

# Deep Learning Seminar

## 6. Training Neural Network

HA SEUNG HYUN

# Contents

## 1. Review

1-1) Classification

1-2) Learning network

## 2. Training neural network - part I

2-1) Overview

2-2) Activation functions

2-3) Data preprocessing

2-4) Weight initialization

2-5) Batch normalization

2-6) Babysitting the learning process

2-7) Hyperparameter optimization

Reference:    lecture note    (Fei-Fei Li)  
                  lecture note    (Andrew Ng)  
                  모두를 위한 머신러닝 (Sung kim)

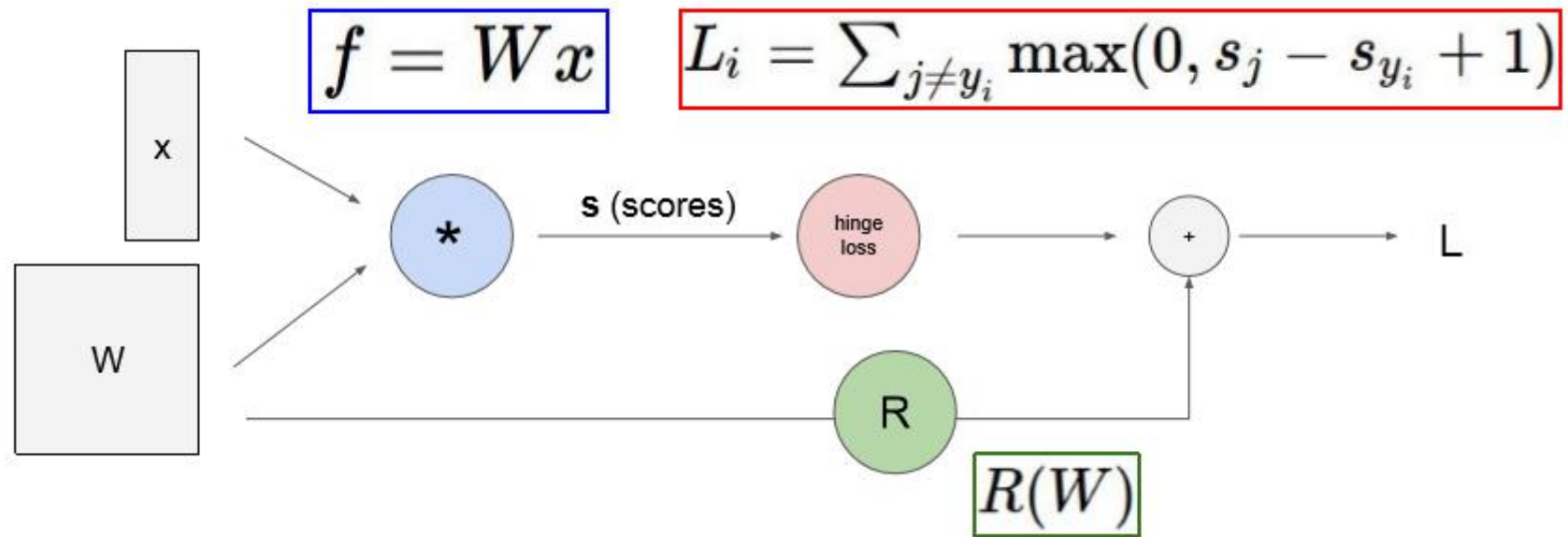
# 1. Review

1-1) Classification

1-2) Learning network

# Classification

## 1. Linear Classifier



# Classification

1. .
2. computation cost ,

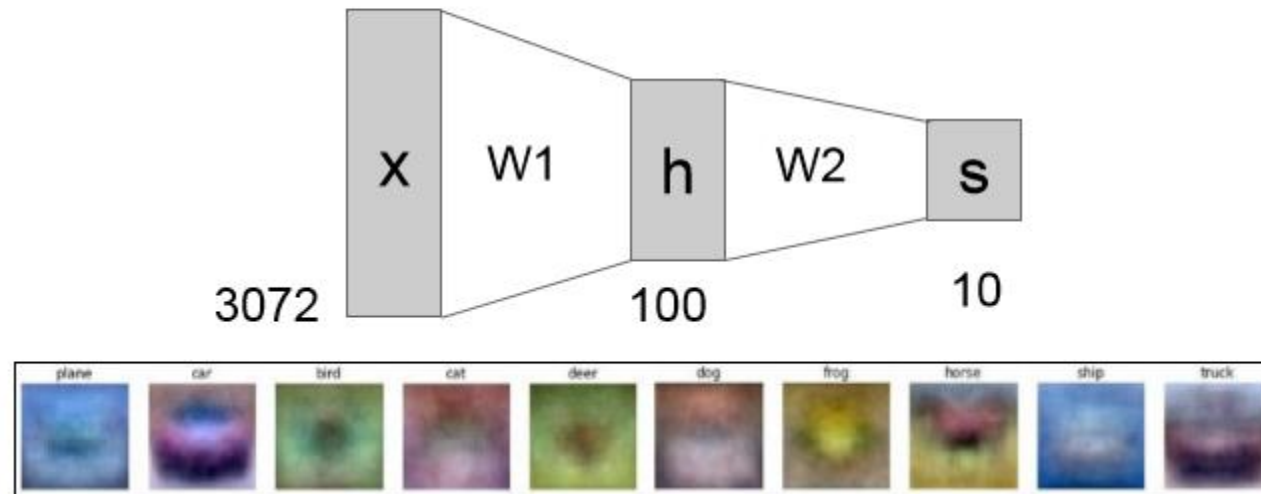
## 2. Fully Connected Network<sup>FCL</sup>

Linear score function:

$$f = Wx$$

2-layer Neural Network

$$f = W_2 \max(0, W_1 x)$$



# Classification

filter

- 1.
- 2.

