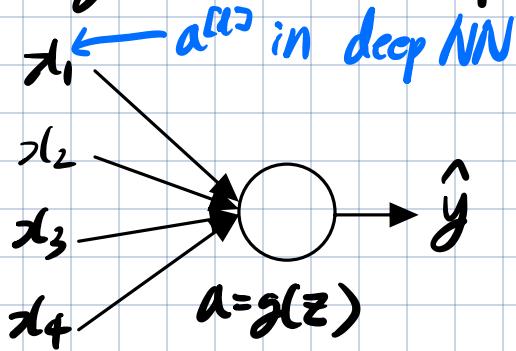
⇒ better choice of random initialization for NN

⇒ Good initialization leads to...

① Speeding up the convergence of gradient descent

② Increase the odds of gradient descent converging to a lower training/generalization error

- Single neuron example



$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b \rightarrow \# \text{ of input features}$$

⇒ to make z not too big/small...

∴ larger n , smaller w_i

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

layer l till ℓ 까지

of input features

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

$$\textcircled{2} \text{ if activation func = ReLU, set } \text{Var}(x) = \frac{2}{n}$$

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

⇒ 각각의 w 를 1보다 너무 크거나 작아지지 않게 설정하여 소실/폭발(X)

③ Other activation Variants

\tanh :

$$\sqrt{\frac{1}{n^{[l-1]}}}$$

or

$$\sqrt{\frac{2}{n^{[l-1]}}}$$

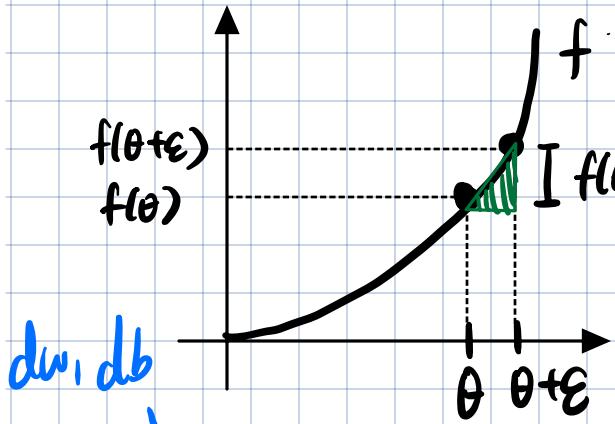
⇒ Xavier initialization

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

⇒ hyperparameter 험수의 의의로 값을 수를 조정 가능

Numerical approximation of gradients

- Check implementation of backprop is correct
- One-side difference

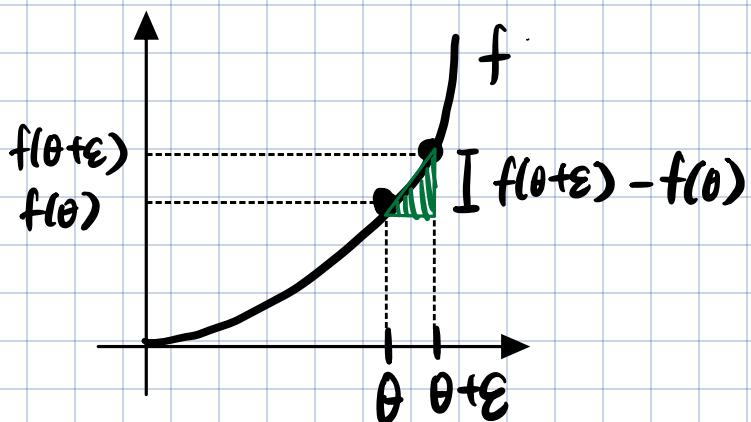


$$g(\theta) = \frac{d}{d\theta} f(\theta) = f'(\theta)$$

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

if $\theta = 1, \epsilon = 0.01,$
 $\Rightarrow \left[\begin{array}{l} g(\theta) = 3 \cdot (1)^3 = 3 \\ \frac{(1.01)^3 - 1^3}{0.01} = 3.0301 \end{array} \right] 0.0301$

- two-side difference



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

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

if $\theta = 1, \epsilon = 0.01,$
 $\Rightarrow \left[\begin{array}{l} g(\theta) = 3 \cdot (1)^3 = 3 \\ \frac{(1.01)^3 - (0.99)^3}{0.02} = 3.000123 \end{array} \right] 0.000123$

⇒ two-side difference vs one-side vs $g(\theta)$ or $f(\theta)$ 의 derivative all
that's correct implementation 대비 confidence 증가.

• Gradient Checking

- Check if you implemented backprop well for NN

① $W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$ ^{가져온} $\xrightarrow{\text{concatenate}}$ ^{reshaped} \Rightarrow θ vector $\in \mathbb{R}^n$
만든다.

⇒ cost function J 또한 w, b all that 가 아니라 θ 에 대한 함수로

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

② $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$ ^{가져온} $\xrightarrow{\text{reshape}} d\theta$ vector

③ Check if $d\theta$ the gradient of $J(\theta)$?

now, $J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots) \rightarrow \theta$ 모든 각각에 대해 check

for each i use two-side difference

$$\Rightarrow d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

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

↳ check by euclidean distance

set $\epsilon = 10^{-7}$, $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} \times \#$

$10^{-9} \Rightarrow \text{great}$

$10^{-5} \Rightarrow \text{Okay, but check that none of the components of vector are too large}$

$-10^{-3} \Rightarrow$ Worry, 특정 i번째 $d\theta$ approx $[i]$ 를 $d\theta[i]$ 의 차이가 크거나
나는 곳을 측정한 필요 있음.

• Gradient checking implementation notes

① Debugging 어떤 사용하고 training에는 사용 X

∴ expensive computation

② algorithmic grad check off 상태 \Rightarrow component 틀림

\rightarrow 특정 i에서 어떤 $db^{(i)}$ 과 $dw^{(i)}$ 에서 θ 를 일어가도록 check

③ Regularization의 J 옆 포함되어 있다는 것을 기억하기

$$J(\theta) = \frac{1}{m} \sum_i h(g^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_i \|w^{(i)}\|^2$$

$d\theta = \text{grad}$ of J w.r.t θ that contains regularization

④ Dropout하고 한 번에는 쓰지 말 것.

∴ Dropout은 node의 subset을 끊어가기 때문에 계산이 어렵다.

keep-prob=1로 설정하여 turn-off $\not\equiv$ grad check

⑤ random initialization $\not\rightarrow$ (몇 번 training 했을 때) 실행

∴ $w, b \approx 0$ in random initialization \Rightarrow network를 잠시 훈련하여 $w & b$ 가 0에서 멀어질 수 있도록