

정규화 모델 (Regularization Model)

2025 Spring

머신러닝1

이 두 호

- 정규화 모델 1
 - 정규화 모델 배경
 - 정규화 모델 개념
 - Ridge Regression (능형회귀)
- 정규화 모델 2
 - Lasso
 - Elastic Net
 - 여러가지 정규화 모델

What is a good model?

현재 데이터(training data)를 잘 설명하는 모델

⇒ **Explanatory modeling**

미래 데이터(testing data)에 대한 예측 성능이 좋은 모델

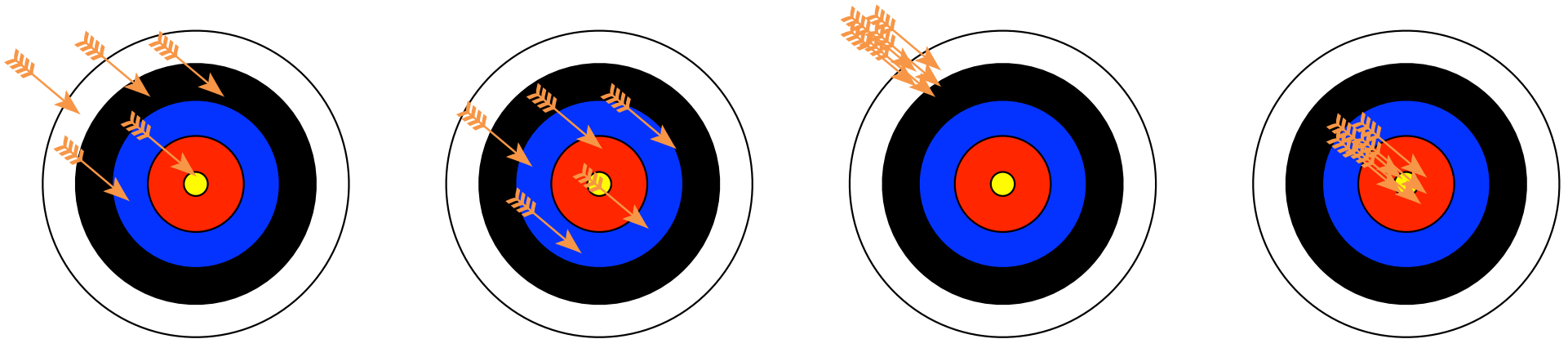
⇒ **Predictive modeling**

현재 데이터(training data)를 잘 설명하는 모델
= 학습오차(training error) 를 최소화하는 모델

$$MSE_{(training)} = \frac{1}{n} \sum_{i=1}^n (Y - \hat{Y})^2$$

Good Predictive Model

$$\begin{aligned} E \left[MSE_{(testing)} \right] &= E \left[(Y - \hat{Y})^2 | X \right] \\ &= \sigma^2 + \left(E[\hat{Y}] - Y \right)^2 + E \left[\left(\hat{Y} - E[\hat{Y}] \right)^2 \right] \\ &= \text{Irreducible error} + \text{Bias}^2 + \text{Variance} \end{aligned}$$



Bias	High	Low	High	Low
Variance	High	High	Low	Low

미래 데이터(testing data)에 대한 예측 성능이 좋은 모델
= 미래 데이터에 대한 expected error 가 낮은 모델

$$E[MSE] = \text{Error} + \text{Bias}^2 + \text{Variance}$$

- Expected MSE를 줄이려면 bias, variance, 혹은 둘다 낮춰야함
- 그렇지 못하다면 둘 중에 하나라도 작으면 좋음
- Bias가 증가되더라도 variance 감소폭이 더 크다면 expected MSE는 감소(예측성능 증가)

Ordinary Linear Regression Model

$$MSE = \sum_{i=1}^n \left\{ Y_i - (w_0 + w_1 x_{1i} + w_2 x_{2i} + \cdots + w_p x_{pi}) \right\}^2$$

$$\min_{w_0, \dots, w_p \in \mathbb{R}} MSE$$

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad \mathbf{w} = \begin{matrix} w_0 \\ w_1 \\ \vdots \\ w_p \end{matrix}$$

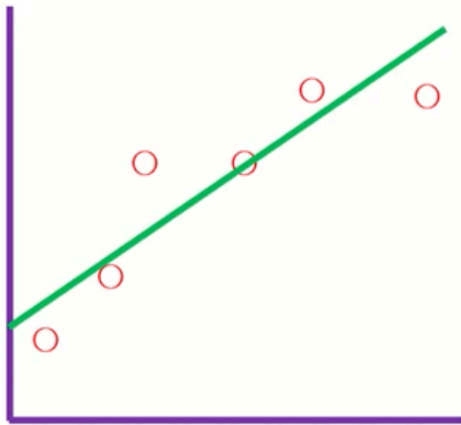
$(1+p) \times 1$

회귀계수 \mathbf{w} 에 대한 unbiased estimator 중 가장 분산이 작은 estimator (**B**est **L**inear and **U**nbiased Estimator: **BLUE**, Gauss-Markov Theorem)

