

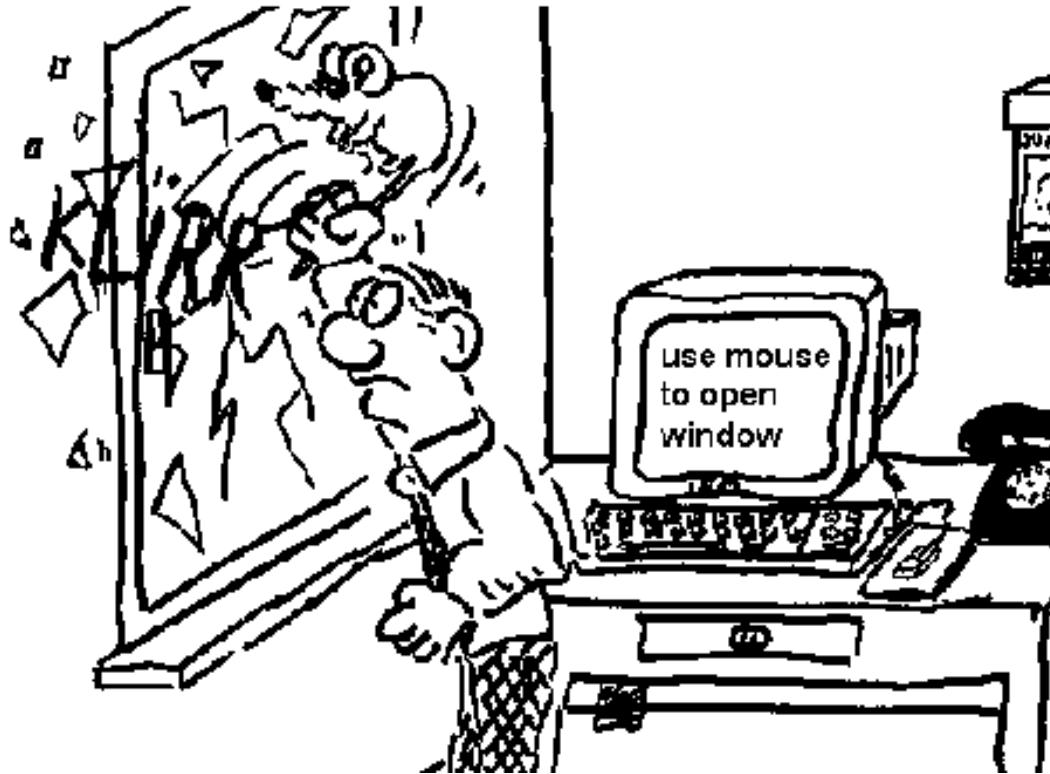
Deep Learning Basics

Lecture 3: Optimization

- Concept → p5~
- Practical Gradient method → p18~
- Regularization → p28~

최성준 (고려대학교 인공지능학과)

Introduction

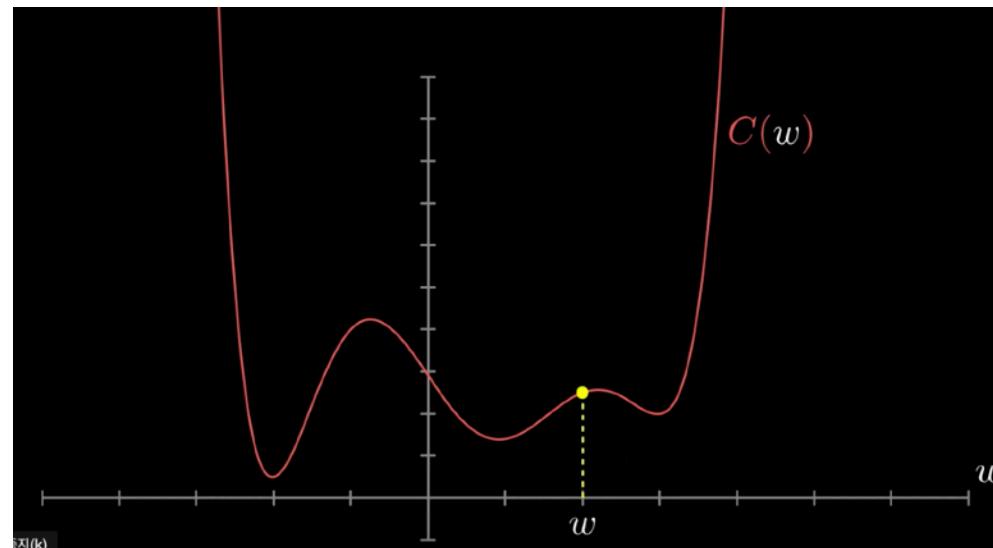
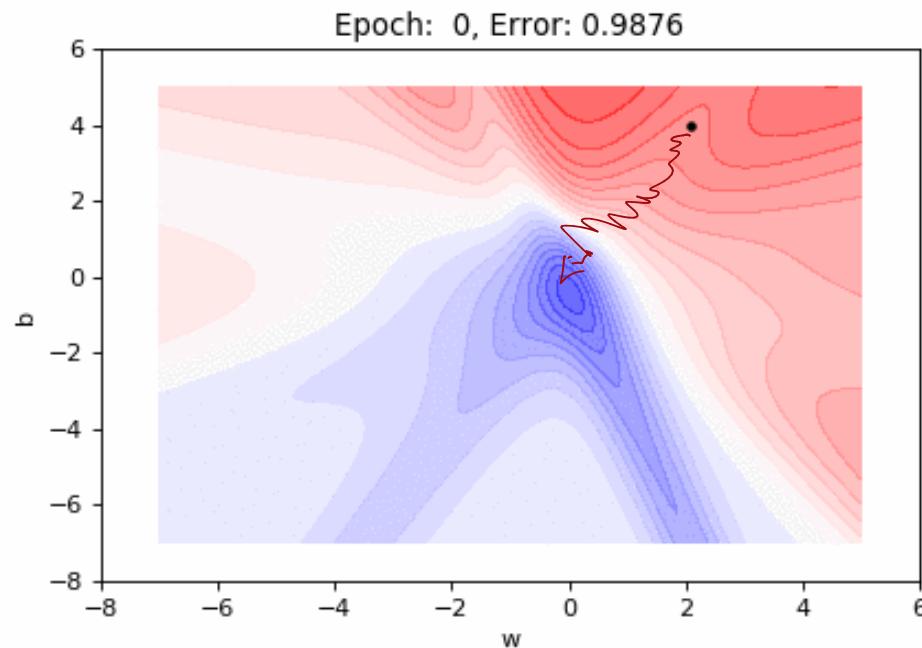


“Language is the source of misunderstandings”
Antoine de Saint-Exupéry (1900-1944)

Introduction

Gradient Descent

- First-order iterative optimization algorithm for finding a local minimum of a differentiable function.



