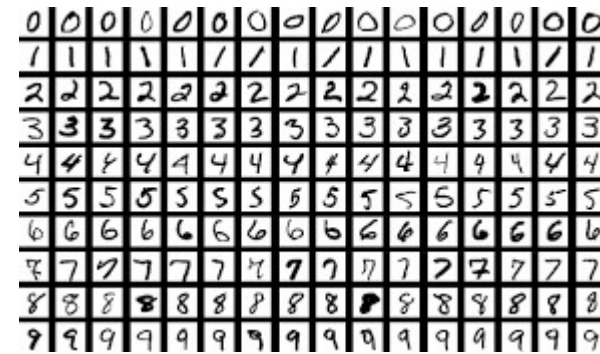# Machine Learning (ML) Training Basics

Saehwa Kim

Information and Communications Engineering
Hankuk University of Foreign Studies

# MNIST Database

- MNIST: Modified National Institute of Standards and Technology



- Database of handwritten digits

- Total 70,000 images

- (cf) Fashion-MNIST: a dataset of clothes images from an article
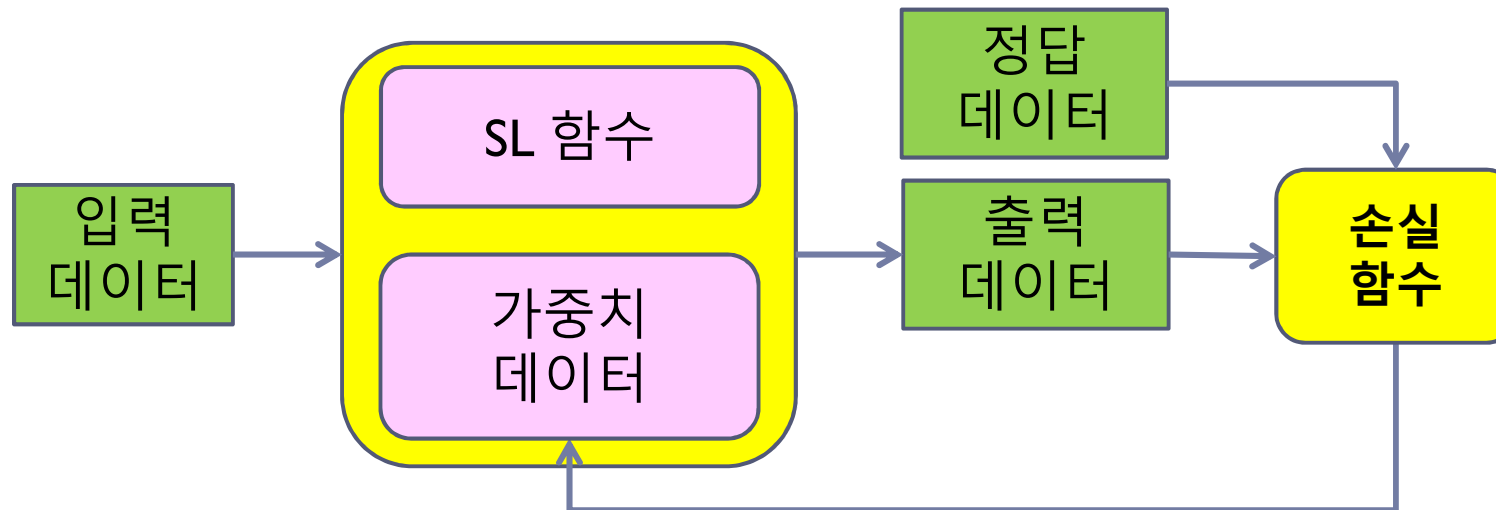
ESE Lab
http://eselab.hufs.ac.kr

# Data Set in Supervised Learning

▶ In supervised learning

▶ Training set and test set

　　▶ Common practice

　　　　▶ 80% data for training (20% data for testing)

▶ Data == Train data + Test data　(훈련, 시험)

▶ Training data == True training data + Validation data (검증)

　　▶ Holdout validation

　　▶ Held-out set ==　Validation set == Development set (dev set)

▶ (K-fold) Cross validation

　　▶ Uses many (K) small validation set

　　▶ Training time is multiplied by K.