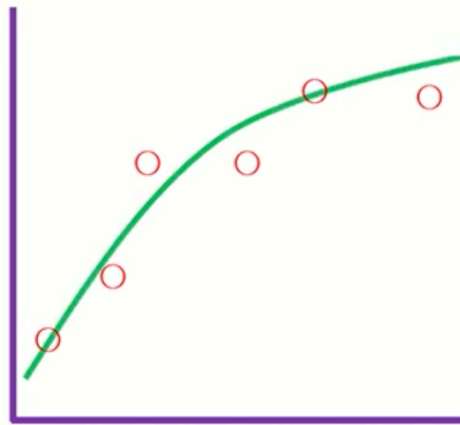
**Unbiased 를 포기 하더라도
(예측) 분산을 획기적으로 줄
일 수 없을까?**

- Subset selection method는 전체 p 개의 입력변수 중 일부 k 개만을 사용하여 회귀계수 \mathbf{w} 를 추정 하는 방법
- 전체 변수 중 일부만을 선택함으로써 bias가 증가할 수 있지만 variance는 감소함
 - Best subset selection
 - Forward stepwise selection
 - Backward stepwise selection
 - Least angle regression
 - Orthogonal matching pursuit

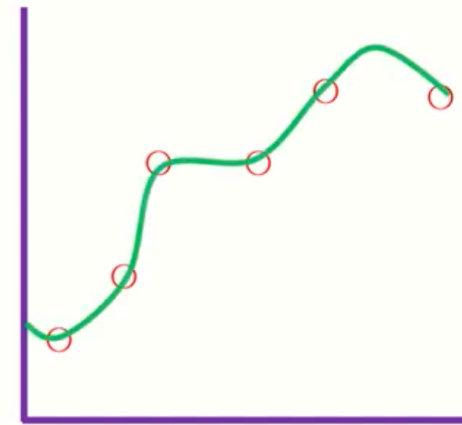
Regularization Concept



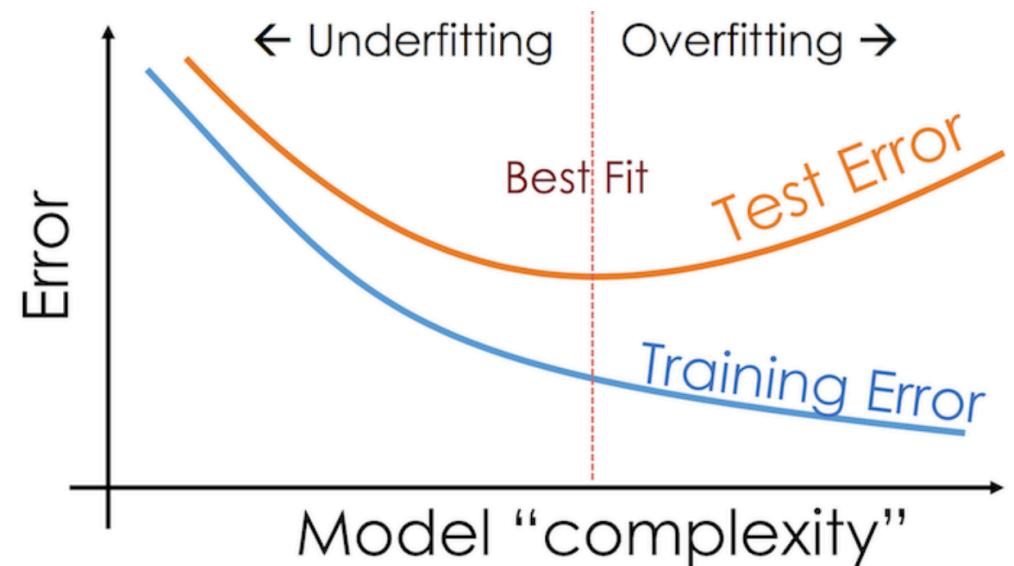
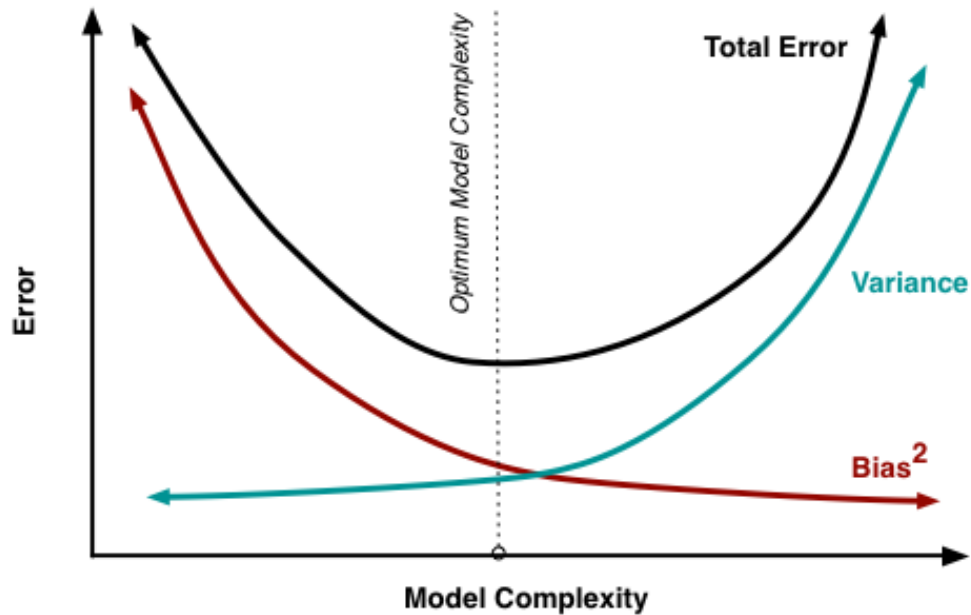
$$w_0 + w_1x$$