Important Concepts in Optimization

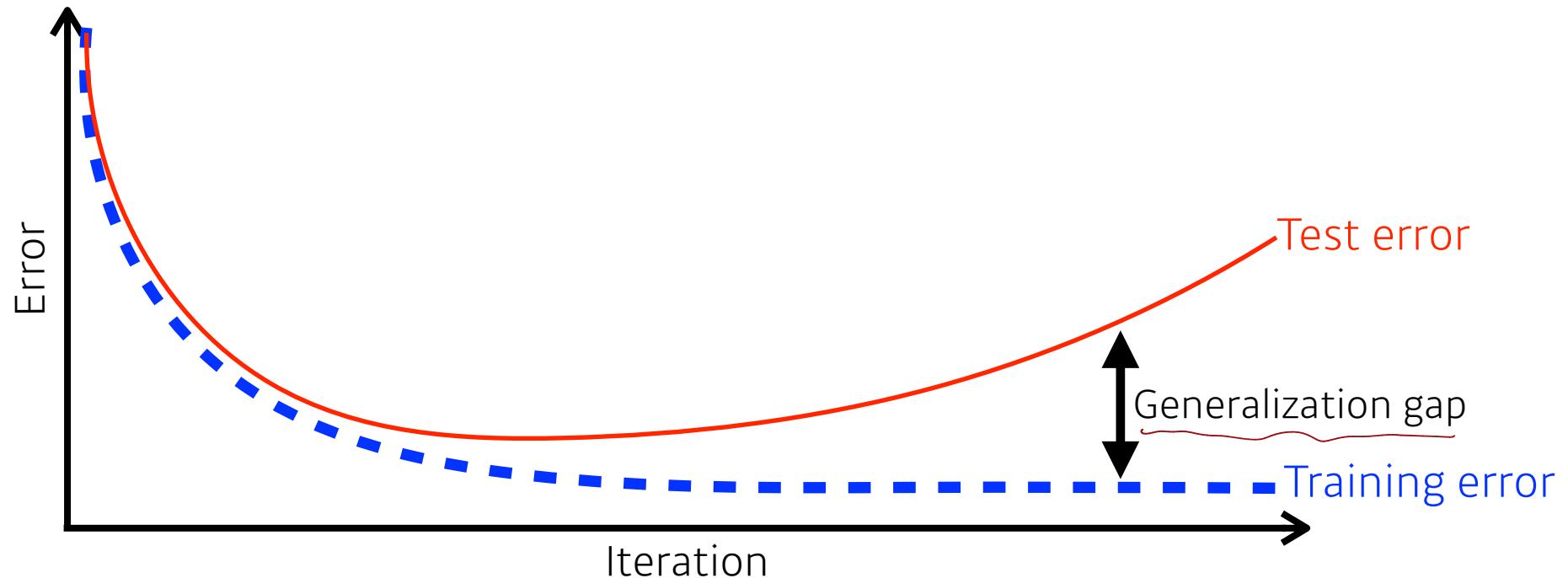
Important Concepts in Optimization

- Generalization (Generalization Gap)
- Under-fitting vs. over-fitting
- Cross validation 교차검증 (k-폴드 CV)
- Bias-variance tradeoff
- Bootstrapping
- Bagging and boosting

Generalization

일반화 성능

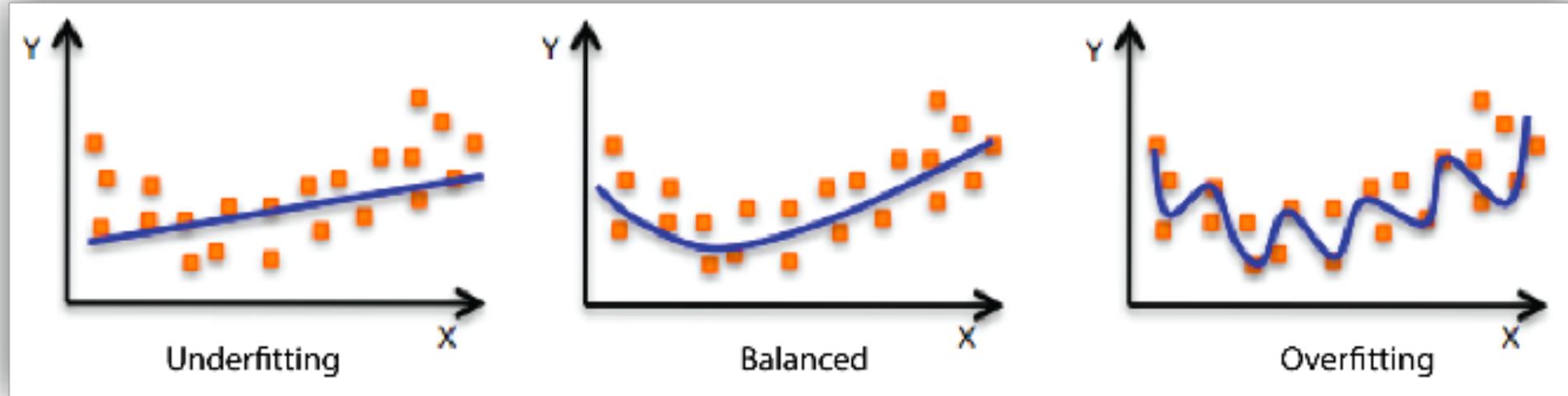
- How well the learned model will behave on unseen data.



Underfitting vs. Overfitting

과소적합

과대적합



<https://docs.aws.amazon.com/machine-learning/latest/dg/model-fit-underfitting-vs-overfitting.html>

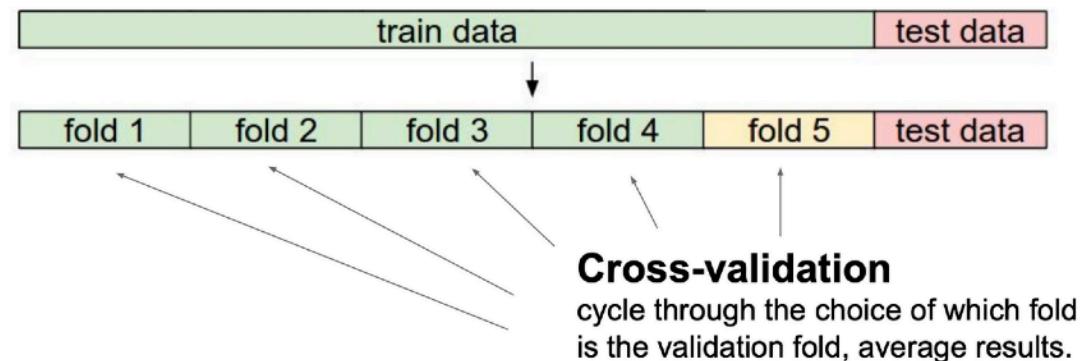
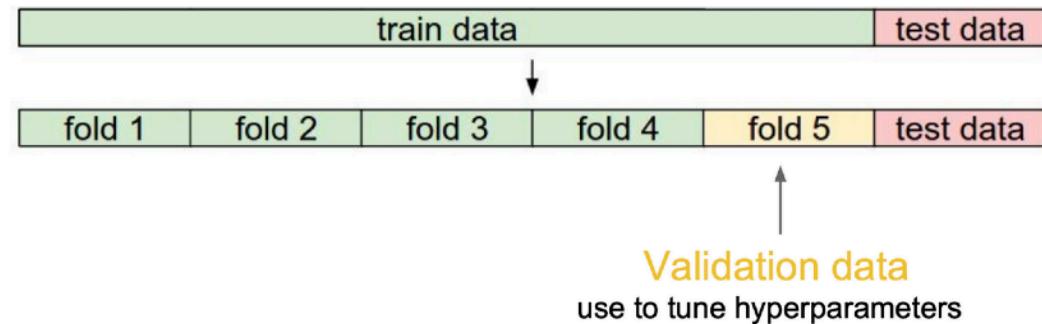
© NAVER Connect Foundation

Cross-validation

교차 검증 (K-fold CV)

→ 허가되는 hyperparameters set을 찾기 위한

- Cross-validation is a model validation technique for assessing how the model will generalize to an independent (test) data set.