computation cost

## 3. Convolutional Neural Network

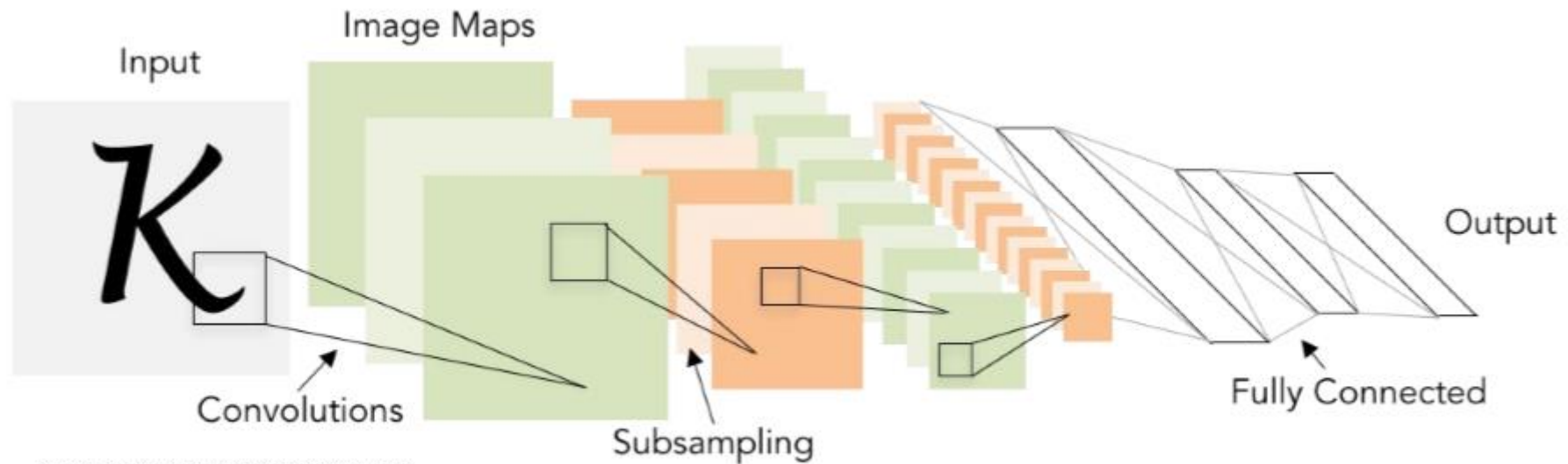
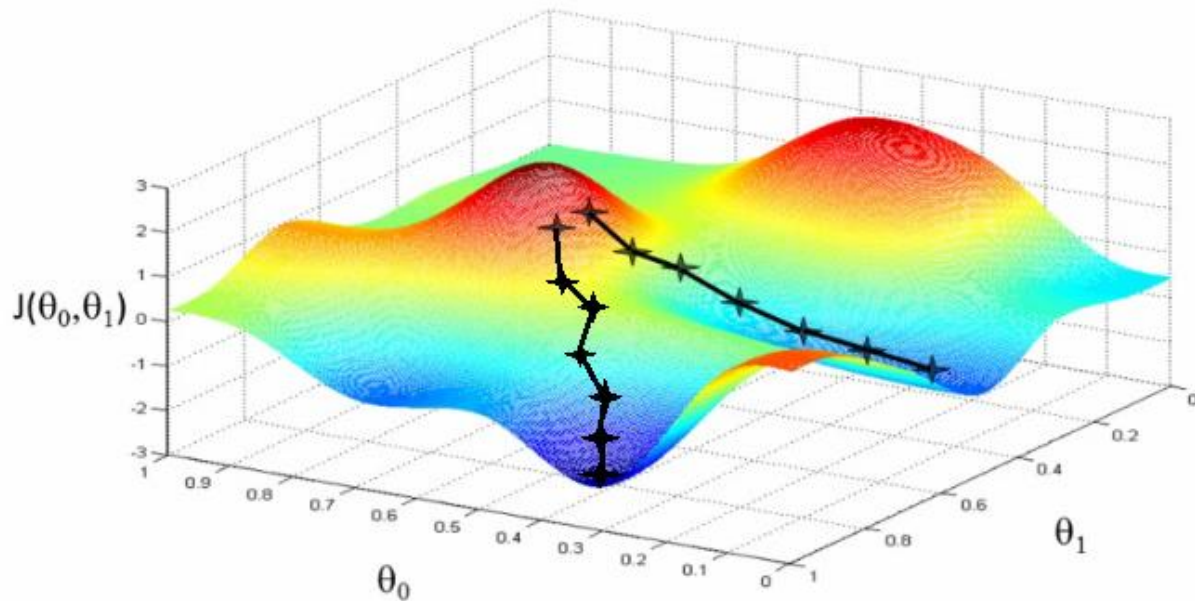


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

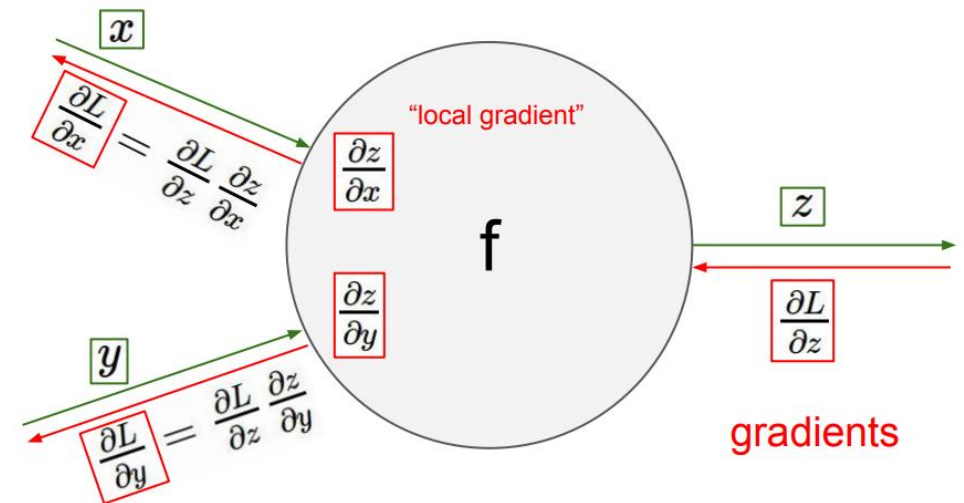
# Learning Network

Back Propagation  
`loss.backward()`  
`model.step()`

SGD  
`S : batch_size`  
`- batch_size loss`



Optimizer



Backpropagation

# Learning Network

## - Deep Learning Pipeline

```
1. Training Data Loading 1_1. . ToTensor()
                        1_2. Loader 100, 64. batch size

2. Training Data Augmentation 1_2. . transforms

3. Deep Neural Network Training with Training Data and Validation Data
   ex) 100 10
      (validation test data)

4. Deep Neural Network Testing with Testing Data
   . validation loss
   backward

5. Inference with verified Deep Neural Network
   test loss 0 acc 100
      (label) test
```

Epoch  
Batch\_size



# Learning Network

1. Training Data Loading
2. Training Data Augmentation

3. Deep Neural Network Training with Training Data

Loop:

1. **Sample** a batch of data
2. **Forward** prop it through the graph (network), get loss
3. **Backprop** to calculate the gradients
4. **Update** the parameters using the gradient

가

가?