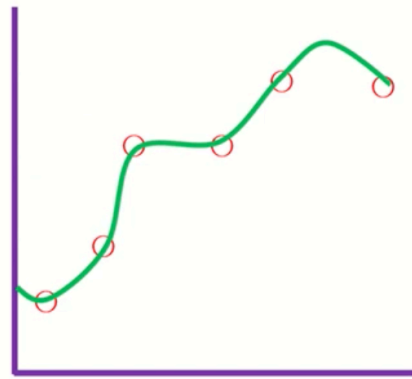
$$w_0 + w_1x + w_2x^2$$



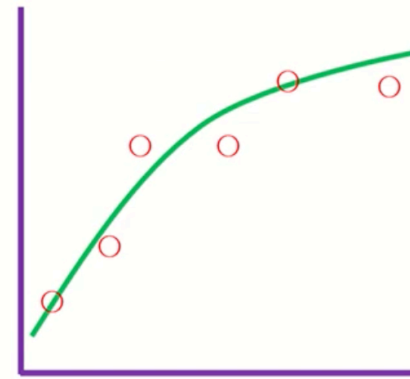
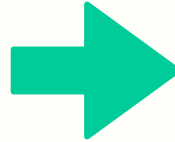
$$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$



Regularization Concept



$$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$



$$w_0 + w_1x + w_2x^2$$

$$\min_{\mathbf{w}} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + 999999w_3^2 + 999999w_4^2$$

$$w_3 \rightarrow 0$$

$$w_4 \rightarrow 0$$

Regularization Concept

$$w_1, w_2, \dots, w_p$$

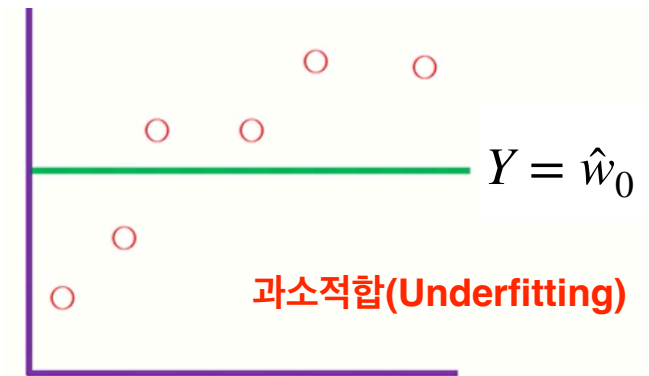
$$C(\mathbf{w}) = \min_{\mathbf{w}} \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{(1) \text{ Training error}} + \underbrace{\lambda \sum_{j=1}^p w_j^2}_{(2) \text{ Generalization}}$$

λ : regularization parameter that controls the tradeoff between (1) and (2)

Regularization Concept

$$C(\mathbf{w}) = \min_{\mathbf{w}} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p w_j^2$$

If λ is very big, then w_i approaches 0.



If λ is very small, then w_i may not 0.



Regularization Method

- Regularization method 는 회귀계수 \mathbf{w} 가 가질 수 있는 값에 제약조건을 부여하는 방법
- 제약조건에 의해 bias는 증가할 수 있지만 variance는 감소함

최소제곱법

$$\min_{w_1, w_2} \sum_{i=1}^n \{Y_i - (w_1 x_{1i} + w_2 x_{2i})\}^2$$

정규화 방법

$$\min_{w_1, w_2} \sum_{i=1}^n \{Y_i - (w_1 x_{1i} + w_2 x_{2i})\}^2$$

st

$$w_1^2 + w_2^2 \leq 30$$

\mathbf{w} 값에 대한 제약조건 추가

Regularization Method

	(w_1, w_2)	$w_1^2 + w_2^2$	MSE
	(4, 5)	41	20
	(3, 5)	34	23
	(4, 4)	32	25
	(2, 5)	27	27
	(2, 4)	18	25
	(2, 3)	13	29

$w_1^2 + w_2^2 \leq 30$

Ridge Regression (능형회귀)

L_2 -norm regularization:

오차제곱합을 최소화하면서 회귀계수 \mathbf{w} 의 L_2 -norm을 제한

$$\mathbf{X} = \begin{array}{|c|c|c|c|} \hline X_{11} & X_{12} & \dots & X_{1p} \\ \hline X_{21} & X_{22} & \dots & X_{2p} \\ \hline \vdots & \vdots & \ddots & \vdots \\ \hline X_{n1} & X_{n2} & \dots & X_{np} \\ \hline \end{array} \quad \mathbf{y} = \begin{array}{|c|} \hline Y_1 \\ \hline Y_2 \\ \hline \vdots \\ \hline Y_n \\ \hline \end{array} \quad \mathbf{w} = \begin{array}{|c|} \hline w_1 \\ \hline \vdots \\ \hline w_p \\ \hline \end{array}$$

$n \times p$ $n \times 1$ $p \times 1$

$$\min_{w_1, \dots, w_p} \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \dots + w_p x_{pi}) \right\}^2$$
$$\text{s.t.} \quad w_1^2 + \dots + w_p^2 \leq t$$

Ridge Regression

$$\begin{aligned} \min_{w_1, \dots, w_p} \quad & \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \dots + w_p x_{pi}) \right\}^2 \\ \text{s.t.} \quad & w_1^2 + \dots + w_p^2 \leq t \end{aligned}$$

\Updownarrow Equivalent (Lagrangian multiplier)

Ridge Regression

$$\min_{w_1, \dots, w_p} \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \dots + w_p x_{pi}) \right\}^2 + \lambda \sum_{j=1}^p w_j^2$$

MSE Contour

$$\begin{aligned}MSE(w_1, w_2) &= \sum_{i=1}^n \{Y_i - (w_1x_{1i} + w_2x_{2i})\}^2 \\&= \left(\sum_{i=1}^n x_{1i}^2\right) w_1^2 + \left(\sum_{i=1}^n x_{2i}^2\right) w_2^2 + \left(2 \sum_{i=1}^n x_{1i}x_{2i}\right) w_1w_2 \\&\quad - \left(2 \sum_{i=1}^n Y_i x_{1i}\right) w_1 - \left(2 \sum_{i=1}^n Y_i x_{2i}\right) w_2 + \sum_{i=1}^n Y_i^2 \\&= Aw_1^2 + Bw_1w_2 + Cw_2^2 + Dw_1 + Ew_2 + F\end{aligned}$$

Conic equation (2차원의 경우)

Conic section

$$Aw_1^2 + Bw_1w_2 + Cw_2^2 + Dw_1 + Ew_2 + F = 0$$

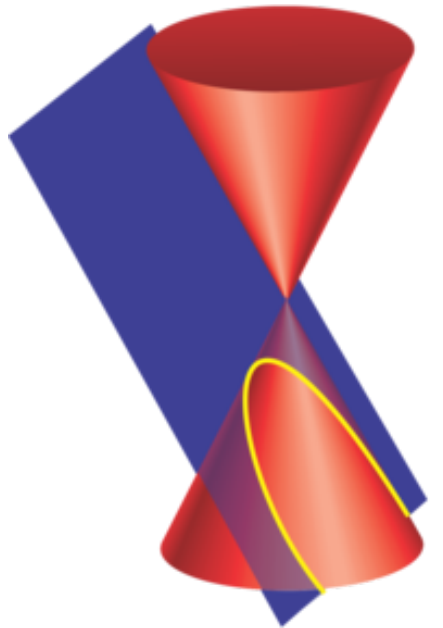
Discriminant of conic equation (판별식): $B^2 - 4AC$

$B^2 - 4AC = 0 \rightarrow$ parabola (포물선)

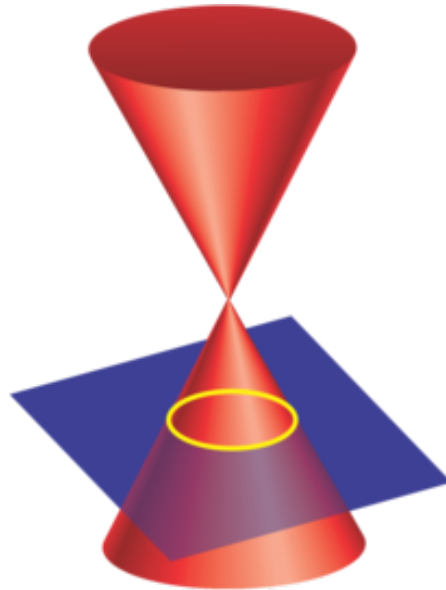
$B^2 - 4AC > 0 \rightarrow$ hyperbola (쌍곡선)

$B^2 - 4AC < 0 \rightarrow$ ellipse (타원)