▶ The case where training set error is low while the test set error is high

　　▶ Overfitting error == Generalization error == Out-of-sample error

ESE Lab
http://eselab.hufs.ac.kr

# Training in Supervised Learning (SL)

지도 학습

입력 데이터 → 프로그램 [에러 계산] → 결과 값

정답 출력 값
**Label**

입력 데이터 → [SL 함수 / 가중치 데이터] → 출력 데이터 → 손실 함수 ← 정답 데이터

손실 == Loss == Error == Cost

손실 함수 == Loss function

# Parameter Learning with Gradient Descent for (Simple Univariate) Linear Regression

- https://wikidocs.net/7635

- Hypothesis function h:  $h_\theta(x)=\theta_0+\theta_1 x$

- Cost function J: mean-squared-error (MSE): LSE (least squared error) criterion
  - 예측값과 실제값의 차이

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( \hat{y}^{(i)} - y^{(i)} \right)^2 = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$
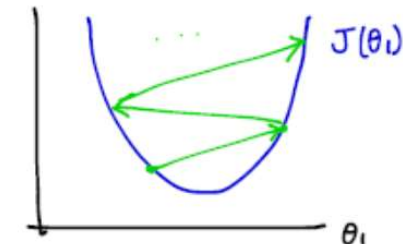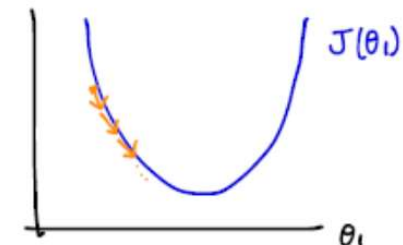
- **Gradient descent algorithm**

repeat until convergence{

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \text{for } j = 0, j = 1$$

}

- $:=$ | "assignment" operator

- $\alpha$ | learning rate

- $j$ | feature index number, should be updated simultaneously

ESE Lab
http://eselab.hufs.ac.kr

$x_j^{(T)} \leftarrow$ instance idx $1 \sim m$

$x_j \leftarrow$ feature idx $0 \sim n$

# Multiple Linear Regression

$$\hat{y}^{(i)} = \theta_0 \underset{=\theta_0}{x_0^{(i)}} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \cdots + \theta_n x_n^{(i)}$$

$$x_0^{(i)} = 1$$

bias (편향, 절편)

$(1,1)$

$$\hat{y}^{(i)} = h_{\bar{\theta}}(\bar{x}^{(i)}) = \bar{\theta}^T \bar{x}^{(i)} \quad {}^{(n+1,1)} \quad {}^{(1,n+1)}$$

$(m,1)$ $(m,n+1)$ $(n+1,1)$

$$\bar{y} = \bar{X}\bar{\theta}$$

$[ 1 ]$

$$\bar{y} = \begin{bmatrix} [y^{(1)}] \\ \vdots \\ [y^{(m)}] \end{bmatrix} \quad \bar{X} = \begin{bmatrix} \bar{x}^{(1)T} \\ \vdots \\ \bar{x}^{(m)T} \end{bmatrix} = \begin{bmatrix} [1, x_1^{(1)}, \cdots, x_n^{(1)}] \\ \vdots \\ [1, x_1^{(m)}, \cdots, x_n^{(m)}] \end{bmatrix} \quad \bar{\theta} = \begin{bmatrix} [\theta_0] \\ [\theta_1] \\ \vdots \\ [\theta_n] \end{bmatrix}$$
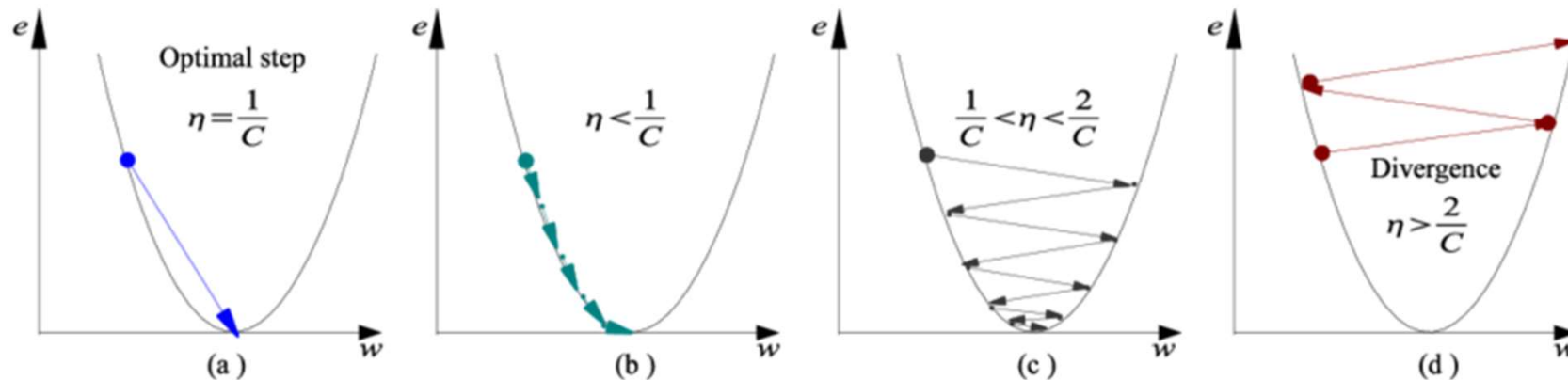
ESE Lab
http://eselab.hufs.ac.kr

