

ITE4052 Computer Vision

Deep Learning Review

Dong-Jin Kim
Spring 2024



Image Classification

The **simplest** task in visual recognition!



(assume given set of discrete labels)
 $\{\text{dog, cat, truck, plane, ...}\} \leftarrow \text{Image Class}$



cat

Result of Image classification

Image Classification Challenges

Viewpoint variation



Scale variation



Deformation



Occlusion



Illumination conditions



Background clutter



Intra-class variation



Data-Driven Approach

→ All kinds of Image classification challenges \Rightarrow 해결할 수 있.

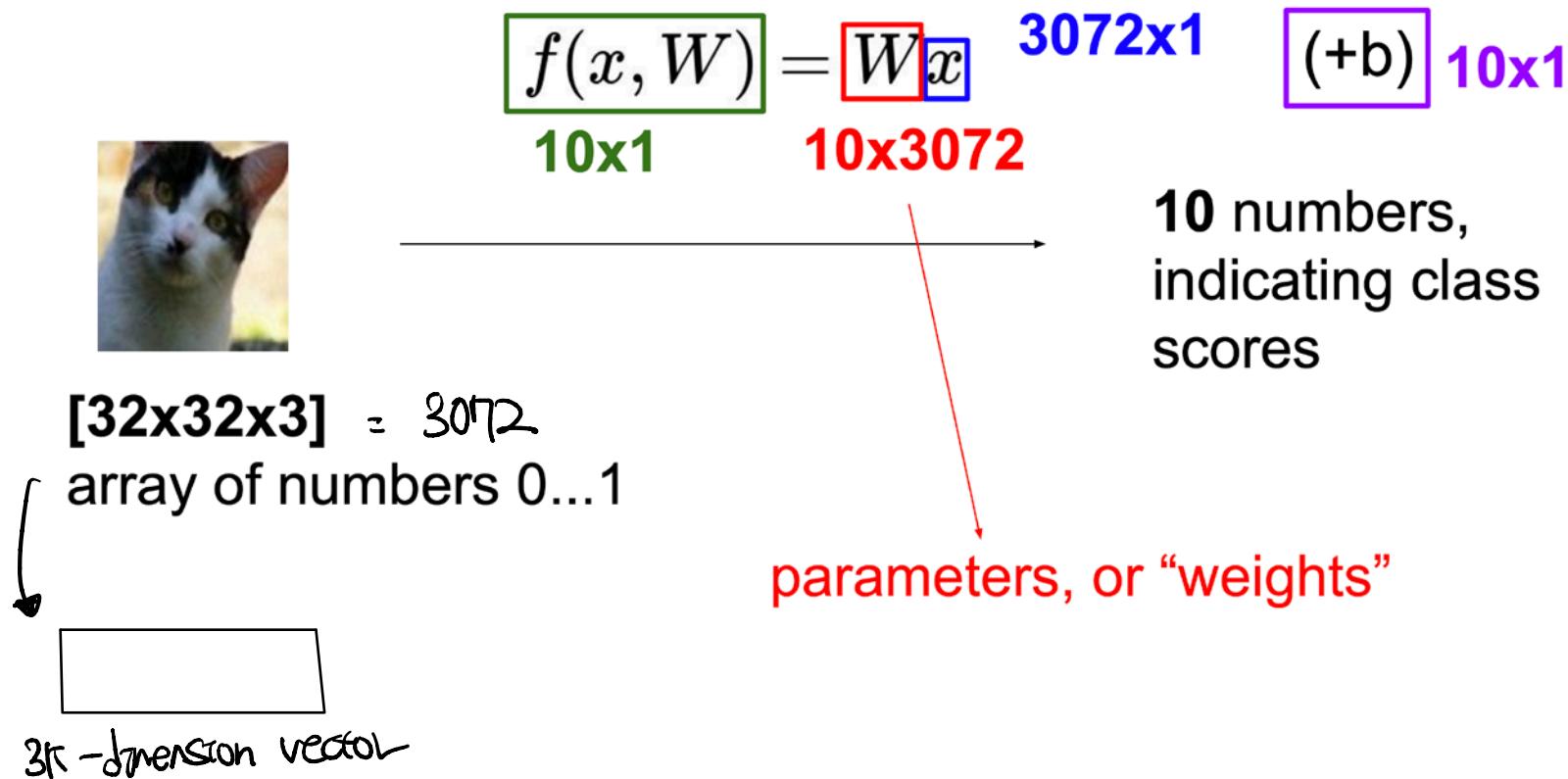
1. Collect a **dataset** of images and labels
2. Use **Machine Learning** to train an image classifier
3. Evaluate the classifier on a withheld set of **test images**

```
def train(train_images, train_labels):  
    # build a model for images -> labels...  
    return model  
  
def predict(model, test_images):  
    # predict test_labels using the model...  
    return test_labels
```

Example training set

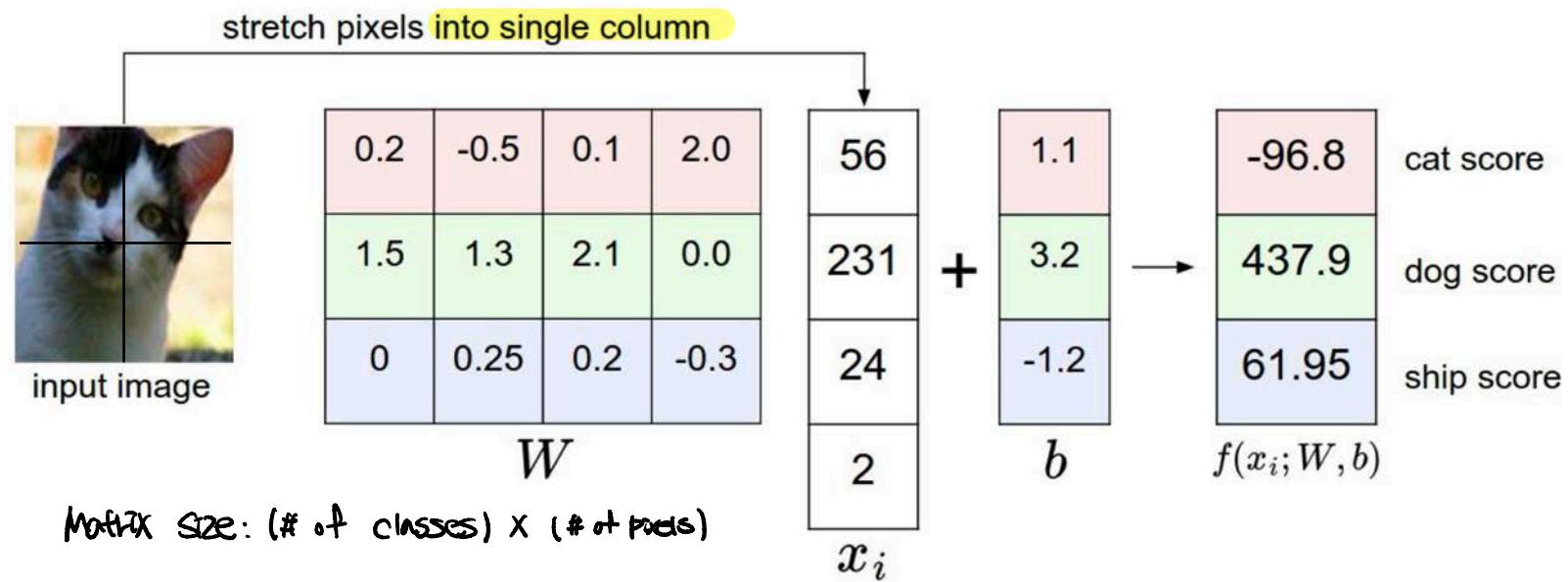


Parametric Approach : Linear Classifier



Parametric Approach : Linear Classifier

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



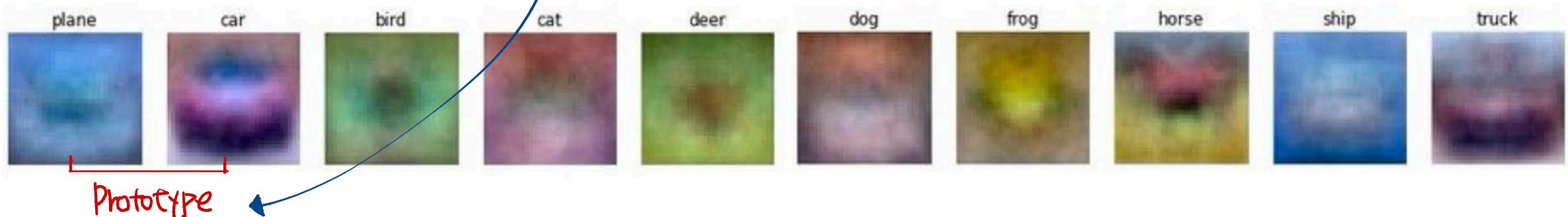
Interpreting a Linear Classifier: Visual

Linear Classifier Formulation

$$f(x_i, W, b) = \boxed{W}x_i + b$$

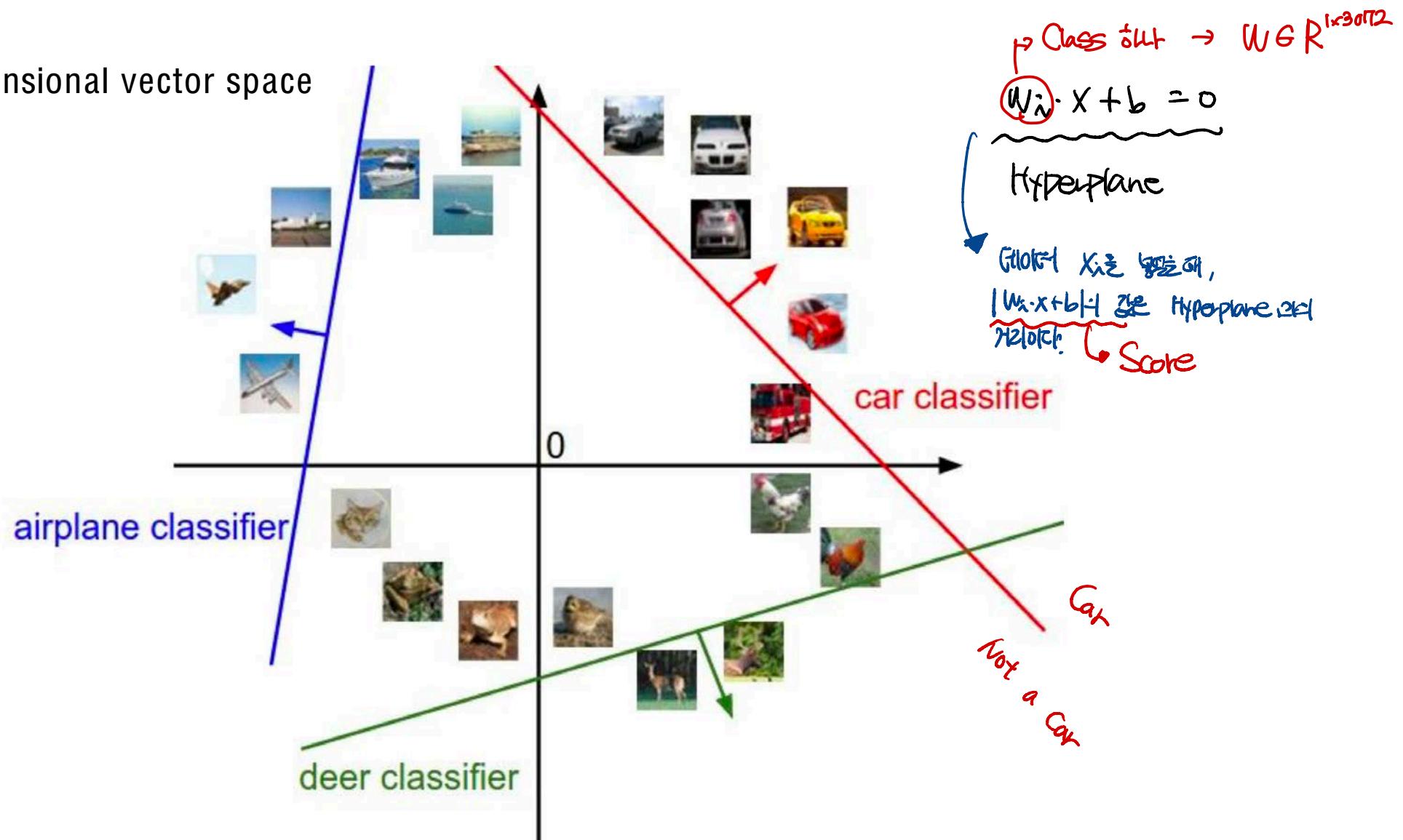
$$\hat{f}_i = \underline{W_i x_i + b_i}$$

Example **trained weights** of a linear classifier trained on CIFAR-10:



Interpreting a Linear Classifier: Geometric

In $32 \times 32 \times 3 = 3072$ dimensional vector space



Training a Linear Classifier

Suppose: 3 training examples, 3 classes.

With some W , the scores are:

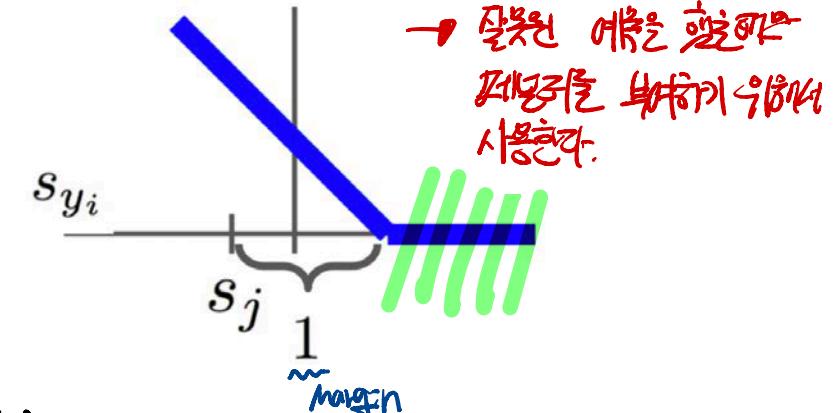
$$f(x, W) = \overbrace{Wx}^{3 \times 3072}$$



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9		

1. Multiclass SVM loss:

“Hinge loss”



2. Softmax Loss

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Wrong *Answer*

Ex) Car-cat : $\max(0, 5.1 - 3.2 + 1) = 2.9$

Softmax Classifier

= Multinomial Logistic Regression

고양이 = 1 이기 때문!