## 2. Training neural network (part I)

" : Accuracy ."

2-1) Training overview

2-2) Activation functions

Network  
가

2-3) Data preprocessing

2-4) Weight initialization

2-5) Batch normalization

2-6) Babysitting the learning process

2-7) Hyperparameter optimization

# Training overview

## 1. One time setup

ReLU()  
sigmoid

!!

!!

**activation functions**, preprocessing, weight  
initialization, regularization, gradient checking

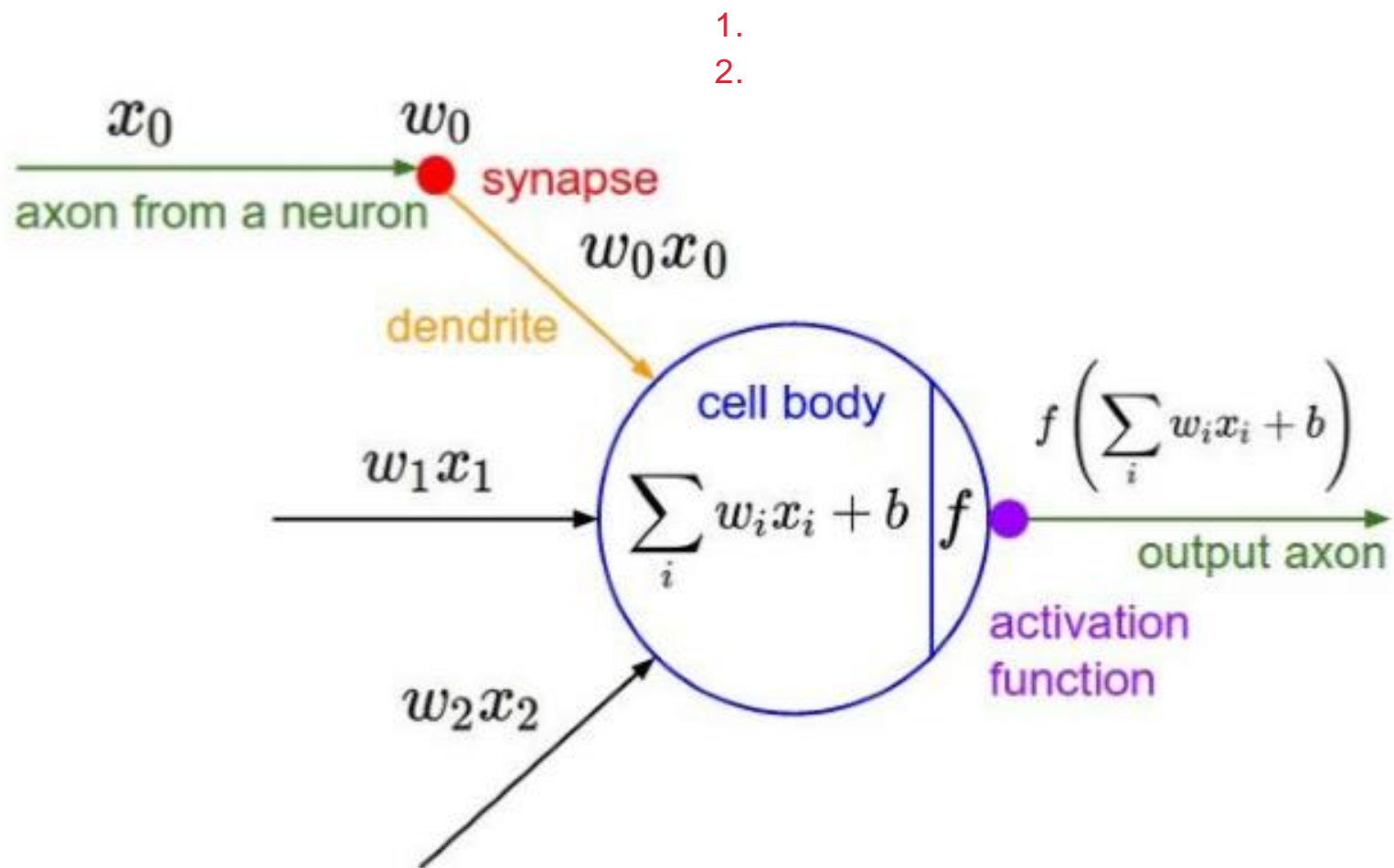
## 2. Training dynamics

*babysitting the learning process,  
parameter updates, hyperparameter optimization*

## 3. Evaluation

*model ensembles*

# Activation Functions



가

# Activation Functions

- Limit of linear classification

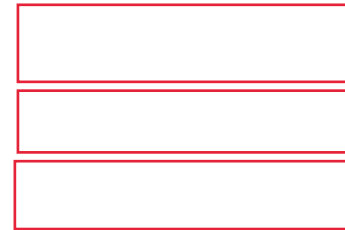
Layer 1:  $y_1 = w_1x + b_1$

Layer 2:  $y_2 = w_2y_1 + b_2$

Layer 3:  $y_3 = w_3y_2 + b_3$

...

Layer n:  $y_n = w_ny_{n-1} + b_n$



linear  
Activation function



가

linear

# Activation Functions

- Limit of linear classification

Layer 1:  $y_1 = w_1x + b_1$

Layer 2:  $y_2 = w_2y_1 + b_2 = w_2(w_1x + b_1) + b_2$

Layer 3:  $y_3 = w_3y_2 + b_3 = w_3(w_2(w_1x + b_1) + b_2) + b_3$

...