# Gradient Descent Algorithm
# for Multiple Linear Regression

$$\nabla_a b = \frac{\partial b}{\partial a}$$

$$\bar{\theta} := \bar{\theta} - \alpha \nabla_{\theta} MSE(\bar{\theta})$$

direction

step size

cf. Gradient Ascent Algorithm for RL. (reward)



(a) Optimal step $\eta = \frac{1}{C}$

(b) $\eta < \frac{1}{C}$

(c) $\frac{1}{C} < \eta < \frac{2}{C}$

(d) Divergence $\eta > \frac{2}{C}$

이미지 출처: ttps://www.researchgate.net/figure/Convergence-Conditions-in-Gradient-Descent-Algorithm_fig2_224324276

ESE Lab
http://eselab.hufs.ac.kr

# Batch Gradient Descent vs. Stochastic Gradient Descent

▸ **Batch Gradient Descent**
  - ▸ 모든 training data 활용(하여 $\theta$ update)

▸ **Stochastic Gradient Descent (SGD)**
  - ▸ "임의의 하나"의 training data 활용(하여 $\theta$ update)
  - ▸ 이걸 True SGD라고 부르기도 함

▸ **Mini-batch Gradient Descent**
  - ▸ "임의의 일부(mini-batch)" training data 활용(하여 $\theta$ update)
  - ▸ 이걸 SGD라고 부르는 사람들도 있음

ESE Lab
http://eselab.hufs.ac.kr

# Batch Gradient Descent

$$\mathrm{Err}^{(i)}$$

$$\mathrm{MSE}(\overline{X}, h_{\overline{\theta}}) = \mathrm{MSE}(\overline{\theta}) = \frac{1}{m} \sum_{i=1}^{m} \left( \overline{\theta}^T \overline{x}^{(i)} - y^{(i)} \right)^2$$

$$\frac{\partial}{\partial \theta_j} \mathrm{MSE}(\overline{\theta}) = \frac{2}{m} \sum_{i=1}^{m} \left( \overline{\theta}^T \overline{x}^{(i)} - y^{(i)} \right) \cdot \overline{x}_j^{(i)}$$

$$\overline{y} = \overline{X}\,\overline{\theta}$$

$$\mathrm{Err}$$

$$\mathrm{MSE}(\overline{\theta}) = \frac{1}{m} \left( \overline{X}\,\overline{\theta} - \overline{y} \right)^T \left( \overline{X}\,\overline{\theta} - \overline{y} \right)$$

$$(n+1)\times m \qquad \mathrm{Err} \qquad m\times 1$$

$$\frac{\partial}{\partial \overline{\theta}} \mathrm{MSE}(\overline{\theta}) = \frac{2}{m} \cdot \overline{X}^T \left( \overline{X}\,\overline{\theta} - \overline{y} \right)$$

$$(n+1)\times 1$$

# BGD vs. (Mini-Batch) SGD (1)



Descending with step coefficient 0.005 (iteration 50)

$f(x) = x^2 * \sin(x)$

Start (2.5, 3.7)

End (4.9, -23.7)