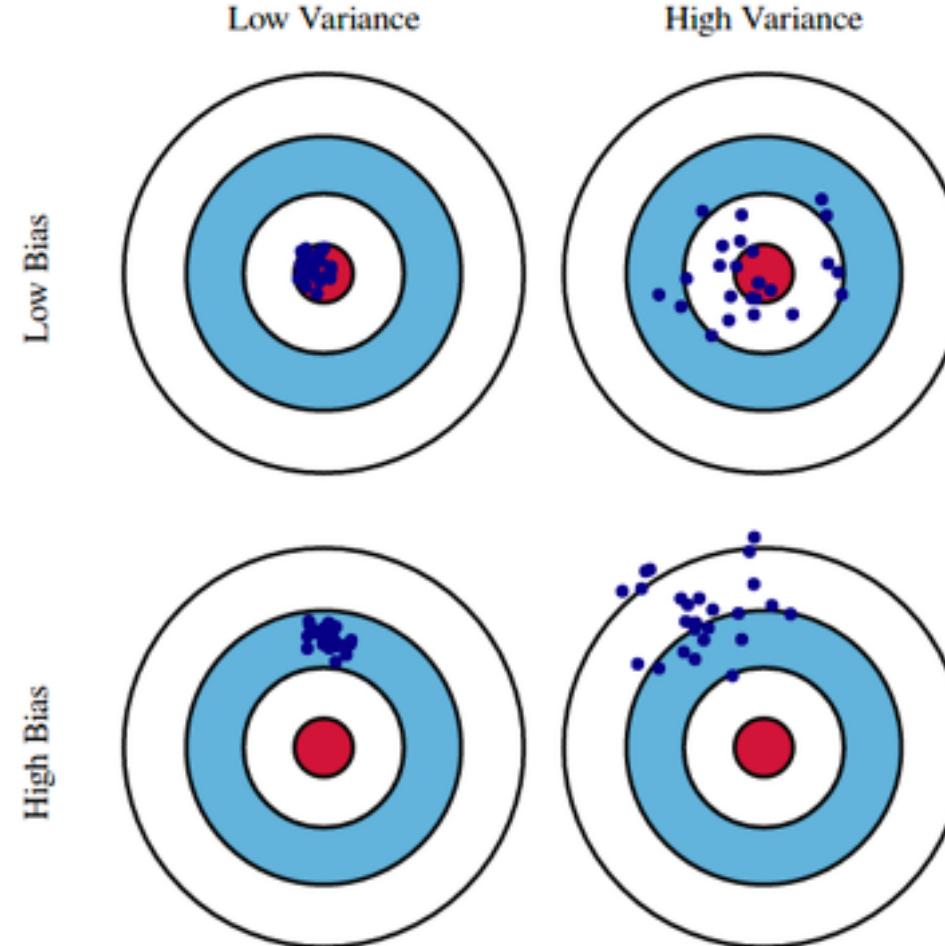
<https://blog.quantinsti.com/cross-validation-machine-learning-trading-models/>

▷ 허락의 험과 더불어 높은 정밀도가 있는지.

(Variance가 크면 Over-fitting의 위험성이 있다)

Bias and Variance

▷ 예측값으로 예상치와 차이가 얼마나 있는지.



Bias and Variance Tradeoff

Given $\mathcal{D} = \{(x_i, t_i)\}_{i=1}^N$, where $t = f(x) + \epsilon$ and $\epsilon \sim \mathcal{N}(0, \sigma^2)$

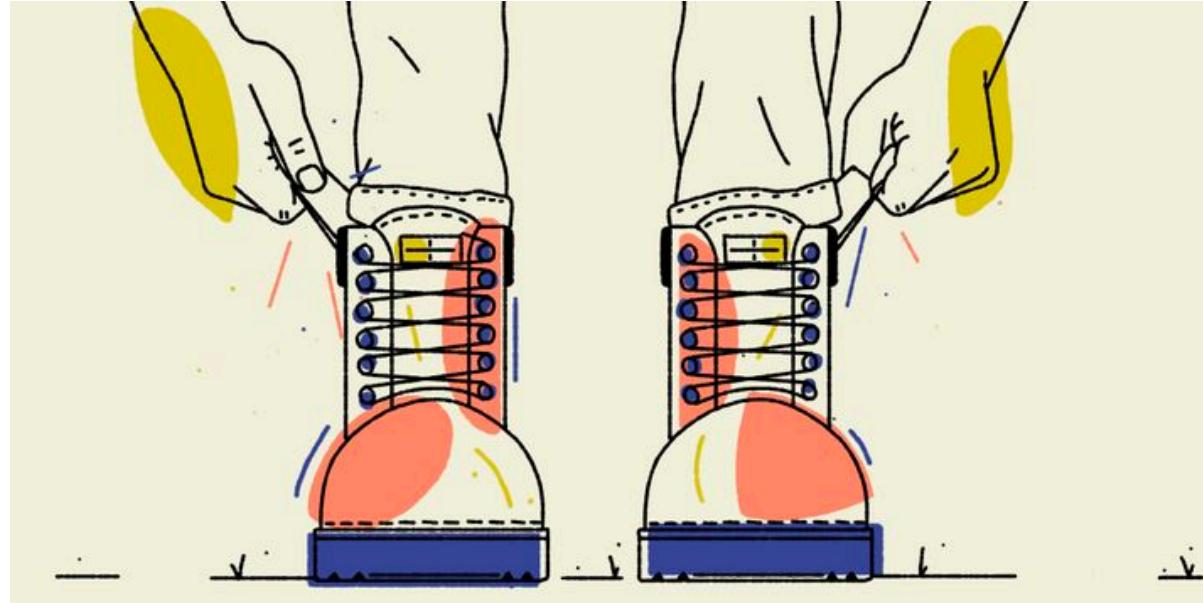
We can derive that what we are minimizing (**cost**) can be decomposed into three different parts: bias², variance, and noise. 

$$\begin{aligned}
 & \text{target} \quad \text{NN} \xrightarrow{\text{fit}} \hat{f} \\
 \mathbb{E} \left[(t - \hat{f})^2 \right] &= \mathbb{E} \left[(t - f + f - \hat{f})^2 \right] \\
 \text{cost} &= \dots \\
 &\qquad\qquad\qquad \underbrace{\text{trade off } (\text{bias} \xrightarrow{\text{small}} \text{variance} \xrightarrow{\text{large}} \text{noise})}_{\text{bias}^2 \quad \text{variance} \quad \text{noise}}
 \end{aligned}$$

Bootstrapping

⇒ 무작위 샘플링을 사용하는 여러 표본을 만들어서 평균화하는 것.
(학습 데이터가 고정되어 있을 때, 그 학습에서 샘플링을 통해)

학습데이터를 여러개 만들고, 이를 통해 모수, 예측: (S를 만드는 것)



Pull yourself up by the **bootstraps**.

- Bootstrapping is any test or metric that uses random sampling with replacement.

Bagging vs. Boosting

Bagging (Bootstrapping aggregating) = bootstrapping (~ampling)

- Multiple models are being trained with bootstrapping.
- ex) Base classifiers are fitted on random subset where individual predictions are aggregated (voting or averaging).