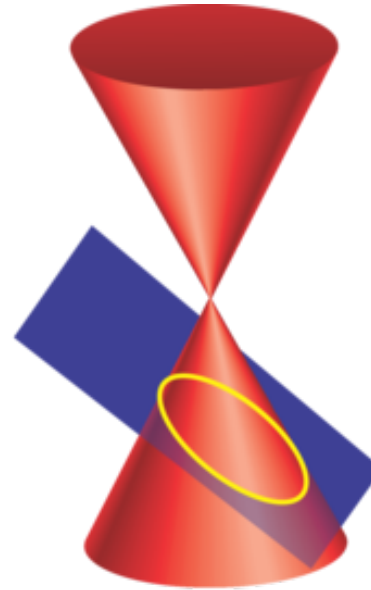
$B = 0$ and $A = C \rightarrow$ circle (원)



parabola



circle



ellipse



hyperbola

MSE Contour

$$\begin{aligned}MSE(w_1, w_2) &= \left(\sum_{i=1}^n x_{1i}^2 \right) w_1^2 + \left(\sum_{i=1}^n x_{2i}^2 \right) w_2^2 + \left(2 \sum_{i=1}^n x_{1i} x_{2i} \right) w_1 w_2 - \left(2 \sum_{i=1}^n Y_i x_{1i} \right) w_1 - \left(2 \sum_{i=1}^n Y_i x_{2i} \right) w_2 + \sum_{i=1}^n Y_i^2 \\&= Aw_1^2 + Bw_1w_2 + Cw_2^2 + Dw_1 + Ew_2 + F\end{aligned}$$

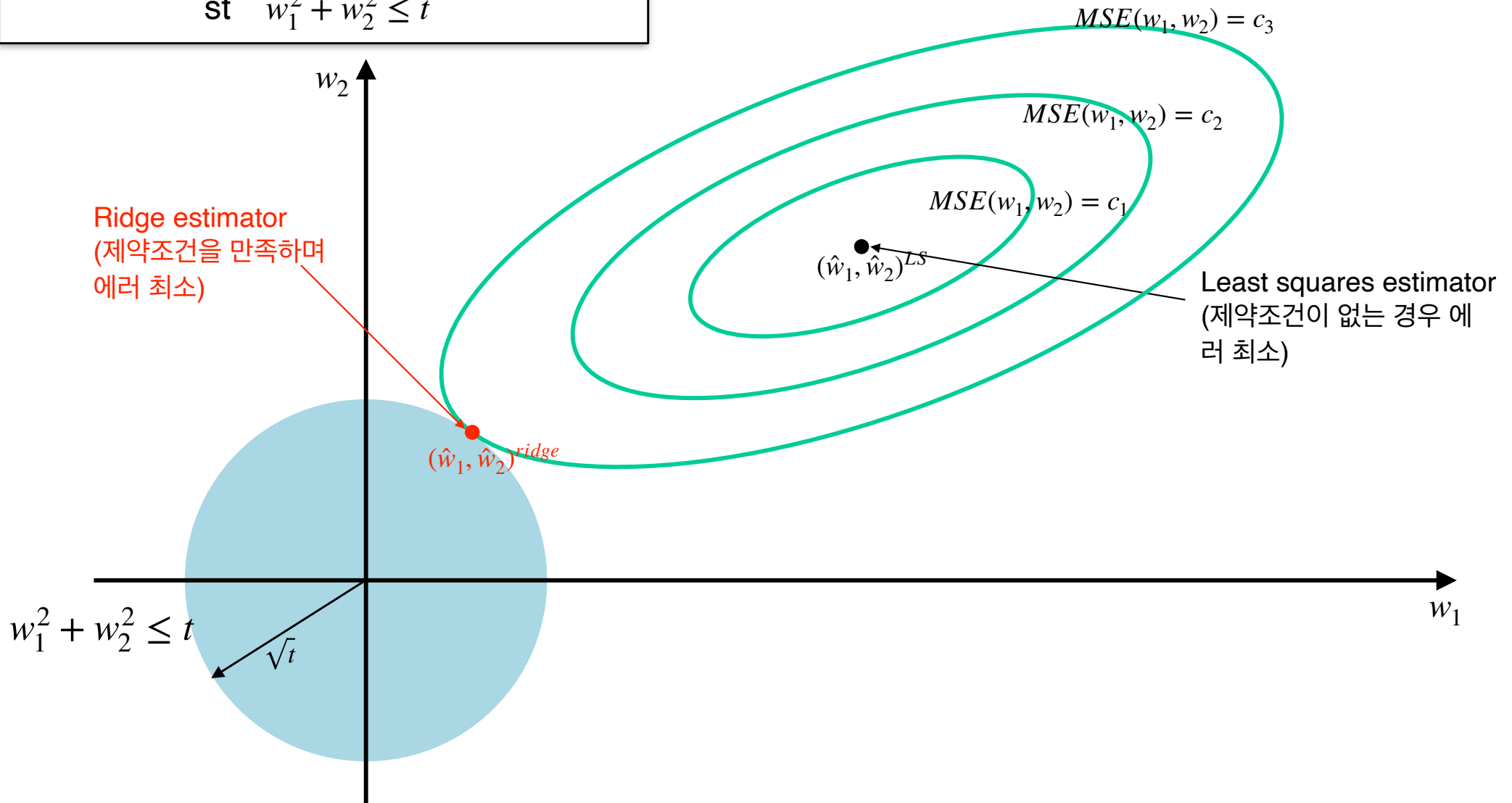
$$\begin{aligned}B^2 - 4AC &= \left(2 \sum_{i=1}^n x_{1i} x_{2i} \right)^2 - 4 \sum_{i=1}^n x_{1i}^2 \sum_{i=1}^n x_{2i}^2 \\&= 4 \left\{ \left(\sum_{i=1}^n x_{1i} x_{2i} \right)^2 - \sum_{i=1}^n x_{1i}^2 \sum_{i=1}^n x_{2i}^2 \right\} < 0\end{aligned}$$

