

기계학습 (Machine Learning)

L10

- Regularization

한밭대학교

정보통신공학과

최 해 철

- ◆ Overfitting & Underfitting
- ◆ Bias and Variance
- ◆ Regularization by Weight Penalty

References

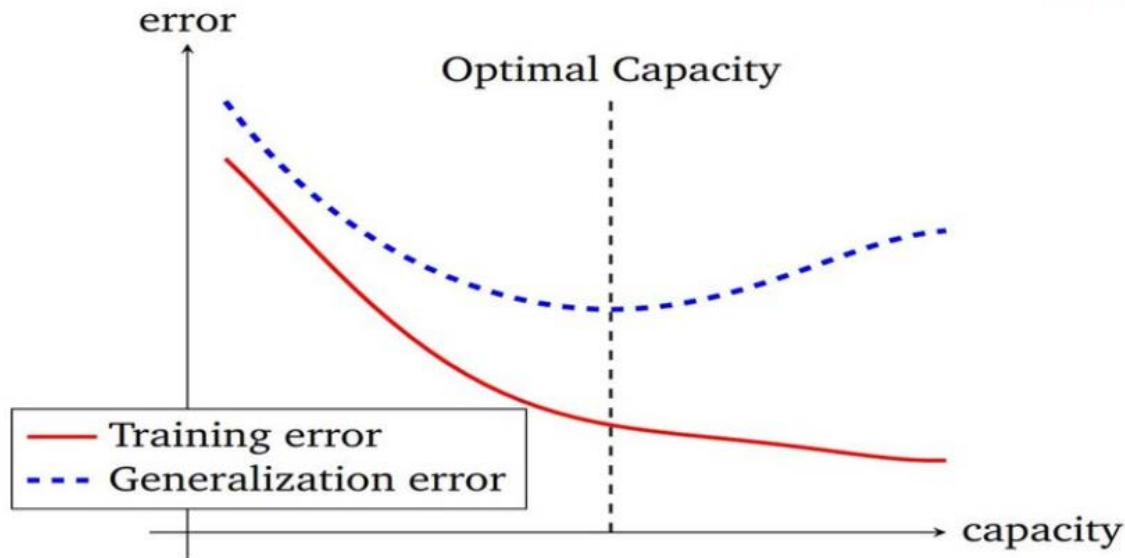
- 기계 학습 “3장 다층 퍼셉트론” by 오일석, 패턴 인식 by 오일석
- 단단한 머신러닝 by 조우쯔와

1. Overfitting & Underfitting

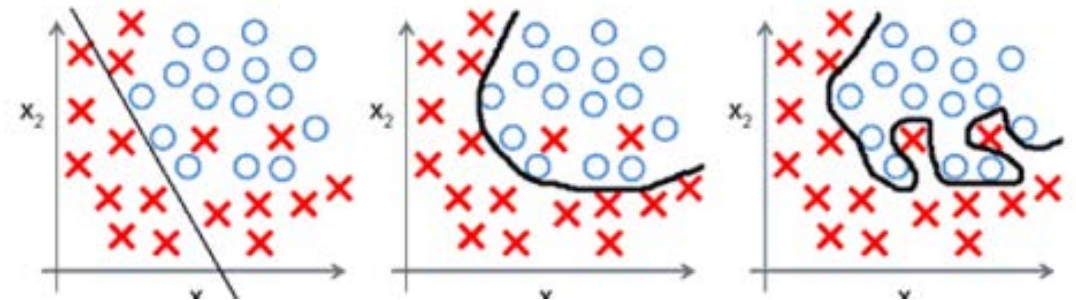
Generalization

일반화

- 모델이 학습 데이터에 대해 학습한 후, 이전에 본 적이 없는 새로운 데이터에 대해 정확하게 예측할 수 있는 능력



• Example: in Classification...



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$

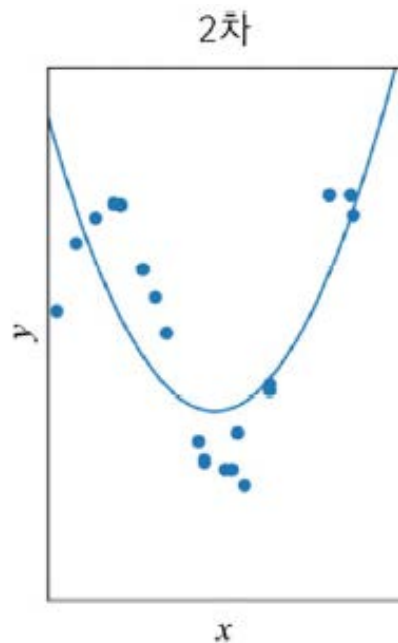
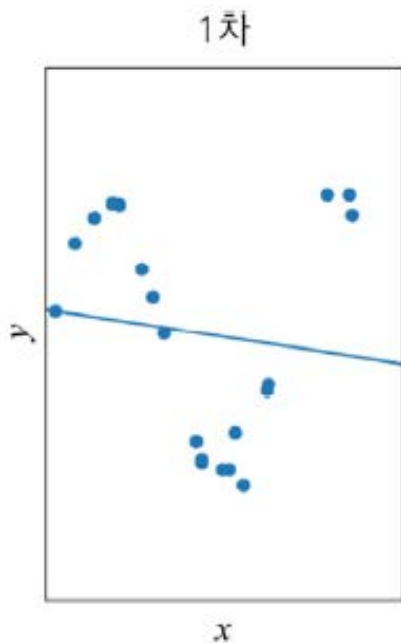
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2^2 + \dots)$$

Underfitting

◆ 과소적합 과 훈련 오차

- '모델의 용량이 너무 작아' or '훈련집합이 너무 작아' 오차가 클 수밖에 없는 현상
 - 예) 아래 그림의 선형(1차 다항식) 또는 2차 다항식 모델을 사용한 경우

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x)$$

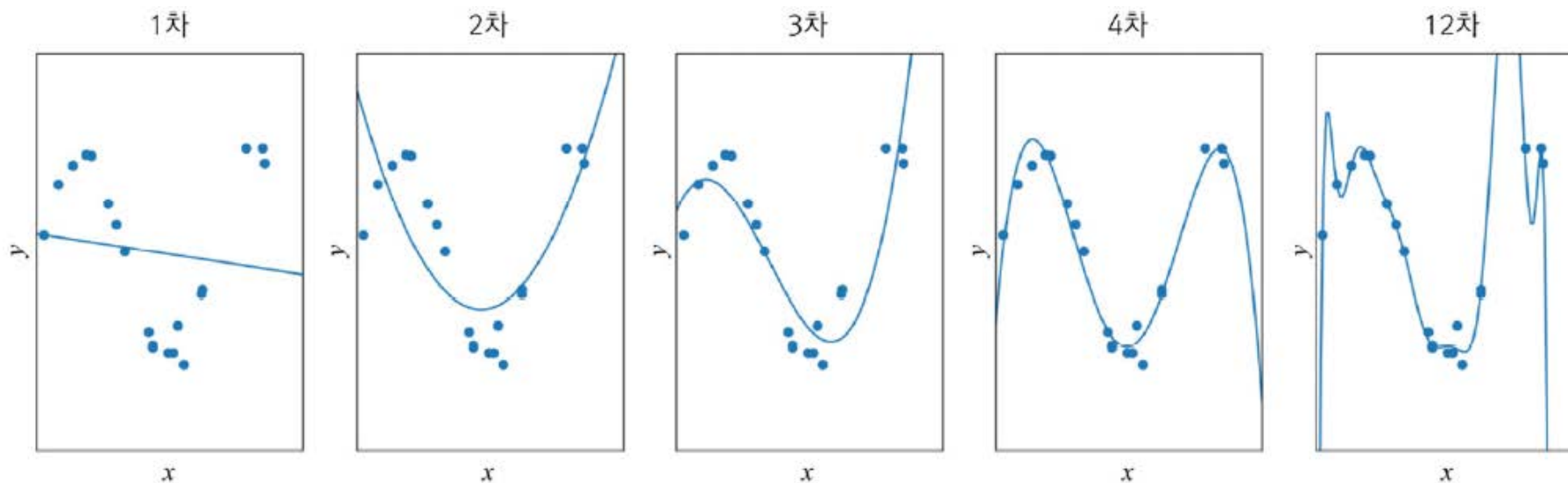


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x + \theta_2 x^2)$$

Underfitting

◆ Underfitting 방지

- 비선형 모델 등과 같이 **용량이 더 큰 모델을 사용**
- 충분한 훈련 집합을 활용
 - 예) 아래 그림의 3차, 4차, 12차 다항식 모델의 경우

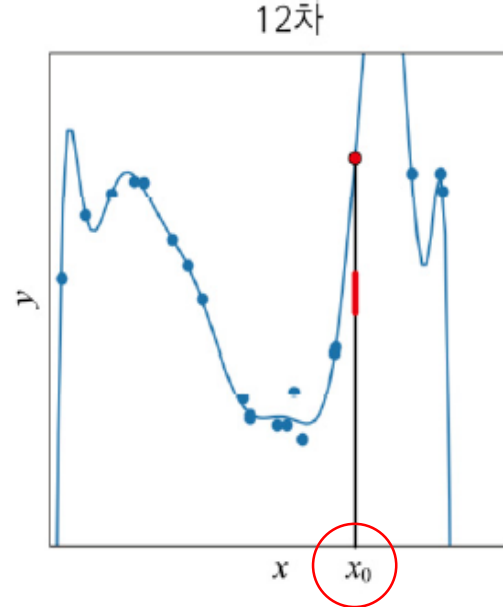
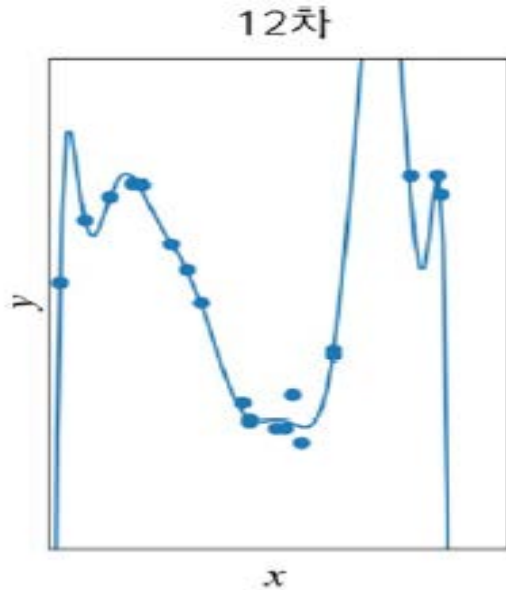


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x + \theta_2 x^2 + \dots \theta_{11} x^{11} + \theta_{12} x^{12})$$

Overfitting

◆ 과적합 (overfitting) 과 예측(시험) 오차

- 12차 다항식 곡선을 채택한다면 훈련집합에 대해 거의 완벽하게 근사화함
- 하지만 '새로운' 데이터를 예측한다면 큰 문제 발생 (빨간 점)
- 이유는 '용량이 너무 크기' 때문에 학습 과정에서 잡음까지 수용 → 과잉적합 현상



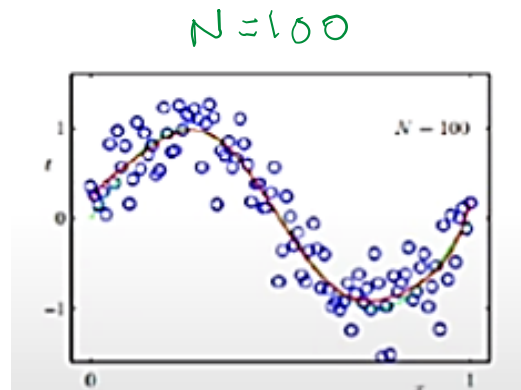
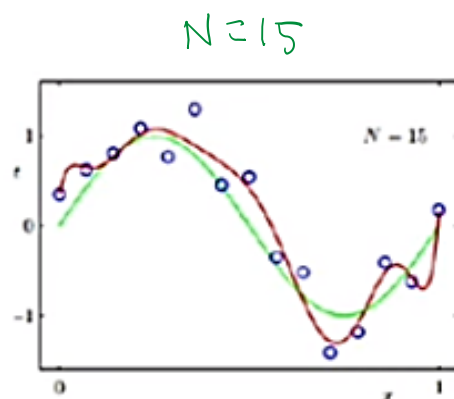
$$\hat{y} = f(x) + \text{noise}$$

Causes of Overfitting – 1. Data

◆ Insufficient # of Training Examples

- the training set may be too sparse or cannot represent the full variety of the data

데이터가 들통들성 있다 (데이터 부족할 때)



N : # of training examples

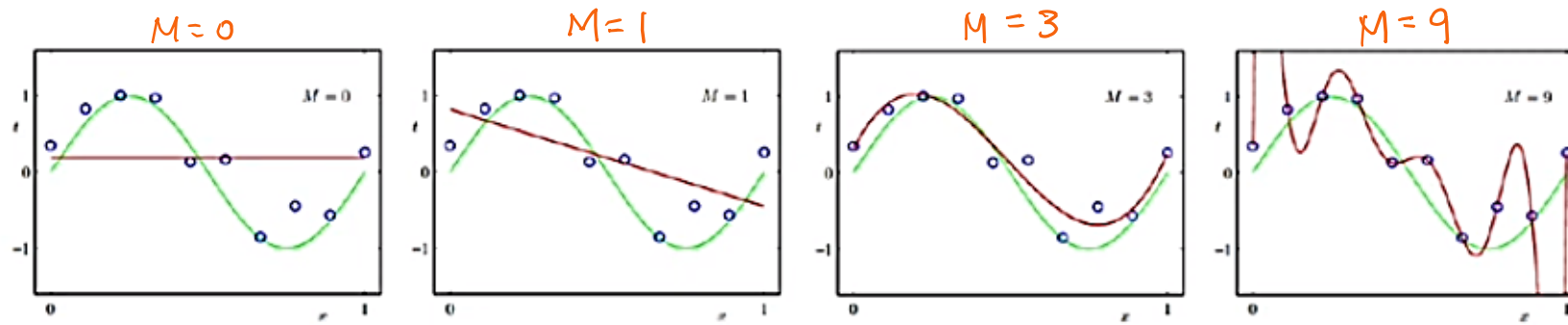
- 해결책: 충분히 많은 Training Data 사용

- Cf.) 데이터 증대(Data Augmentation) 기법 등을 통해 기존 Training Data 증대 가능

Causes of Overfitting – 2. Model

◆ Too Large # of Parameters (Model Capacity)

- the model is relatively too flexible for the dataset
- the resulting parameters tend to have large values



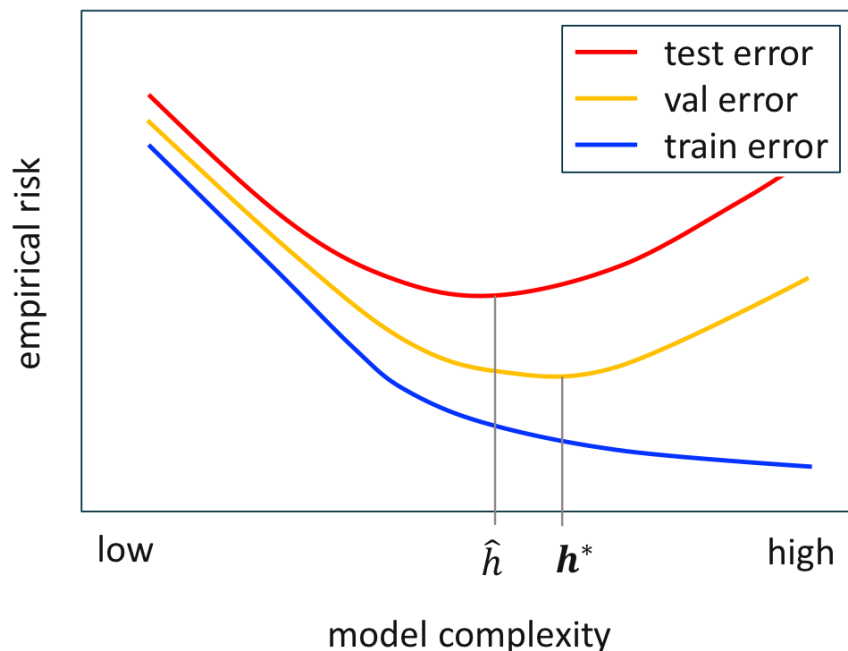
	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