Cat probability ↑

→ Cat probability ↓ and frog probability ↓

cat

3.2

car

5.1

frog

-1.7



$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

예측 Data given Softmax

where $s = f(x_i; W)$

Want to maximize the **log likelihood** of the correct class:

$$L_i = -\log P(Y = y_i | X = x_i)$$

→ minimize negative log likelihood

In summary :

$$L_i = -\log \left(\frac{e^{s_i}}{\sum_j e^{s_j}} \right)$$

→ Softmax는 확률 분포다.

Scores as **probabilities** (≥ 0 , sum to 1) of the classes.

$$\text{Softmax} = \frac{e^s}{\sum_j e^{s_j}}$$

\leftarrow exponential
Normalization

Softmax Classifier

= Multinomial Logistic Regression



unnormalized probabilities

cat
car
frog

3.2
5.1
-1.7

exp

24.5
164.0
0.18

normalize

0.13
0.87
0.00

$$L_i = -\log(0.13) = 0.89$$

Comparing with **ground truth** (one-hot vector)
:Kullback–Leibler (**KL**) Divergence

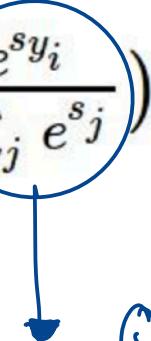
$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

ex: cat one-hot $P = [1, 0, 0]$

$D_{\text{KL}}(P \parallel Q)$

$$= 1 \cdot \log\left(\frac{1}{0.13}\right) + 0 + 0$$

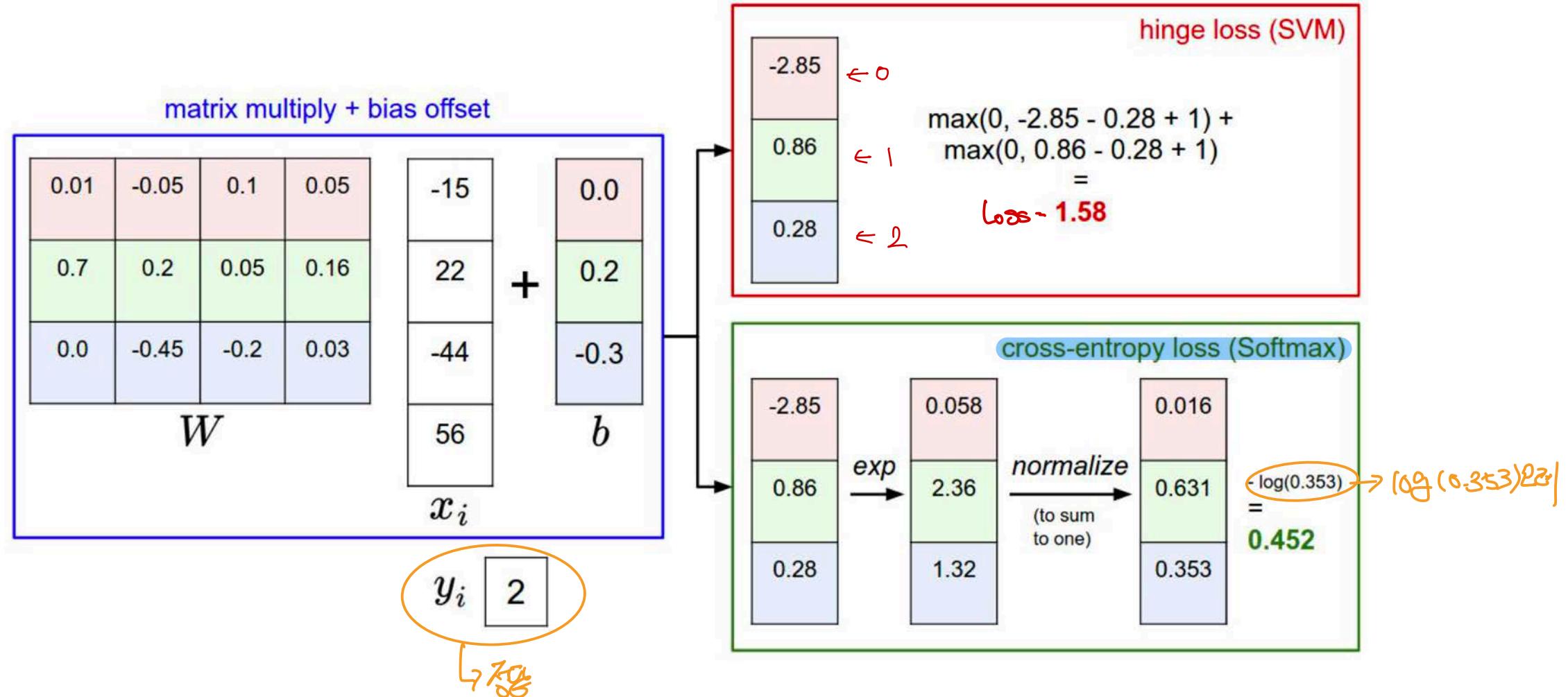
$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$



unnormalized log probabilities

probabilities

SVM vs Softmax (Cross Entropy Loss)



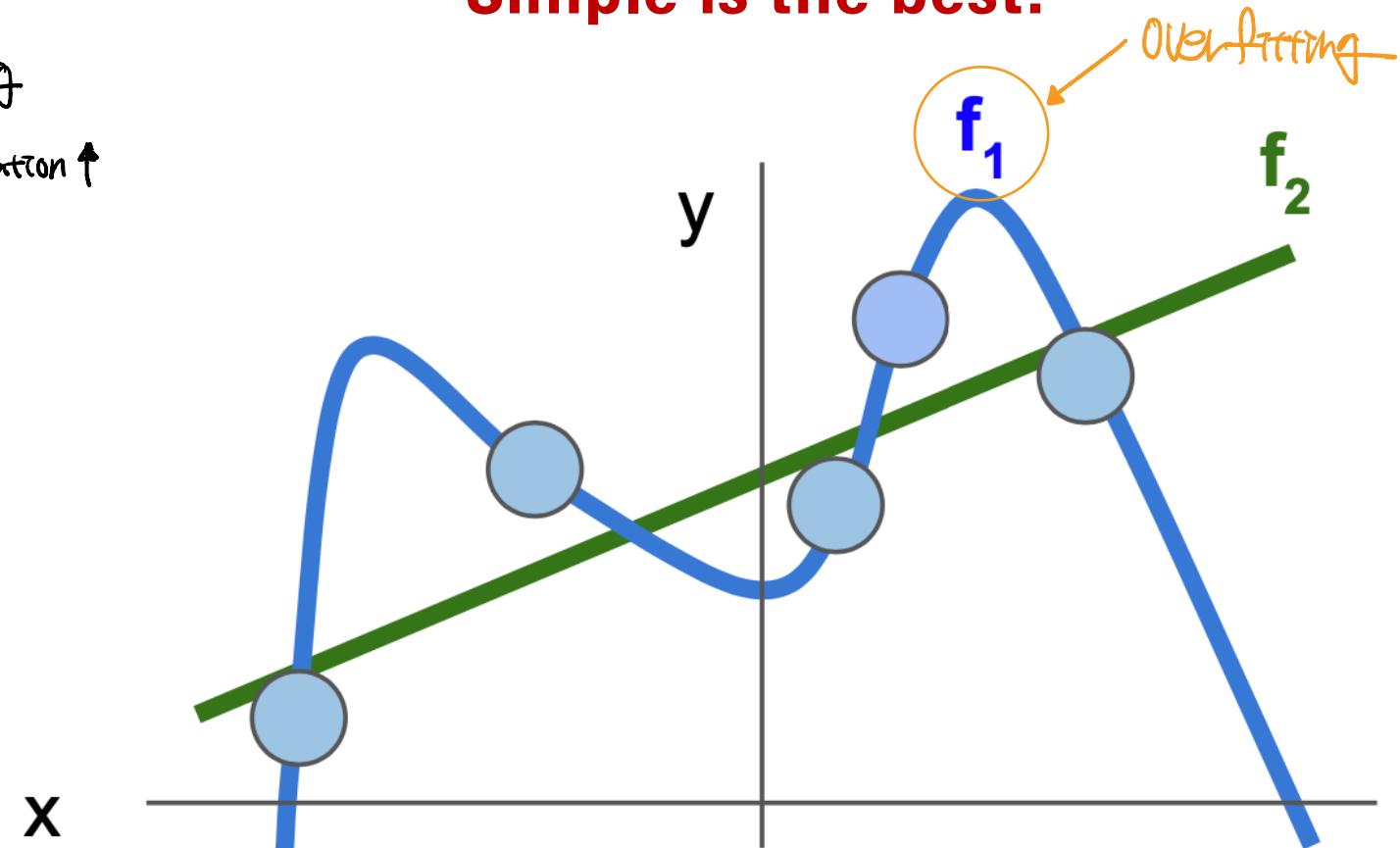
Regularization Intuition

: Simpler Models

If not overfitting

Loss↑ but regularization↑

“Simple is the best!”



Regularization avoids overfitting on the training data

Too much fitting on training data

Weight Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss : Following data distribution.}} - \lambda R(W)$$

: Regularization Strength (hyperparameter)

Data loss :
Following data distribution.

Regularization :
Preventing **overfitting**.

In common use:

L2 regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

L1 regularization

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

Elastic net (L1 + L2)

$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

L2, L1 etc

More complex :

Dropout

Batch normalization

Stochastic depth, etc



Optimization

Follow the Slope – Analytic Gradient

Analytic Gradient

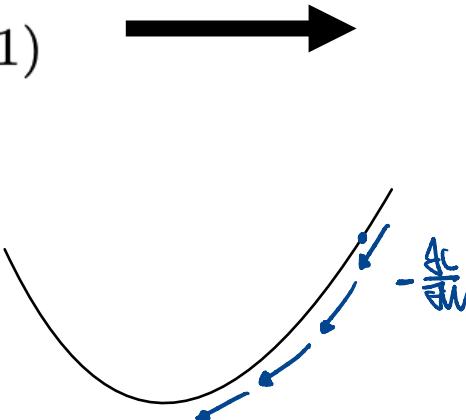
The loss is just a function of W :

$$\text{minimize} \rightarrow L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$s = f(x; W) = Wx$$

want : $\nabla_W L$ $\frac{\partial L}{\partial W}$



Use **calculus** to compute an analytic gradient

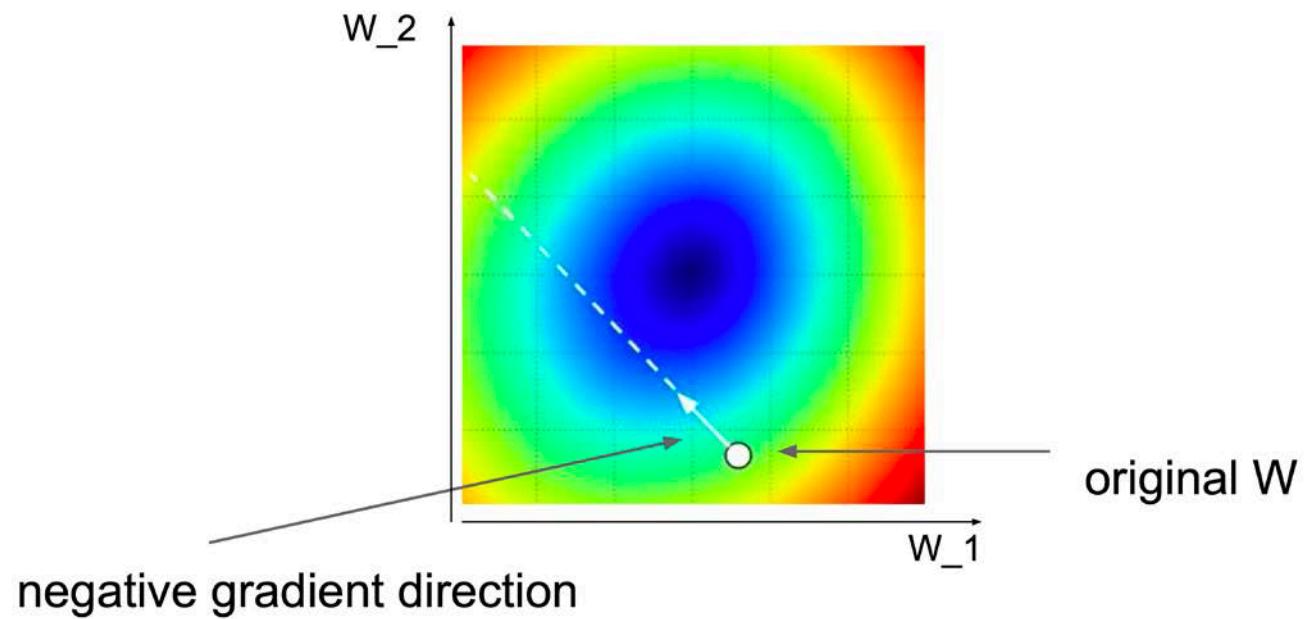
→ **partial derivative!**

Gradient Descent

Full - batch

$$\hookrightarrow \frac{1}{N} \sum_{i=1}^N \frac{\partial L}{\partial w_i}$$

```
● ● ●  
# Vanilla Gradient Descent  
  
while True:  
    weights_grad = evaluate_gradient(loss_fn, data, weights)  
    weights += - step_size * weights_grad # perform parameter  
update
```