⇒ 허깅 샘플을 여러 개 만들고, 여러 모델을 만들어 out+put+의
합을 출력하는 것.
즉각적인 모델을

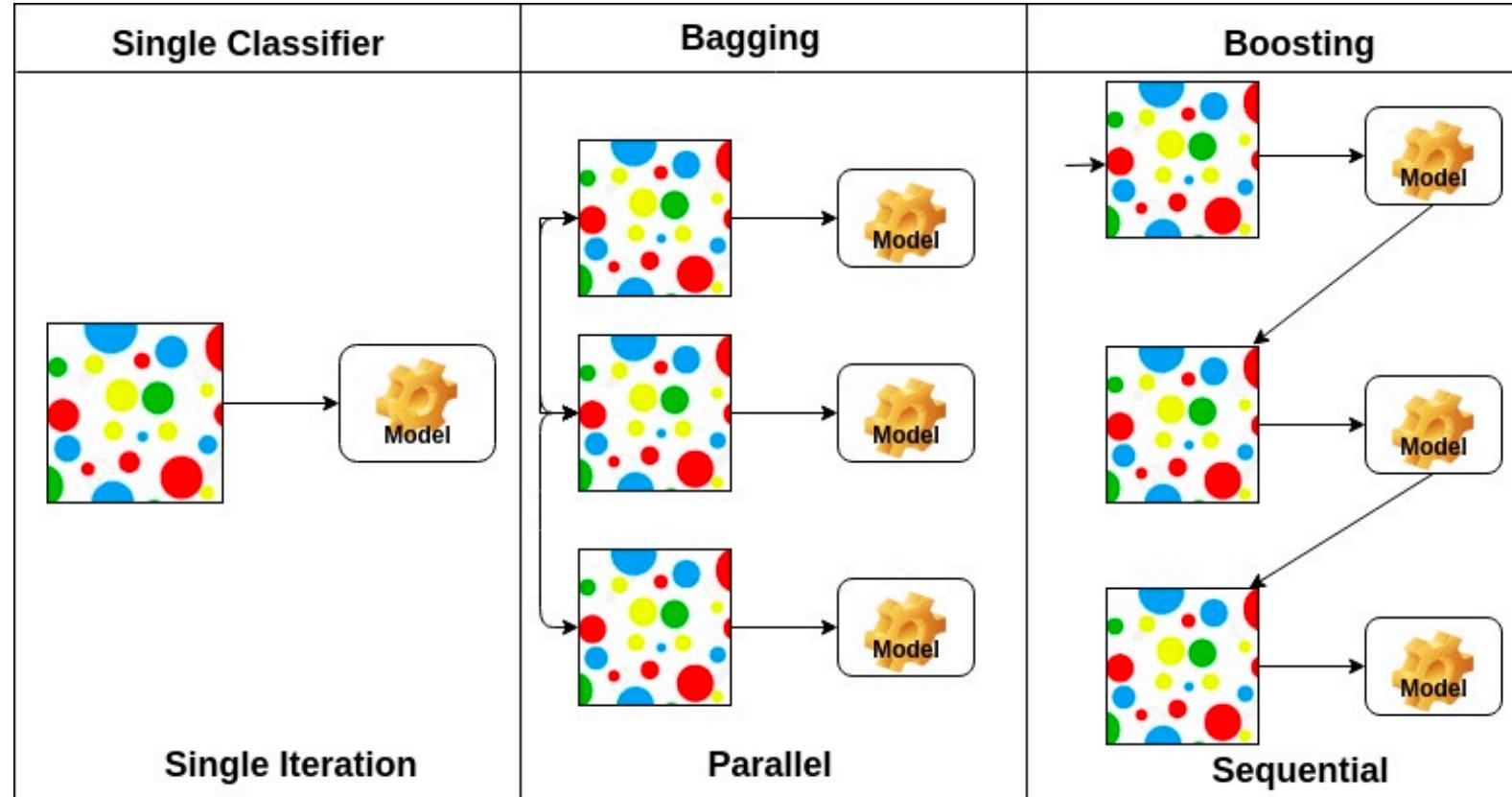
Boosting

- It focuses on those specific training samples that are hard to classify.
- A strong model is built by combining weak learners in sequence where each learner learns from the mistakes of the previous weak learner.

⇒ 여러 개의 모델을 만들어서 합친다. (여러 개의 weak learner을 합쳐서 strong model을 만드는 것)

↳ 같은 weak learner의 수가 많아
내수는 많다.
(sequential 학습이나)

Bagging vs. Boosting



<https://www.datacamp.com/community/tutorials/adaboost-classifier-python>

Practical Gradient Descent Methods

Gradient Descent Methods

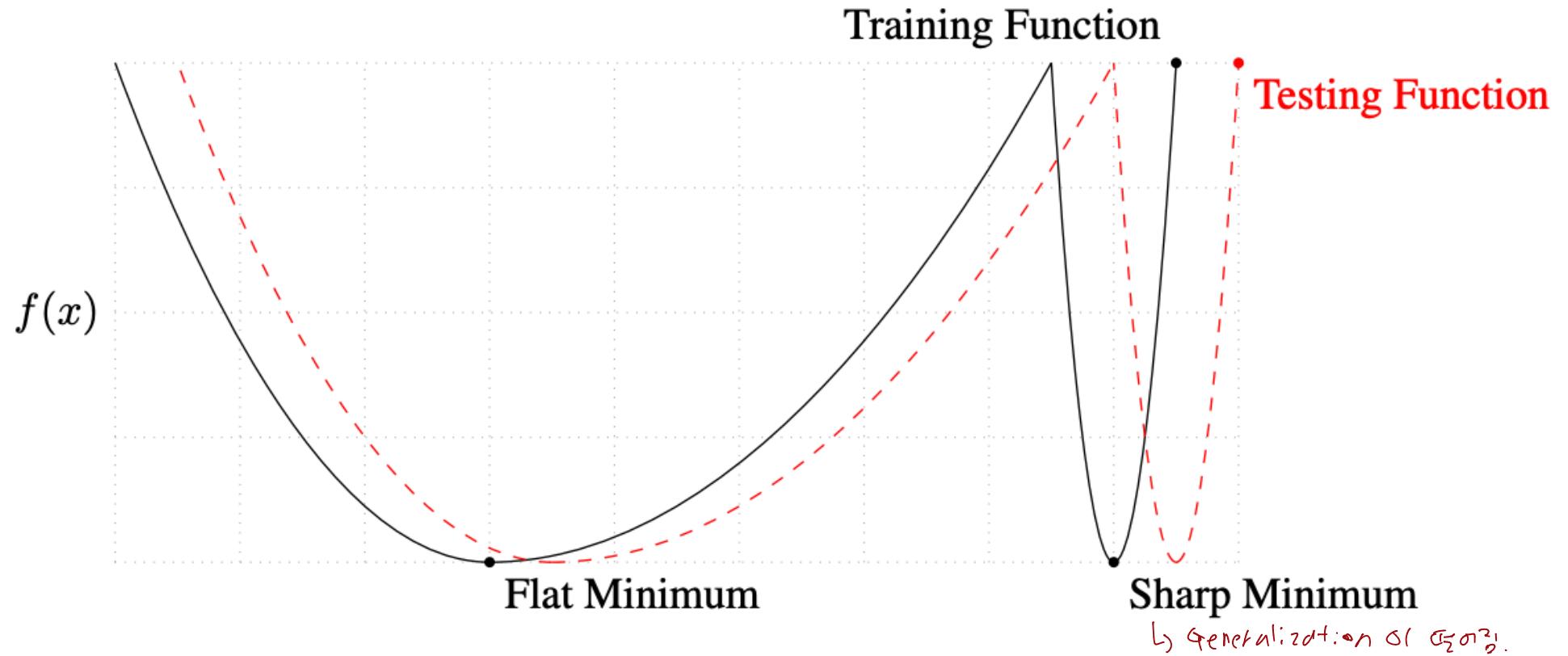
- ➊ Stochastic gradient descent
 - ➌ Update with the gradient computed from a single sample.
- ✓ ➋ Mini-batch gradient descent
 - ➌ Update with the gradient computed from a subset of data.
- ➌ Batch gradient descent
 - ➌ Update with the gradient computed from the whole data.

Batch-size Matters

- "It has been observed in practice that when using a larger batch there is a degradation in the quality of the model, as measured by its ability to generalize."
- "We ... present numerical evidence that supports the view that large batch methods tend to converge to sharp minimizers of the training and testing functions. In contrast, small-batch methods consistently converge to flat minimizers... this is due to the inherent noise in the gradient estimation."