뉴럴넷은 loss(or cost) function을 가지고 있습니다. 쉽게 말하면 "틀린 정도"
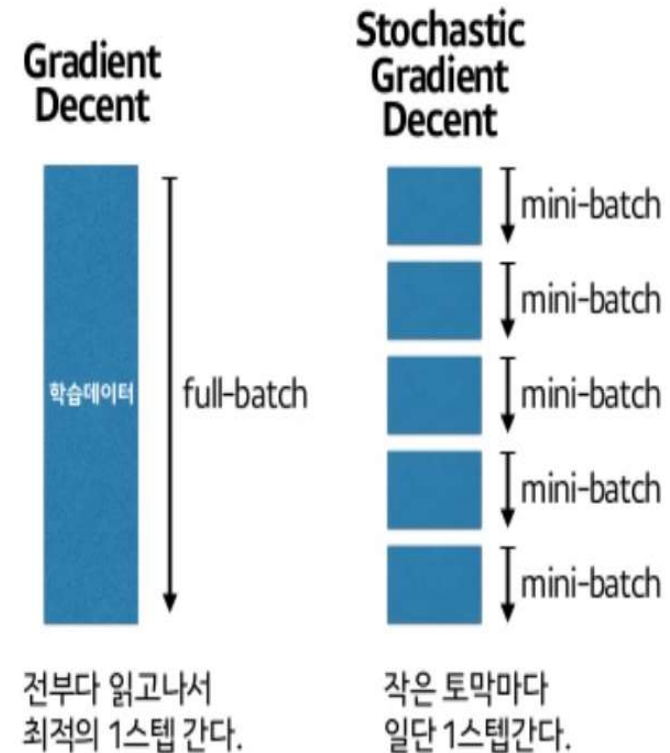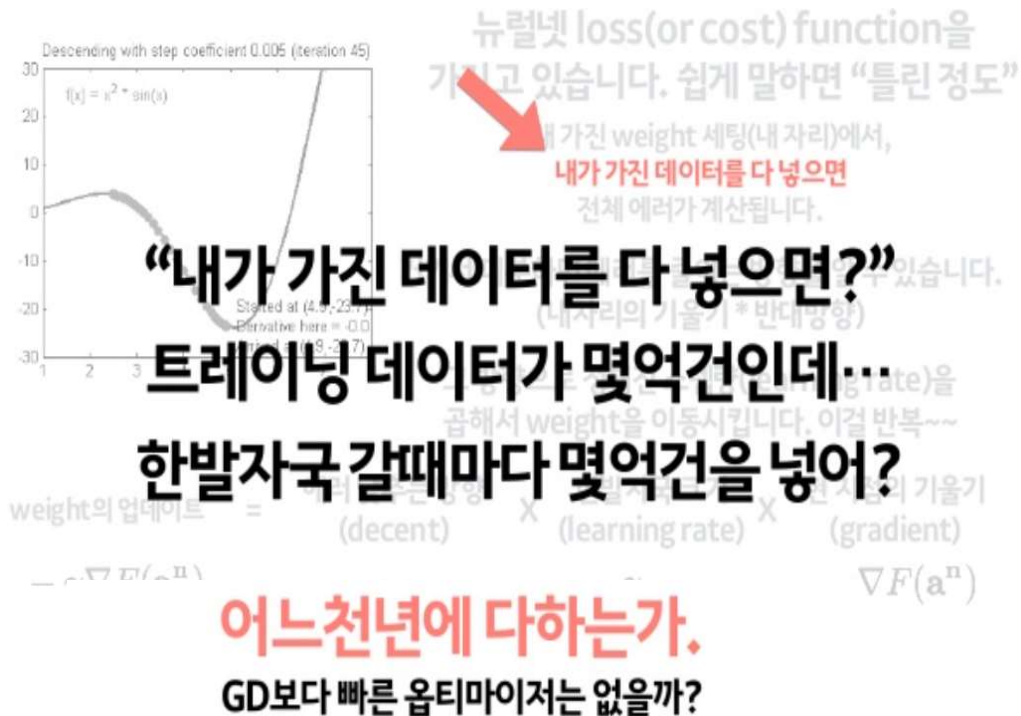
현재 가진 weight 세팅(내 자리)에서, 내가 가진 데이터를 다 넣으면 전체 에러가 계산됩니다.

거기서 미분하면 에러를 줄이는 방향을 알 수 있습니다. (내 자리의 기울기 * 반대방향)

그 방향으로 정해진 스텝량(learning rate)을 곱해서 weight을 이동시킵니다. 이걸 반복~~

weight의 업데이트 = 에러 낮추는 방향 (decent) X 한발자국 크기 (learning rate) X 현 지점의 기울기 (gradient)

$$-\gamma \nabla F(\mathbf{a}^n) \qquad - \qquad \gamma \qquad \nabla F(\mathbf{a}^n)$$

ESE Lab
http://eselab.hufs.ac.kr

# BGD vs. (Mini-Batch) SGD (2)



출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

# BGD vs. (Mini-Batch) SGD (3)



출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

# Gradient Descent Pitfalls



선형회기에서는 **MSE cost function**이 볼록함수여서 항상 최솟값 찾는 것이 가능

ESE Lab
http://eselab.hufs.ac.kr

# Feature Scaling

- **Min-max scaling**
  - Makes values ranging from 0 to 1
  - V' = (V-min)/(max-min)

- **Standardization**
  - Achieves zero mean and unit variance. (0, 1)
  - V' = (V-mean)/standard_dev
  - Much less effected by outliers (이상치)