Mini-batch Gradient Descent

Only use a **small portion** of the training set to compute the gradient

$$\frac{1}{B} \sum_{i=1}^B \frac{\partial C}{\partial W_i}$$

Randomly selected Batch Sample

```

● ● ●

# Vanilla Minibatch Gradient Descent

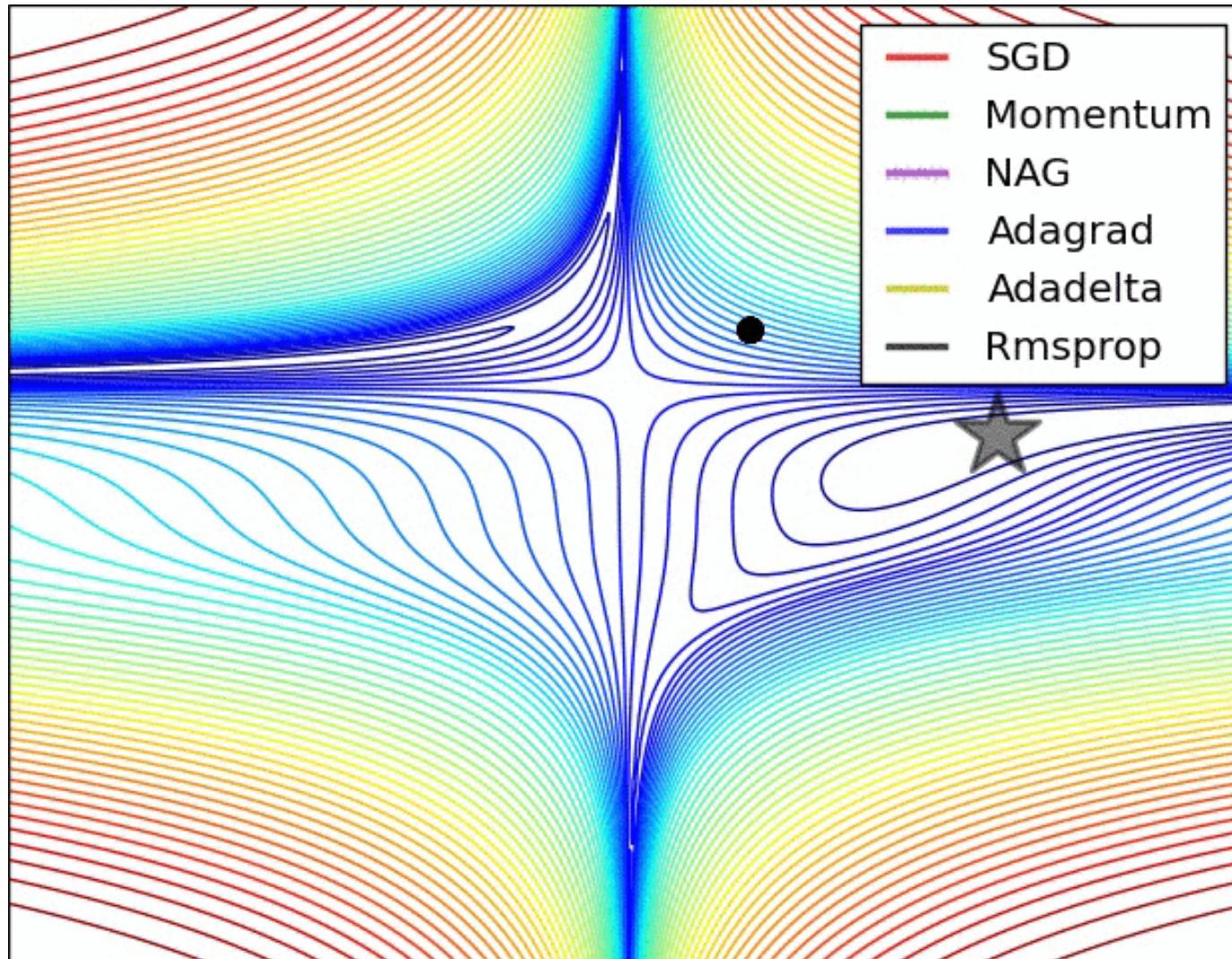
while True:
    data_batch = sample_training_data(data, 256) # sample 256
    examples
    weights_grad = evaluate_gradient(loss_fn, data_batch,
    weights)
    weights += - step_size * weights_grad # perform parameter
    update
  
```

Common mini-batch sizes are 32/64/128/... examples

$\frac{N}{B}$ 개의 Gradient descent 실행
 → 뚱뚱개처럼 Batch GD보다 빠름
 → 일정한 흐름

- Various Optimization Methods
- Momentum
 - Adagrad
 - RMSProp
 - Adam
 - :

Gradient Descent



Stochastic Gradient Descent (SGD)

Mini-Batch Size \rightarrow Sampling Error \rightarrow 흐름적

$$w \leftarrow w - \eta \frac{1}{N} \sum_{i=1}^N \nabla L_i(x_i, y_i, w)$$

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

Full sum expensive
when N is large!

→ Reason why use
small number of batch



```
# Vanilla Minibatch Gradient Descent
```

```
while True:
    data_batch = sample_training_data(data, 256) # sample 256
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```

Learning Rate Schedules

$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

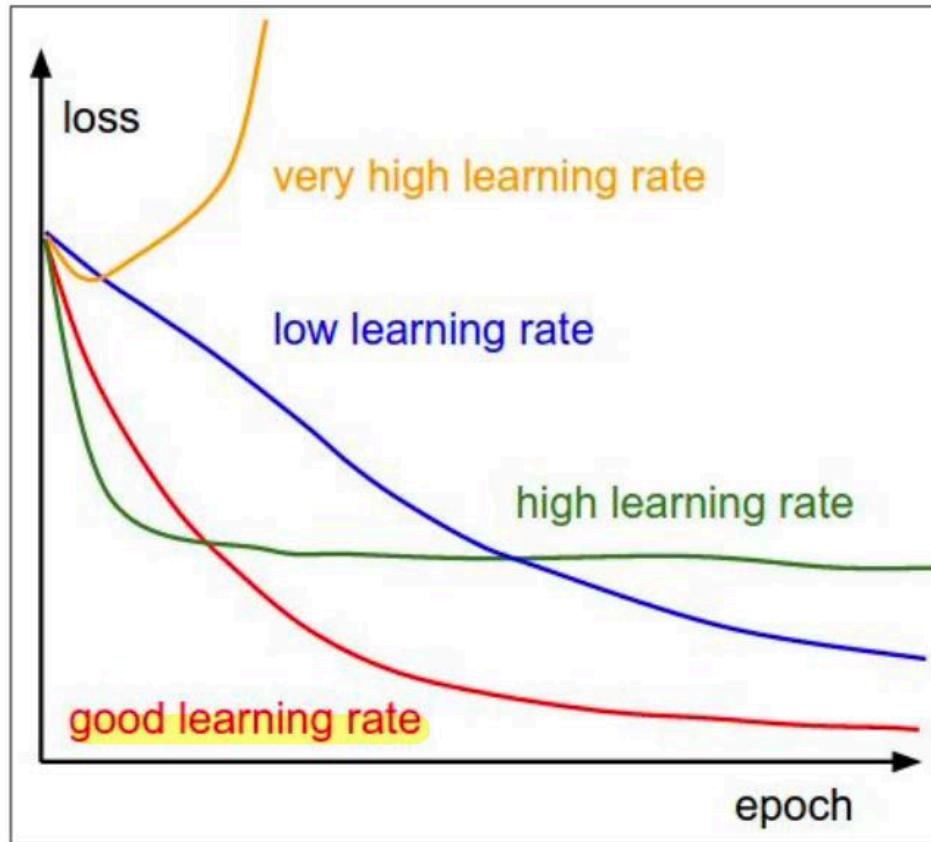
Learning rate (step size)