3차 예, $h_{\theta}(x) = g(\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3)$

Causes of Overfitting – 2. Model (cont'd)

◆ 해결책 1: 검증집합 을 이용한 모델 선택

- 훈련집합과 테스트집합과 다른 별도의 **검증집합**을 준비한다.
- **모델집합**에 속한 각각의 모델에 대해 **훈련집합**으로 학습시킨다. (훈련 성능)
 - 앞의 예에서는 서로 다른 차수의 다항식의 집합(서로 다른 용량)이 모델집합인 셈
- **검증집합**에 대해 최고의 성능을 보인 모델을 선택한다. (검증 성능) → *Overfitting 방지*



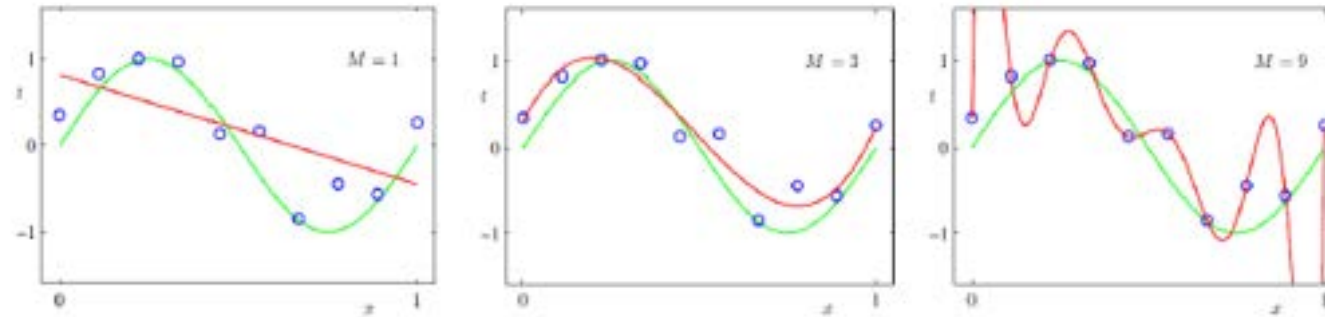
h^* : model with lowest validation error

\hat{h} : model with lowest test error

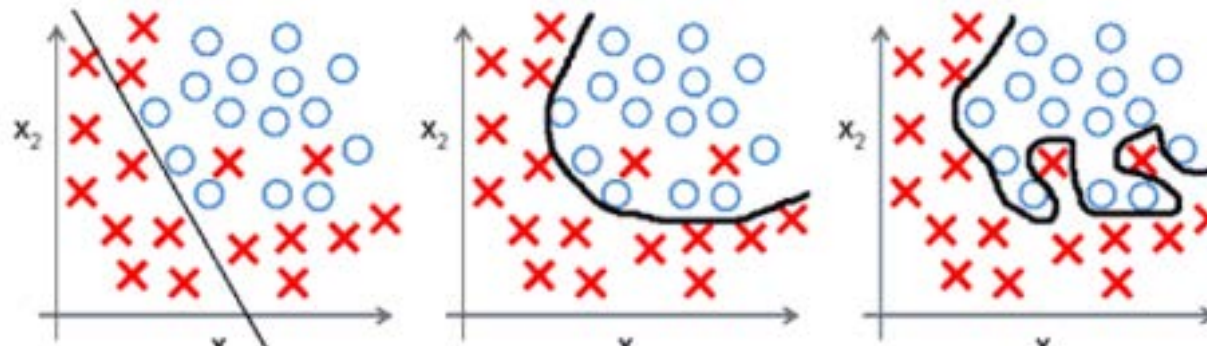
Overfitting

- ◆ 해결책 2: 적당한 용량의 모델을 선택
 - Model selection, model evaluation 작업을 수행

- *Example: in Regression...*



- *Example: in Classification...*



Causes of Overfitting – 2. Model (cont'd)

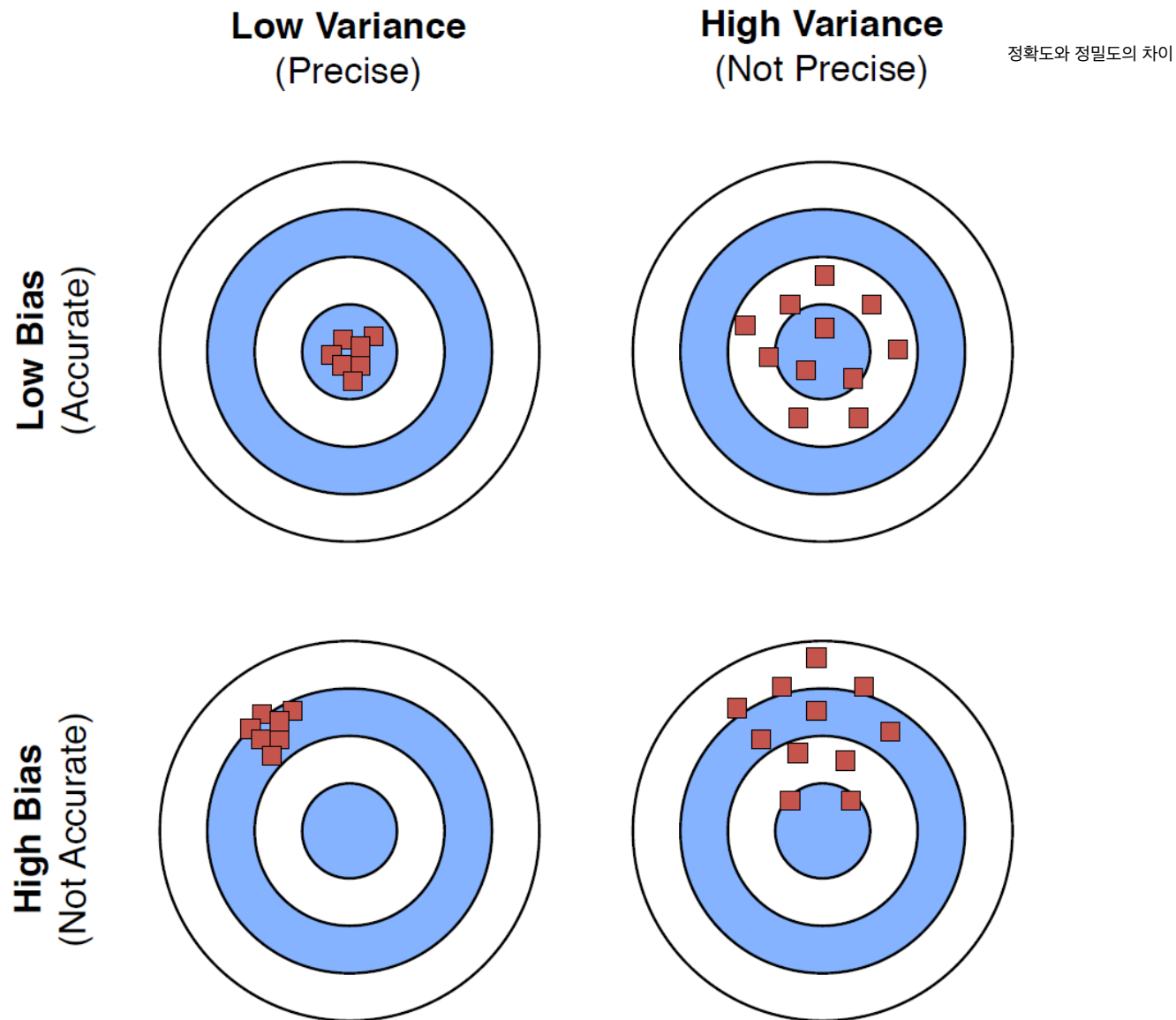
규제

◆ 해결책 3: _____

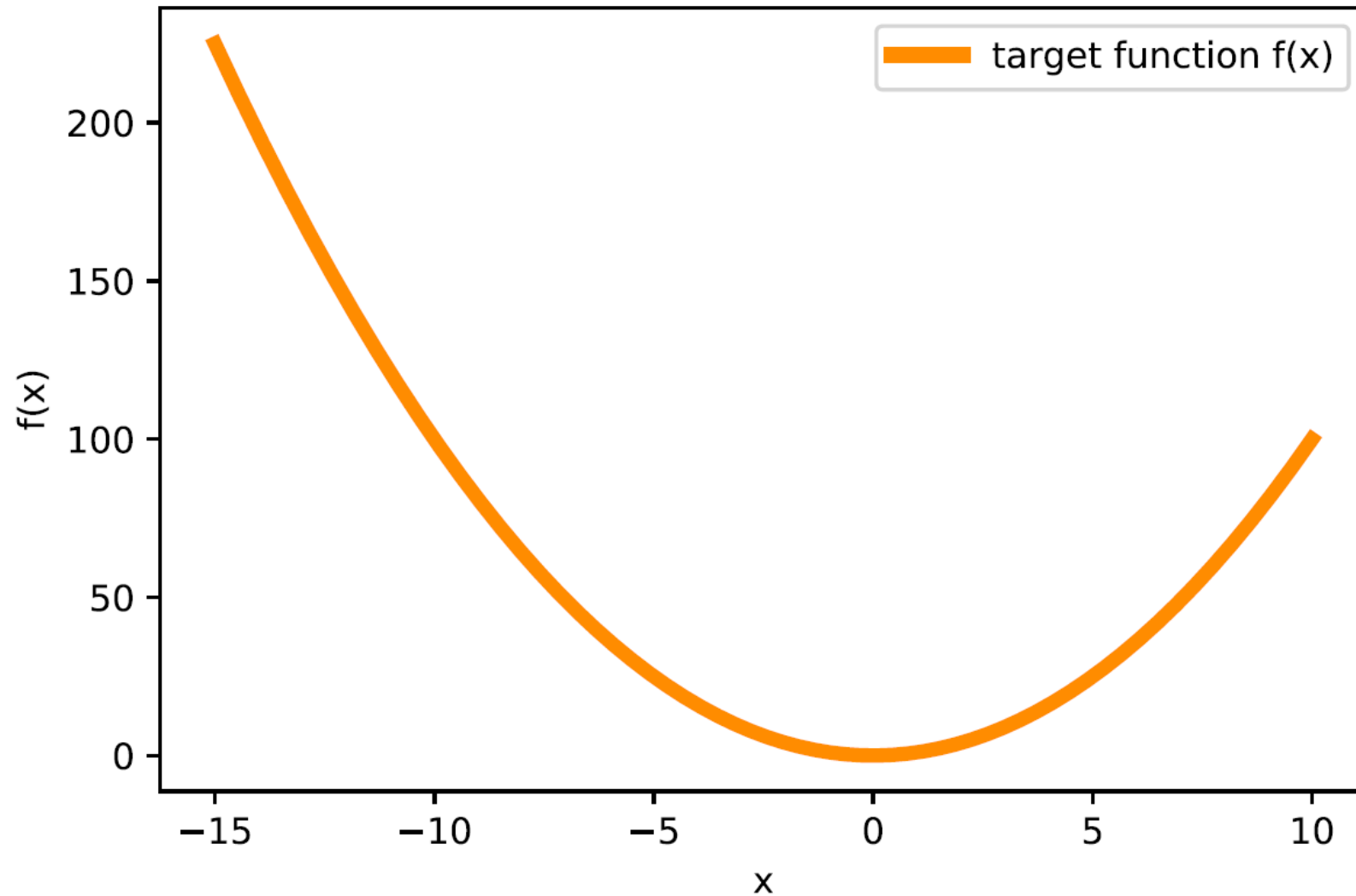
- 용량이 충분히 큰 모델 + 다양한 **규제**(*Regularization*) 기법을 적용
 - 예시: *Weight Penalty*, Drop-out...
 - Overfitting을 방지하기 위한 기술을 통칭하여 '규제'라고 부르기도 함
 - **Regularization Parameter**들은 *Validation Set* 을 이용하여 결정 가능 (model selection)