Layer n:  $y_n = w_ny_{n-1} + b_n = w_n(w_{n-1}(w_{n-2}(\dots w_1x + b_1) + \dots) + b_n$

# Activation Functions

- Limit of linear classification

Layer 1:  $y_1 = w_1x + b_1$

Layer 2:  $y_2 = w_2y_1 + b_2 = w_2(w_1x + b_1) + b_2 = c_2x + d_2$

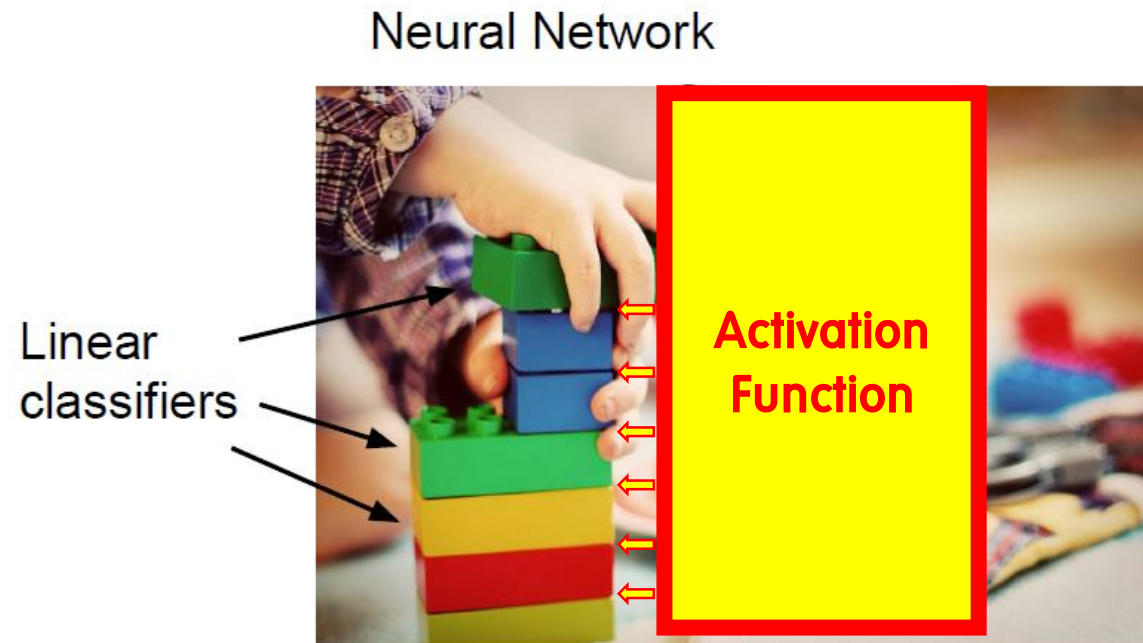
Layer 3:  $y_3 = w_3y_2 + b_3 = w_3(w_2(w_1x + b_1) + b_2) + b_3 = c_3x + d_3$

...

Layer n:  $y_n = w_ny_{n-1} + b_n = w_n(w_{n-1}(w_{n-2}(\dots w_1x + b_1) + \dots) + b_n = c_nx + d_n$

# Activation Functions

Generate **non-linear mappings** from inputs to outputs





# Activation Functions

Layer 1:

$$y_1 = w_1x + b_1$$

Act.

$$\hat{y}_1 = \tanh(y_1) = \text{ReLU}(\tanh(w_1x + b_1))$$

Layer 2:

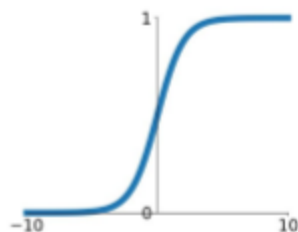
$$y_2 = w_2\hat{y}_1 + b_2 = w_2(\tanh(w_1x + b_1)) + b_2 \neq c_2x + d_2$$

(Activation function: tanh)

# Activation Functions

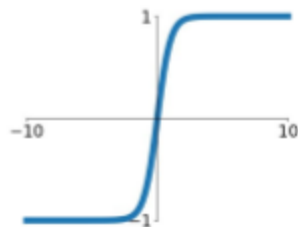
**Sigmoid**

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



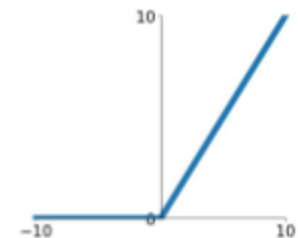
**tanh**

$$\tanh(x)$$



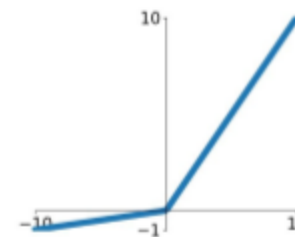
**ReLU**

$$\max(0, x)$$



**Leaky ReLU**

$$\max(0.1x, x)$$

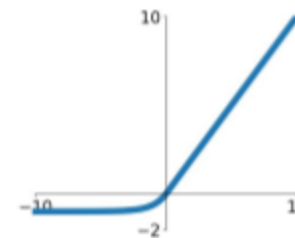


**Maxout**

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



acc 가 90%~95%

97% 98%

!!

# Activation Functions

## - Summary

1. Use **ReLU**. Be careful with your learning rates
2. Don't use **Sigmoid**

Diagram illustrating the choice of Sigmoid activation function for binary classification:

- Input:  $x = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$
- Weights:  $W = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$
- Bias:  $b = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$
- Linear combination:  $z = Wx + b = \begin{bmatrix} 30 \\ 7 \\ 9 \end{bmatrix}$
- Activation: Sigmoid function applied to each element of  $z$ .
- Output:  $\begin{bmatrix} 0.96 \\ 0.05 \\ 0.09 \end{bmatrix}$
- Decision: The highest probability is 0.96, which is rounded to 1.0, indicating the class 'cat'.

# Data Preprocessing

- RGB Image range:
- Network input range:

0~255

0~1

weight  
overfitting  
high frequency

. noise

!!

/255 (scaling)  
transforms ToTensor()

0~1

!!