```
● ● ●  
# Vanilla Gradient Descent  
  
while True:  
    weights_grad = evaluate_gradient(loss_fun, data, weights)  
    weights += - step_size * weights_grad # perform parameter update
```

Learning rate

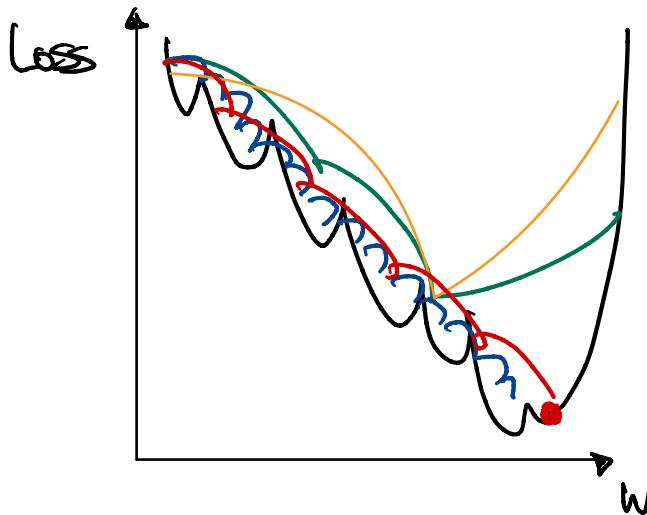
Learning Rate Schedules

Learning rate as a hyperparameter.

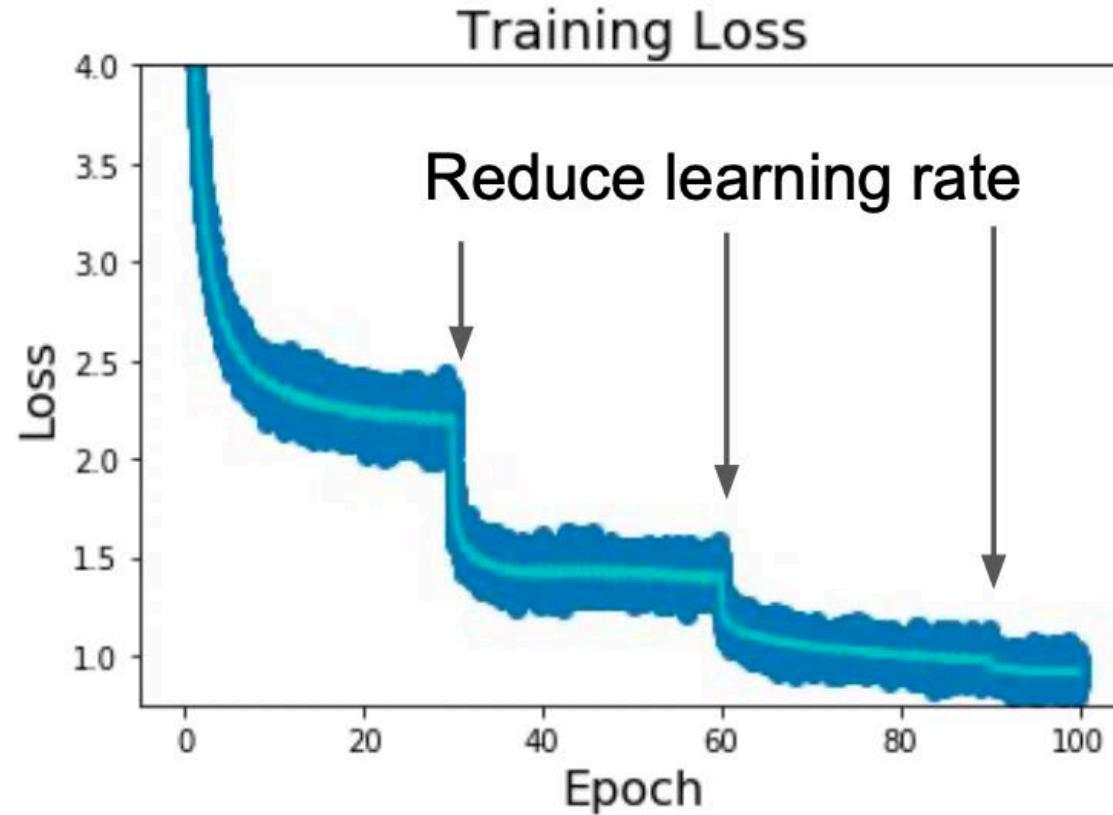


Q. Which one of these learning rates is best to use?

→ ↗ Learning rate가 낮을 때 Loss가 빠르게 줄어들고
Learning rate가 높을 때 Loss가 빠르게 증가함.



Learning Rate Decay



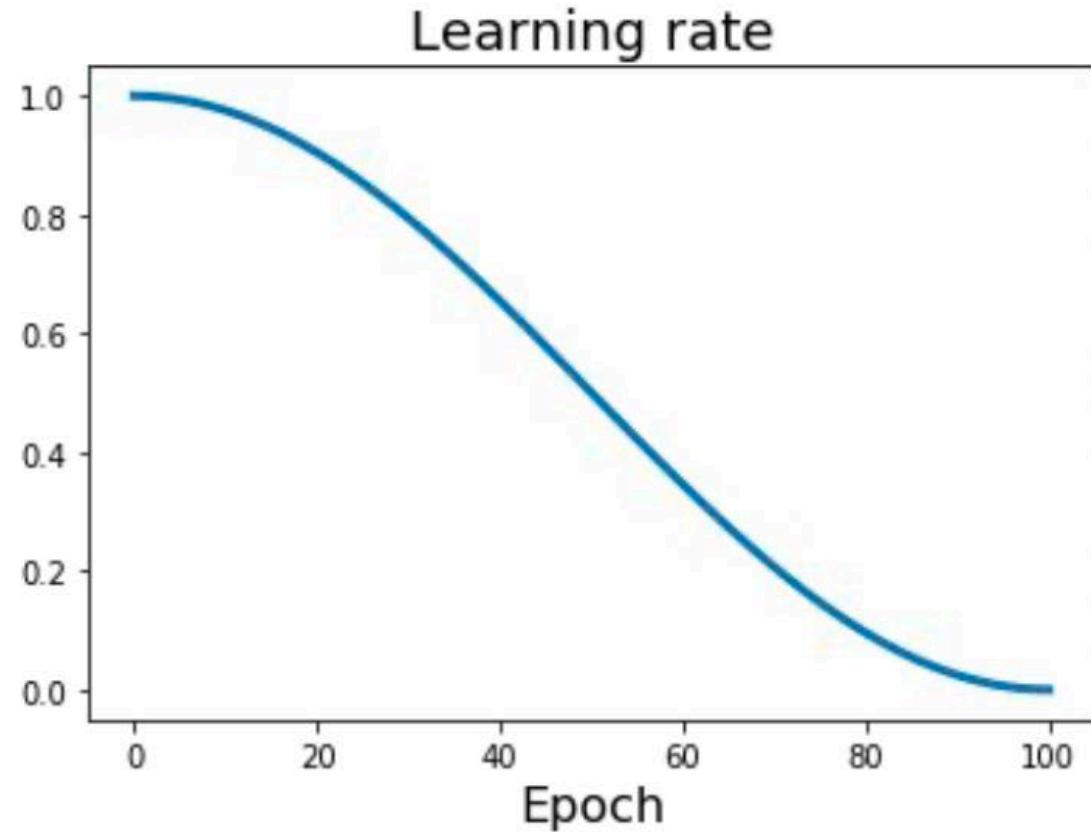
Step: Reduce LR at a few fixed points
(E.g. multiply LR by 0.1 after epochs 30, 60, 90)

Epoch마다 LR을 줄임

학습률: W 가 커지면 $\rightarrow LR \uparrow$

학습률: W 가 작아지면 $\rightarrow LR \downarrow$

Learning Rate Decay



Step: Reduce LR at a few fixed points
(E.g. multiply LR by 0.1 after epochs 30, 60, 90)

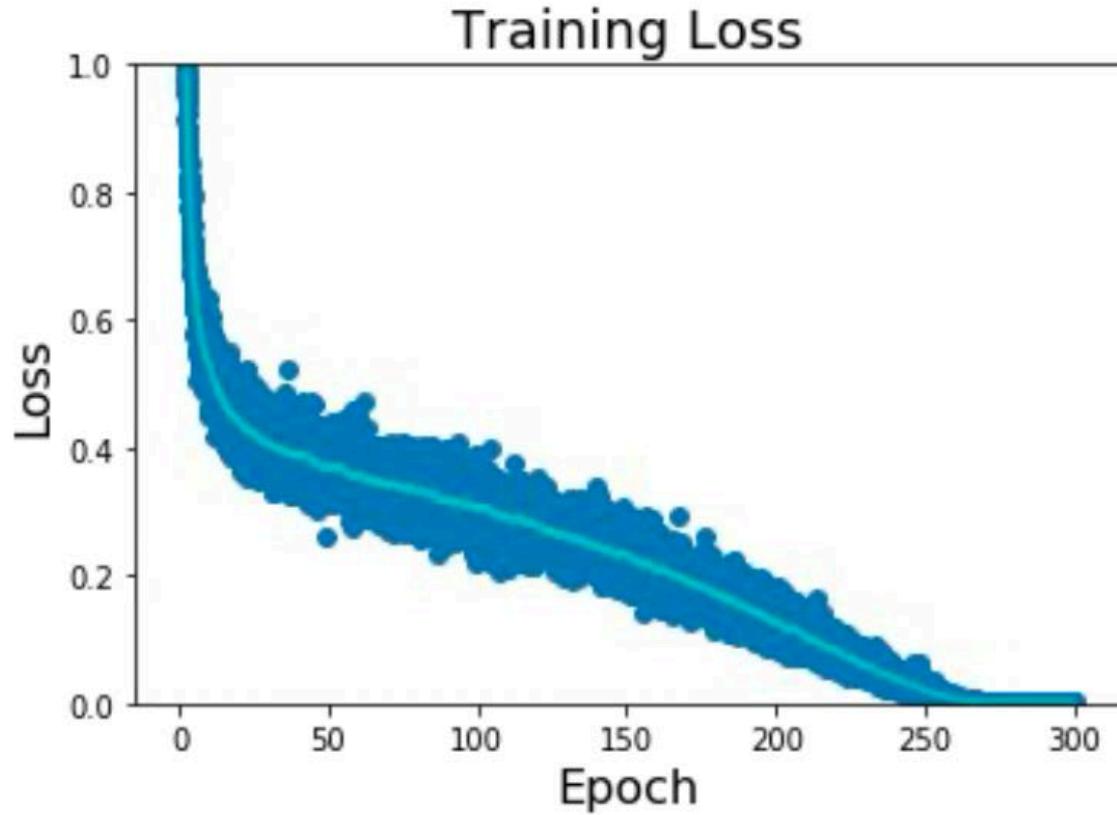
Cosine: $\alpha_t = \frac{1}{2}\alpha_0(1 + \cos(t\pi/T))$

α_0 : Initial learning rate

α_t : Learning rate at epoch t

T : Total number of epochs

Learning Rate Decay



Step: Reduce LR at a few fixed points
(E.g. multiply LR by 0.1 after epochs 30, 60, 90)

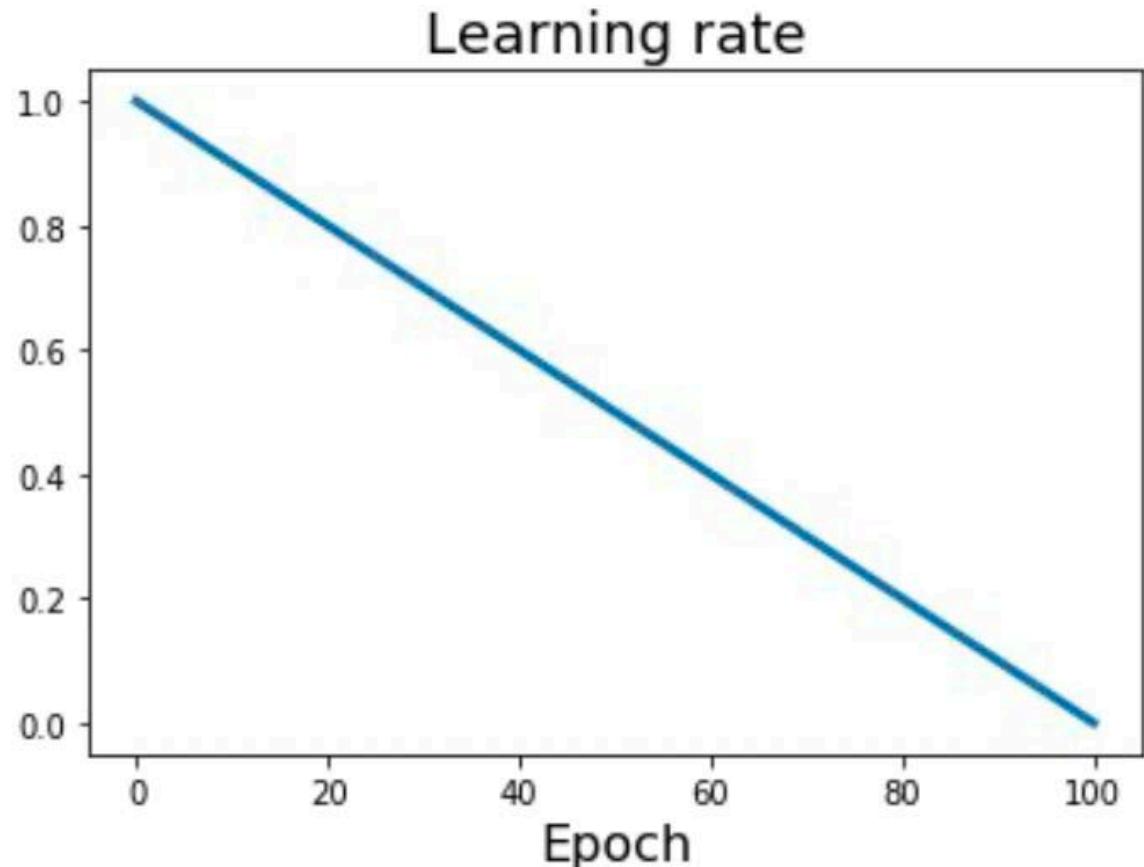
Cosine: $\alpha_t = \frac{1}{2}\alpha_0(1 + \cos(t\pi/T))$

α_0 : Initial learning rate

α_t : Learning rate at epoch t

T : Total number of epochs

Learning Rate Decay



Step: Reduce LR at a few fixed points
(E.g. multiply LR by 0.1 after epochs 30, 60, 90)

Cosine: $\alpha_t = \frac{1}{2}\alpha_0(1 + \cos(t\pi/T))$

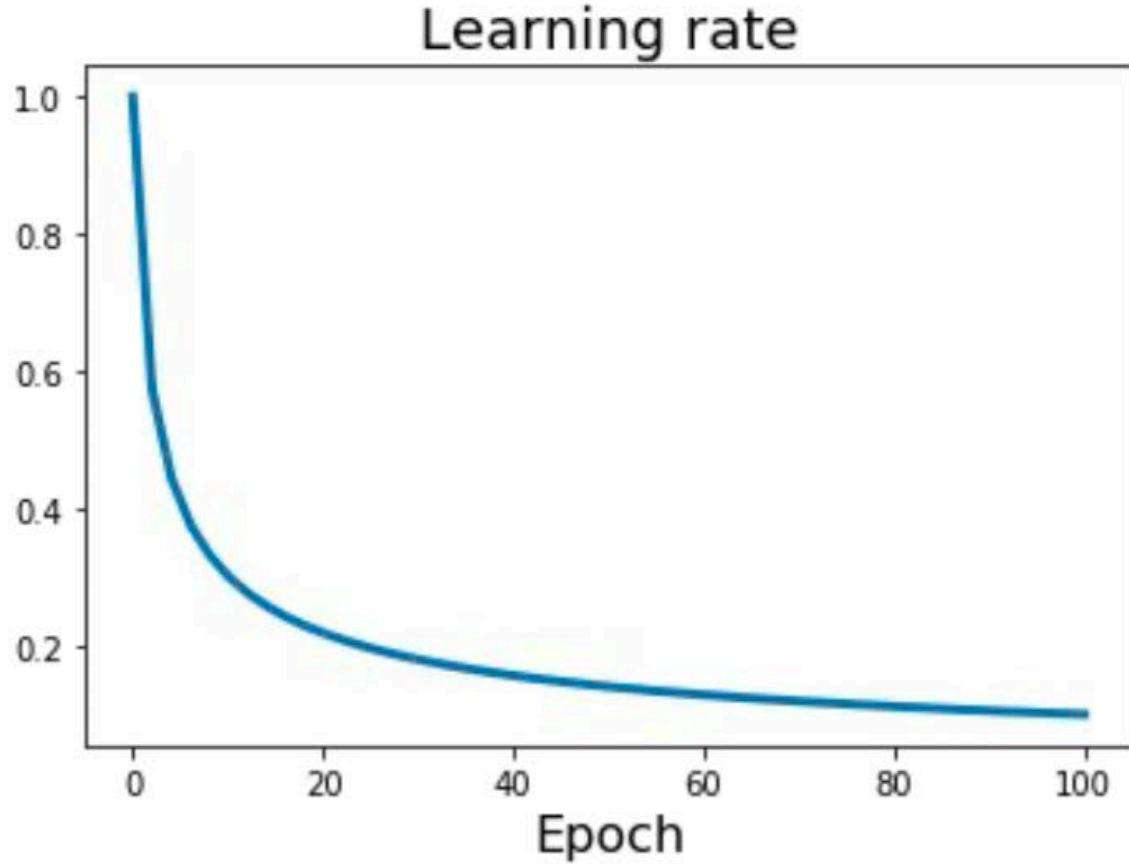
Linear: $\alpha_t = \alpha_0(1 - t/T)$

α_0 : Initial learning rate

α_t : Learning rate at epoch t

T : Total number of epochs

Learning Rate Decay



Epoch마다 끝나는/ 고정