2. Bias and Variance

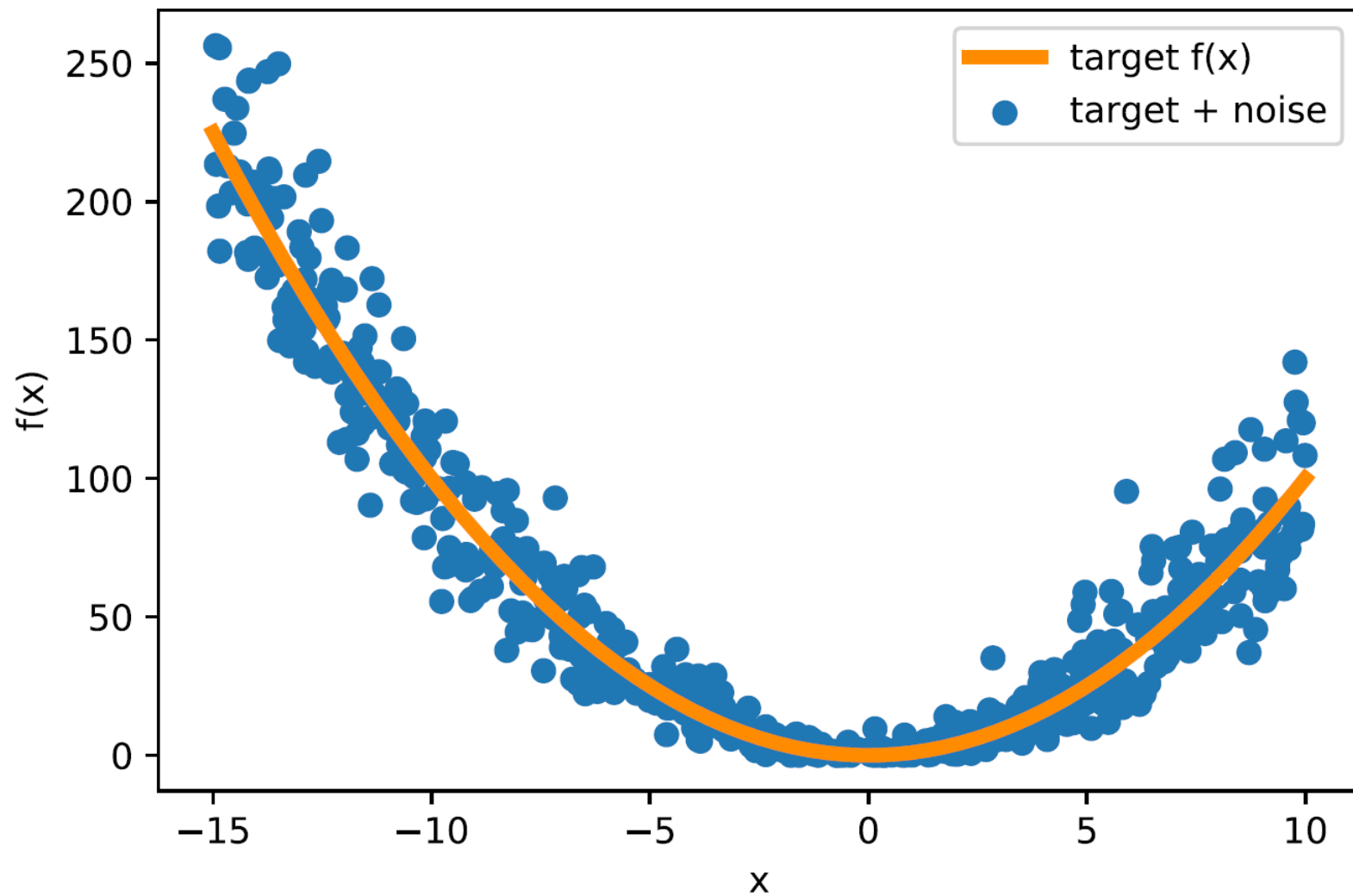
Bias-Variance 직관



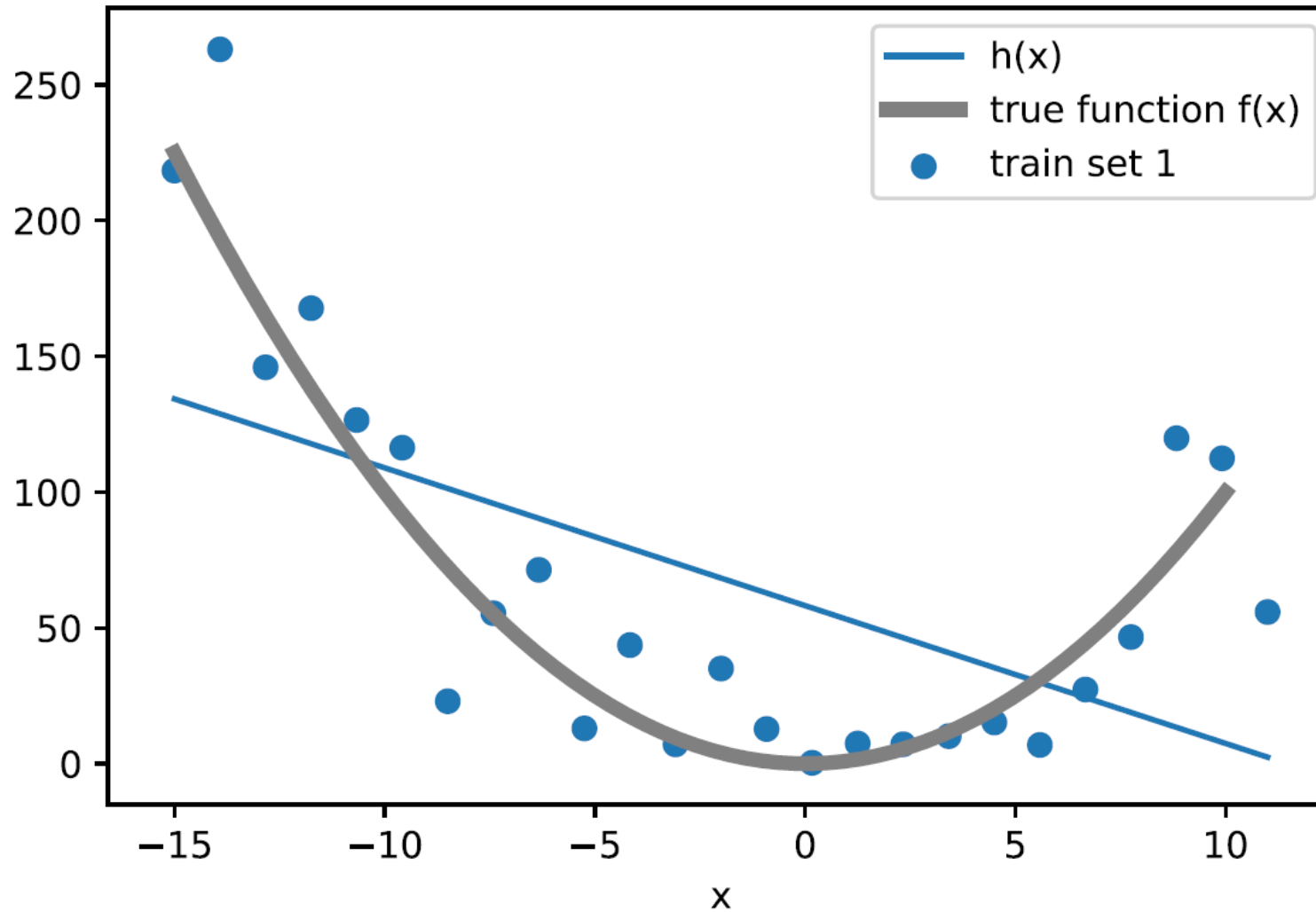
Bias-Variance 직관



Bias-Variance 직관

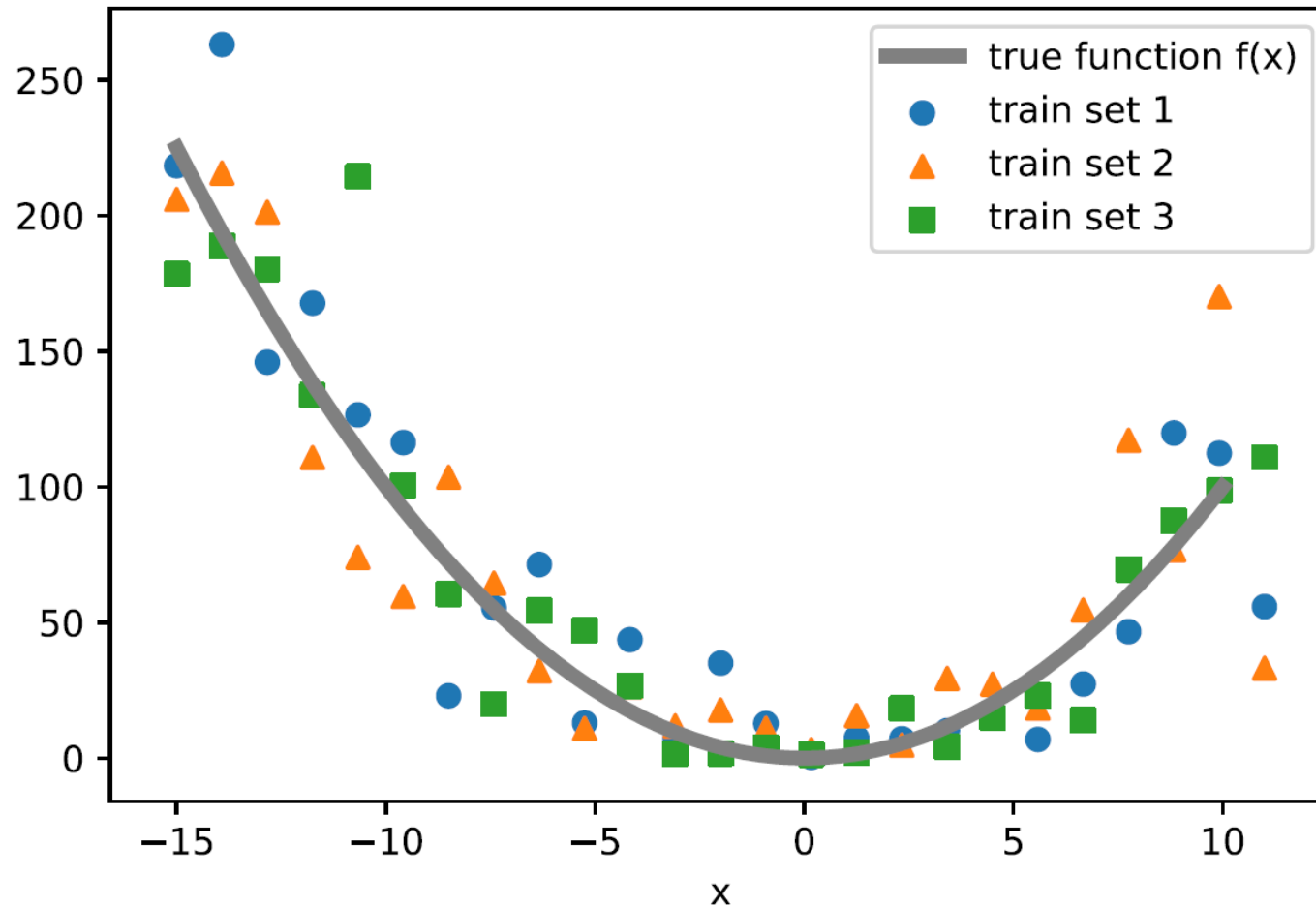


Bias-Variance 직관

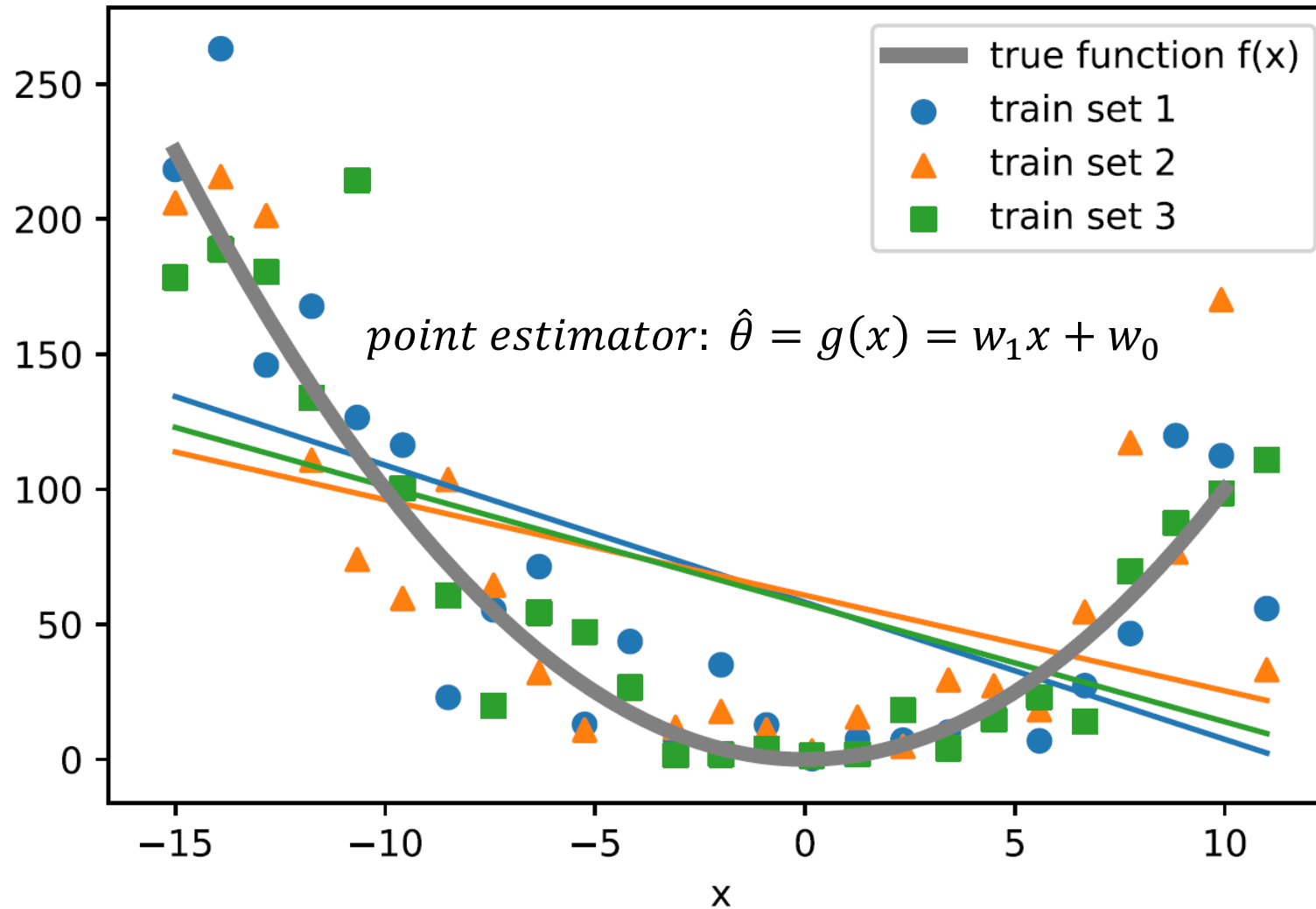


Bias-Variance 직관

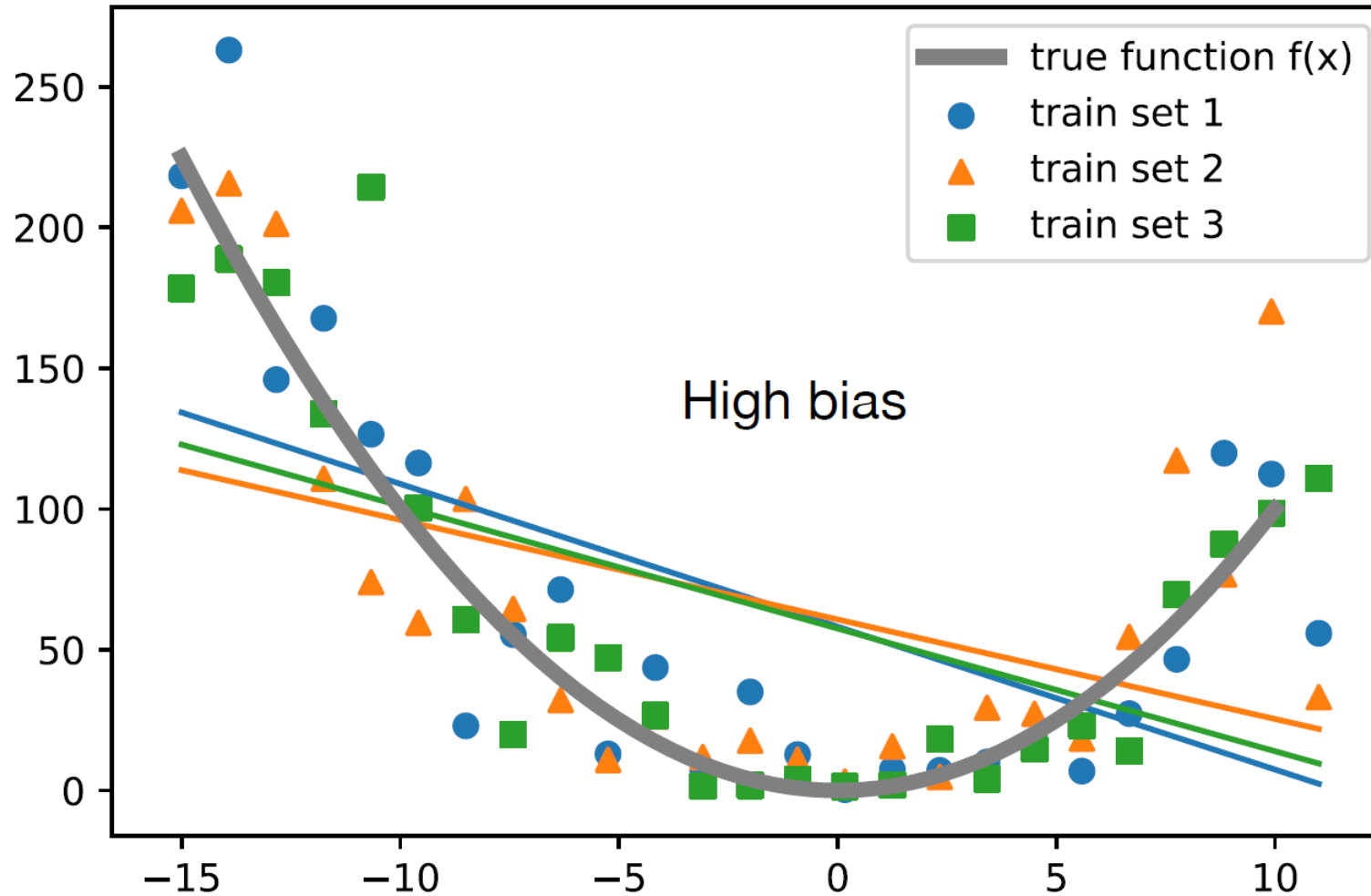
suppose we have multiple training sets



Bias-Variance 직관

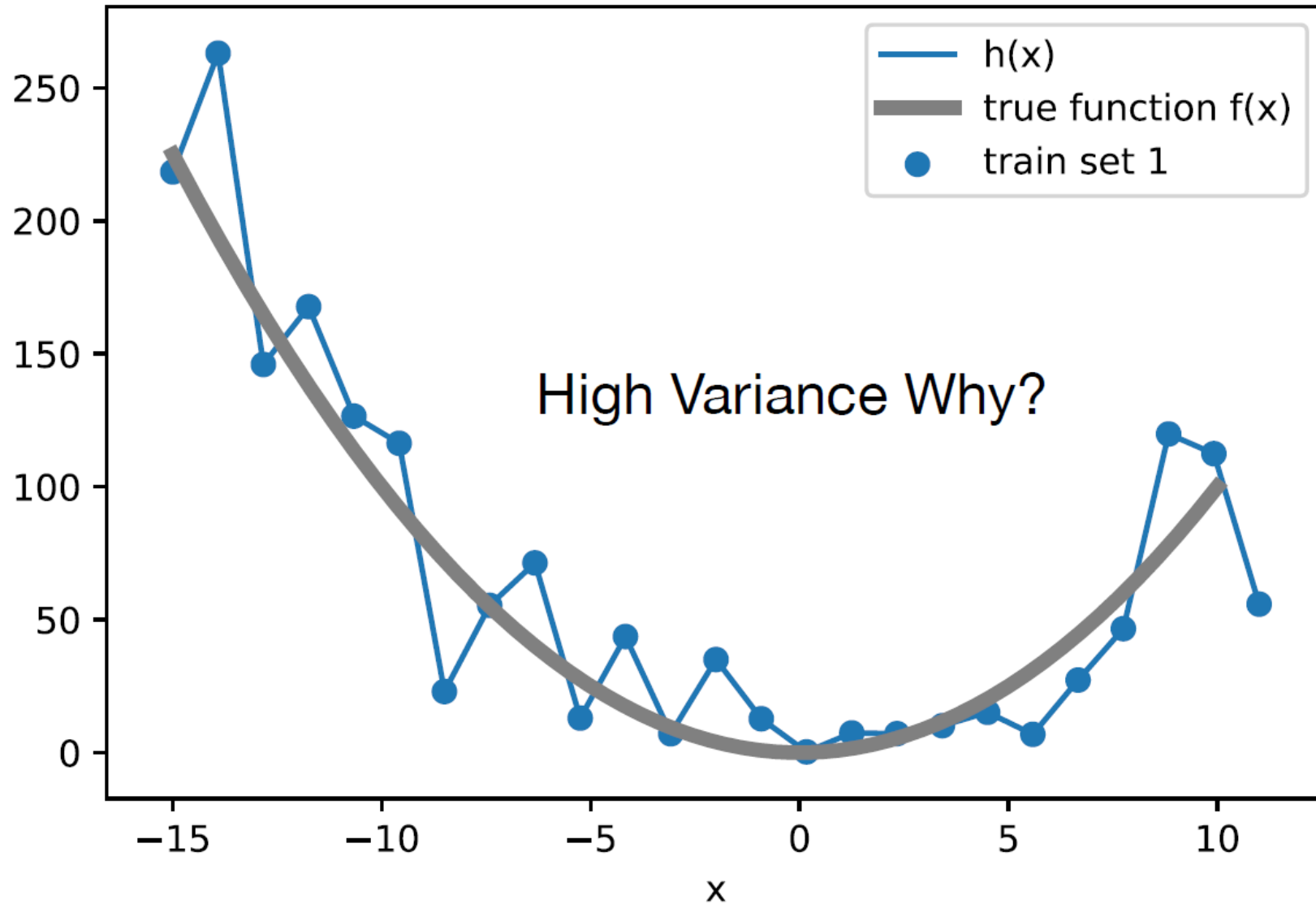


Bias-Variance 직관



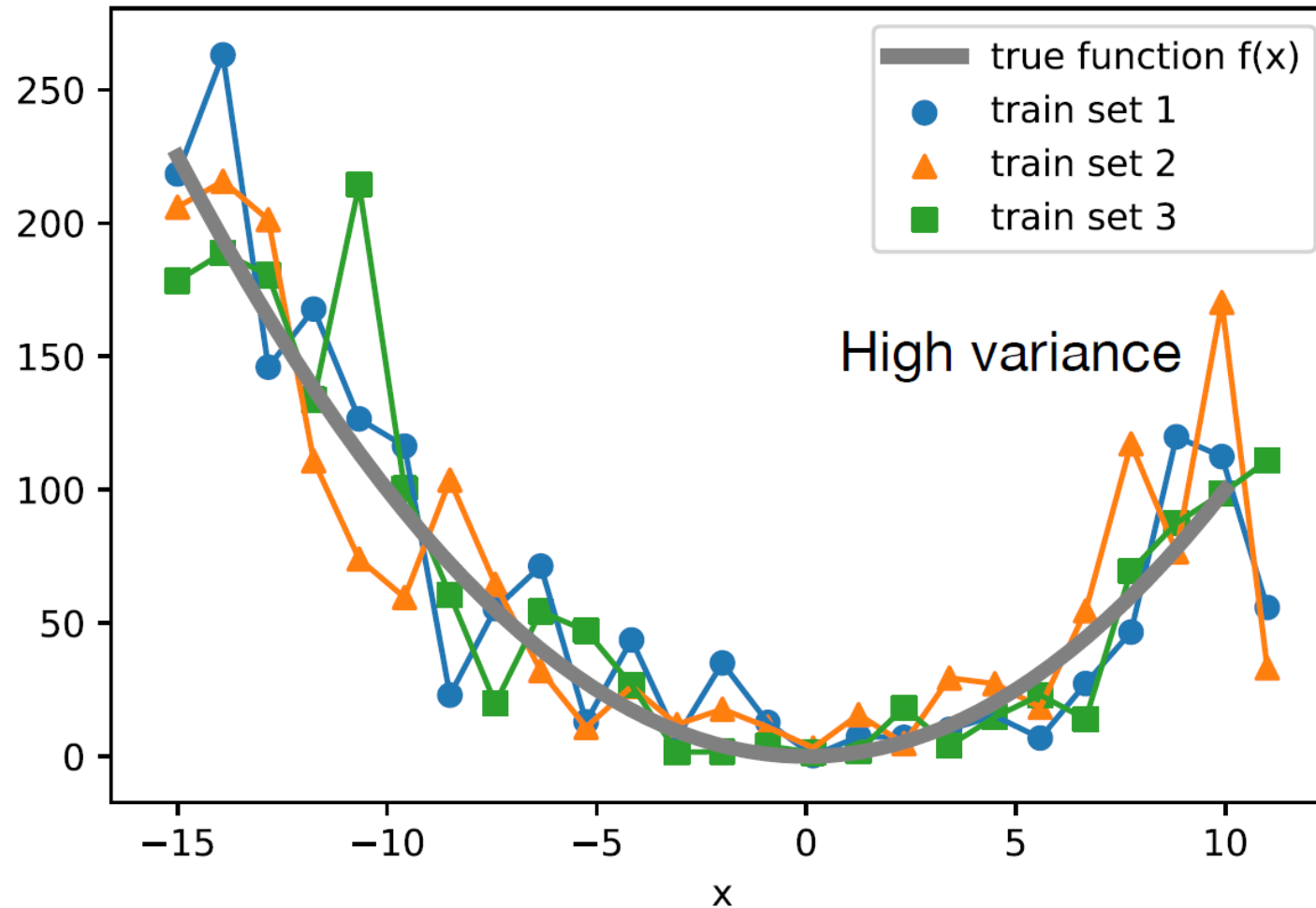
✓ ↓
B ↑

Bias-Variance 직관



Bias-Variance 직관

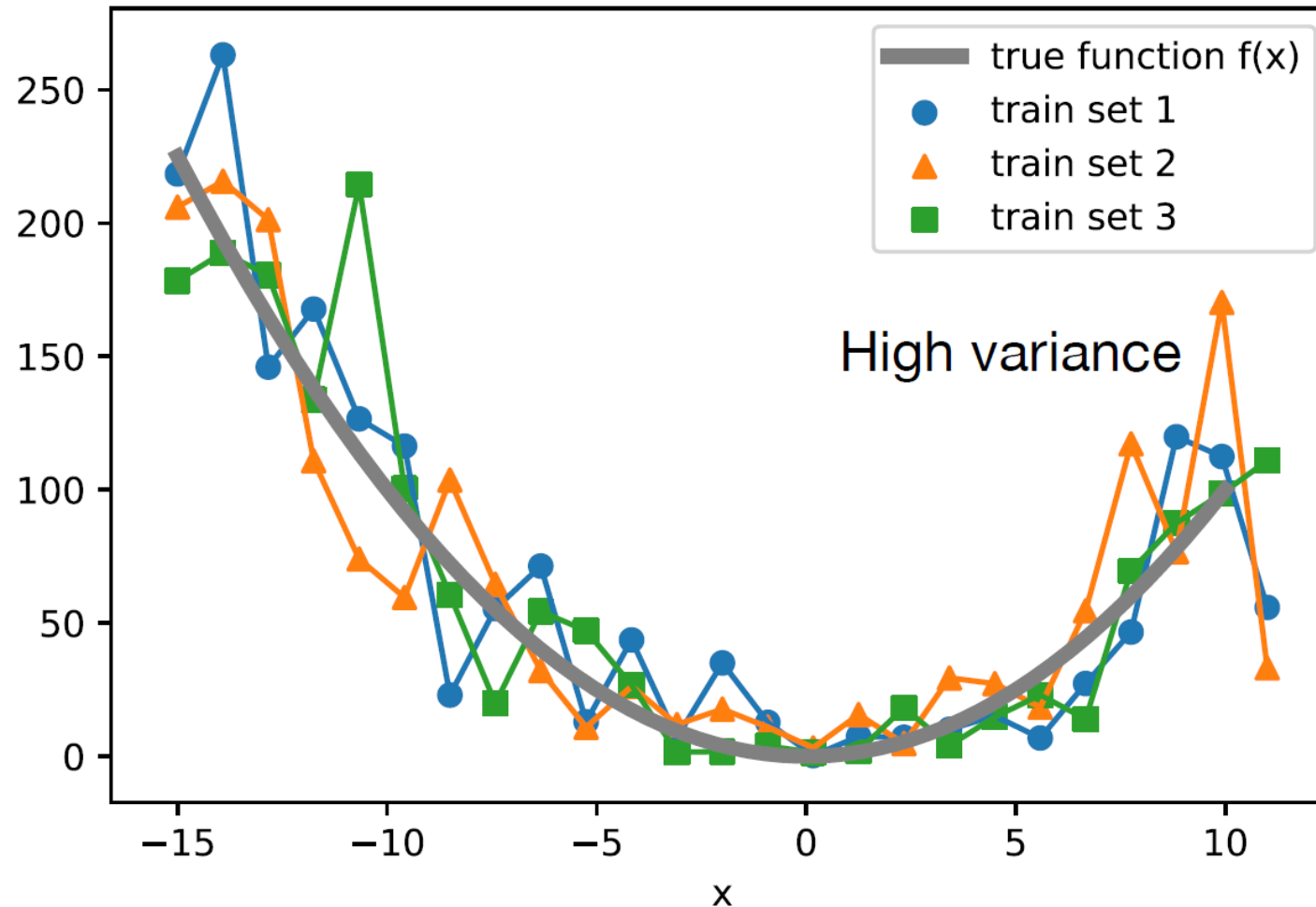
suppose we have multiple training sets



Bias-Variance 직관

평균의 취하면?

suppose we have multiple training sets



Bias-Variance 정의

◆ A point estimator $\hat{\theta}$ of some parameter or function θ