\Rightarrow flat min: 미끄러지기
flat: 미끄러운, 미끄는.

Batch-size Matters

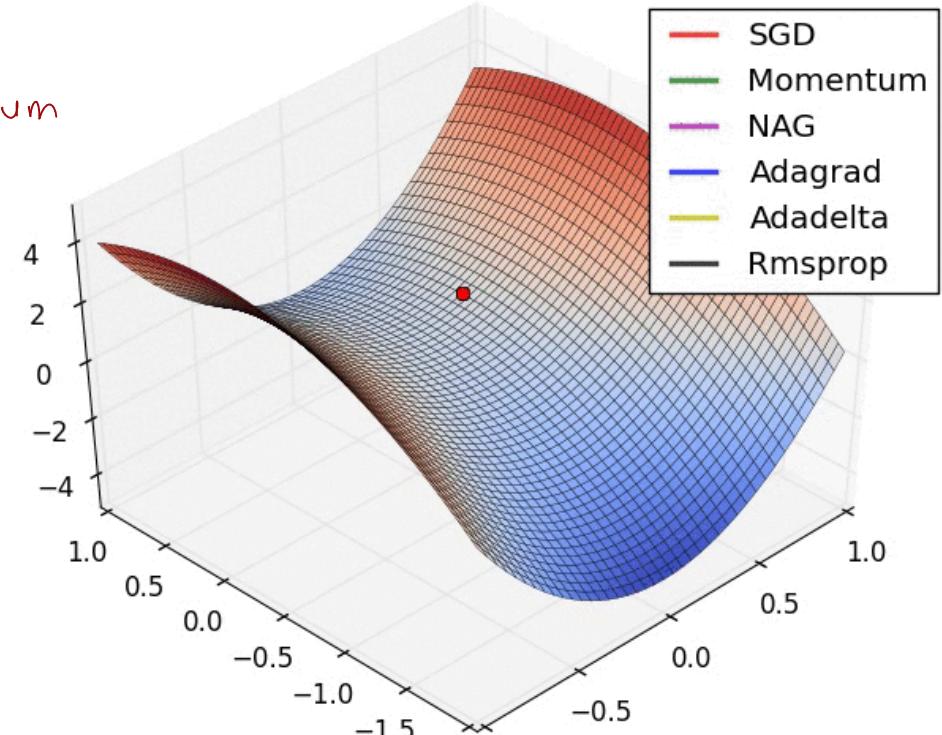


Gradient Descent Methods

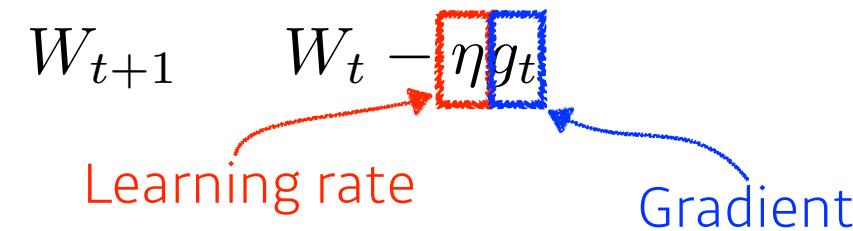
- Stochastic gradient descent
- Momentum
- Nesterov accelerated gradient
- Adagrad
- Adadelta
- RMSprop
- Adam \rightarrow Adaptive + momentum

momentum

Adaptive

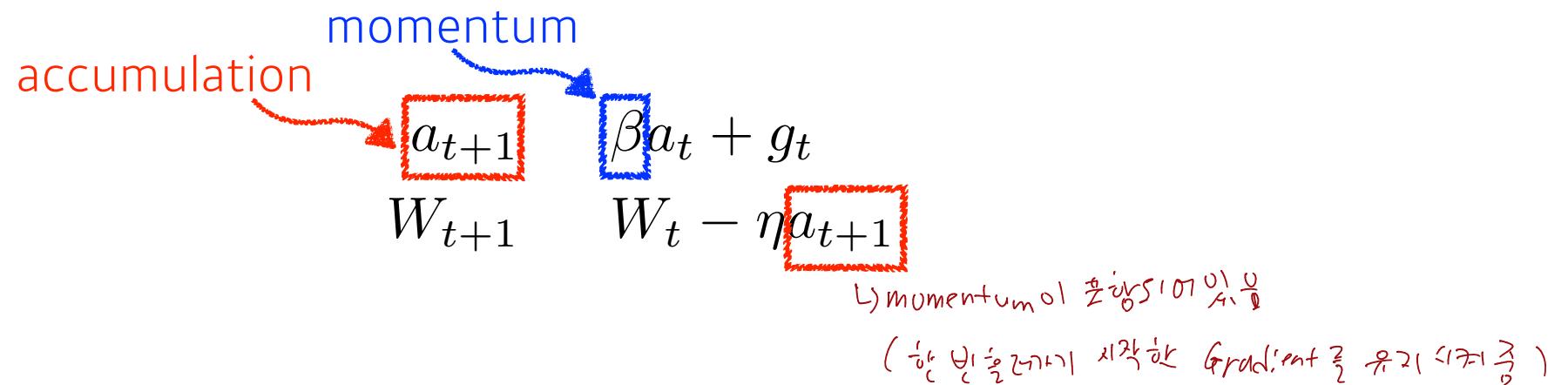


Stochastic - Gradient Descent



Momentum

(모멘텀)

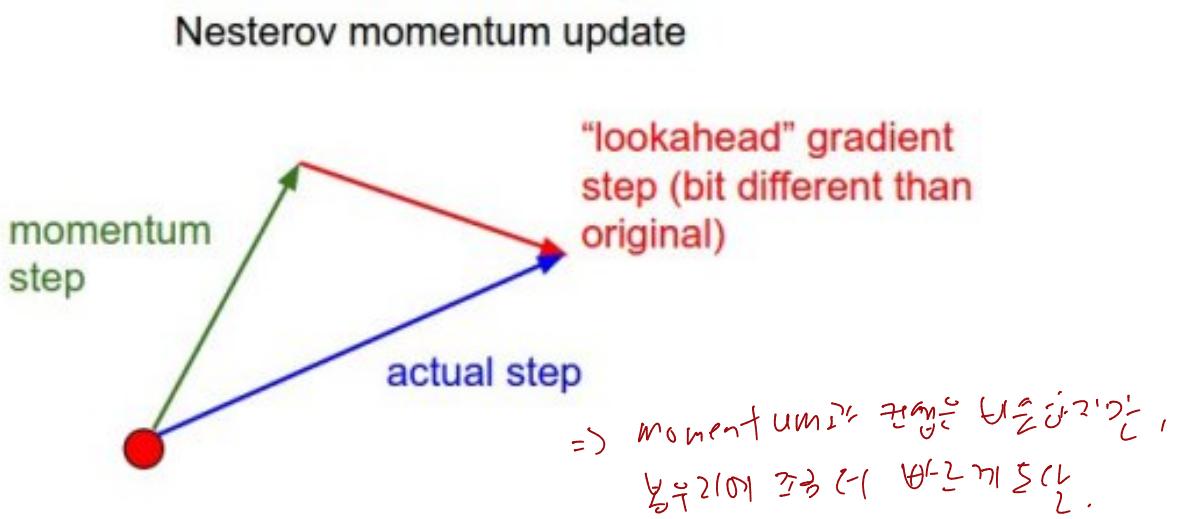
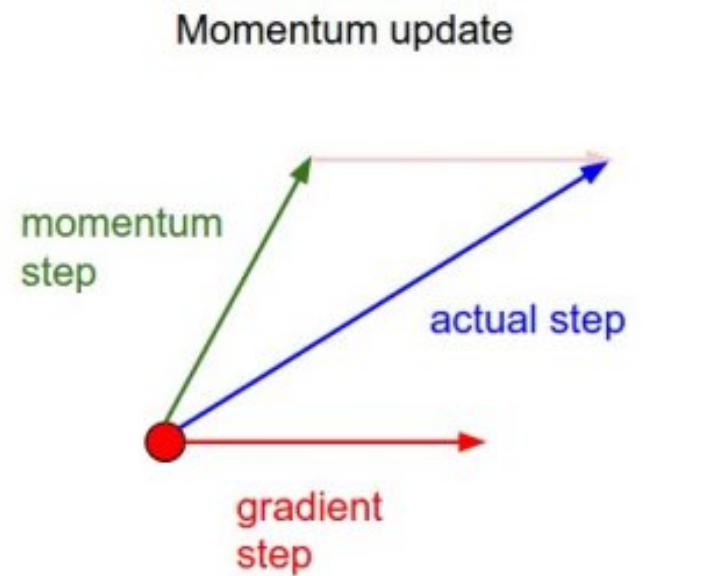