ESE Lab
http://eselab.hufs.ac.kr

# Gradient Descent with and without Feature Scaling



그림 출처: https://medium.com/@mlgomez230/optimization-techniques-in-machine-learning-5d06725942b
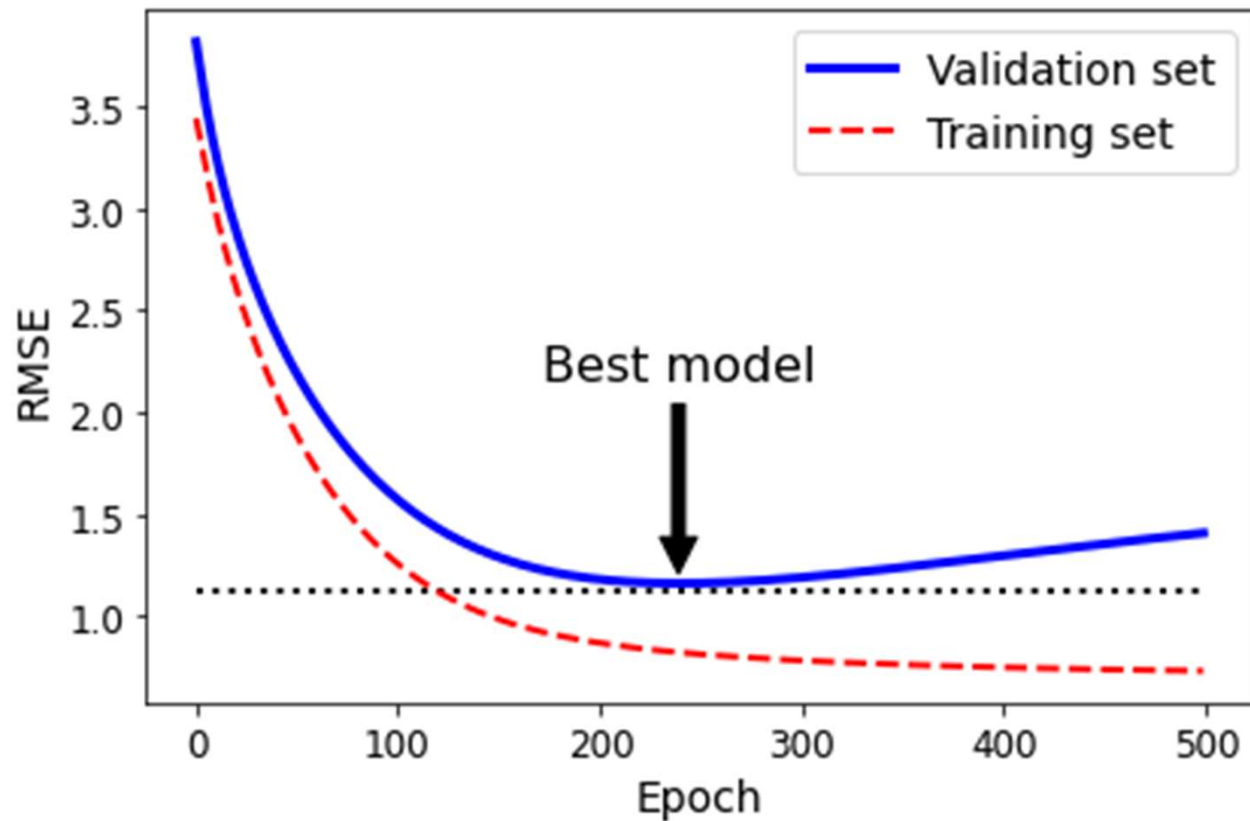
ESE Lab
http://eselab.hufs.ac.kr

# Path of Gradient Descent Parameters



출처: Aurélien Géron, Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow, 2nd Edition, O'Reilly Media

ESE Lab
http://eselab.hufs.ac.kr

# Early Stopping



출처: Aurélien Géron, Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow, 2nd Edition, O'Reilly Media

ESE Lab
http://eselab.hufs.ac.kr

# Vanishing Gradients Problem (1)



뉴럴넷의 학습방법 **Back propagation**
(사실 별거 없고 그냥 "뒤로 전달")

뭐를 전달하는가?
현재 내가 틀린정도를 '미분(**기울기**)' 한 거

입력 → 출력

미분하고, **곱하고**, 더하고를 역방향으로 반복하며 업데이트한다.

출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

## 근데 문제는?

우리가 activation 함수로 sigmoid 를 썼다는 것

여기의 미분(기울기)는 뭐라도 있다. 다행

근데 여기는 기울기 0.. 이런거 중간에 곱하면 **뭔가 뒤로 전달할게 없다?!**

그런 상황에서 이걸 반복하면??????

미분하고, 곱하고, 더하고를 역방향으로 반복

출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

# Vanishing Gradients Problem (3)

**Vanishing gradient 현상 :** 레이어가 깊을 수록 업데이트가 사라져간다.
그래서 fitting이 잘 안됨(underfitting)



출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

# Vanishing Gradients Problem (4)

사그라드는 sigmoid대신
죽지않는 activation func을 쓰자!

→ **ReLU** (Rectified Linear Units)

이녀석은 양의구간에서 전부 미분 값(1)이 있다!

끝줄학생 ← 줄좀맞추자   줄좀맞추자   줄좀맞추자   줄좀맞추자   줄좀맞추자 → 교장샘

끝 줄 학생까지 이야기가 전달이 잘 되고 위치를 고친다!