Bias-Variance of the Squared Error

$$\text{Bias}[\hat{\theta}] = \underbrace{E[\hat{\theta}]}_{\text{추정치}} - \underbrace{\theta}_{\text{정답}}$$

$$\text{Var}[\hat{\theta}] = E[\hat{\theta}^2] - (E[\hat{\theta}])^2$$

$$\text{Var}[\hat{\theta}] = E \left[(E[\hat{\theta}] - \hat{\theta})^2 \right]$$

"ML Notation" for Squared Error Loss

$$y = f(x) \text{ target} \leftarrow \text{For simplicity, we ignore the noise term}$$

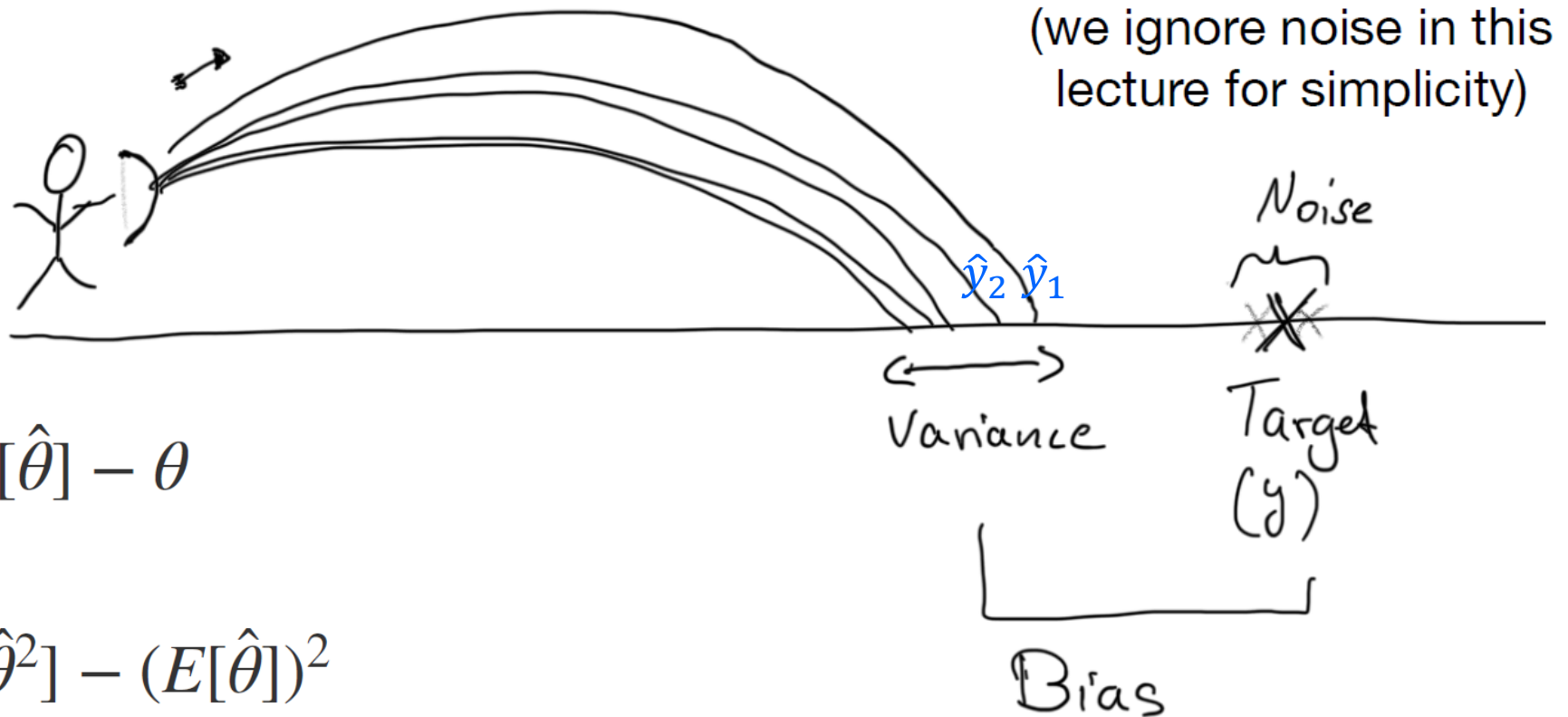
$$\hat{y} = \hat{f}(x) = h(x) \text{ prediction}$$

$$S = (y - \hat{y})^2 \text{ squared error}$$

- Expectation: 확률 변수의 평균값 또는 평균적인 결과를 나타내는 개념
- The **expectation** is over the training data, i.e, the average estimator from different training samples