Step: Reduce LR at a few fixed points
(E.g. multiply LR by 0.1 after epochs 30, 60, 90)

Cosine: $\alpha_t = \frac{1}{2}\alpha_0(1 + \cos(t\pi/T))$

Linear: $\alpha_t = \alpha_0(1 - t/T)$

Inverse sqrt : $\alpha_t = \alpha_0/\sqrt{t}$

α_0 : Initial learning rate

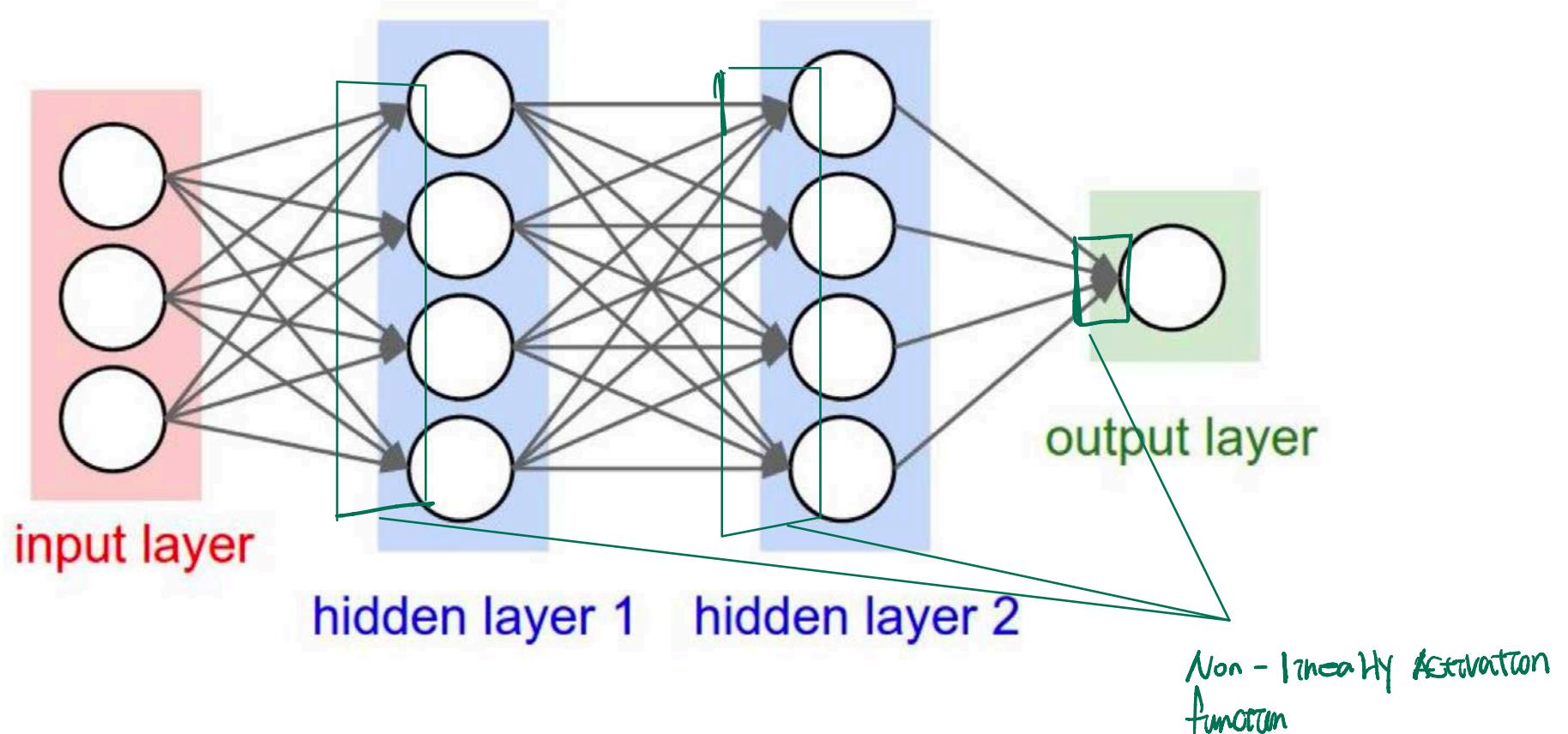
α_t : Learning rate at epoch t

T : Total number of epochs



Neural Networks

Neural Networks: Architectures



“3-layer Neural Net” or “2-hidden-layer Neural Net”

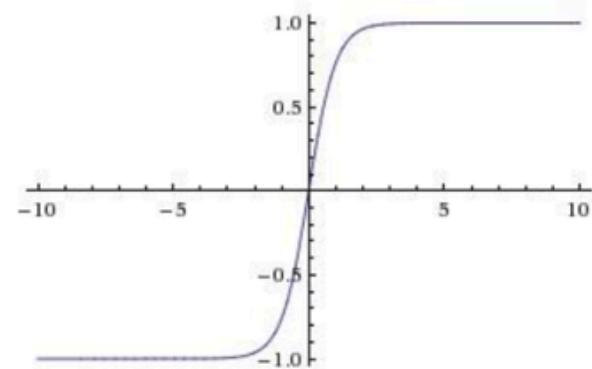
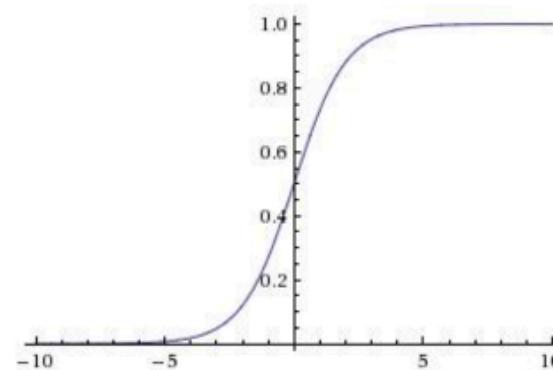
This Neural Networks is “Fully-connected” layers

모든 Input Neurons → Output Neurons

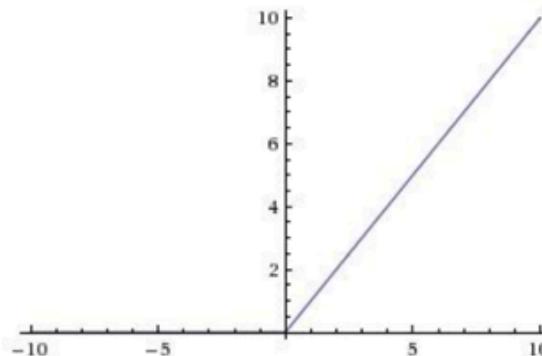
Activation Functions

: Non-linear function

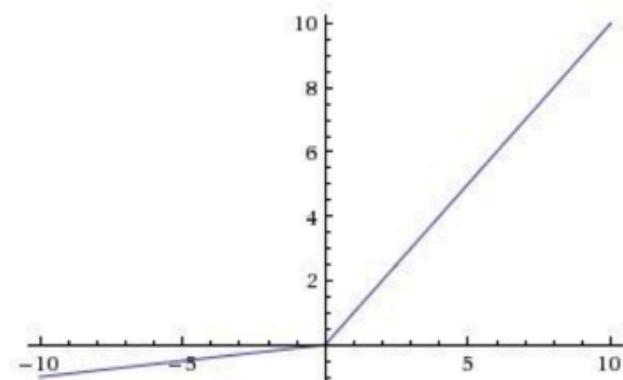
Sigmoid $\sigma(x) = 1 / (1 + e^{-x})$ tanh $\tanh(x)$



ReLU $\max(0, x)$



LeakyReLU $\max(0.1x, x)$



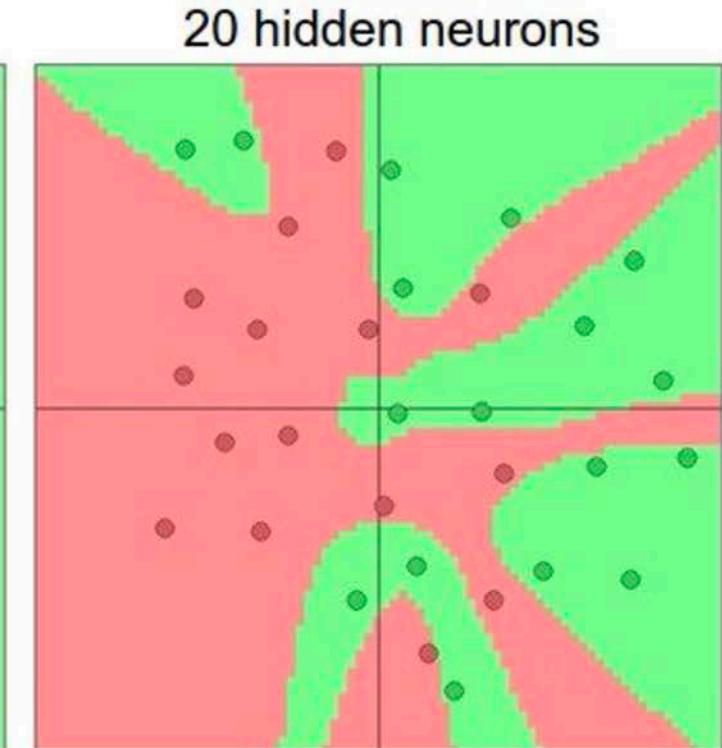
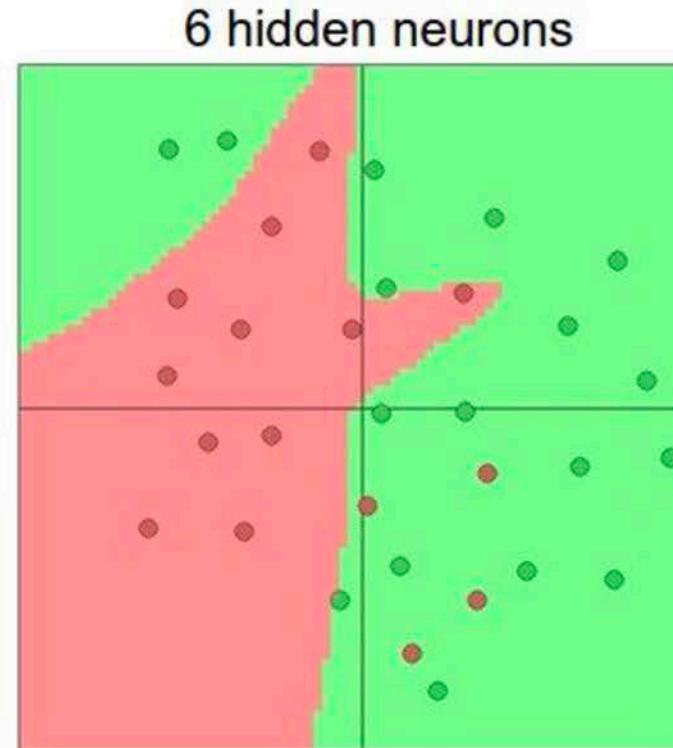
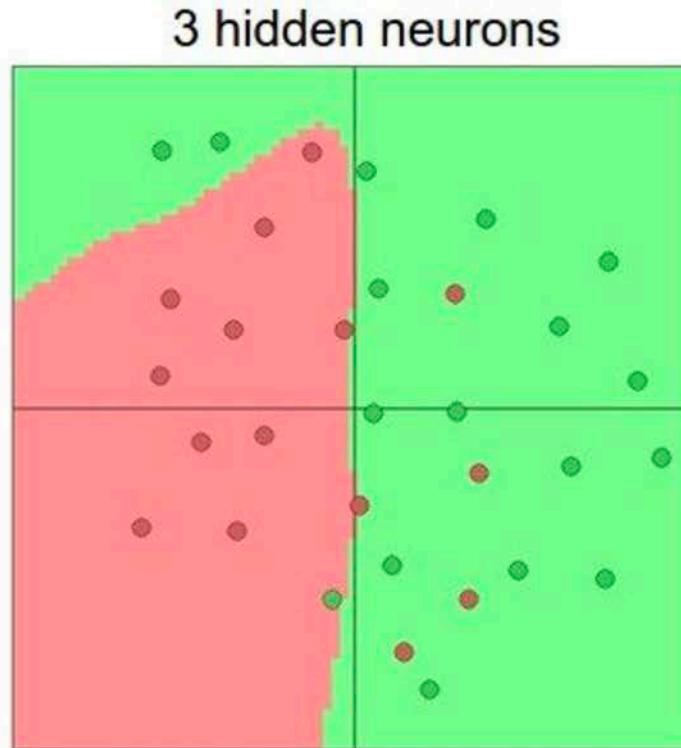
Maxout

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

ELU

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

Setting the Number of Layers and Sizes



Layer ↑ → 비선형성 증가, more complex decision boundary
→ Deep learning의 핵심 원리

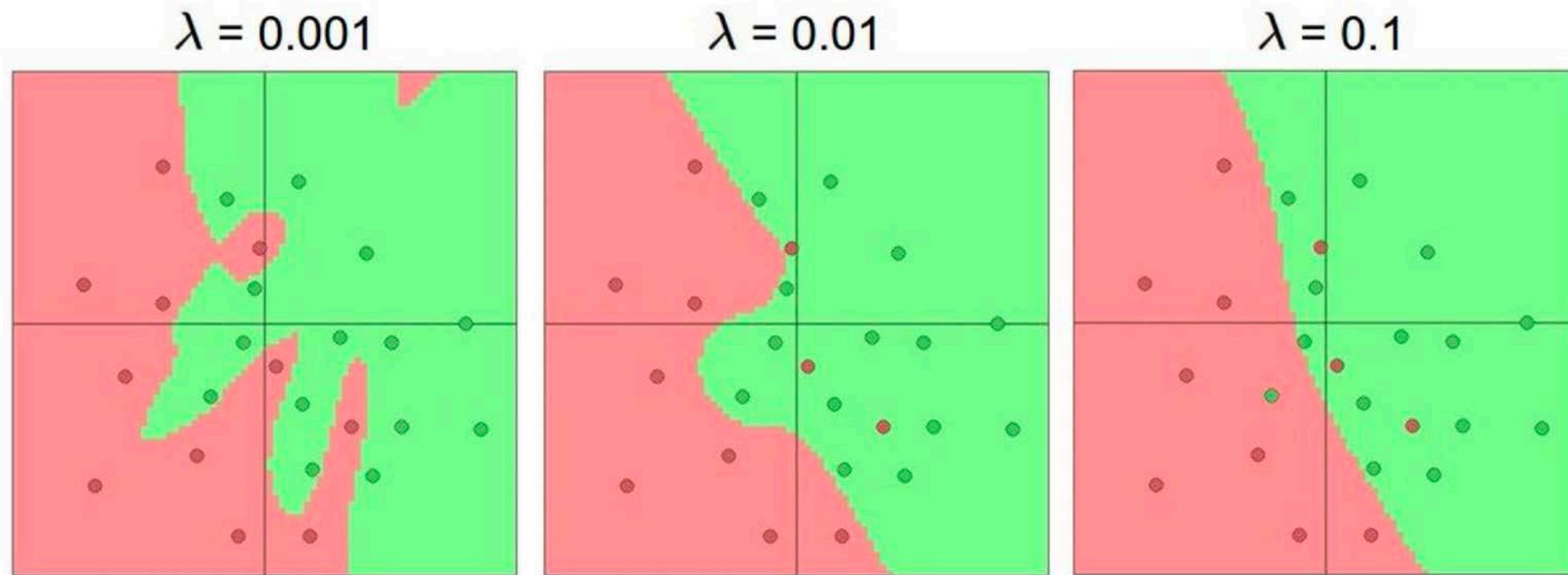
↑
more neurons = more capacity of data

Setting the Number of Layers and Sizes

* More layers & neurons but, add some regularization

$$\frac{1}{n} \sum L + \lambda R(W)$$

Do not use size of neural network as a regularizer. Use **stronger regularization** instead :

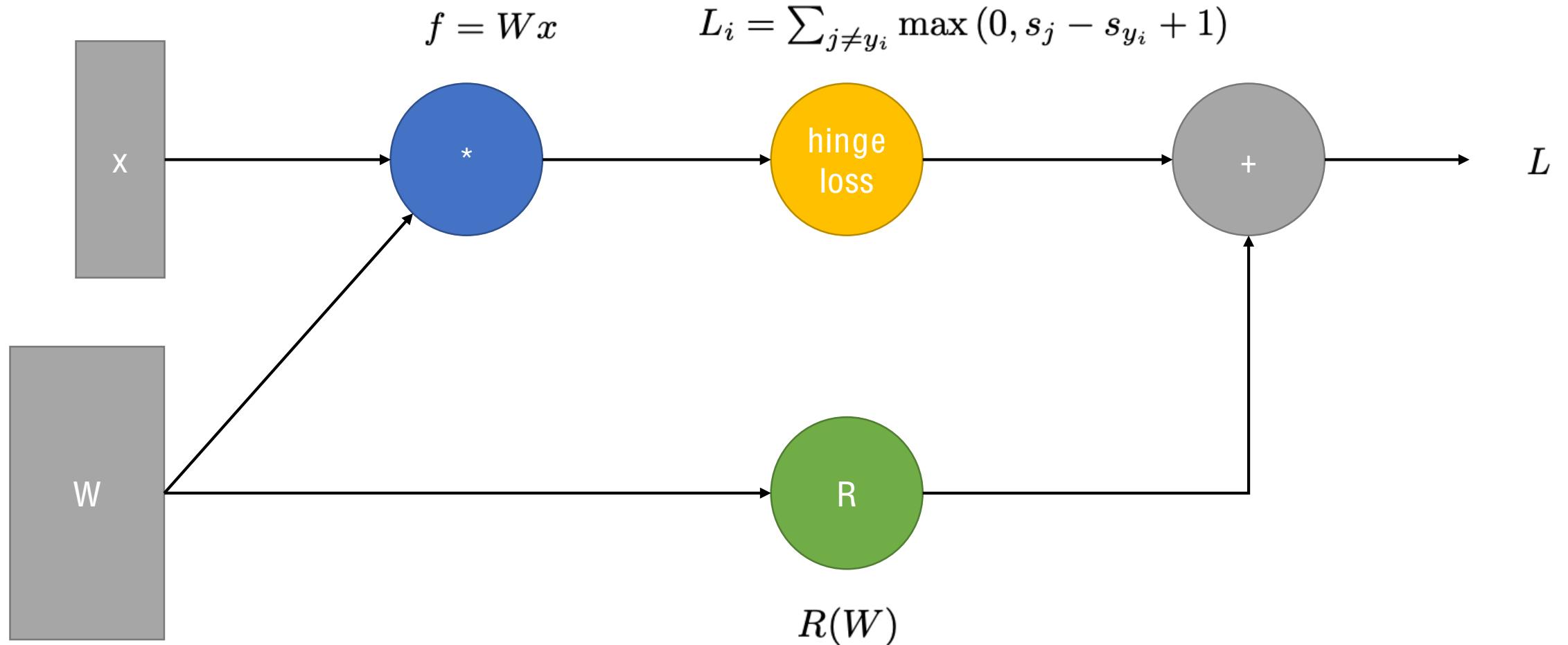


$\lambda \uparrow \rightarrow$ Overfitting ↓

* Overfitting \Leftrightarrow Model Size
Large, Regularization small!