Nesterov Accelerated Gradient

$$\begin{aligned} a_{t+1} &= \beta a_t + \boxed{\nabla \mathcal{L}(W_t - \eta \beta a_t)} \\ W_{t+1} &= W_t - \eta a_{t+1} \end{aligned}$$

Lookahead gradient

Nesterov Accelerated Gradient



Adagrad

- Adagrad adapts the learning rate, performing larger updates for infrequent and smaller updates for frequent parameters.

- 누적 미트워크의 parameter가 그립가지 뜯만큼 많이 변하는지를 볼.
· 초기 변화한 파라미터에 대해서는 크기 변화시키고, 많이 변화한 파라미터는 초기 변화시킬.

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{G_t} + \epsilon} g_t$$

for numerical stability
(0으로 안 나눠지게 만드는 값)

Sum of gradient squares

What will happen if the training occurs for a long period?

• 가 학습률이 적으면 학습이 중장 멈추게 되는 경향 발생
(분수가 학습률, W 가 무의 안정화)

Adadelta

- Adadelta extends Adagrad to reduce its monotonically decreasing the learning rate by restricting the accumulation window.

EMA of gradient squares

$$G_t = \gamma G_{t-1} + (1 - \gamma) g_t^2$$

EMA of difference squares

$$W_{t+1} = W_t - \frac{\sqrt{H_{t-1} + \epsilon}}{\sqrt{G_t + \epsilon}} g_t$$

$$H_t = \gamma H_{t-1} + (1 - \gamma)(\Delta W_t)^2$$

There is no learning rate in Adadelta.

↳ ∵ 학습률은 고정된 값이 아닙니다.

↳ 학습률은

RMSprop

- RMSprop is an unpublished, adaptive learning rate method proposed by Geoff Hinton in his lecture.

EMA of gradient squares

$$G_t = \gamma G_{t-1} + (1 - \gamma) g_t^2$$
$$W_{t+1} = W_t - \frac{\eta}{\sqrt{G_t + \epsilon}} g_t$$

stepsize η \leftarrow ?

Adam

(가장 빠른 학습률 사용)

- Adaptive Moment Estimation (Adam) leverages both past gradients and squared gradients.

$$\begin{aligned} \text{Momentum} & \rightarrow m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ \text{EMA of gradient squares} & \rightarrow v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \\ W_{t+1} &= W_t - \frac{\eta}{\sqrt{v_t + \epsilon}} \frac{\sqrt{1 - \beta_2^t}}{1 - \beta_1^t} m_t \end{aligned}$$

Adam effectively combines momentum with adaptive learning rate approach.

Regularization

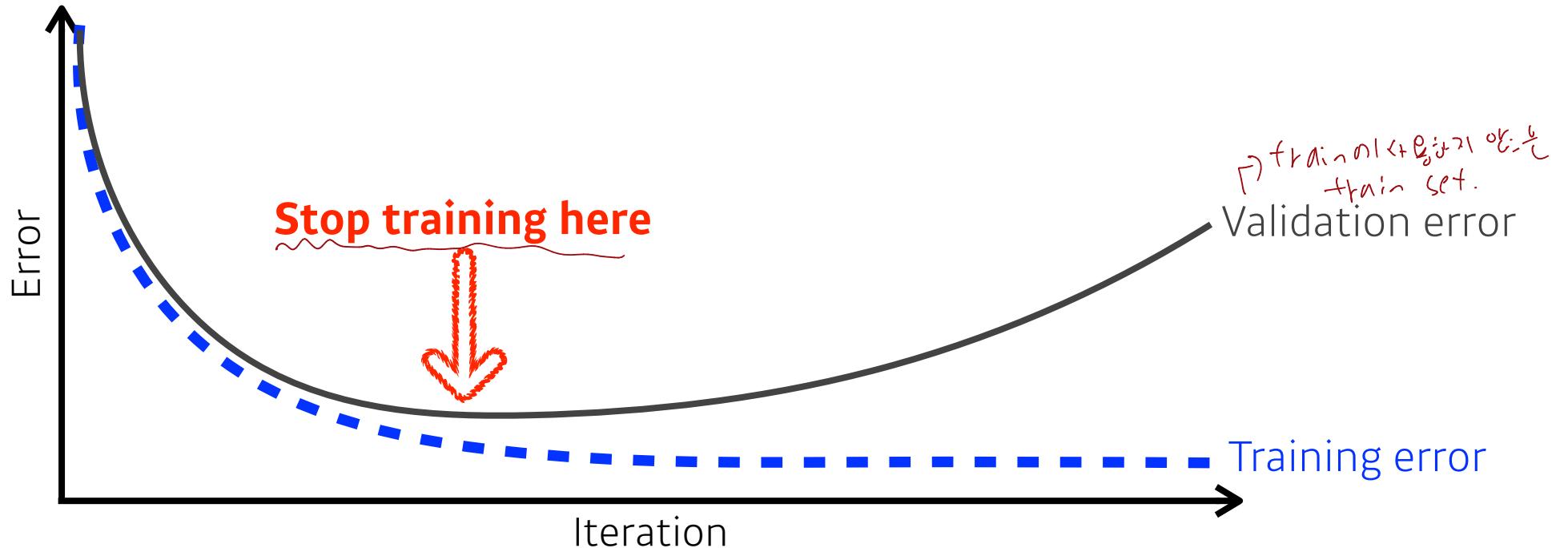
Regularization

: 규제(적합을 방해하지 Generalization 할 수 있도록 하는 것.)

- Early stopping
- Parameter norm penalty
- Data augmentation
- Noise robustness
- Label smoothing
- Dropout
- Batch normalization

설명

Early Stopping



Note that we need **additional validation data** to do early stopping.

k-fold, ...

Parameter Norm Penalty

- It adds smoothness to the function space.