Bias-Variance 정의

Intuition

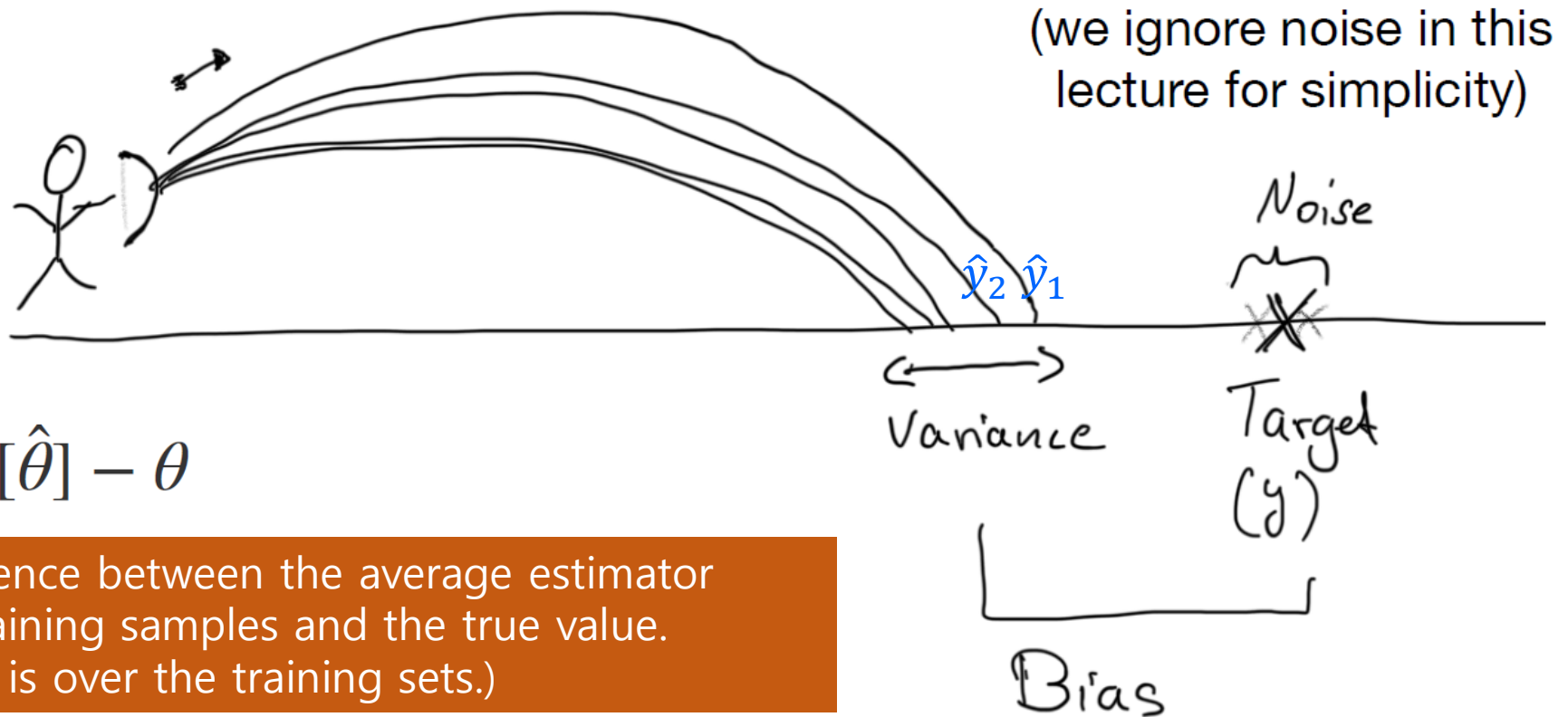


$$\text{Bias}[\hat{\theta}] = E[\hat{\theta}] - \theta$$

$$\text{Var}[\hat{\theta}] = E[\hat{\theta}^2] - (E[\hat{\theta}])^2$$

Bias-Variance 정의

Intuition

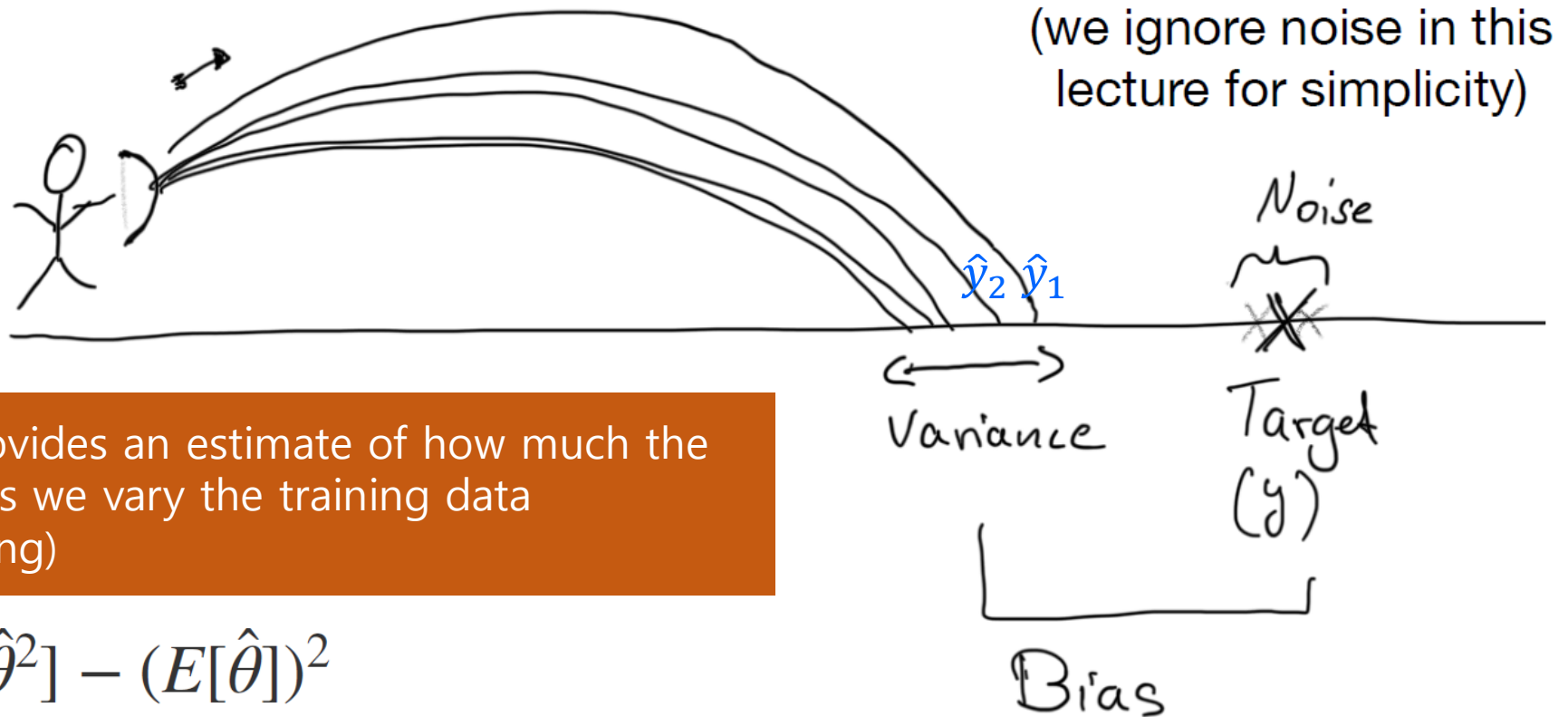


$$\text{Bias}[\hat{\theta}] = E[\hat{\theta}] - \theta$$

Bias is the difference between the average estimator from different training samples and the true value. (The expectation is over the training sets.)

Bias-Variance 정의

Intuition



The **variance** provides an estimate of how much the estimate varies as we vary the training data (e.g. by resampling)

$$\text{Var}[\hat{\theta}] = E[\hat{\theta}^2] - (E[\hat{\theta}])^2$$

Bias-Variance Decomposition

$$\text{Loss} = \text{Bias} + \text{Variance} + \text{Noise}$$

Bias-Variance of the Squared Error

$$\text{Bias}[\hat{\theta}] = E[\hat{\theta}] - \theta$$

$$y = f(x) \text{ target}$$

$$\text{Var}[\hat{\theta}] = E[\hat{\theta}^2] - (E[\hat{\theta}])^2$$

$$\hat{y} = \hat{f}(x) = h(x) \text{ prediction}$$

$$\text{Var}[\hat{\theta}] = E \left[(E[\hat{\theta}] - \hat{\theta})^2 \right]$$

$$S = (y - \hat{y})^2 \text{ squared error}$$

$$\begin{aligned} (y - \hat{y})^2 &= (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2 \\ &= (y - E[\hat{y}])^2 + (E[\hat{y}] - \hat{y})^2 + 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y}) \end{aligned}$$

Bias-Variance of the Squared Error

$$\text{Bias}[\hat{\theta}] = E[\hat{\theta}] - \theta$$

$$y = f(x) \text{ target}$$

$$\text{Var}[\hat{\theta}] = E[\hat{\theta}^2] - (E[\hat{\theta}])^2$$

$$\hat{y} = \hat{f}(x) = h(x) \text{ prediction}$$

$$\text{Var}[\hat{\theta}] = E \left[(E[\hat{\theta}] - \hat{\theta})^2 \right]$$

$$S = (y - \hat{y})^2 \text{ squared error}$$

$$\begin{aligned} (y - \hat{y})^2 &= (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2 \\ &= (y - E[\hat{y}])^2 + (E[\hat{y}] - \hat{y})^2 + 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y}) \end{aligned}$$

$$E[S] = E \left[(y - \hat{y})^2 \right]$$

Bias-Variance of the Squared Error

$$\text{Bias}[\hat{\theta}] = E[\hat{\theta}] - \theta$$

$y = f(x)$ target

$$\text{Var}[\hat{\theta}] = E[\hat{\theta}^2] - (E[\hat{\theta}])^2$$

$\hat{y} = \hat{f}(x) = h(x)$ prediction

$$\text{Var}[\hat{\theta}] = E \left[(E[\hat{\theta}] - \hat{\theta})^2 \right]$$

$S = (y - \hat{y})^2$ squared error

$$\begin{aligned} (y - \hat{y})^2 &= (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2 \\ &= (y - E[\hat{y}])^2 + (E[\hat{y}] - \hat{y})^2 + 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y}) \end{aligned} \quad E[\] = 0$$

$$E[S] = E \left[(y - \hat{y})^2 \right] = (y - E[\hat{y}])^2 + E \left[(E[\hat{y}] - \hat{y})^2 \right]$$

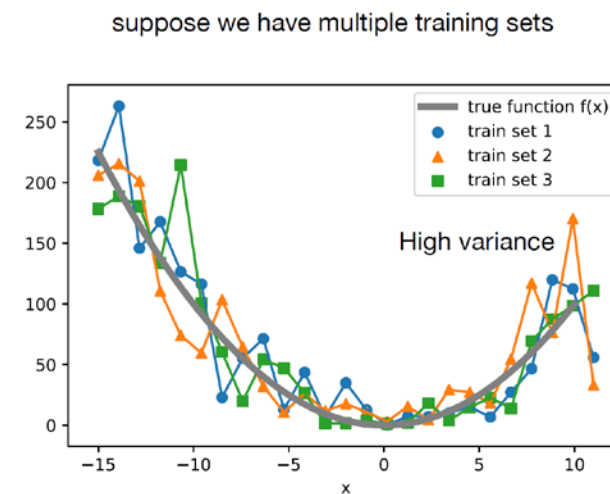
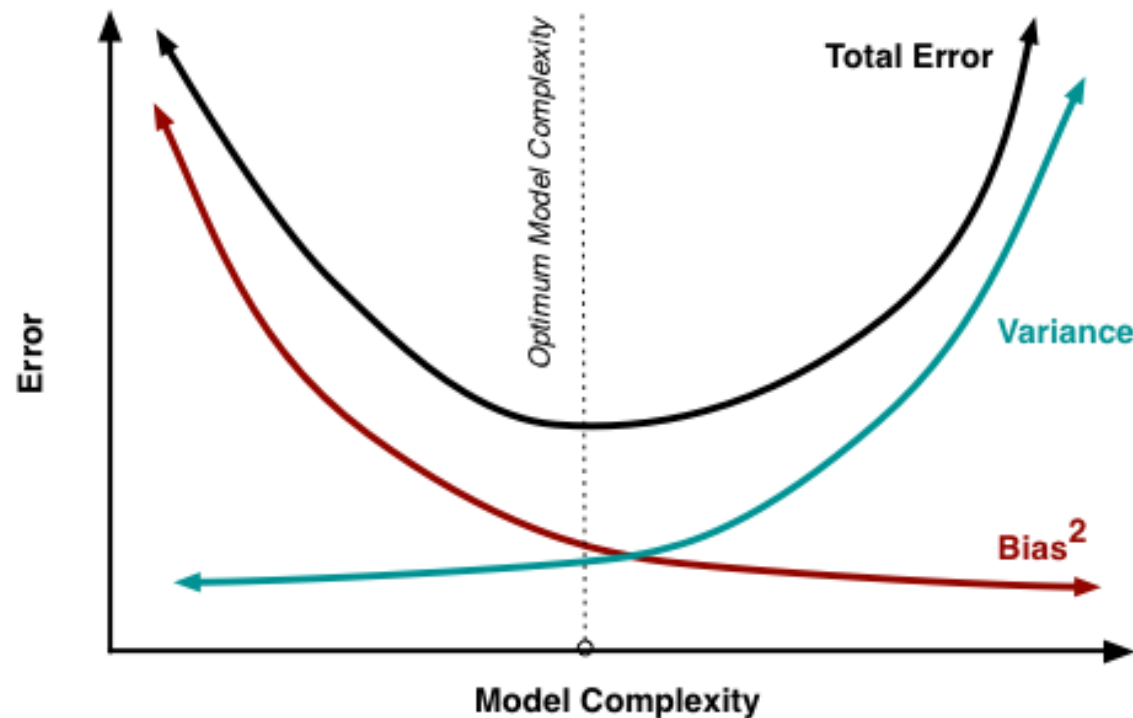
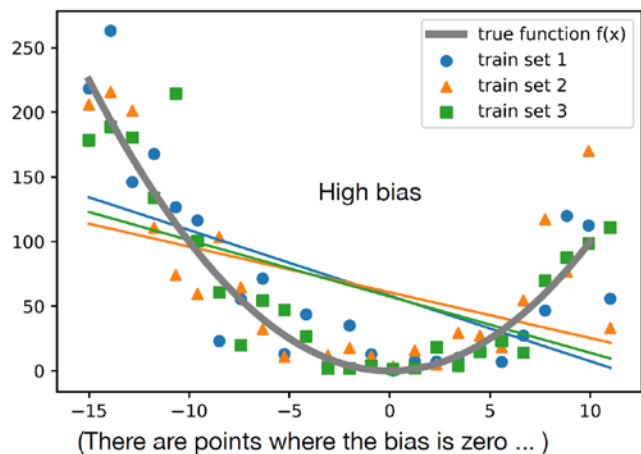
$$= \text{Bias}^2 + \text{Var}$$