출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

# 기본 Optimizer: (Mini-Batch) SGD (1)
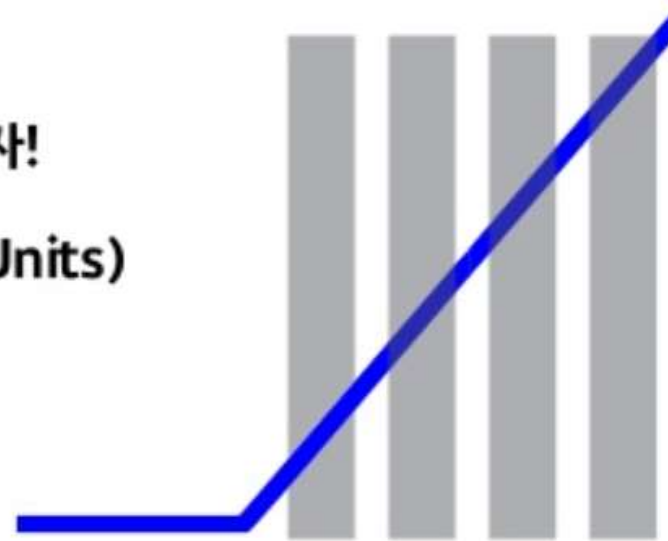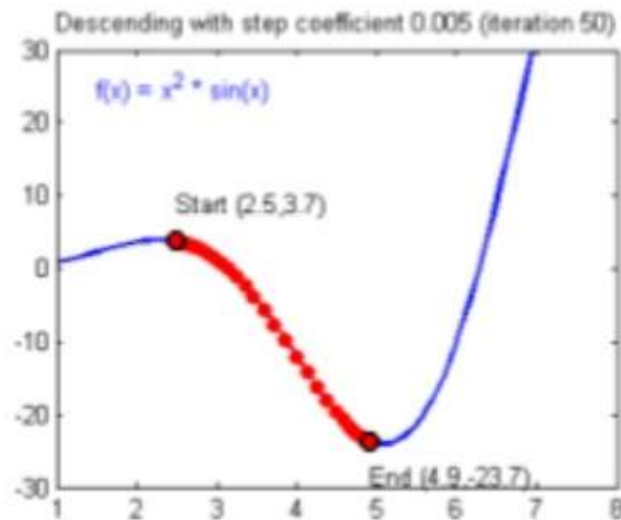
Descending with step coefficient 0.005 (iteration 50)

$f(x) = x^2 \ast \sin(x)$

Start (2.5,3.7)

End (4.9,-23.7)

뉴럴넷은 loss(or cost) function을 가지고 있습니다. 쉽게 말하면 "틀린 정도"

현재 가진 weight 세팅(내 자리)에서, 내가 가진 데이터를 다 넣으면 전체 에러가 계산됩니다.

거기서 미분하면 에러를 줄이는 방향을 알 수 있습니다. (내자리의 기울기 * 반대방향)

그 방향으로 정해진 스텝량(learning rate)을 곱해서 weight을 이동시킵니다. 이걸 반복~~

| weight의 업데이트 = | 에러 낮추는 방향 (decent) | X | 한발자국 크기 (learning rate) | X | 현 지점의 기울기 (gradient) |
|---|---|---|---|---|---|
| $-\gamma \nabla F(\mathbf{a}^n)$ | $-$ | | $\gamma$ | | $\nabla F(\mathbf{a}^n)$ |

출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

# 기본 Optimizer: (Mini-Batch) SGD (2)

근데 미니 배치를 하다보니 와리가리(?) 방향 문제가 있다.

딱 봐도 더 잘 갈 수 있는데
훨씬 더 헤매면서 간다.

훑기도 잘 훑으면서,
좀 더 휙휙 더 좋은 방향으로 갈 순 없을까?

출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

# 기본 Optimizer: (Mini-Batch) SGD (3)

## 스텝사이즈(learning rate)도 문제가 된다.



보폭이 너무 작으면 오래 헤매고(파란라인)

보폭이 너무 크면, 오솔길을 지나친다(녹색라인)

출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

# Faster Optimizers (1)

$$-\gamma\nabla F(\mathbf{a^n}) \qquad \text{산을 잘 타고 내려오는 것은}$$

$$\nabla F(\mathbf{a^n}) \qquad \text{어느 방향으로 발을 디딜지}$$

$$\gamma \qquad \text{얼마 보폭으로 발을 디딜지}$$

**두가지를 잘잡아야 빠르게 타고 내려온다.**

**SGD를 더 개선한 멋진 optimizer가 많다!**
**SGD의 개선된 후계자들**

출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

# Faster Optimizers (2)



그래프 출처: https://imgur.com/NKsFHJb

# Faster Optimizers (3)



산 내려오는 작은 오솔길 잘찾기(Optimizer)의 발달 계보

모든 자료를 다 검토해서
내 위치의 산기울기를 계산해서
갈 방향을 찾겠다.