By Cauchy-Schwartz inequality

Ridge Regression

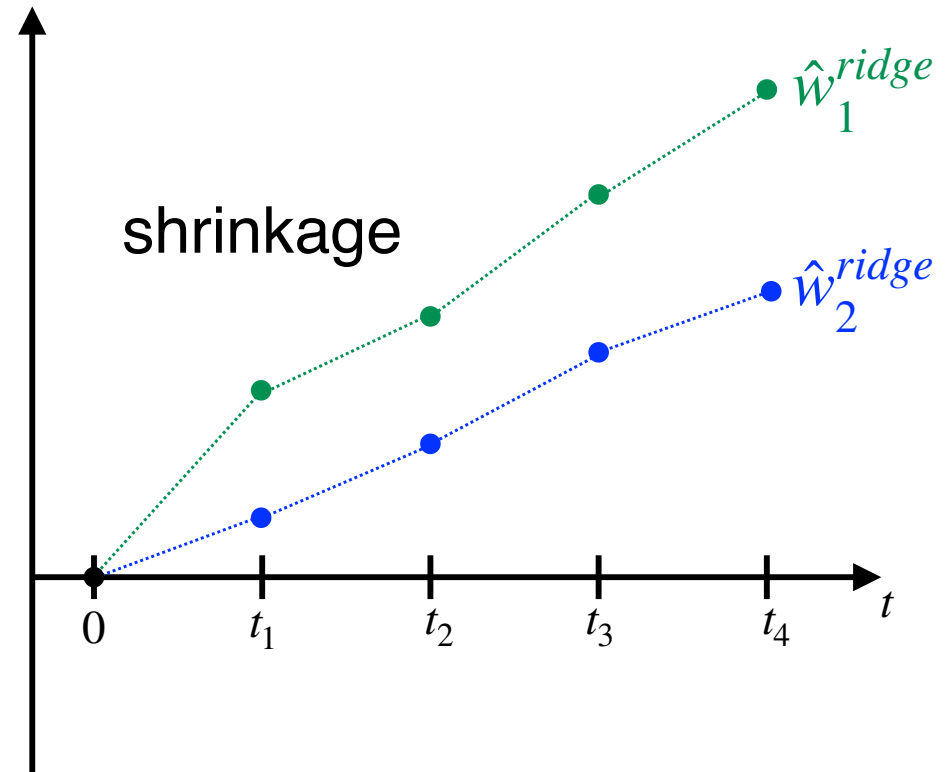
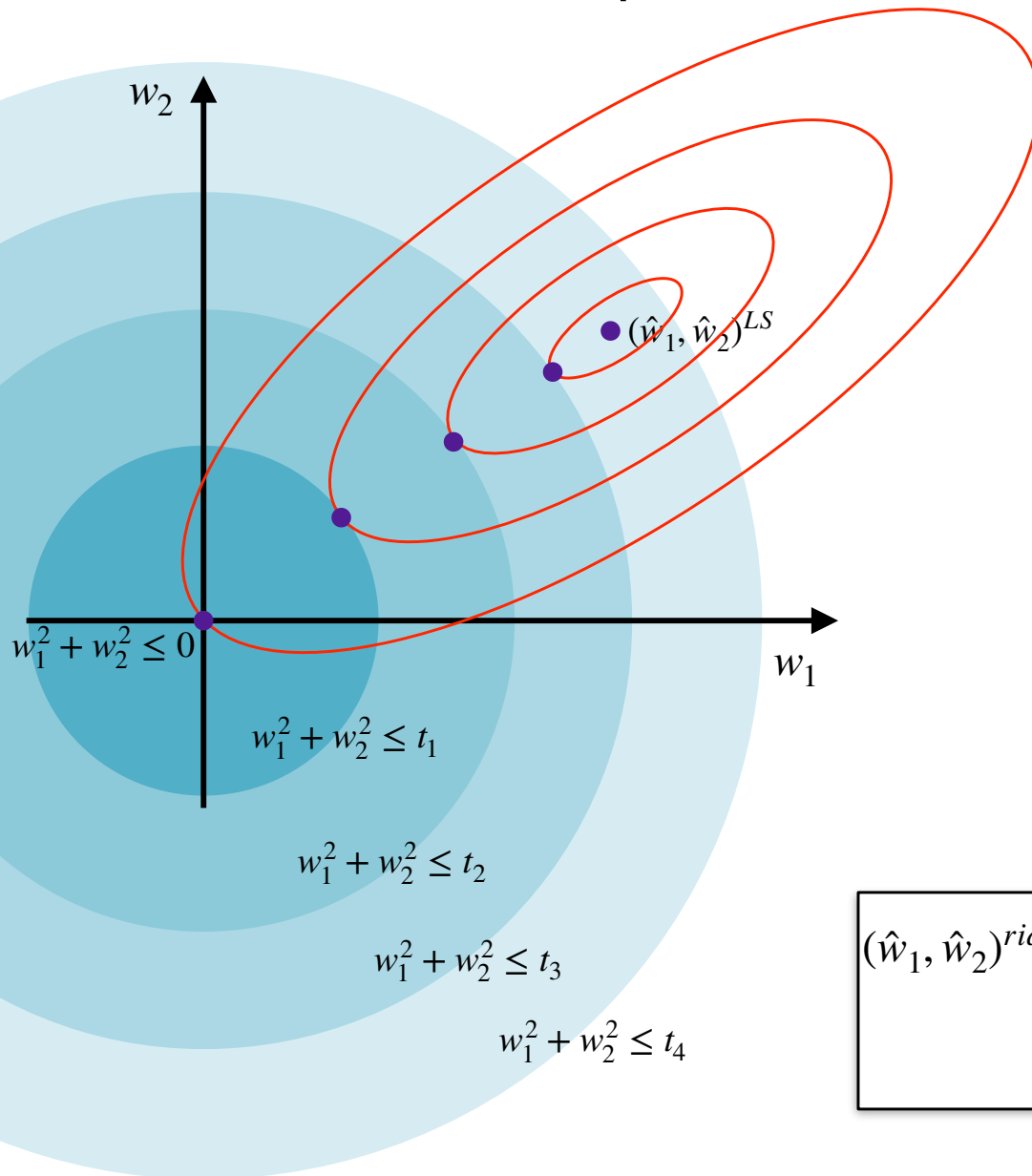
$$(\hat{w}_1, \hat{w}_2)^{ridge} = \arg \min_{w_1, w_2} \sum_{i=1}^n (Y_i - x_{1i}w_1 - x_{2i}w_2)^2$$
$$\text{st } w_1^2 + w_2^2 \leq t$$