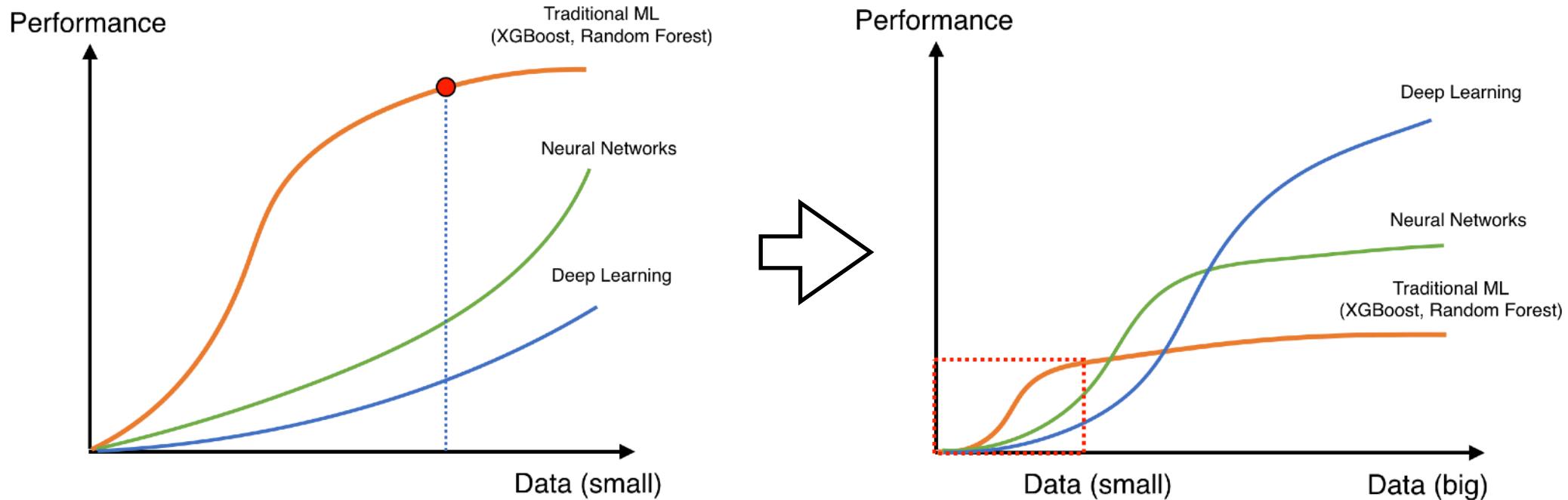
- 파라미터가 넓은 범위를 탐색하는 것.
- 파라미터를 제한하여 계산량을 줄여 cost function이 계산하는데 있어
(부드러운 학습으로 만족하기 위해)

Parameter Norm Penalty

$$\text{total cost} = \text{loss}(\mathcal{D}; W) + \frac{\alpha}{2} \|W\|_2^2$$

Data Augmentation

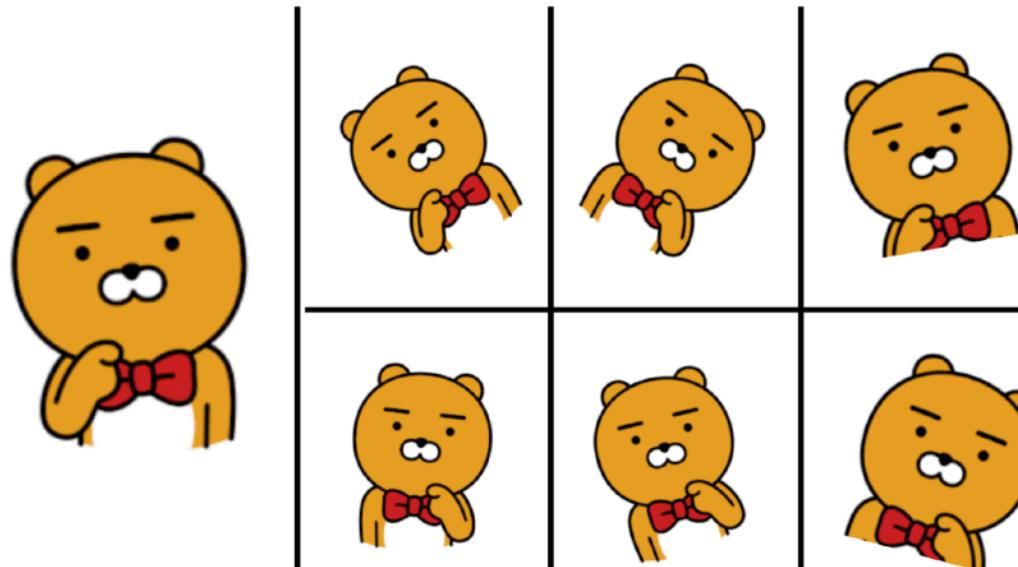
- More data are always welcomed.



Data Augmentation

- More data are always welcomed.
- However, in most cases, training data are given in advance.
- In such cases, we need data augmentation.

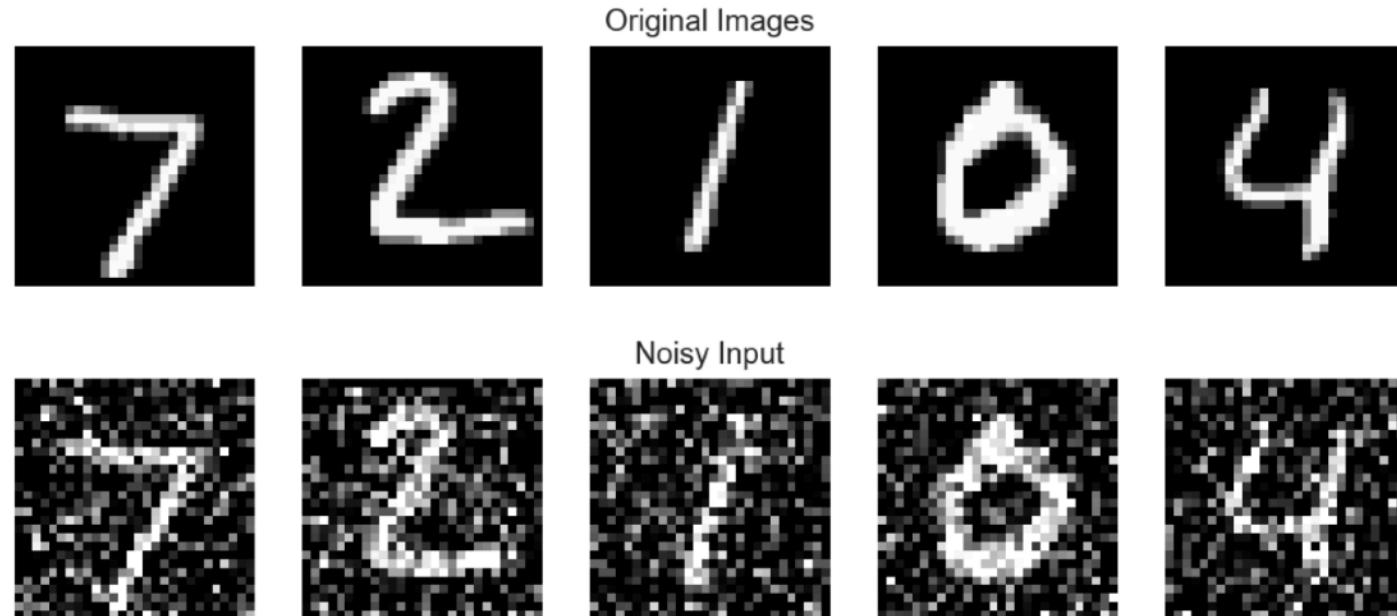
=> 기존 데이터를 토대로 신작 데이터를 만들 수 있는 것.