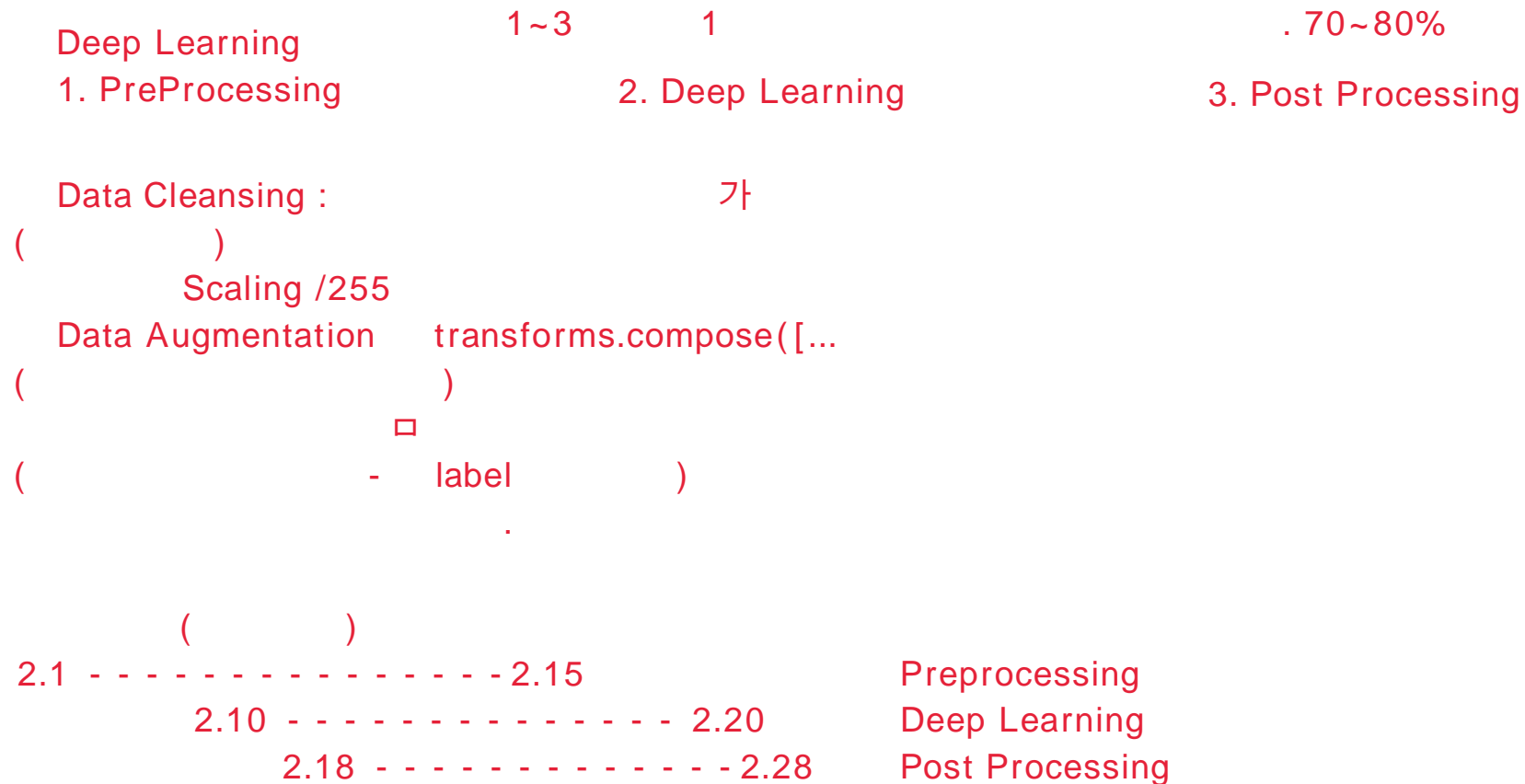
(Recommandable)

In practice,

`network_input = image_matrix / 255`

# Data Preprocessing

- Data augmentation



# Batch Normalization

- Batch (=mini batch)

ex) 7 : 6 | 1

batch size 100

1 Epoch 600

100  
loss가

100  
backward가

loss가

100

W

loss

600

Batch Normalization

layer

sampling

가

Dataset

batch size



1 Epoch



Normalization

x

(

)

(

stable

.)

Batch Normalization

batch size

stable

stable ...

# Batch Normalization

variance가 (Distribution)  
internal covariate shift  
Batch Normalization

Data가 (= )

## • Impact of batch size on error

Data가 (= )  
Data Dense  
Scaling normalization

$(x - \text{mean}) / \text{std}$   
2, 100

noise(stable)  
stable .  
batch size 16  
batch size 16

batch 200 ? !!

computation  
!!

가 ...