Bias-Variance of the Squared Error

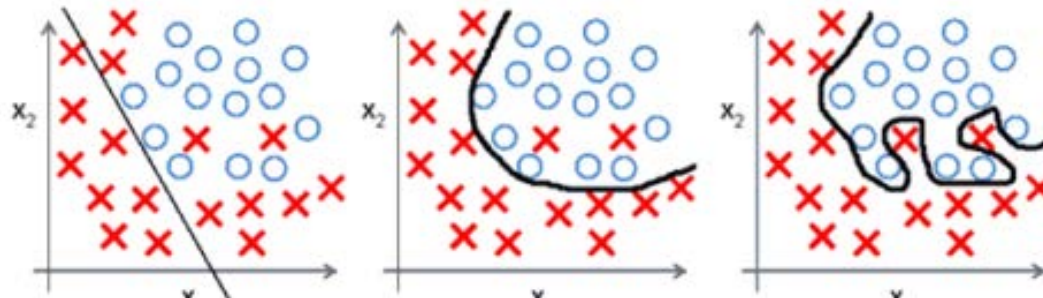
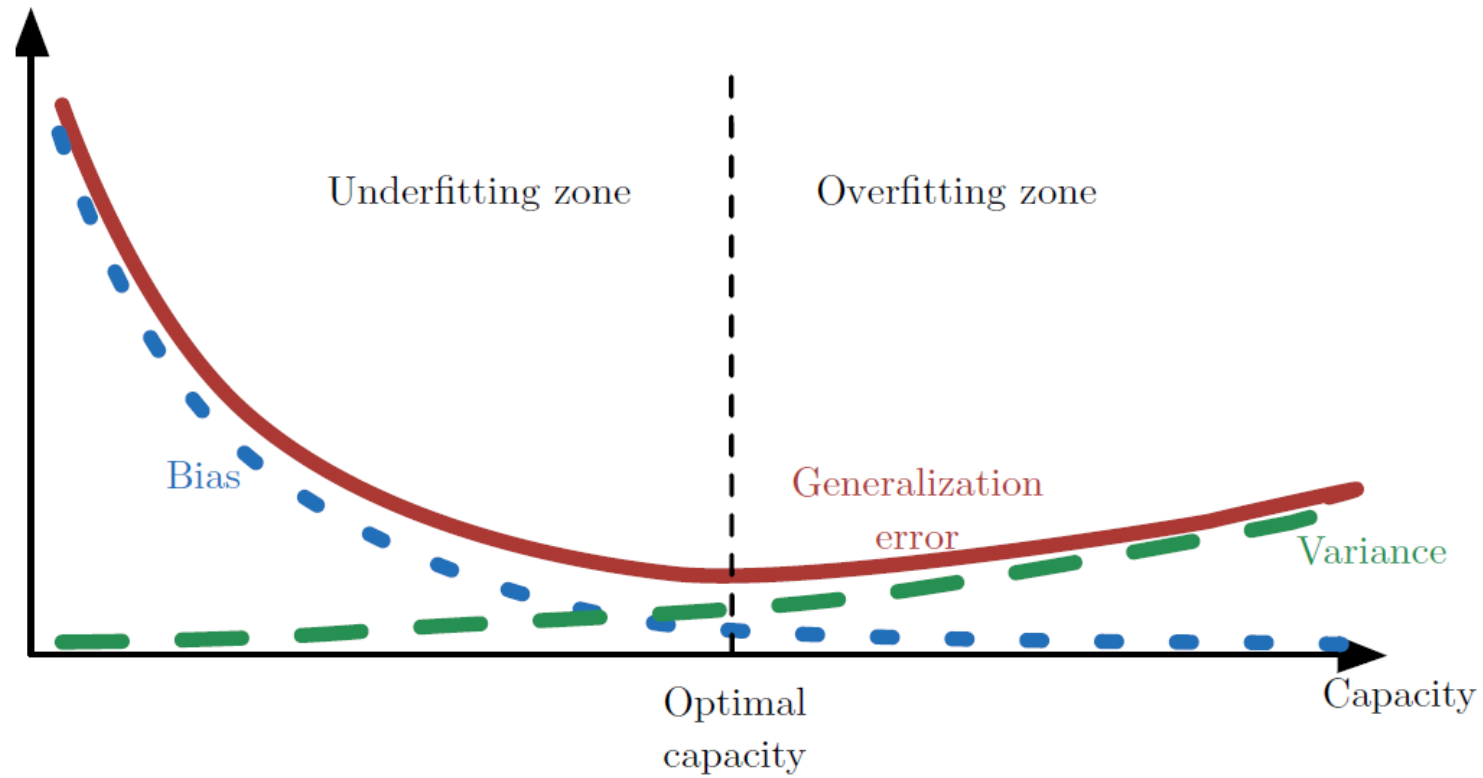
$$\begin{aligned} E[2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})] &= 2E[(y - E[\hat{y}])(E[\hat{y}] - \hat{y})] \\ &= 2(y - E[\hat{y}])E[(E[\hat{y}] - \hat{y})] \\ &= 2(y - E[\hat{y}])(E[E[\hat{y}]] - E[\hat{y}]) \\ &= 2(y - E[\hat{y}])(E[\hat{y}] - E[\hat{y}]) \\ &= 0 \end{aligned}$$

Bias-Variance Tradeoff

- ◆ 모델의 복잡도 관점에서 봤을 때 분산과 편향이 트레이드 오프(trade-off) 관계



Generalization Error



3. Regularization by Weight Penalty

Regularization (규제)

◆ 『Deep Learning』 책의 규제 정의

- "...any modification we make to a learning algorithm that is intended to *reduce its generalization error* ..."
- (일반화 오류를 줄이려는 의도를 가지고 학습 알고리즘을 수정하는 방법 모두)

◆ 규제는 오래 전부터 수학과 통계학에서 연구해온 주제

- **모델 용량에 비해 데이터가 부족한 경우**의 불량 문제를 ill-posed problem 푸는 데 사용
- 현대 기계학습도 규제를 널리 사용

◆ 명시적 규제와 암시적 규제

- **명시적 규제**: 가중치 감쇠나 드롭아웃처럼 **목적함수나 신경망 구조를 직접 수정**하는 방식
- **암시적 규제**: 조기 멈춤, 데이터 증대, 잡음 추가, 앙상블처럼 **간접적으로 영향을 미치는** 방식

Regularization by Weight Penalty (가중치 감쇠)

◆ Regularized Cost Function

$$\underbrace{J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \lambda \underbrace{R(\Theta)}_{\text{규제 항}}$$

- **규제항**은 훈련집합과 무관하며, 데이터 생성 과정에 내재한 **사전 지식**에 해당

Regularization by Weight Penalty (가중치 감쇠)

$$\underbrace{J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \lambda \underbrace{R(\Theta)}_{\text{규제 항}}$$

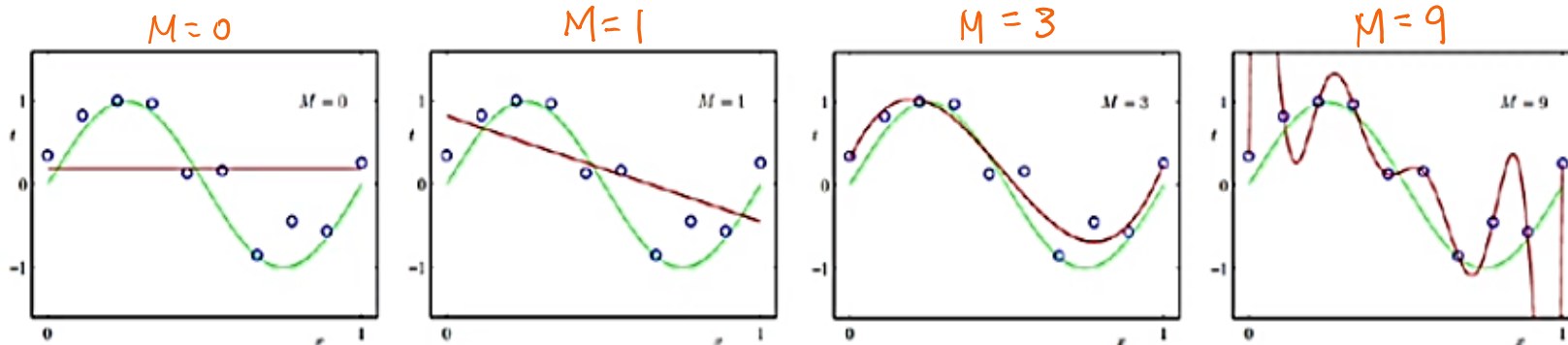
◆ 규제항 $R(\Theta)$ 로 무엇을 사용할 것인가? → **가중치 감쇠 (가중치 벌칙)**

● 큰 가중치(Θ)에 벌칙을 가해 작은 가중치를 유지 → 최종해를 원점 가까이 당기는 효과

● L2 norm 사용: $R(\Theta) = \|\Theta\|_2^2$

● L1 norm 사용: $R(\Theta) = \|\Theta\|_1$

● 가중치 감쇠는 모델의 구조적 용량을 충분히 크게 하고 모델의 수치적 용량을 제한하는 규제 기법



	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

Regularization – L2 Norm,

◆ Regularized Cost & Gradient

$$\underbrace{J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \underbrace{\lambda \|\Theta\|_2^2}_{\text{규제 항}}$$

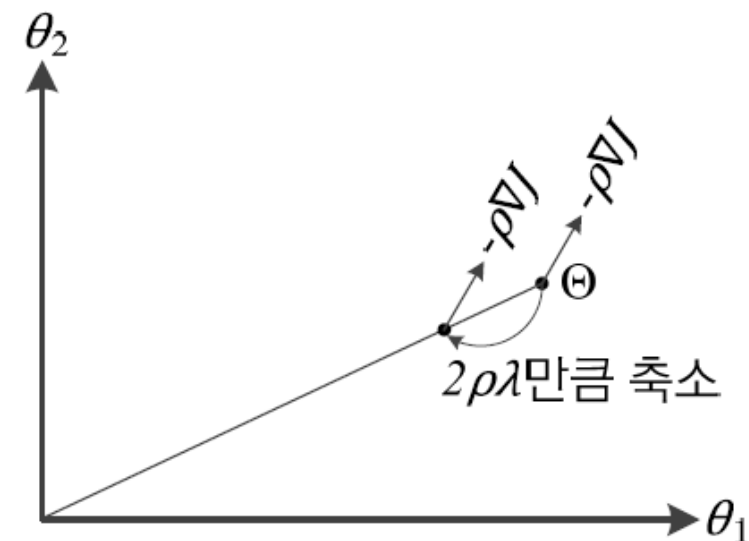
$$\nabla J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y}) = \nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + 2\lambda\Theta$$

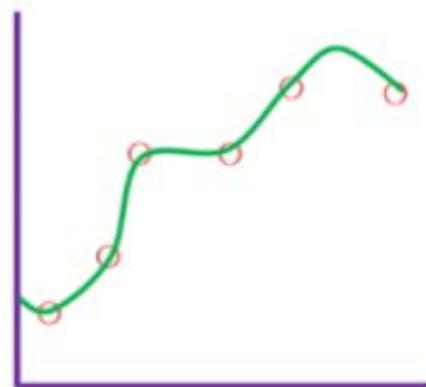
◆ Parameter Update

$$\begin{aligned}\Theta &= \Theta - \rho \nabla J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y}) \\ &= \Theta - \rho (\nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + 2\lambda\Theta) \\ &= (1 - 2\rho\lambda)\Theta - \rho \nabla J(\Theta; \mathbb{X}, \mathbb{Y})\end{aligned}$$

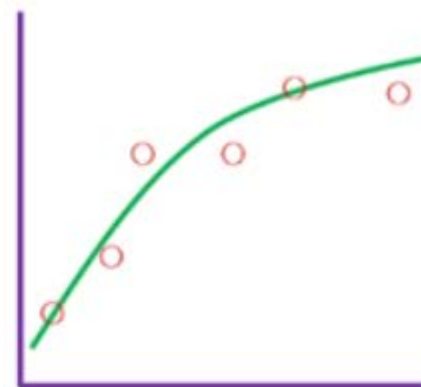
- L2 규제는 Θ 를 $2\rho\lambda$ 의 비율로 줄인 후 업데이트 하는 셈
 - 즉, 가중치 감소 정도가 현재 가중치 크기에 비례함 함

Weight decay (가중치 감소)





$$\beta_0 + \beta_1 x + \beta_2 x^2 + \cancel{\beta_3 x^3} + \cancel{\beta_4 x^4}$$



$$\beta_0 + \beta_1 x + \beta_2 x^2$$

$$\underbrace{J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \lambda \underbrace{\|\Theta\|_2^2}_{\text{규제 항}}$$

$$\min_{\beta} \sum_{i=1} (y_i - \hat{y}_i)^2 + 5000\beta_3^2 + 5000\beta_4^2$$



$$\beta_3 \approx 0 \quad \beta_4 \approx 0$$

Regularization – L1 Norm

◆ Regularized Cost & Gradient

$$\underbrace{J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \underbrace{\lambda \|\Theta\|_1}_{\text{규제 항}}$$

$$\nabla J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y}) = \nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + \lambda \text{sign}(\Theta)$$

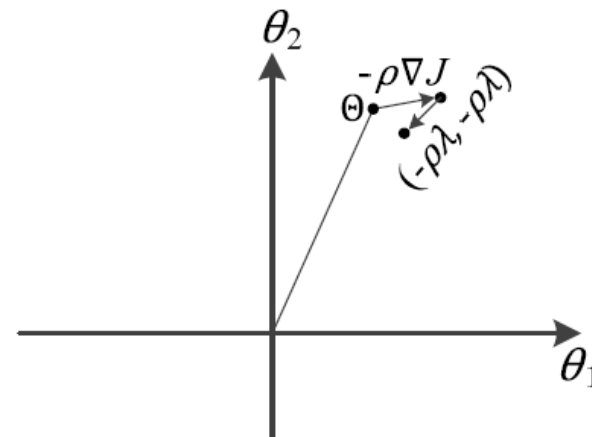
$\text{sign}(\Theta)$: Θ 의 부호 벡터 (1, -1)