**GD**

**Nesterov Accelerated Gradient**
**NAG**
일단 관성 방향 먼저 움직이고,
움직인 자리에 스텝을 계산하니
더 빠르더라

**Momentum**
스텝 계산해서 움직인 후,
아까 내려 오던 관성 방향 또 가자

**Nadam**
Adam에 Momentum
대신 NAG를 붙이자.

스텝방향

**SGD**
전부 다봐야 한걸음은
너무 오래 걸리니까
조금만 보고 빨리 판단한다
같은 시간에 더 많이 간다

스텝사이즈

**Adam**
RMSProp + Momentum
방향도 스텝사이즈도 적절하게!

**RMSProp**
보폭을 줄이는 건 좋은데
이전 맥락 상황봐가며 하자.

**Adagrad**
안 가본곳은 성큼 빠르게 걸어 훑고
많이 가본 곳은 잘 아니까
갈수록 보폭을 줄여 세밀히 탐색

**AdaDelta**
종종걸음 너무 작아져서
정지하는걸 막아보자.

출처: https://www.slideshare.net/yongho/ss-79607172, 하용호@kakao

ESE Lab
http://eselab.hufs.ac.kr

# Bias and Variance Tradeoff (1)

▶ http://scott.fortmann-roe.com/docs/BiasVariance.html



The center of the target is a model that perfectly predicts the correct values. (As we move away from the bulls-eye, our predictions get worse.)

Fig. 1 Graphical illustration of bias and variance.

ESE Lab
http://eselab.hufs.ac.kr

# Bias and Variance Tradeoff (2)

- https://towardsdatascience.com/the-bias-variance-tradeoff-8818f41e39e9

- Suppose that we have independent variables $x$ that affect the value of a dependent variable $y$

$$y = f(x) + \epsilon$$

- Noise is modeled by random variable $\epsilon$ with zero mean and variance $\sigma\epsilon^2$

$$\mathbb{E}[\epsilon] = 0, \text{var}(\epsilon) = \mathbb{E}[\epsilon^2] = \sigma_\epsilon^2$$

- In the linear regression, mean square error

$$\text{MSE} = \mathbb{E}[(y - \hat{f}(x))^2]$$

$$\text{bias}[\hat{f}(x)] = \mathbb{E}[\hat{f}(x)] - f(x)$$

$$\text{var}(\hat{f}(x)) = \mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2]$$

$$\mathbb{E}[\mathbb{E}[(y - \hat{f}(x))^2]] = \mathbb{E}[\text{bias}[\hat{f}(x)]^2] + \mathbb{E}[\text{var}(\hat{f}(x))] + \sigma_\epsilon^2$$

$$Err(x) = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

$$\mathbb{E}[(y - \hat{f}(x))^2] = \mathbb{E}[(f(x) + \epsilon - \hat{f}(x))^2] \tag{1}$$

$$= \mathbb{E}[(f(x) - \hat{f}(x))^2] + \mathbb{E}[\epsilon^2] + 2\mathbb{E}[(f(x) - \hat{f}(x))\epsilon]$$

$$= \mathbb{E}[(f(x) - \hat{f}(x))^2] + \underbrace{\mathbb{E}[\epsilon^2]}_{=\sigma_\epsilon^2} + 2\mathbb{E}[(f(x) - \hat{f}(x))] \underbrace{\mathbb{E}[\epsilon]}_{=0} \tag{2}$$

$$= \mathbb{E}[(f(x) - \hat{f}(x))^2] + \sigma_\epsilon^2 \tag{3}$$

$$\mathbb{E}[(f(x) - \hat{f}(x))^2] = \mathbb{E}\left[\left((f(x) - \mathbb{E}[\hat{f}(x)]) - (\hat{f}(x) - \mathbb{E}[\hat{f}(x)])\right)^2\right] \tag{4}$$

$$= \mathbb{E}\left[\left(\mathbb{E}[\hat{f}(x)] - f(x)\right)^2\right] + \mathbb{E}\left[\left(\hat{f}(x) - \mathbb{E}[\hat{f}(x)]\right)^2\right]$$