MSE contour
(중심에서 멀어질수록 에러 증가)



Ridge Solution Path

Solution path : t 값에 따른 $(\hat{w}_1, \hat{w}_2)^{ridge}$ 의 변화



$$(\hat{w}_1, \hat{w}_2)^{ridge} = \arg \min_{w_1, w_2} \sum_{i=1}^n (Y_i - x_{1i}w_1 - x_{2i}w_2)^2$$

$$\text{st } w_1^2 + w_2^2 \leq t$$

Least Squares Solutions

일반선형회귀분석의 회귀계수 \mathbf{w} 는 행렬 연산을 통해 구할 수 있음

$$\mathbf{X} = \begin{array}{c|c|c|c|c} 1 & X_{11} & X_{12} & \dots & X_{1p} \\ 1 & X_{21} & X_{22} & \dots & X_{2p} \\ 1 & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{np} \end{array} \quad \mathbf{y} = \begin{array}{c} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{array} \quad \mathbf{w} = \begin{array}{c} w_0 \\ w_1 \\ \vdots \\ w_p \end{array}$$

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

$$C(\mathbf{w})^{LS} = (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\frac{\partial C(\mathbf{w})^{LS}}{\partial \mathbf{w}} = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{0}$$

$$\hat{\mathbf{w}}^{LS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Ridge Solutions

능형회귀분석의 회귀계수 \mathbf{w} 도 행렬 연산을 통해 구할 수 있음

$$\mathbf{X} = \begin{array}{|c|c|c|c|} \hline X_{11} & X_{12} & \dots & X_{1p} \\ \hline X_{21} & X_{22} & \dots & X_{2p} \\ \hline \vdots & \vdots & \ddots & \vdots \\ \hline X_{n1} & X_{n2} & \dots & X_{np} \\ \hline \end{array} \quad \begin{array}{c} \\ \\ \\ n \times p \end{array} \quad \mathbf{y} = \begin{array}{|c|} \hline Y_1 \\ \hline Y_2 \\ \hline \vdots \\ \hline Y_n \\ \hline \end{array} \quad \begin{array}{c} \\ \\ n \times 1 \end{array} \quad \mathbf{w} = \begin{array}{|c|} \hline w_1 \\ \hline \vdots \\ \hline w_p \\ \hline \end{array} \quad \begin{array}{c} \\ \\ p \times 1 \end{array}$$

$$C(\mathbf{w})^{ridge} = (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w}$$

$$\frac{\partial C(\mathbf{w})^{ridge}}{\partial \mathbf{w}} = -2\mathbf{X}^\top \mathbf{y} + 2 \left(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p \right) \mathbf{w} = \mathbf{0}$$

$$\hat{\mathbf{w}}^{ridge} = \left(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