◆ Parameter Update

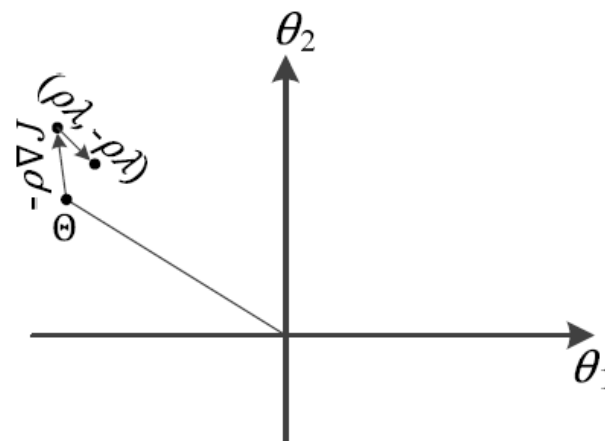
$$\begin{aligned} \Theta &= \Theta - \rho \nabla J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y}) \\ &= \Theta - \rho (\nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + \lambda \text{sign}(\Theta)) \\ &= \Theta - \rho \nabla J(\Theta; \mathbb{X}, \mathbb{Y}) - \rho \lambda \text{sign}(\Theta) \end{aligned}$$

원점 방향으로

- L1 규제는 Θ 를 $\rho\lambda$ (고정값)만큼 줄인 후 업데이트 하는 셈
- L1 규제의 희소성(Sparse) 효과: 0이 되는 가중치가 많이 발생
- 선형 회귀에 적용하면 특징 선택 효과



(a) $\text{sign}(\Theta) = (1, 1)^T$ 인 경우

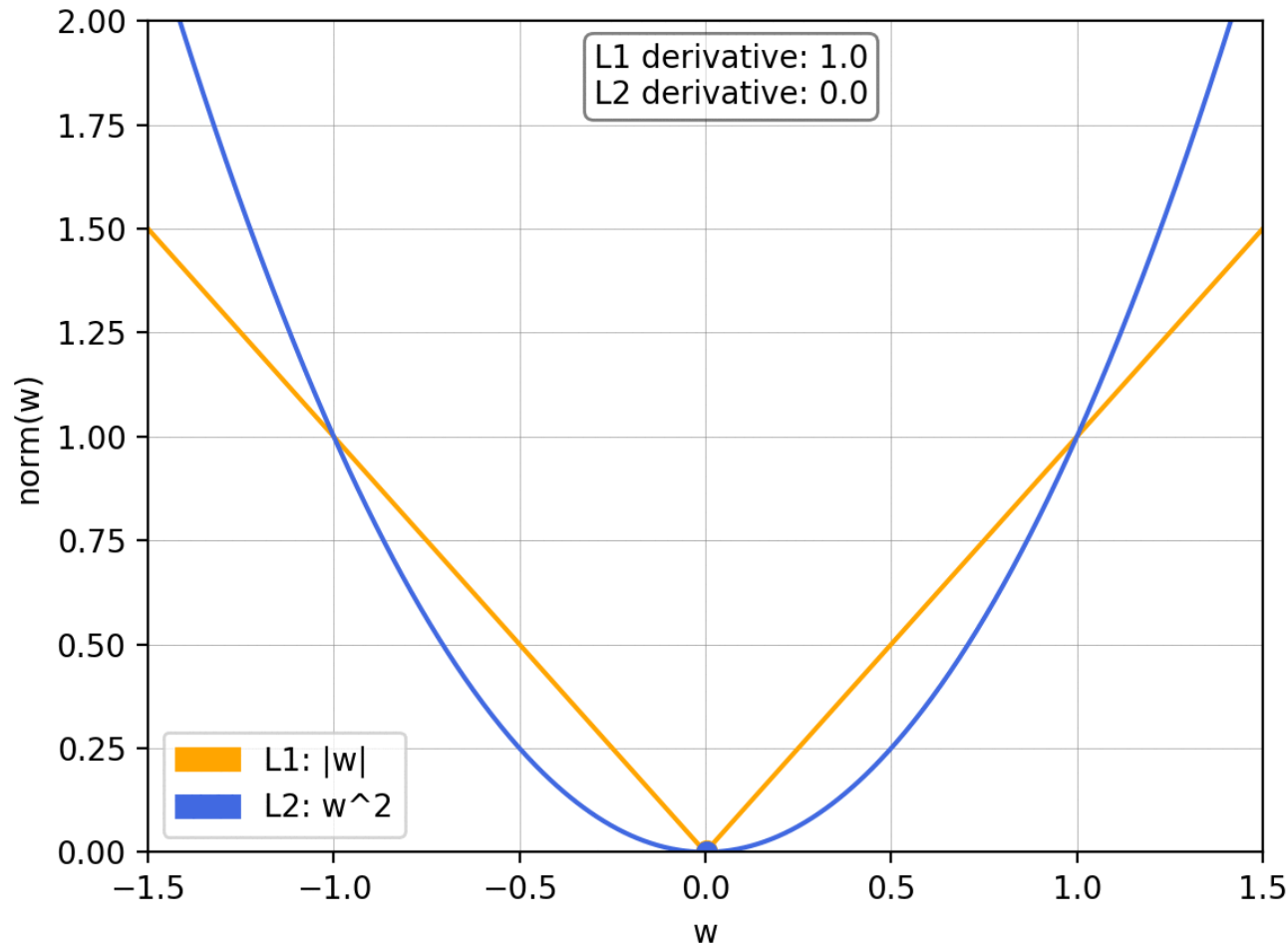


(b) $\text{sign}(\Theta) = (-1, 1)^T$ 인 경우

Regularization – L1 norm vs. L2 norm

L1 norm이 0이 되는 가중치가 많이 발생하는 이유 :

1. 업데이트 속도차이 L1은 고정적인 비율로 갱신 L2는 가중치 크기에 비례
- 2.



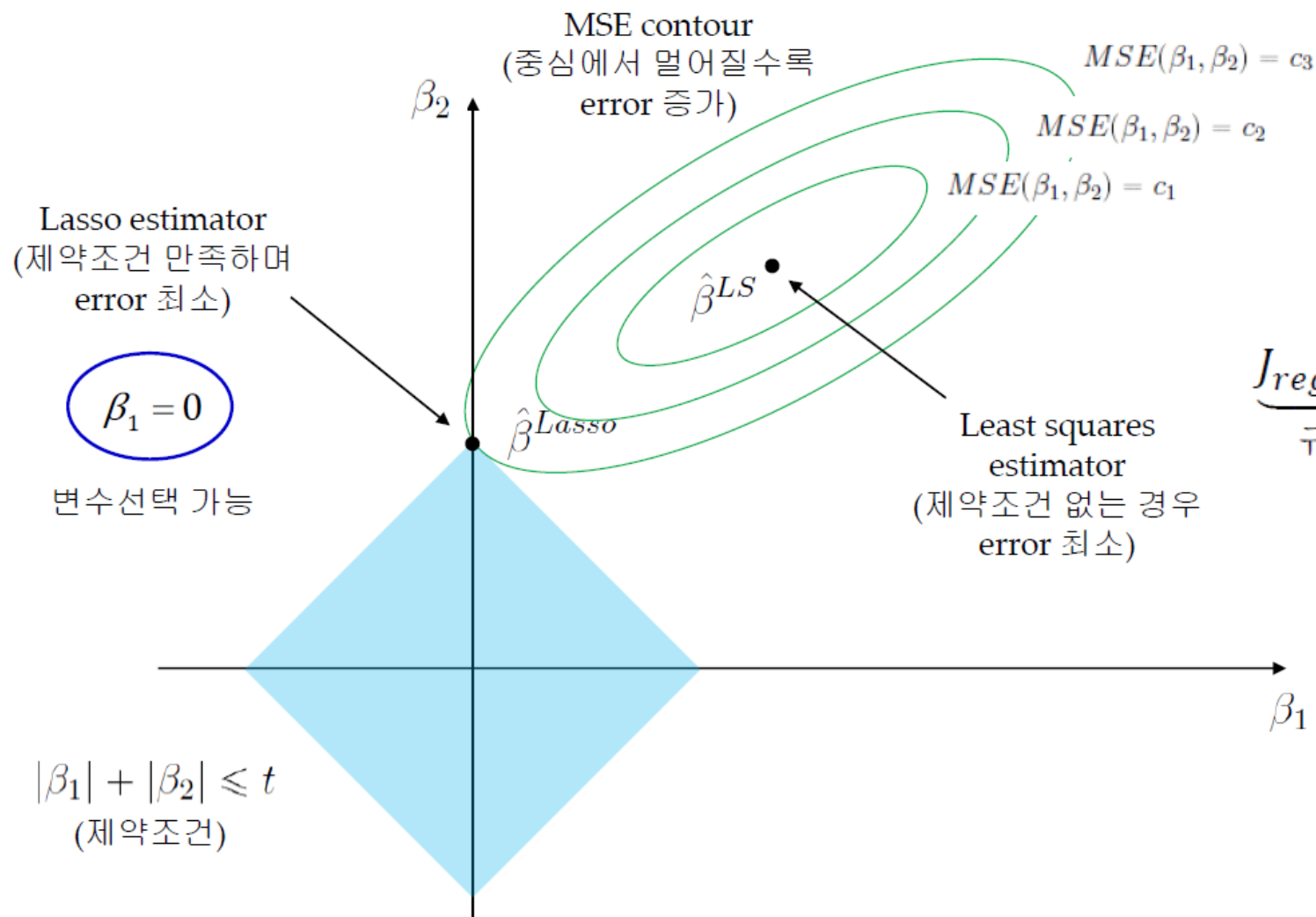
◆ L1 norm update

$$\begin{aligned}\Theta &= \Theta - \rho \nabla J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y}) \\ &= \Theta - \rho (\nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + \lambda \text{sign}(\Theta)) \\ &= \Theta - \rho \nabla J(\Theta; \mathbb{X}, \mathbb{Y}) - \rho \lambda \text{sign}(\Theta)\end{aligned}$$

◆ L2 norm update

$$\begin{aligned}\Theta &= \Theta - \rho \nabla J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y}) \\ &= \Theta - \rho (\nabla J(\Theta; \mathbb{X}, \mathbb{Y}) + 2\lambda \Theta) \\ &= (1 - 2\rho\lambda) \Theta - \rho \nabla J(\Theta; \mathbb{X}, \mathbb{Y})\end{aligned}$$

Lasso Regression의 기하학적 이해

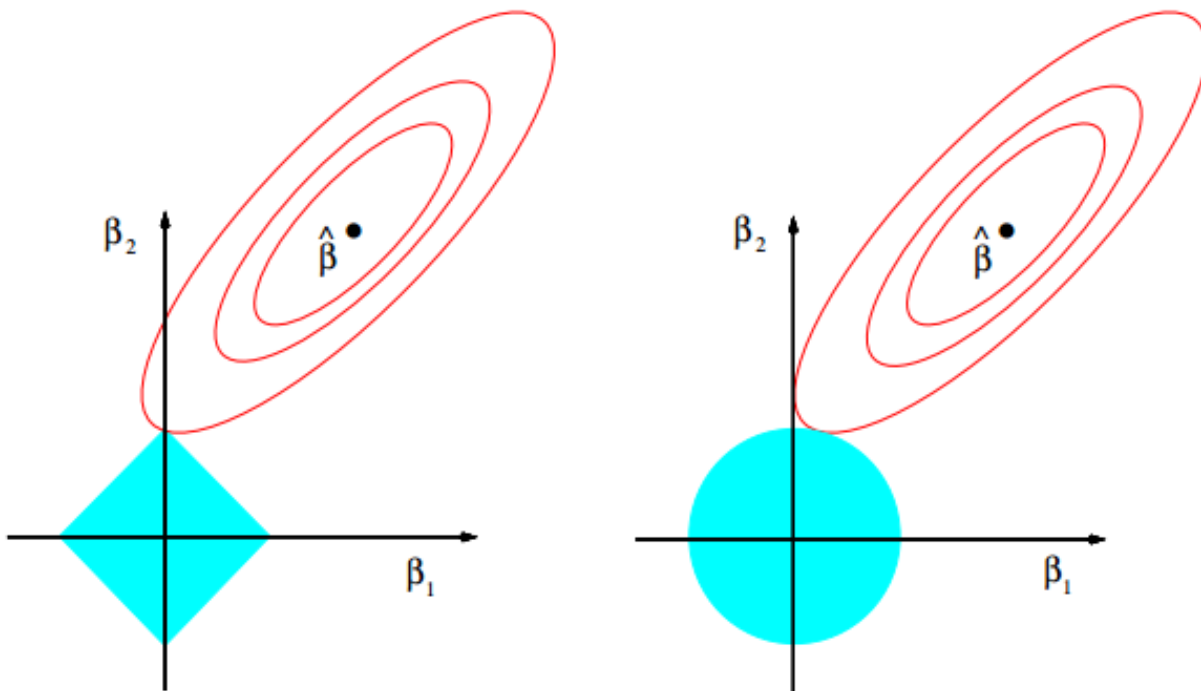


$$\underbrace{J_{regularized}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \lambda \underbrace{\|\Theta\|_1}_{\text{규제 항}}$$

Regularization – L1 norm vs. L2 norm

◆ Lasso (L1) vs. Ridge (L2) Regression

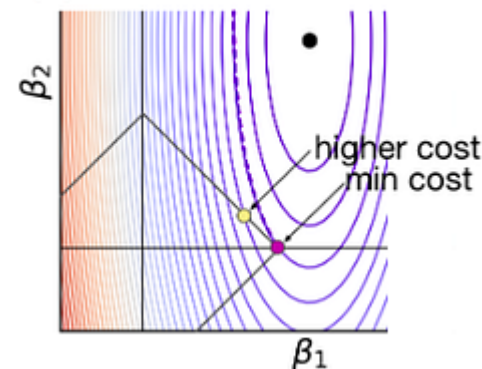
$$\underbrace{J_{\text{regularized}}(\boldsymbol{\Theta}; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\boldsymbol{\Theta}; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \lambda \underbrace{R(\boldsymbol{\Theta})}_{\text{규제 항}}$$



$$|\beta_1| + |\beta_2| < t$$

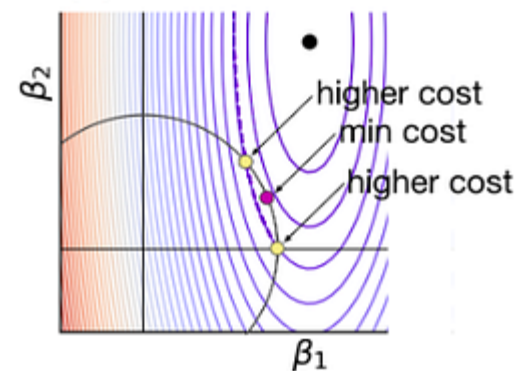
$$\beta_1^2 + \beta_2^2 < t$$

(a) L1 Constraint Diamond



0이 되는 가중치가
많이 발생

(b) L2 Constraint Circle

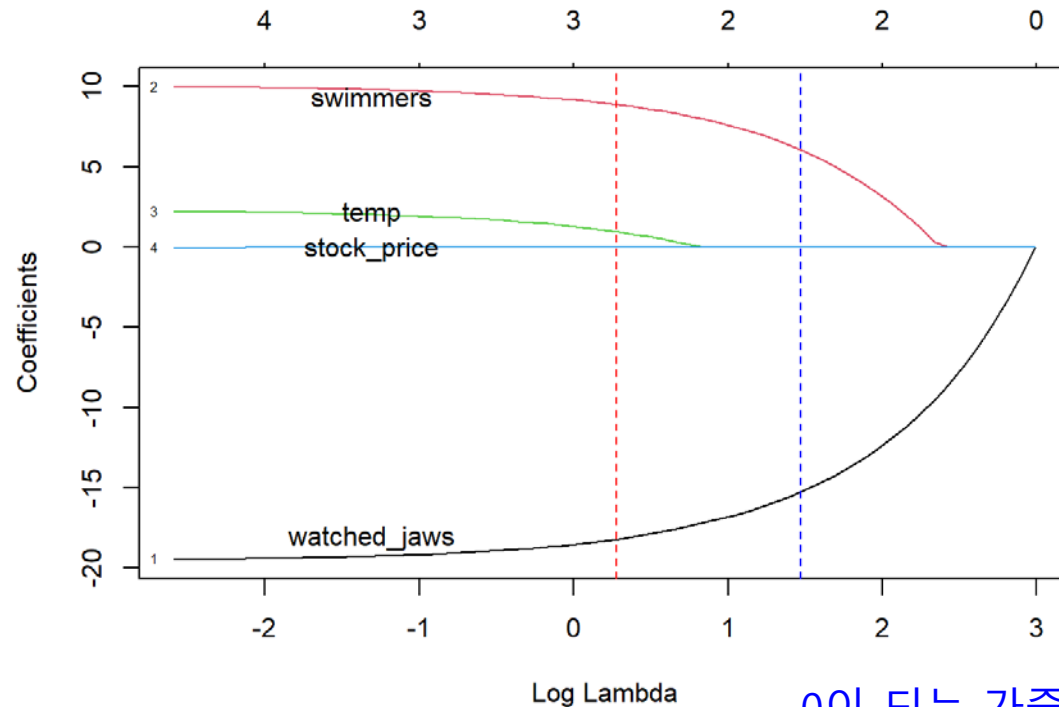


<https://medium.com/@mukulranjan/how-does-lasso-regression-l1-encourage-zero-coefficients-but-not-the-l2-20e4893cba5d>