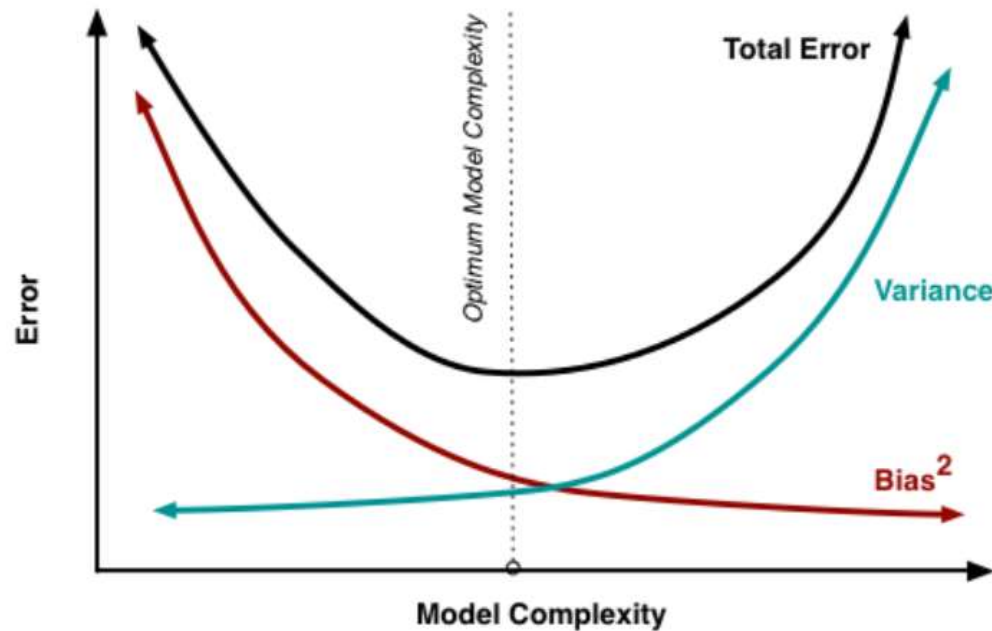$$- 2\mathbb{E}\left[\left(f(x) - \mathbb{E}[\hat{f}(x)]\right)\left(\hat{f}(x) - \mathbb{E}[\hat{f}(x)]\right)\right] \tag{5}$$

$$= \underbrace{(\mathbb{E}[\hat{f}(x)] - f(x))^2}_{=\text{bias}[\hat{f}(x)]} + \underbrace{\mathbb{E}\left[\left(\hat{f}(x) - \mathbb{E}[\hat{f}(x)]\right)^2\right]}_{=\text{var}(\hat{f}(x))}$$

$$- 2\left(f(x) - \mathbb{E}[\hat{f}(x)]\right)\mathbb{E}\left[\left(\hat{f}(x) - \mathbb{E}[\hat{f}(x)]\right)\right] \tag{6}$$

$$= \text{bias}[\hat{f}(x)]^2 + \text{var}(\hat{f}(x))$$

$$- 2\left(f(x) - \mathbb{E}[\hat{f}(x)]\right)\left(\mathbb{E}[\hat{f}(x)] - \mathbb{E}[\hat{f}(x)]\right) \tag{7}$$

$$= \text{bias}[\hat{f}(x)]^2 + \text{var}(\hat{f}(x)) \tag{8}$$

ESE Lab
http://eselab.hufs.ac.kr

# Bias and Variance Tradeoff (3)

▸ http://scott.fortmann-roe.com/docs/BiasVariance.html



Fig. 6 Bias and variance contributing to total error.

$$\frac{dBias}{dComplexity} = -\frac{dVariance}{dComplexity}$$

ESE Lab
http://eselab.hufs.ac.kr

혼용 분산

# Bias-Variance Tradeoff (4)

▸ In parameter estimation of a model https://en.wikipedia.org/wiki/Bias%E2%80%93variance_tradeoff

  ▸ A lower bias ←→ A higher variance

▸ Higher bias → Underfitting (too simple)

  ▸ Miss the relevant relations between features and target outputs

▸ Higher variance → Overfitting (too complex)

  ▸ Sensitivity to small fluctuations in the training set

  ▸ Do not ignore the random noise in the training data

▸ Both bias and variance decrease when increasing the width of a neural network.

▸ Bias-variance decomposition

  ▸ Regularization with the expected generalization error

\* ksaehwa: [ Bias : similarity를 찾음
            [ Variance : difference를 찾음

ESE Lab
http://eselab.hufs.ac.kr

# Bias-Variance Tradeoff (5)

- [https://medium.com/@mp32445/understanding-bias-variance-tradeoff-ca59a22e2a83](https://medium.com/@mp32445/understanding-bias-variance-tradeoff-ca59a22e2a83)

- **Examples of low-bias and high-variance machine learning algorithms**
  - Support Vector Machines
  - Decision Trees
  - k-Nearest Neighbors

- **Examples of high-bias and low-variance machine learning algorithms**
  - Linear Regression
  - Logistic Regression

ESE Lab
http://eselab.hufs.ac.kr