Backpropagation

Computational Graph



Backpropagation

$$\frac{\Delta L}{\Delta z} = \frac{\Delta f}{\Delta q} \cdot \frac{\Delta q}{\Delta z}$$

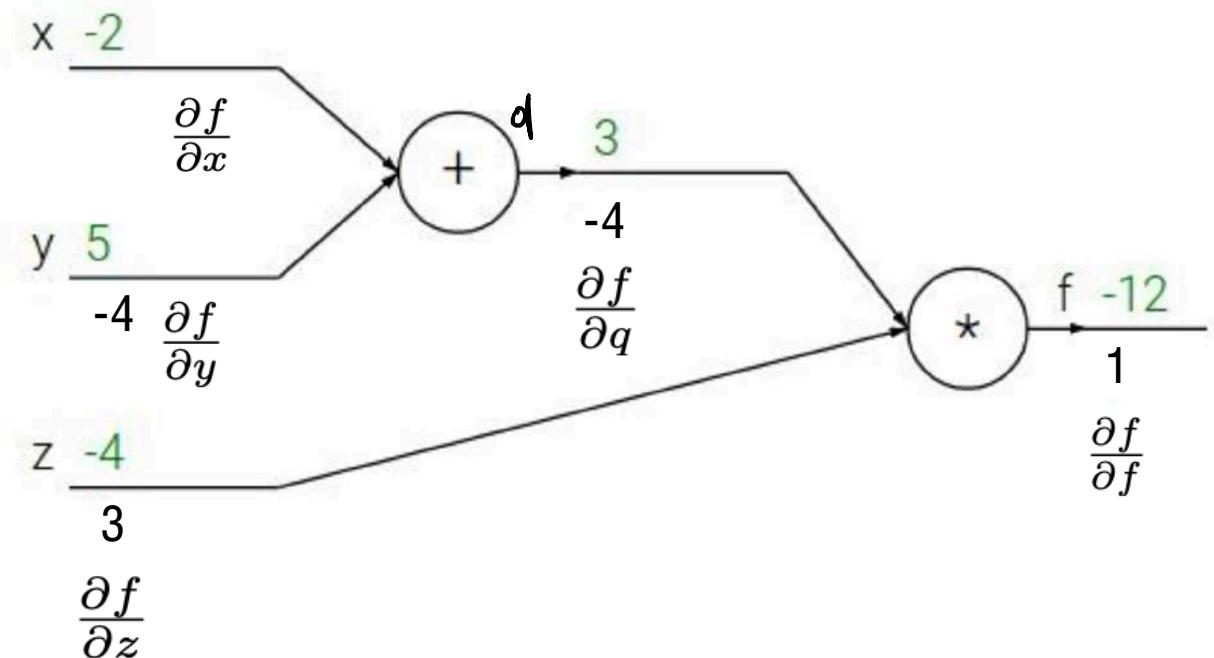
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want : $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

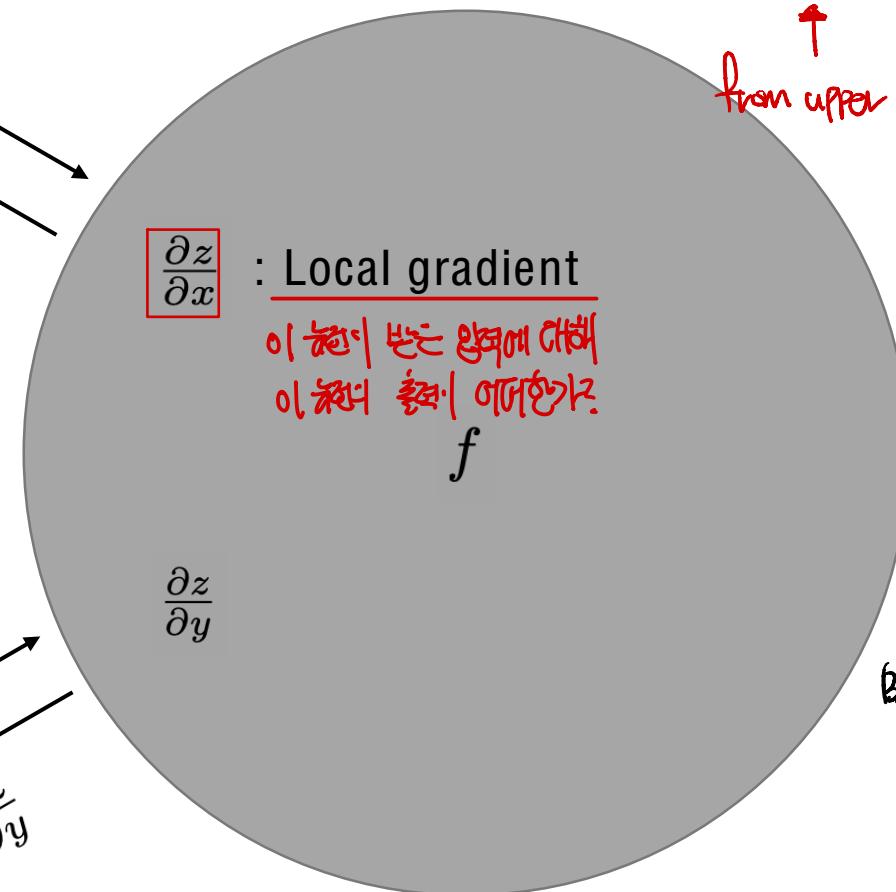
$$\text{Chain Rule : } \frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

$$z \cdot 1 = z$$

x : Activations

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y}$$



Passed to the next layer

$$\frac{\Delta L}{\Delta z} = \frac{1}{\Delta z} \cdot \frac{\Delta z}{\Delta a} \leftarrow \text{Local gradient}$$

from upper layer

(Inside layer, 뉴런 간의 연산.)

$$z = f(\theta, w)$$

z

$\frac{\partial L}{\partial z}$ gradients

Backpropagation

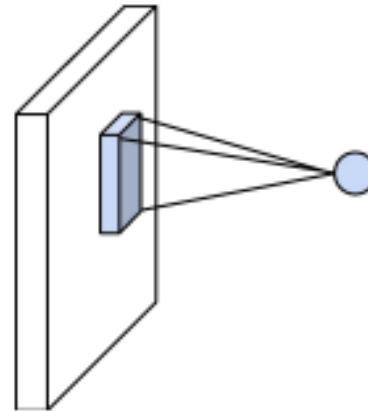
상위 뉴런에서 오는 그라디언트는 뉴런 간의 연산에 따라
증가하고 하위 뉴런으로 전달된다.



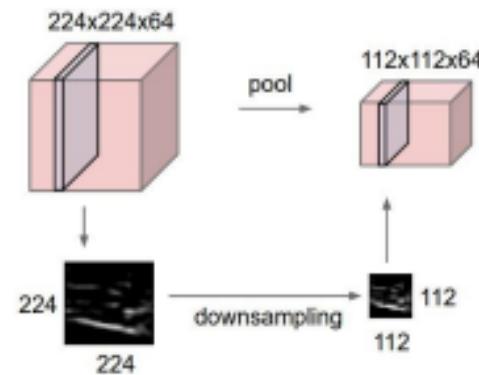
Convolutional Neural Networks

Components of CNNs

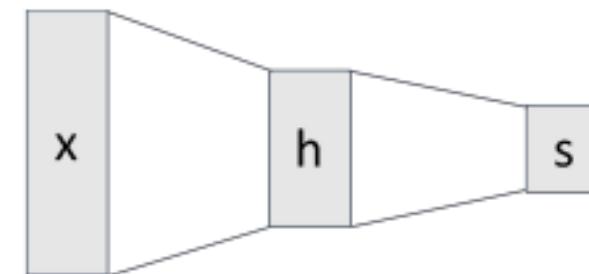
Convolution Layers



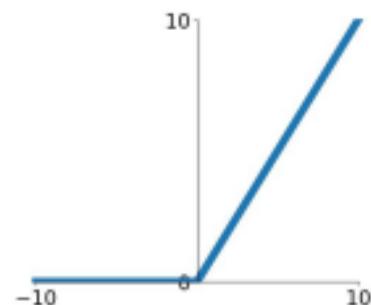
Pooling Layers



Fully-Connected Layers



Activation Function



Normalization

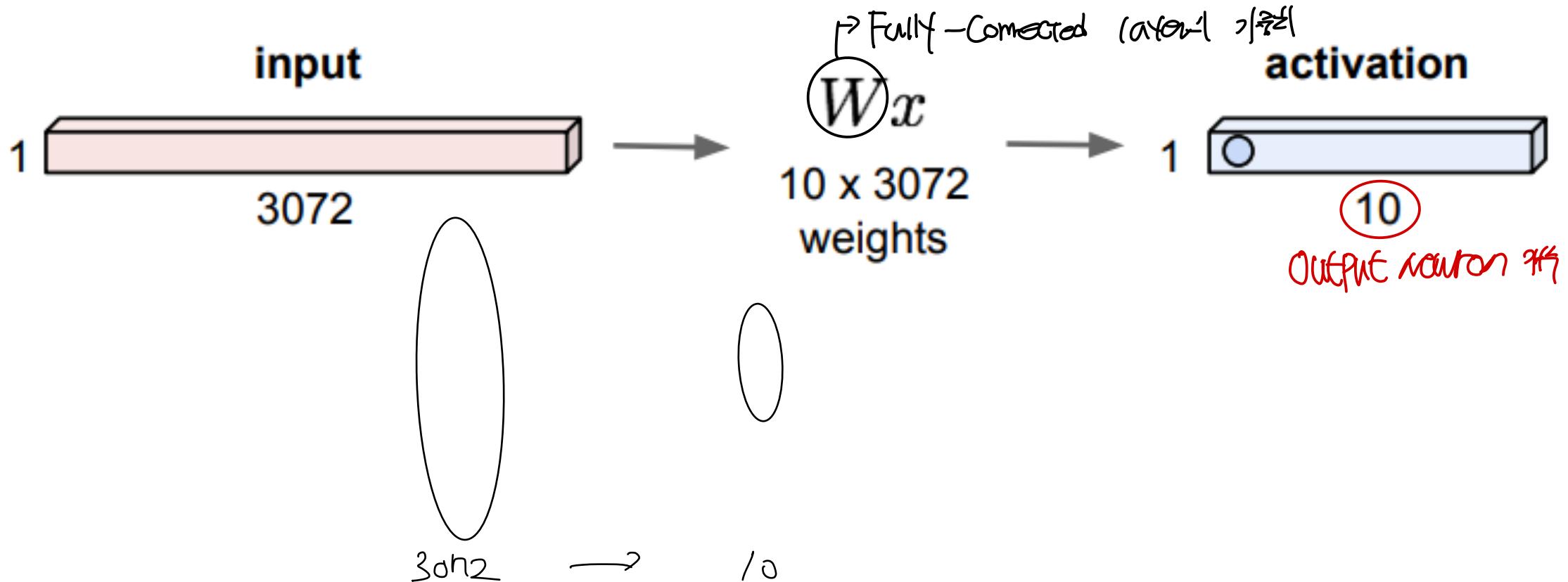
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Question: How should we put them together?

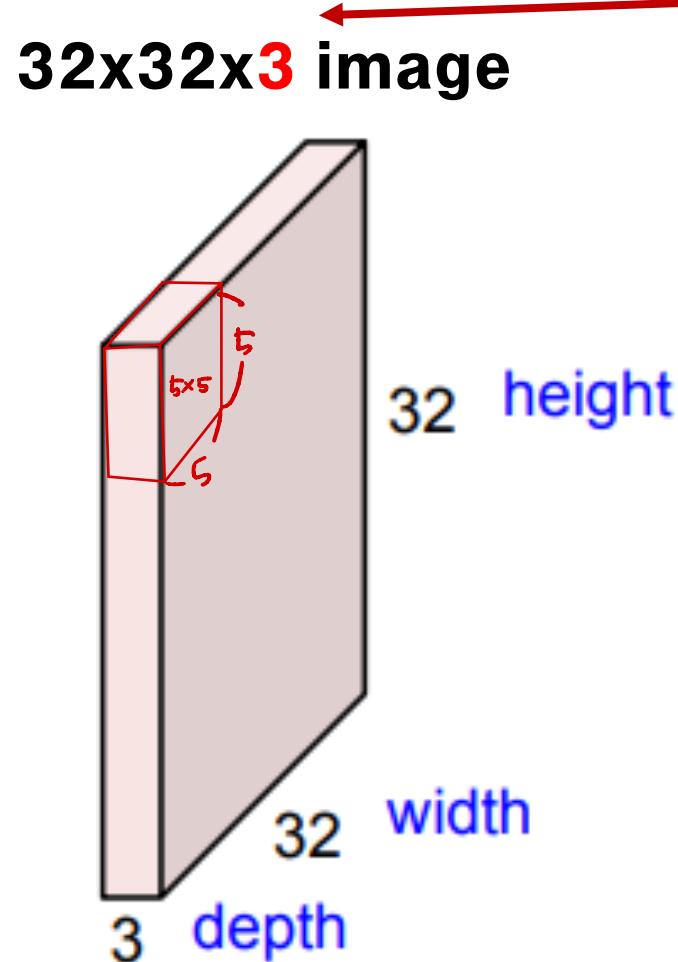
Recap: Fully Connected Layer

Fully Connected 1-dim Vector^을 Input^을 한 10개의 neuron^을

32x32x3 image → stretch to 3072x1



Convolution Layer



Filter always extend the full depth of the input volume

Filter Depth는 Image 깊이와 동일해야 함

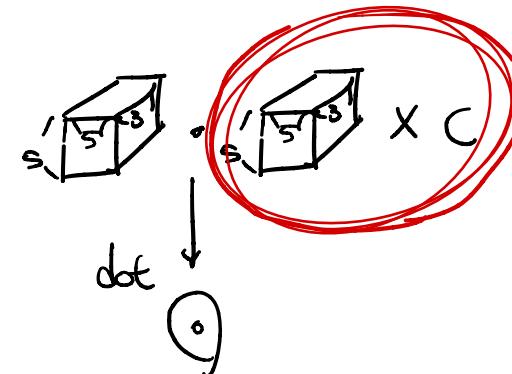
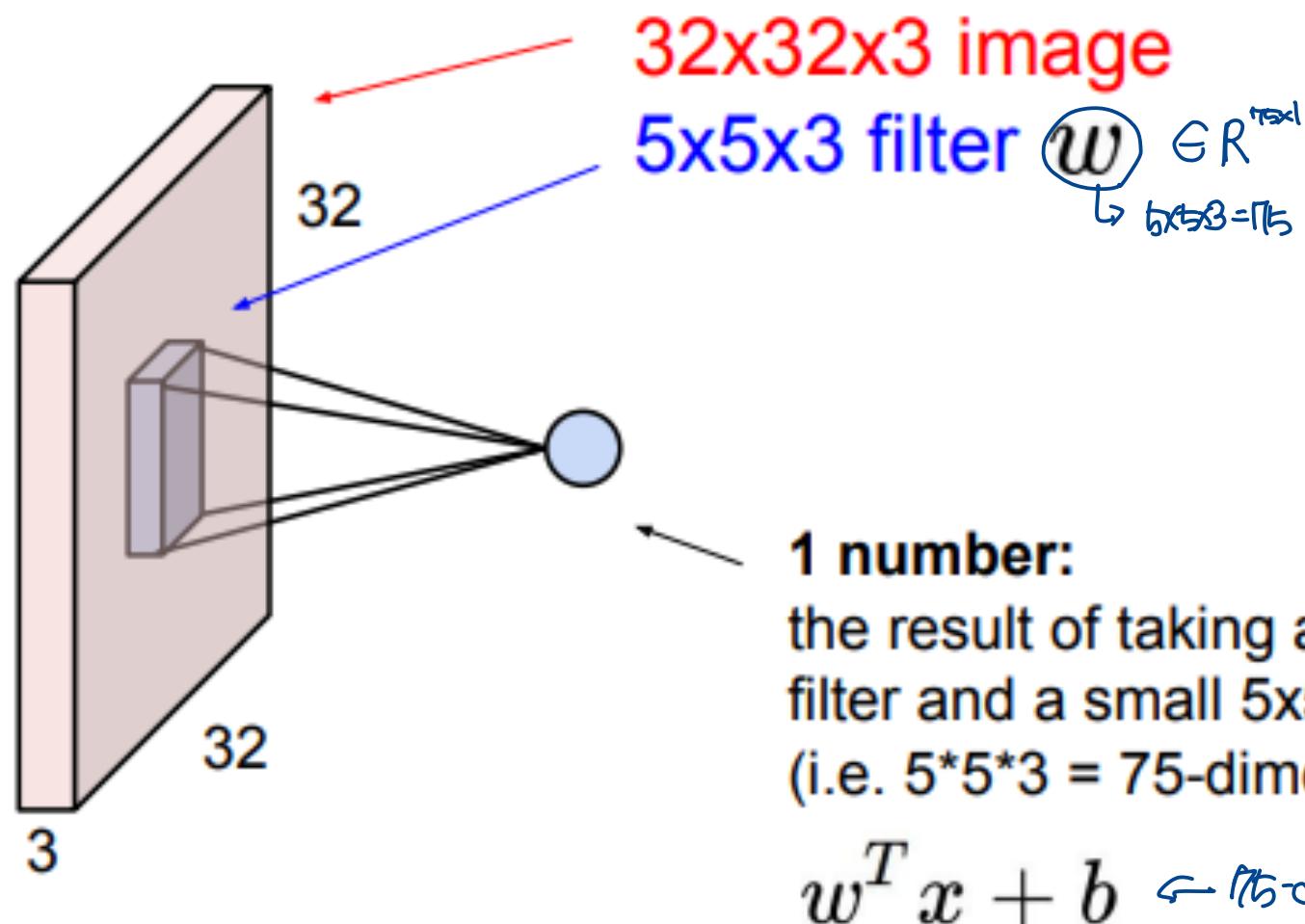
5x5x3 filter



→ Scalar Value

Convolve the filter with the image