100

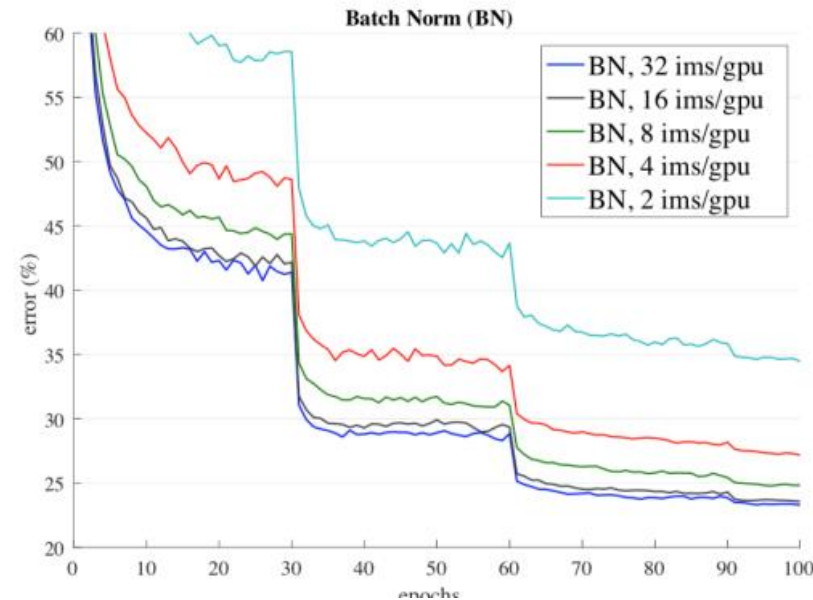
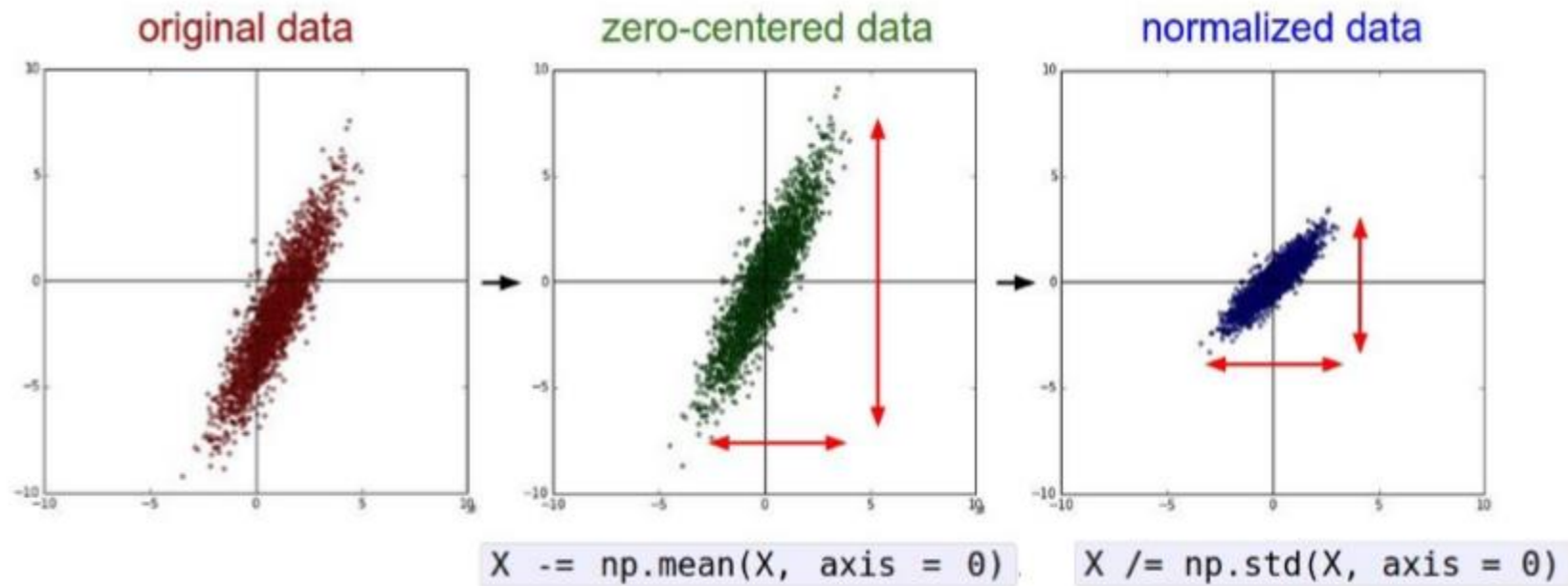


Figure 5. **Sensitivity to batch sizes:** ResNet-50's validation error of BN (left) and GN (right) trained with 32, 16, 8, 4, and 2 images/GPU

Batch size 가

# Babysitting the Learning Process

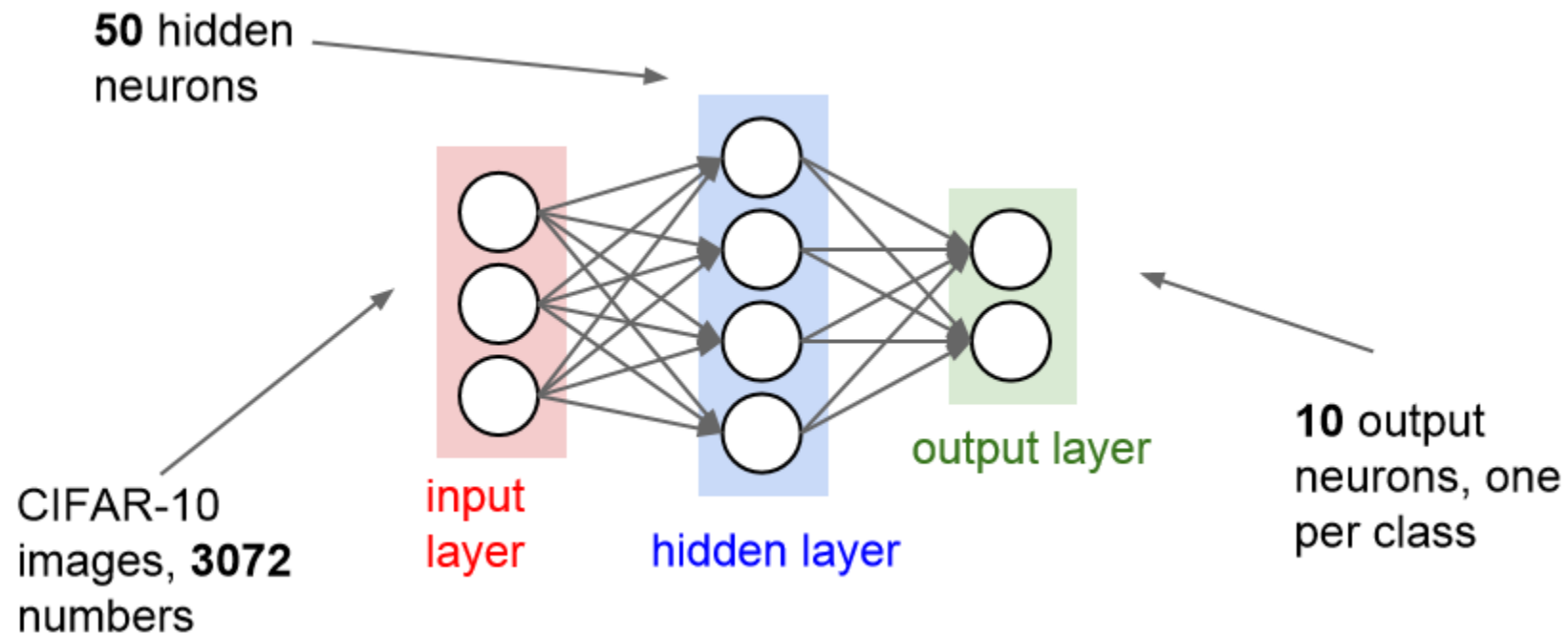
- Step1: Preprocessing the data (+ data augmentation)