Noise Robustness

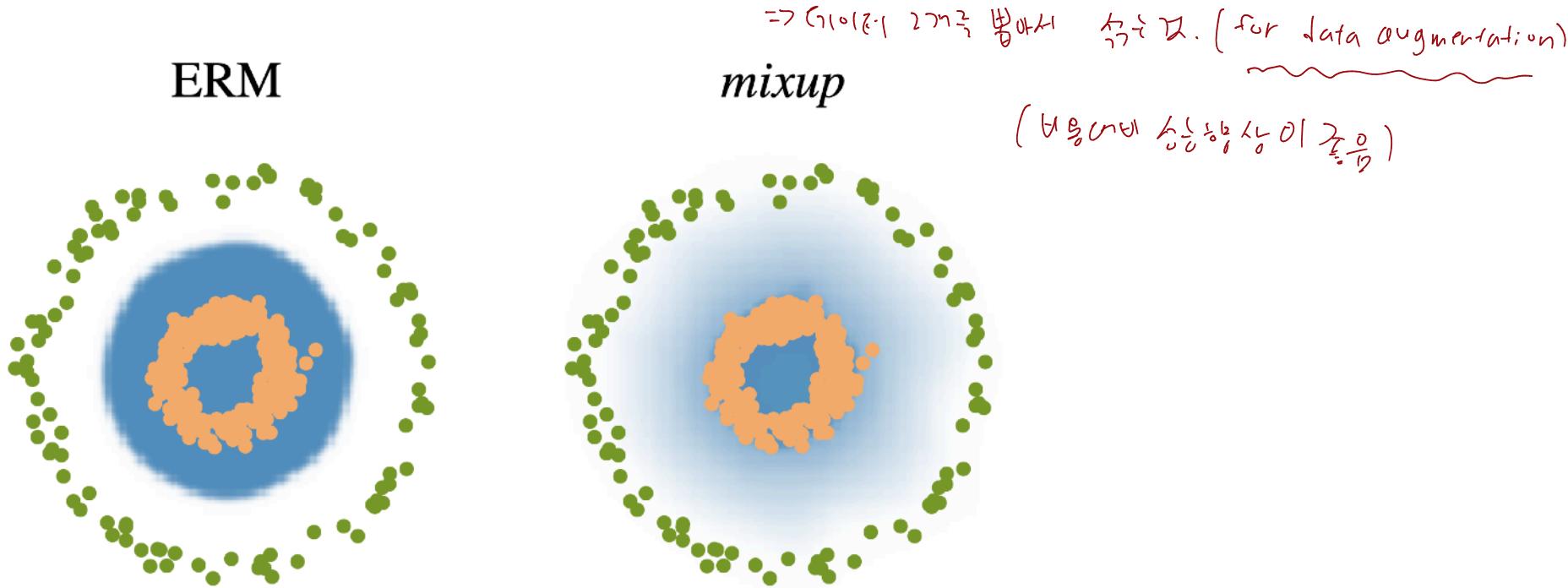
- Add random noises inputs or weights.

random noise는 input or weight에 주입
Test 결과는 72% 정확률로 4를 예측함



Label Smoothing

- Mix-up constructs augmented training examples by mixing both input and output of two randomly selected training data.



mixup: Beyond Empirical Risk Minimization, 2018

Label Smoothing

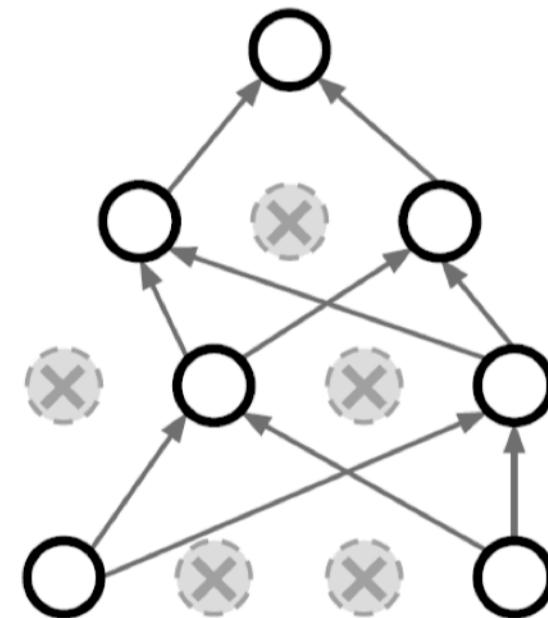
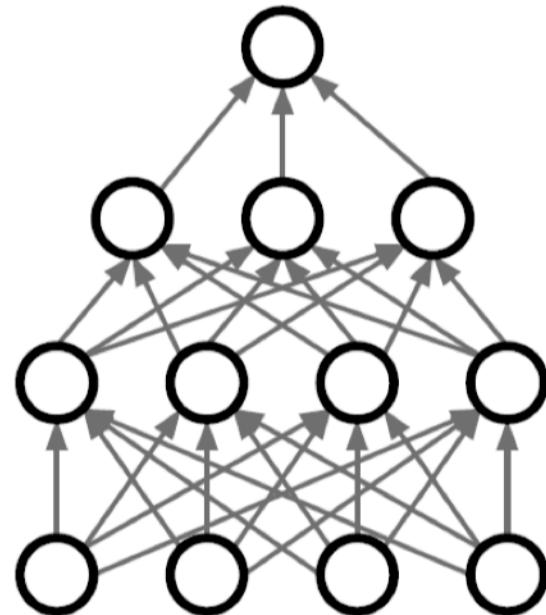
- CutMix constructs augmented training examples by mixing inputs with cut and paste and outputs with soft labels of two randomly selected training data.

	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

CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features, 2019

Dropout

- In each forward pass, randomly set some neurons to zero.



Batch Normalization

- Batch normalization compute the empirical mean and variance independently for each dimension (layers) and normalize.

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

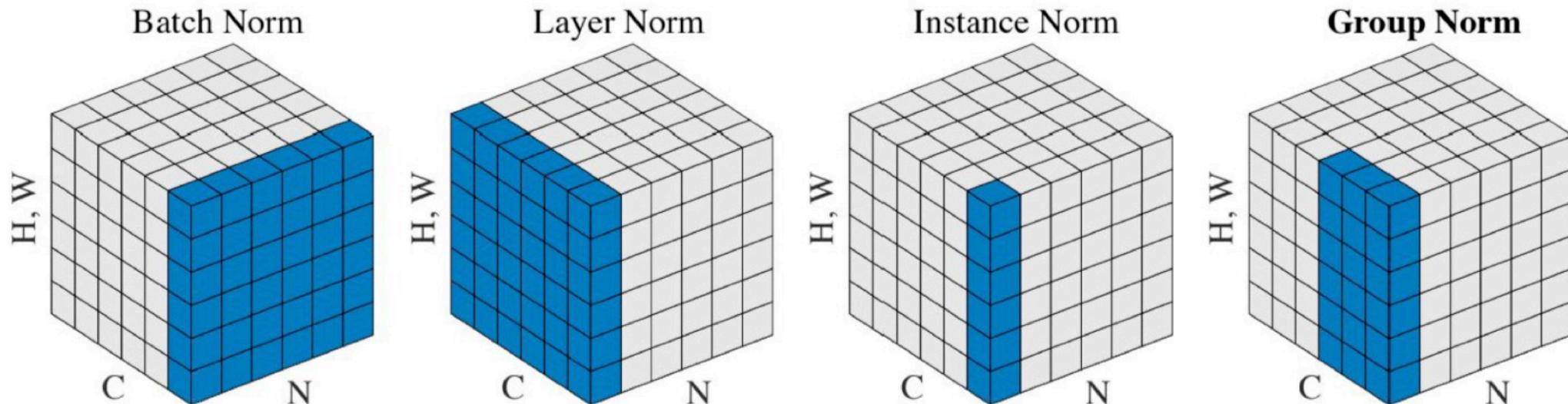
Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift, 2015

© NAVER Connect Foundation

Batch Normalization

- There are different variances of normalizations.

각 차원(2차원)에 대해 독립적으로 정규화 하는 것.



Group Normalization, 2018

Thank you for listening