i.e. "slide over the image spatially, computing dot products"

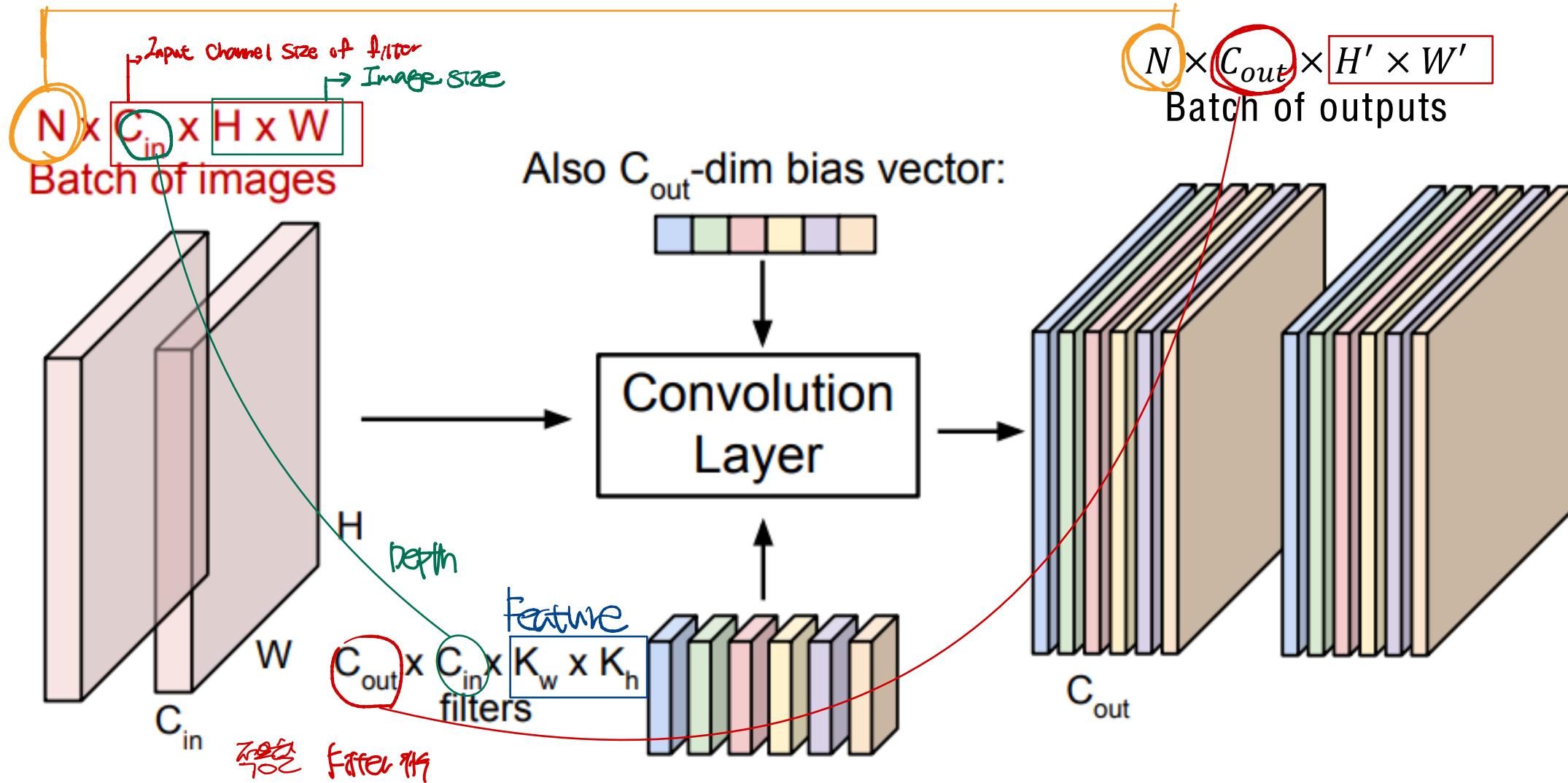
Convolution Layer



Convolution Layer

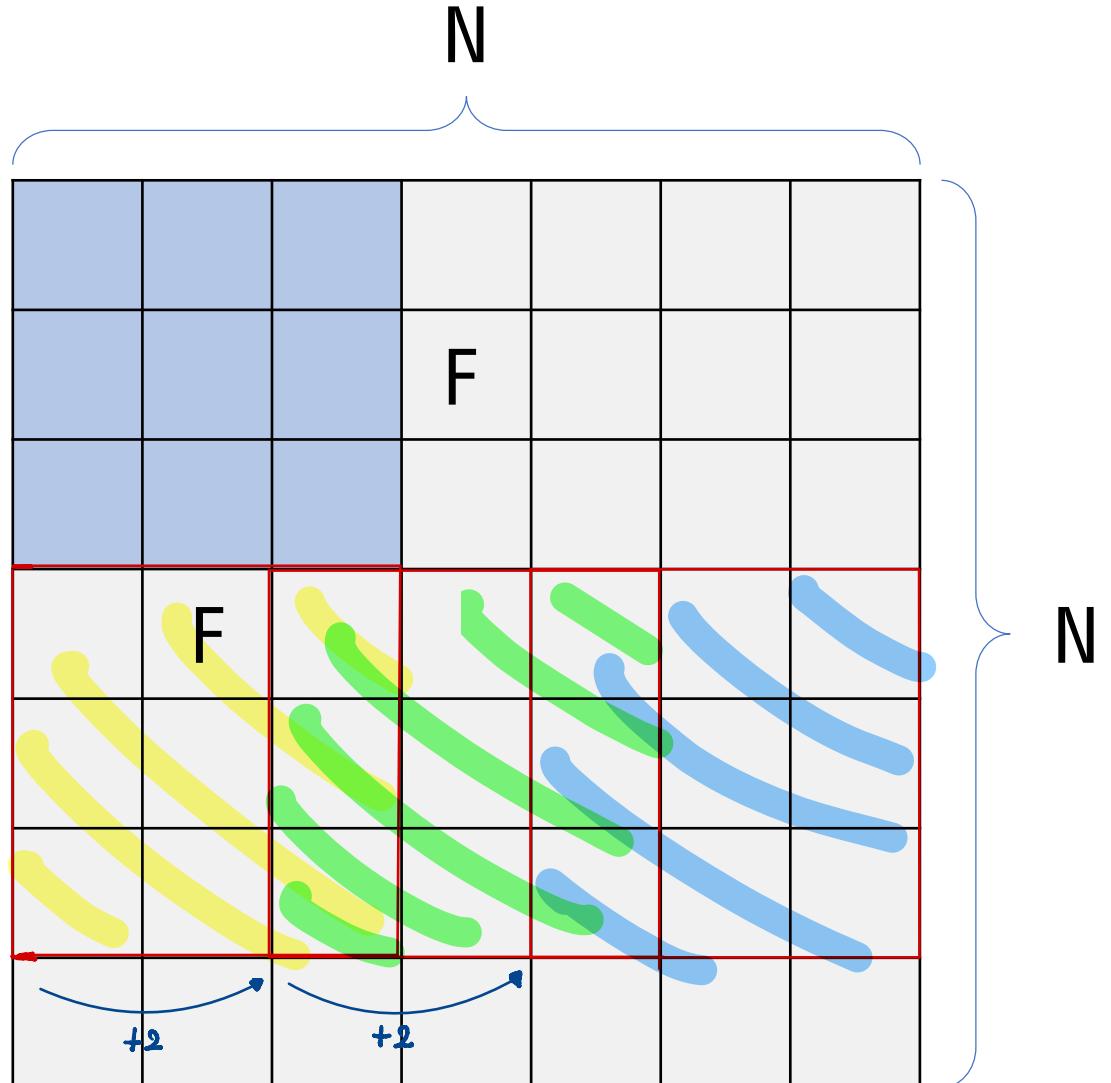
$N \times 1$ ($H \times W \times C_{in}$) 이미지를 $C_{out} \times 1$ (K_w, K_h, C_{in}) filters

Convolution 끝면 $N \times C_{out} \times H' \times W'$



A closer look at spatial dimensions

$$H \times W \times D \leftarrow \text{kernel} \Rightarrow H' \times W' \\ (x, y, D)$$



Output size

$$: (N - F) / \text{stride} + 1$$

e.g. $N=7, F=3$

stride 1 $\rightarrow (7 - 3) / 1 + 1 = 5$

stride 2 $\rightarrow (7 - 3) / 2 + 1 = 3$

stride 3 $\rightarrow (7 - 3) / 3 + 1 = 2.33$

卷积步长
Convolution stride
跳跃步长跳跃步长

→ Input size는 필터의 양이 3으로 되었다.

In practice: Common to zero pad the border

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

$$\text{Output size} = \frac{N+2P-f}{S} + 1$$

→ 여기서 가장자리인 차이 $\rightarrow P=1$ Output size

$$\text{Output Size } \frac{N+2P-f+1}{S} = \frac{7+2-3+1}{1} = 7$$

e.g. input 7x7

3x3 filter, applied with stride 1
pad with 1 pixel border

Q. what is the output?

(recall)

$(N - F) / \text{stride} + 1$

In practice: Common to zero pad the border

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

e.g. **input 7x7**
3x3 filter, applied with **stride 1**
pad with 1 pixel border

Q. what is the output?
A. 7x7 output

Receptive Fields

• Output의 특정 Pixel을 결정하는 Input

For convolution with kernel size K, each element in the output depends on a $K \times K$ receptive field in the input

Convolution의 결과는 Output의 어느 영역에 의존하는가?



Input

Output

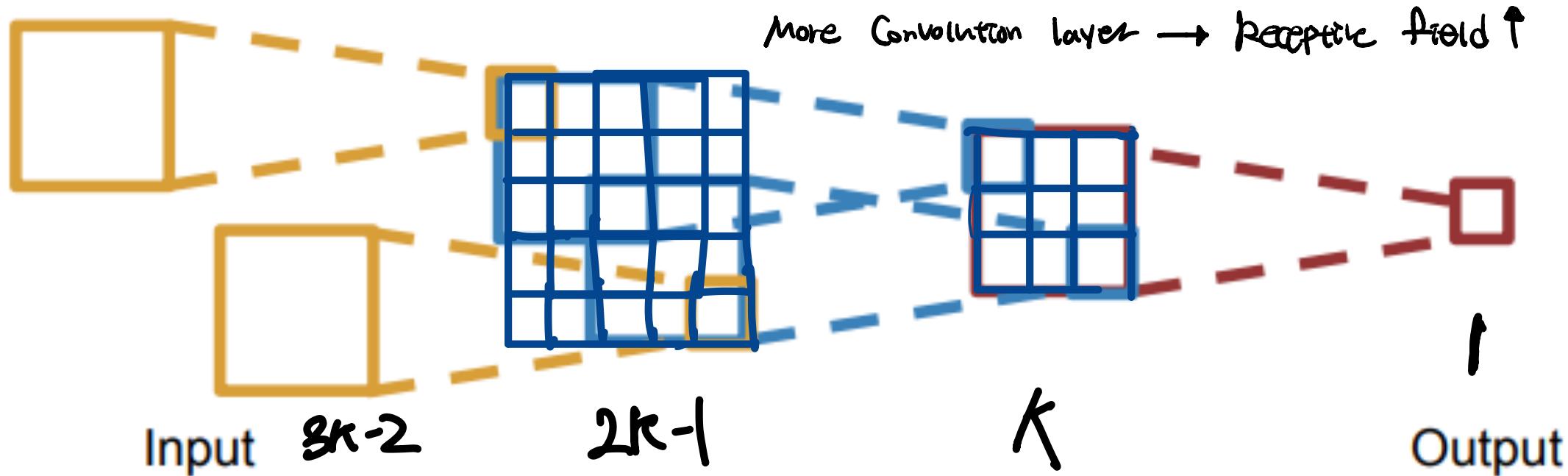
Receptive Fields

kernel size = k

$$I \rightarrow k \times k \rightarrow k \times k - 1 \rightarrow$$