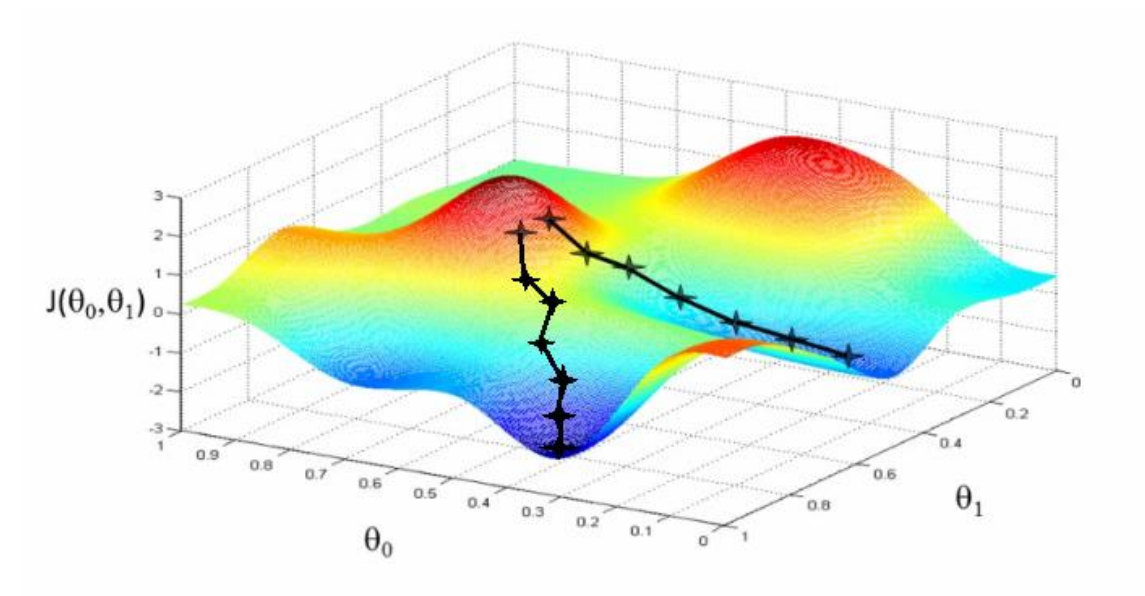
# Babysitting the Learning Process

- Step2: Design the network architecture



# Babysitting the Learning Process

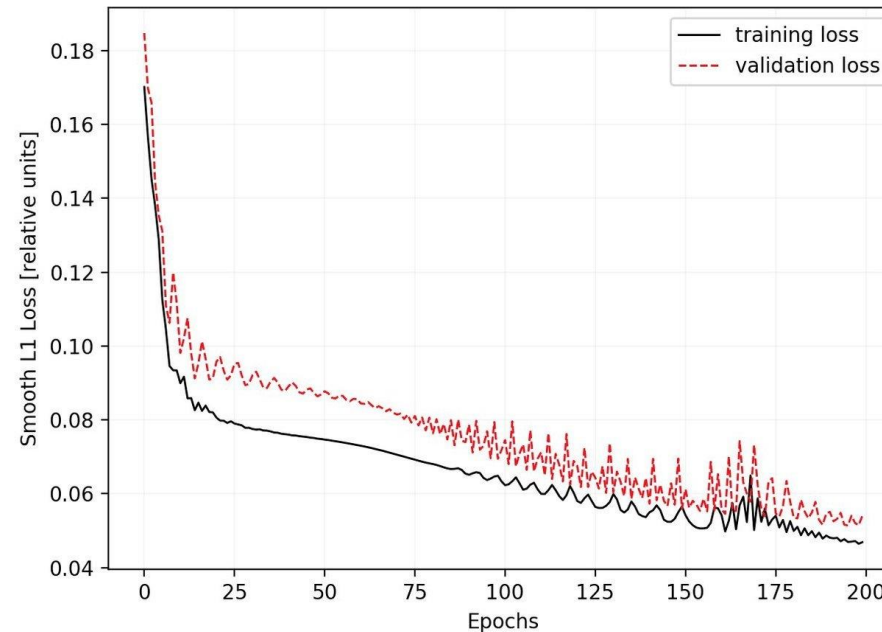
- Step3: Design loss function / Optimizer (lr)



# Babysitting the Learning Process

- Step4: Train model and analyze results

## Visualize train & validation loss



```
Finished epoch 1 / 200: cost 2.302603, train: 0.400000, val 0.400000, lr 1.000000e-03  
Finished epoch 2 / 200: cost 2.302258, train: 0.450000, val 0.450000, lr 1.000000e-03  
Finished epoch 3 / 200: cost 2.301849, train: 0.600000, val 0.600000, lr 1.000000e-03
```

!!

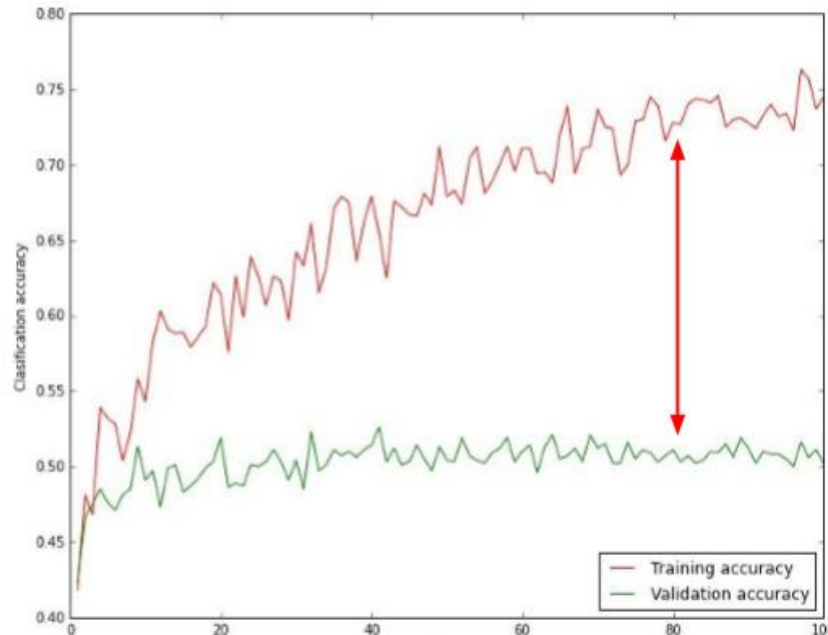
가 loss가

# Babysitting the Learning Process

- Step4: Train model and analyze results

Check **overfitting**

Learning Rate



**big gap = overfitting**

=> increase regularization strength?