Regularization – L1 norm vs. L2 norm

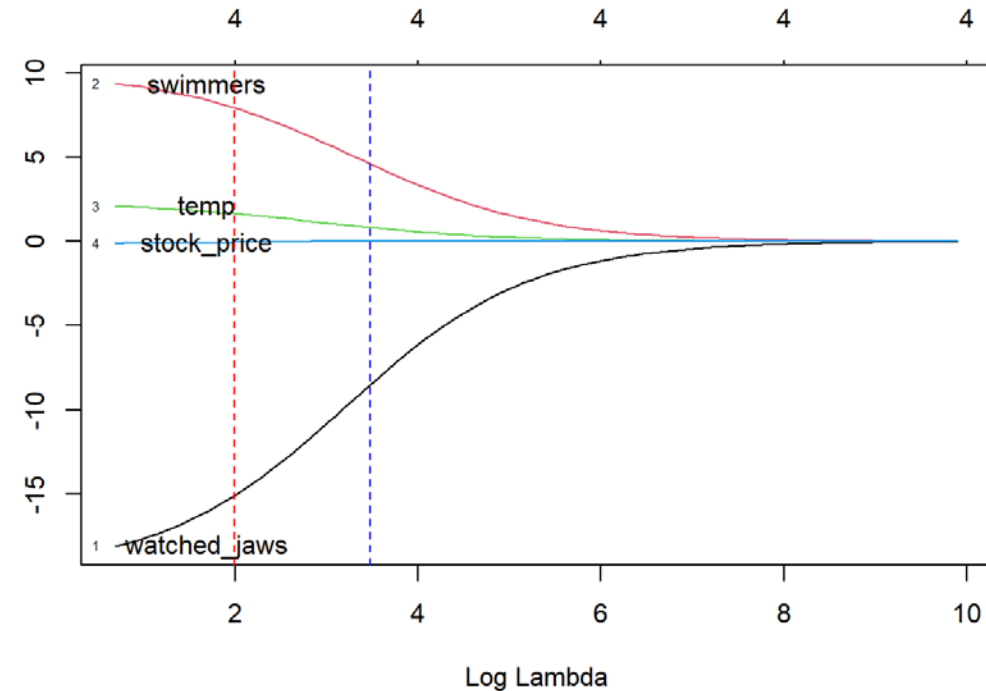
◆ Lasso (L1) vs. Ridge (L2) Regression

$$\underbrace{J_{\text{regularized}}(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{규제를 적용한 목적함수}} = \underbrace{J(\Theta; \mathbb{X}, \mathbb{Y})}_{\text{목적함수}} + \underbrace{\lambda R(\Theta)}_{\text{규제 항}}$$



학습속도가 빠름 정확도는 낮은편
Lasso (L1 norm)

0이 되는 가중치가
많이 발생

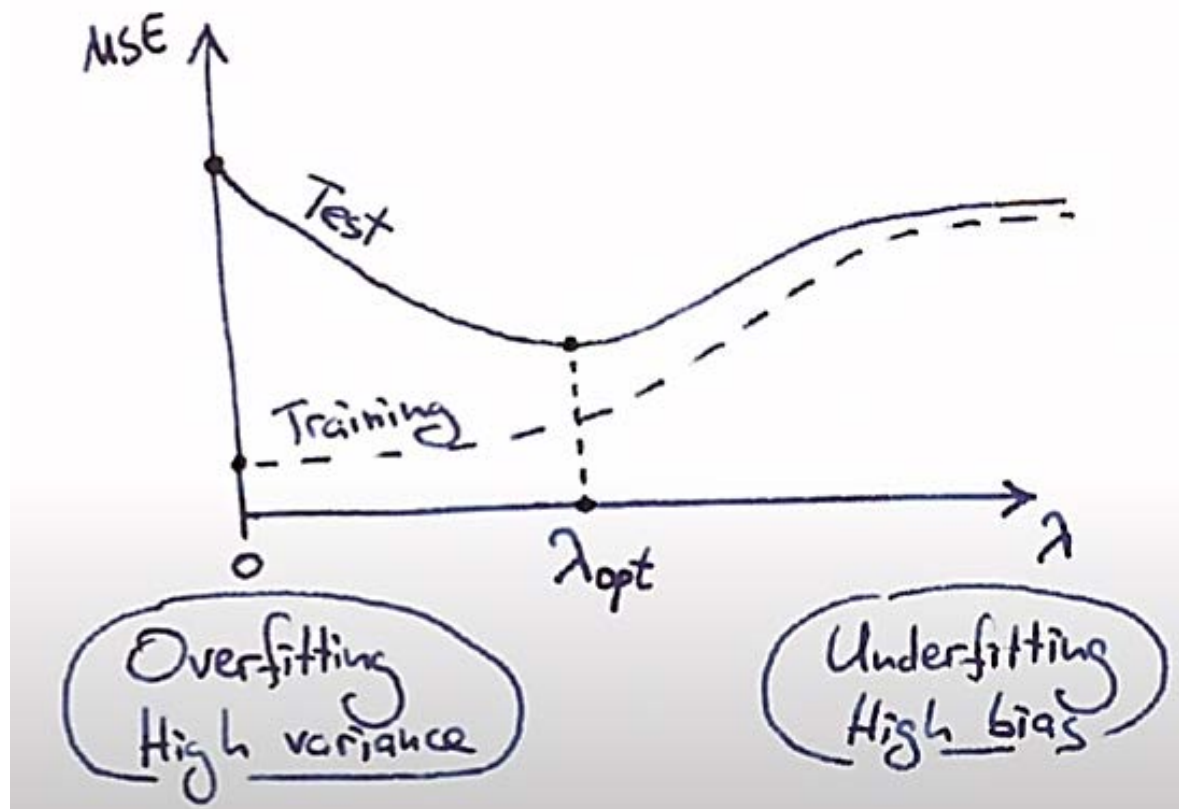


Ridge (L2 norm)

Regularization – Selecting Lambda

◆ Test Error가 가장 작게 되는 λ 가 최적

- 그러나 학습시에는 test set에 접근할 수 없으므로, validation set을 이용하여 최적의 λ 를 선택함



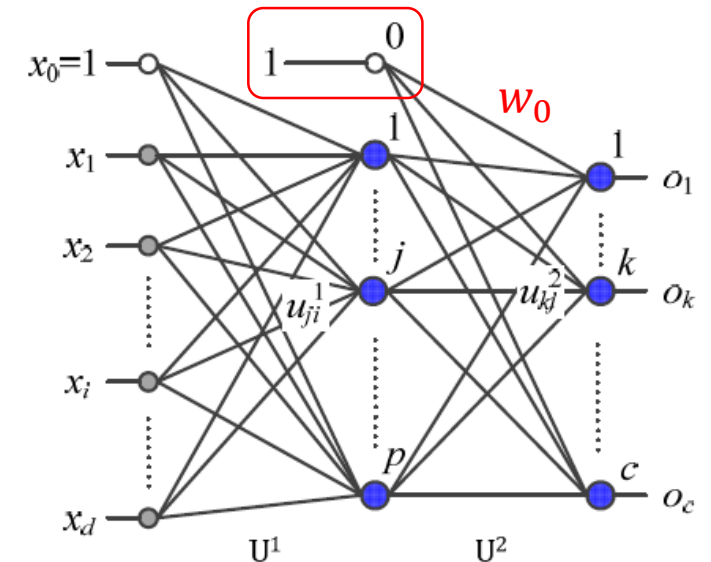
Regularization – Do Not Penalize Bias!

- ◆ For Centered Dataset (when both x and y have zero mean)
 - No problem even if we have zero bias (i.e., $w_0 = 0$).

$$J(\Theta) = \frac{1}{n} \|\mathbf{y} - \mathbf{x}^T \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 = \frac{1}{n} \sum_{i=1}^n \left(y_i - w_0 - \sum_{j=1}^d w_j x_{ij} \right)^2 + \lambda \sum_{j=0}^d w_j^2$$

- ◆ For Non-centered Dataset (the general case)
 - Penalizing bias often leads to bad performance.
 - Thus we need to *exclude the bias (w_0) from the regularization term*:

$$J(\Theta) = \frac{1}{n} \sum_{i=1}^n \left(y_i - w_0 - \sum_{j=1}^d w_j x_{ij} \right)^2 + \lambda \sum_{j=1}^d w_j^2$$



Regularization – Example: Linear Regression

■ 선형 회귀에 적용

- 선형 회귀는 훈련집합 $\mathbb{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, $\mathbb{Y} = \{y_1, y_2, \dots, y_n\}$ 이 주어지면, 식 (5.24)를 풀어 $\mathbf{w} = (w_1, w_2, \dots, w_d)^T$ 를 구하는 문제. 이때 $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$

$$w_1 x_{i1} + w_2 x_{i2} \dots + w_d x_{id} = \mathbf{x}_i^T \mathbf{w} = y_i, \quad i = 1, 2, \dots, n \quad (5.24)$$

- 식 (5.24)를 행렬식으로 바꿔 쓰면,

$$\mathbf{X}\mathbf{w} = \mathbf{y}$$

(5.25)

$$\mathbf{X} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(n)} & x_2^{(n)} & \dots & x_d^{(n)} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

- 가중치 감소를 적용한 목적함수

$$J_{\text{regularized}}(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \|\mathbf{w}\|_2^2 \quad (5.27)$$

Regularization – Example: Linear Regression (cont'd)

- 식 (5.27)을 미분하여 0으로 놓으면,

$$\frac{\partial J_{regularized}}{\partial \mathbf{w}} = \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} + 2\lambda \mathbf{w} = \mathbf{0} \Rightarrow (\mathbf{X}^T \mathbf{X} + 2\lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (5.28)$$

- 식 (5.28)을 정리하면,

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + 2\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (5.29)$$


- 공분산 행렬 $\mathbf{X}^T \mathbf{X}$ 의 대각 요소가 2λ 만큼씩 증가 \rightarrow 역행렬을 곱하므로 가중치를 축소하여 원점으로 당기는 효과 ([그림 5-21])

- 예측 단계에서는.

$$y = \mathbf{x}^T \hat{\mathbf{w}} \quad (5.30)$$

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{\partial}{\partial \mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \quad \text{풀이}$$

$$\begin{aligned} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 &= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \stackrel{\textcircled{1}}{=} (\mathbf{w}^T \mathbf{X}^T - \mathbf{y}^T) (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{y}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y} \\ &\stackrel{\textcircled{2}}{=} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 &= \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}) - \frac{\partial}{\partial \mathbf{w}} (2\mathbf{y}^T \mathbf{X} \mathbf{w}) + \frac{\partial}{\partial \mathbf{w}} (\mathbf{y}^T \mathbf{y}) \\ &\stackrel{\textcircled{3}}{=} 2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y} + 0 \\ &= 2\mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \end{aligned}$$

$$\begin{aligned} \textcircled{1} \quad &(\mathbf{A} - \mathbf{B})^T = \mathbf{A}^T - \mathbf{B}^T \\ &(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad &\mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x} \\ &\mathbf{w}^T \mathbf{X}^T \mathbf{y} = (\mathbf{X}\mathbf{w})^T \mathbf{y} = \mathbf{y}^T (\mathbf{X}\mathbf{w}) \end{aligned}$$

$$\textcircled{3} \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{A} \mathbf{x}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} \quad \begin{array}{l} \mathbf{X}^T \mathbf{X} : \text{symmetric} \\ \mathbf{X}^T \mathbf{X} = (\mathbf{X}^T \mathbf{X})^T \end{array}$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}) &= (\mathbf{X}^T \mathbf{X} + (\mathbf{X}^T \mathbf{X})^T) \mathbf{w} \\ &= (\mathbf{X}^T \mathbf{X} + \mathbf{X}^T \mathbf{X}) \mathbf{w} \\ &= 2\mathbf{X}^T \mathbf{X} \mathbf{w} \end{aligned}$$

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{y}^T \mathbf{x}) = \mathbf{y}$$

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{y}^T \mathbf{X} \mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} ((\mathbf{X}^T \mathbf{y})^T \mathbf{w}) = \mathbf{X}^T \mathbf{y}$$

$\mathbf{y}^T \mathbf{X}$: a row vector

Regularization – Example: Linear Regression (cont'd)

예제 5-1 리지 회귀

훈련집합 $\mathbb{X} = \{\mathbf{x}_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} 2 \\ 3 \end{pmatrix}, \mathbf{x}_3 = \begin{pmatrix} 3 \\ 3 \end{pmatrix}\}$, $\mathbb{Y} = \{y_1 = 3.0, y_2 = 7.0, y_3 = 8.8\}$ 이 주어졌다고 가정하자. 특징 벡터가 2차원이므로 $d=2$ 이고 샘플이 3개이므로 $n=3$ 이다. 훈련집합으로 설계행렬 \mathbf{X} 와 레이블 행렬 \mathbf{y} 를 다음과 같이 쓸 수 있다.

$$\mathbf{X} = \begin{pmatrix} 1 & 1 \\ 2 & 3 \\ 3 & 3 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} 3.0 \\ 7.0 \\ 8.8 \end{pmatrix}$$

이 값들을 식 (5.29)에 대입하여 다음과 같이 $\hat{\mathbf{w}}$ 을 구할 수 있다. 이때 $\lambda = 0.25$ 라 가정하자.

$$\hat{\mathbf{w}} = \left(\begin{pmatrix} 1 & 2 & 3 \\ 1 & 3 & 3 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 2 & 3 \\ 3 & 3 \end{pmatrix} + \begin{pmatrix} 0.5 & 0 \\ 0 & 0.5 \end{pmatrix} \right)^{-1} \begin{pmatrix} 1 & 2 & 3 \\ 1 & 3 & 3 \end{pmatrix} \begin{pmatrix} 3.0 \\ 7.0 \\ 8.8 \end{pmatrix} = \begin{pmatrix} 1.4916 \\ 1.3607 \end{pmatrix}$$

따라서 하이퍼 평면은 $y = 1.4916x_1 + 1.3607x_2$ 이다. 새로운 샘플로 $\mathbf{x} = (5 \ 4)^T$ 가 입력되면 식 (5.30)을 이용하여 12.9009를 예측한다.

감사합니다.