Receptive field는 높은 높은 높은 등이다.

Each successive convolution adds $K - 1$ to the receptive field size
With L layers the receptive field size is $1 + L * (K - 1)$

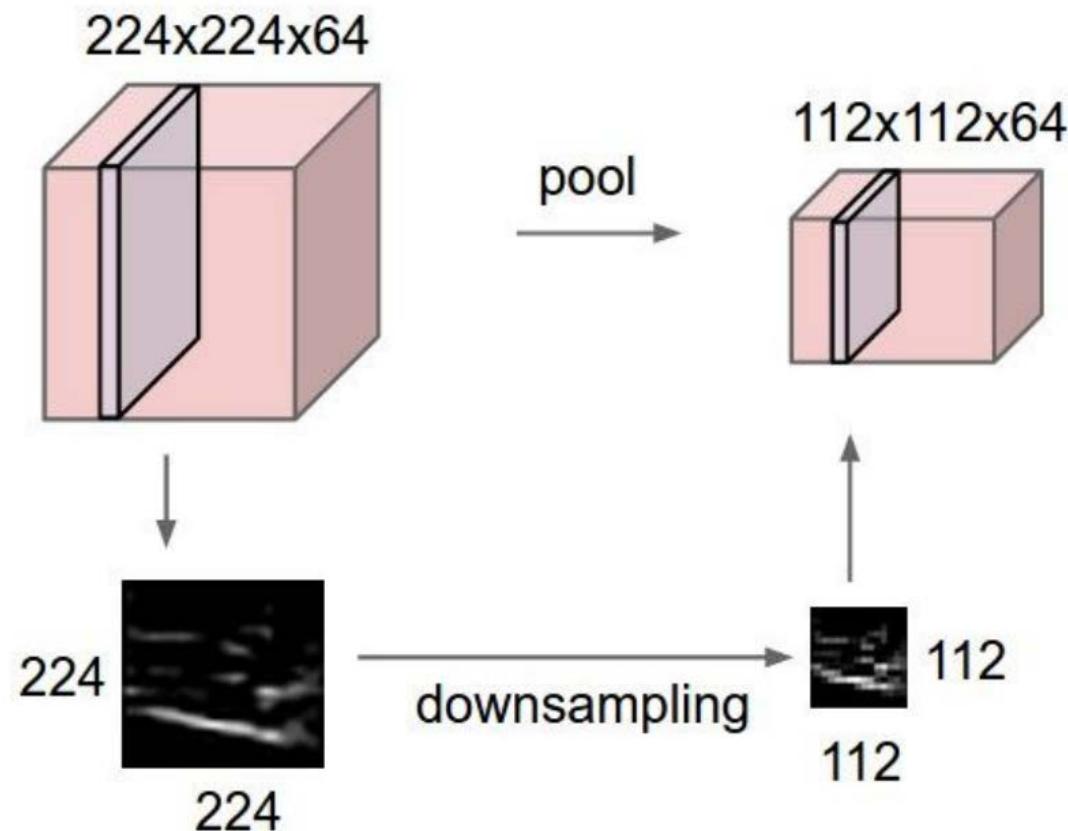


Be careful – "receptive field in the input" vs. "receptive field in the previous layer"

2방 퍼런 문제!!

Pooling Layer

Shrinks the representations **smaller**.

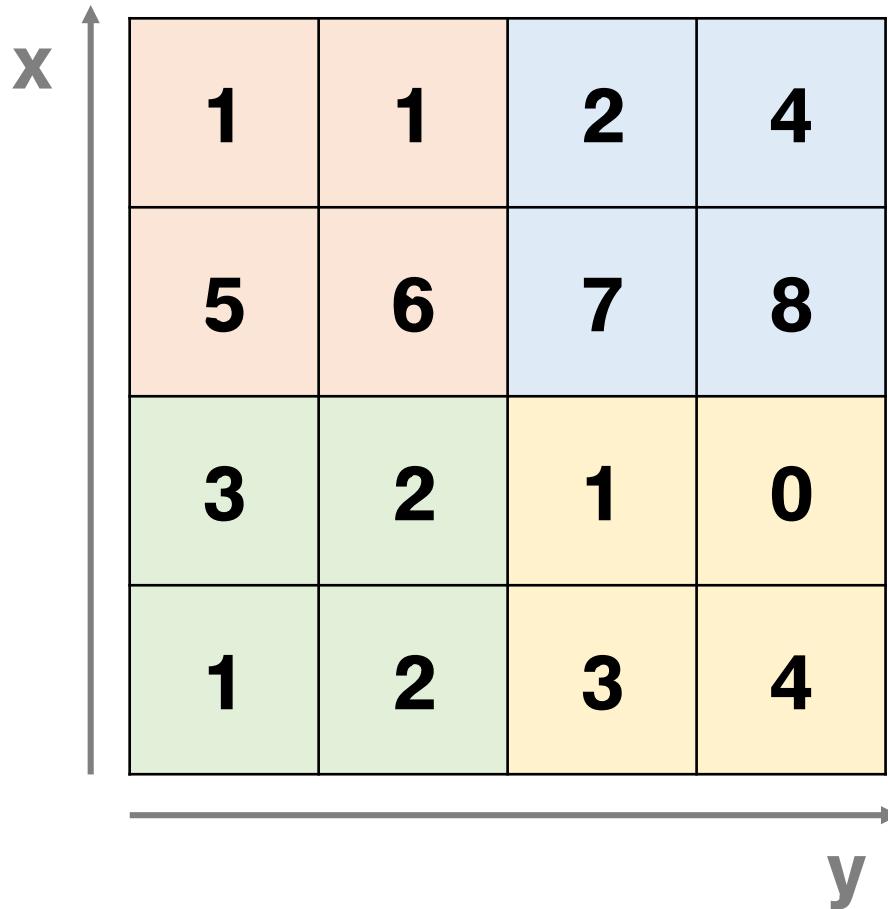


Max Pooling

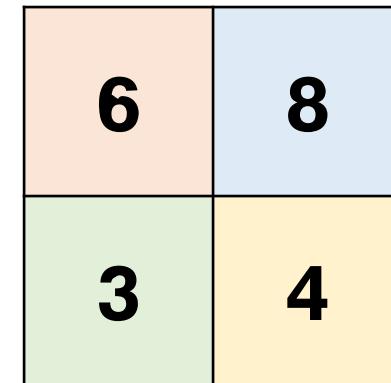
: Choose largest value
in each filter

Activation ↑ → More informative
Max ← 톤막의 가장 잘되는

Single depth slice



max pool with **2x2** filters
and stride 2



6	8
3	4

$$\frac{N-F}{S} + 1 = \frac{N}{\Sigma}$$

Stride = Filter size
→ 원래의 $\frac{1}{2}$ 배

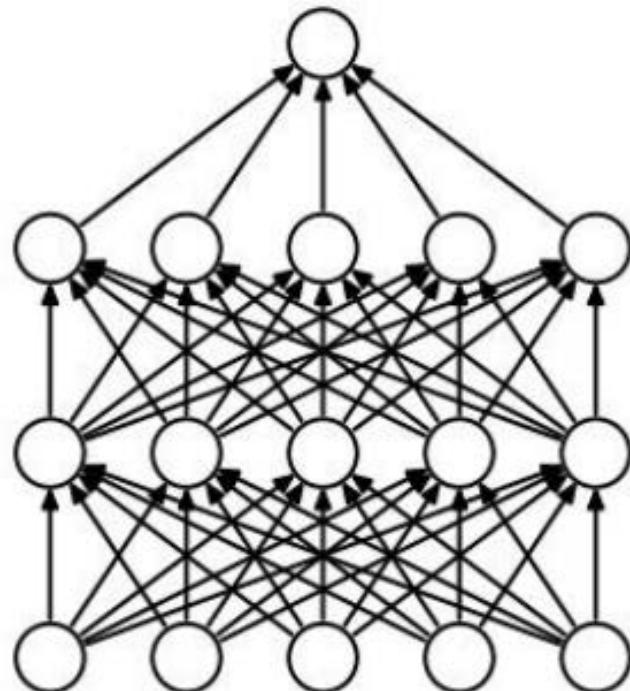


Regularization

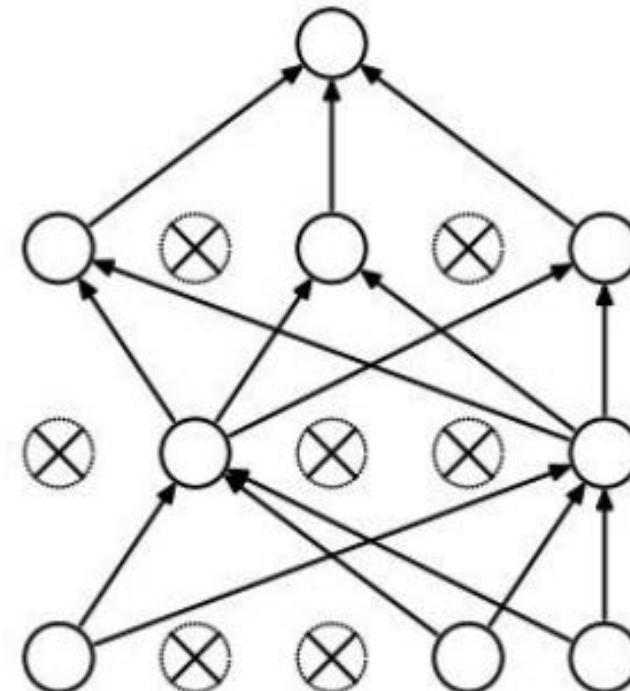
Regularization: Dropout

중요하지 않음 = 0

- Randomly set some **neurons to zero** (commonly use 0.5 probability)
- Can be seen as an **ensembling** of sub-models.



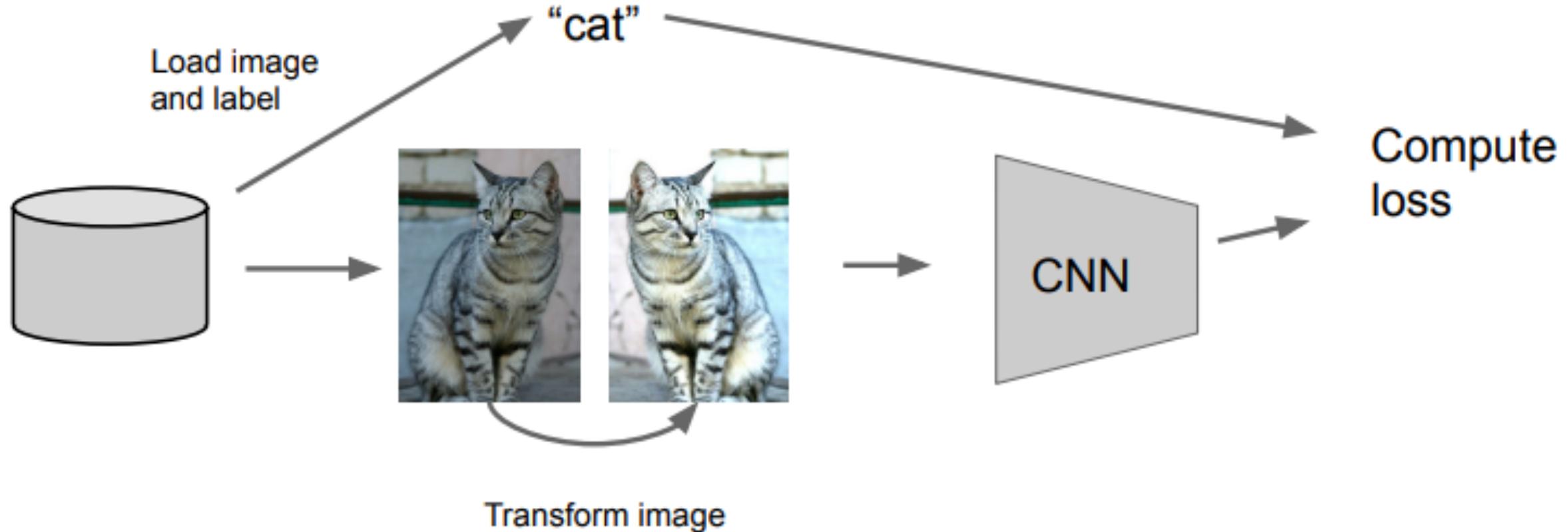
(a) Standard Neural Net



(b) After applying dropout.

Regularization: Data Augmentation

Model 성능↑ : Training data를 늘리는 것 (Increasing Samples)
→ 하나의 사진의 일부 또는 변경을 예상 사용하는 것 | 도움이 된다.



Data Augmentation



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)

(f) Rotate $\{90^\circ, 180^\circ, 270^\circ\}$ 

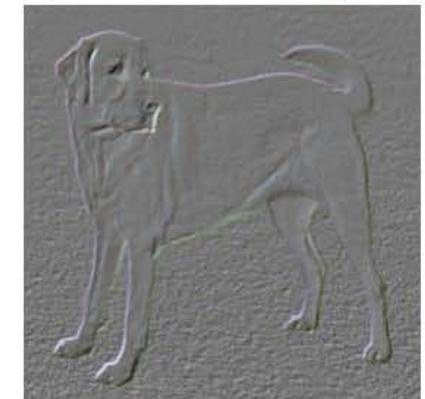
(g) Cutout



(h) Gaussian noise



(i) Gaussian blur



(j) Sobel filtering

Data Augmentation

Mixup and CutMix

같이지 사용되는 고양이
사진을 합쳐도록 한다.

ResNet-50



Mixup [48]



Cutout [3]



CutMix



Image

Label

Dog 1.0

Dog 0.5
Cat 0.5

Dog 1.0

Dog 0.6
Cat 0.4

ITE4052 Computer Vision

Camera & Image Formation

Dong-Jin Kim
Spring 2025