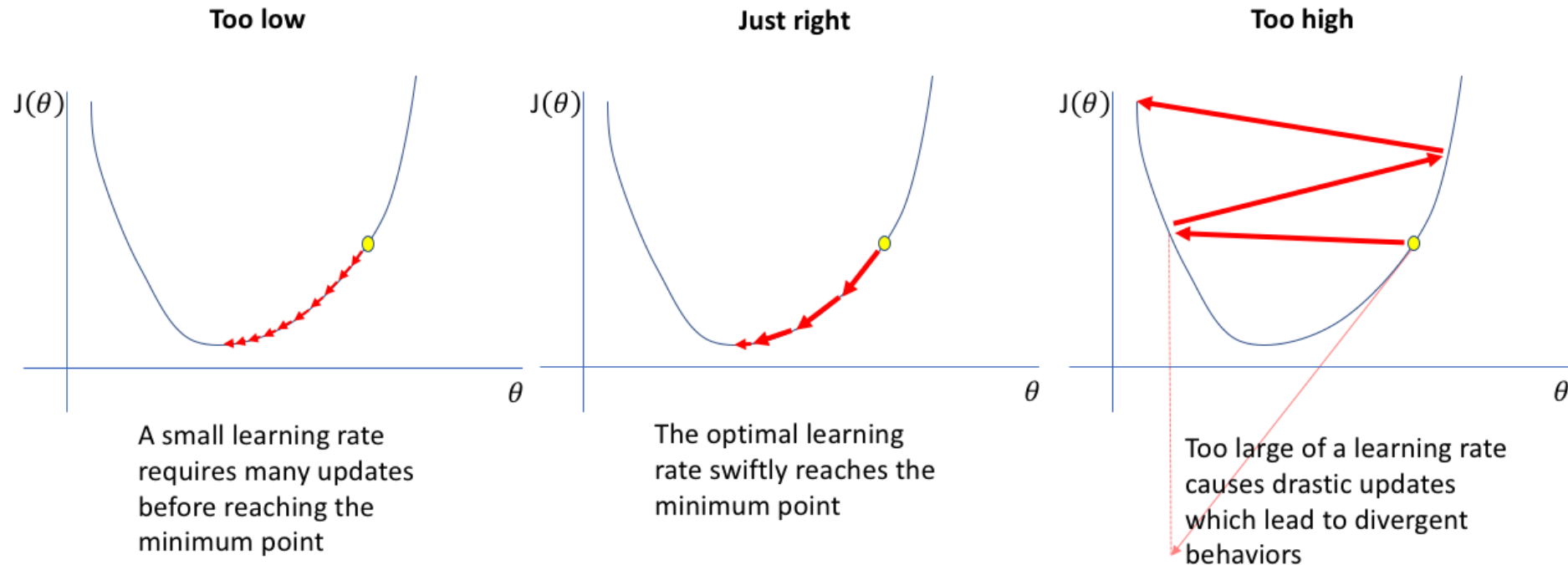
**no gap**

=> increase model capacity?

# Babysitting the Learning Process

- Step4: Train model and analyze results

## Check learning rate



# Babysitting the Learning Process

- Step4: Train model and analyze results



Epoch

Loss가

가 가 X

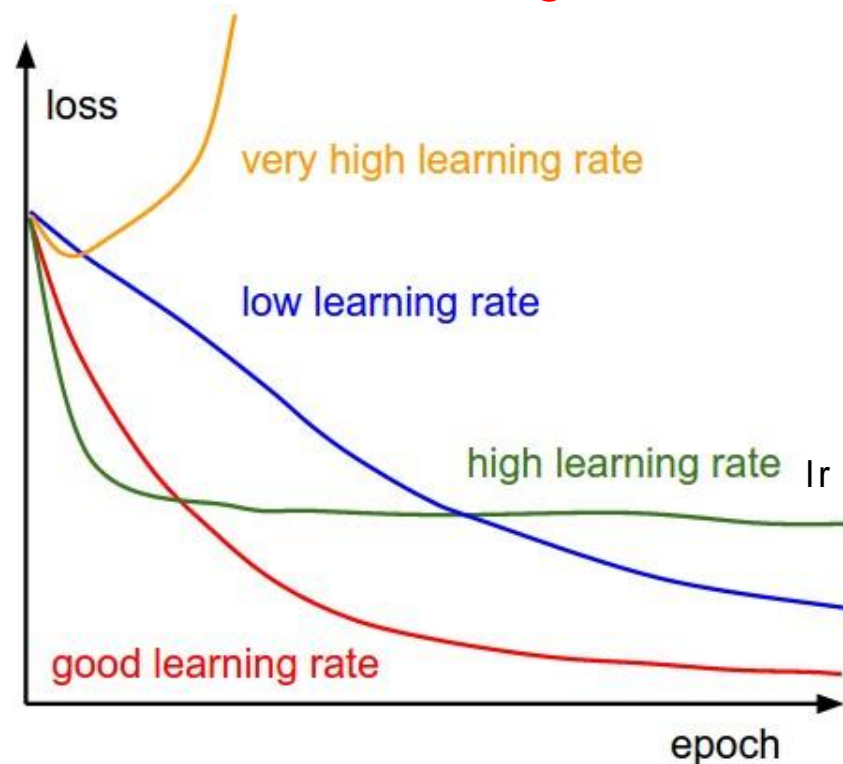
가 .

acc가

...

...

Check learning rate



# Hyperparameter Optimization

Choices about the algorithm that we set rather than learn  
( $\approx$  Heuristic Values)

ex)

- Initial learning rate
- learning rate decay
- batch size 16 32
- epoch 100
- number of layer
- convolutional kernel size 3x3
- pooling type

- activation function
- upsampling method, max pooling
- optimizer Adam .
- data augmentation parameters
- over-sampling ratio
- k-fold cross validation
- etc...

GPU

가 Test

# Summary (Training neural network)

- Activation func. (Use ReLU)
  - Data normalization (Divide 255)
  - Data augmentation (Must do it)
  - Weight Initialization (Xavier init.) Random
  - Batch Normalization (Must do it)
- 
- Babysitting the Learning process
  - Hyperparameter optimization