Least Absolute Shrinkage and Selection Operator

변수 선택 가능

L₁-norm regularization: 회귀계수 \mathbf{w} 의 L₁-norm을 제한

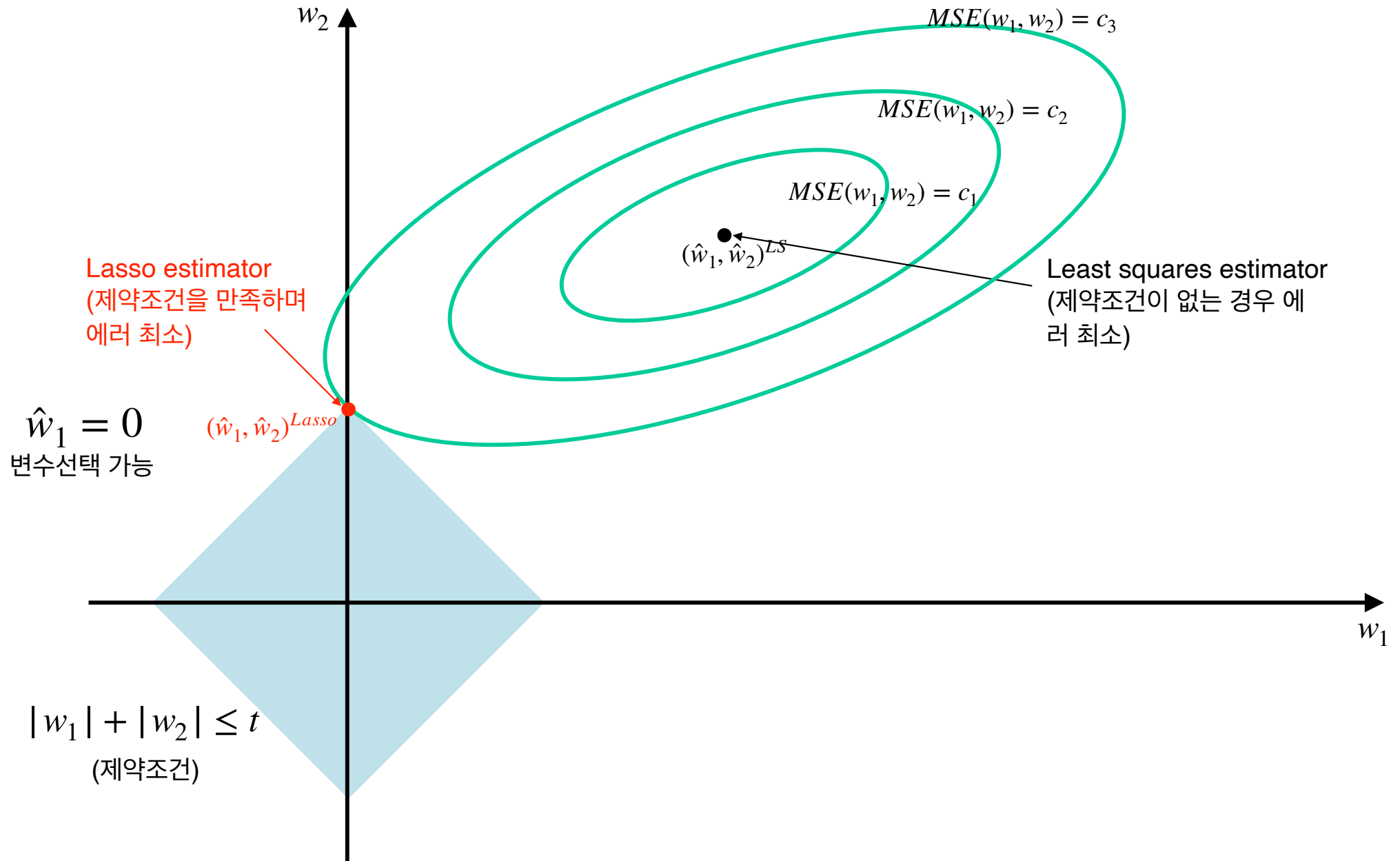
Lasso

$$\begin{aligned} \min_{w_1, \dots, w_p} \quad & \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \dots + w_p x_{pi}) \right\}^2 \\ \text{s.t.} \quad & |w_1| + \dots + |w_p| \leq t \end{aligned}$$

↕ Equivalent (Lagrangian multiplier)

$$\min_{w_1, \dots, w_p} \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \dots + w_p x_{pi}) \right\}^2 + \lambda \sum_{j=1}^p |w_j|$$

Lasso Solution path



Lasso Solutions

$$\min \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \cdots + w_p x_{pi}) \right\}^2 + \lambda \sum_{j=1}^p w_j^2 \quad \Rightarrow \quad \hat{\mathbf{w}}^{ridge} = \left(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\min \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \cdots + w_p x_{pi}) \right\}^2 + \lambda \sum_{j=1}^p |w_j| \quad \Rightarrow \quad \hat{\mathbf{w}}^{lasso} = ?$$

- Ridge와 달리 Lasso는 closed form solution 을 구하는 것이 불가능 (L₁-norm 미분 불가능)
- Numerical optimization methods:
 - Quadratic programming techniques (1996, Tibshirani)
 - LARS algorithm (2004, Efron et al.)
 - Coordinate descent algorithm (2007, Friedman et al.)

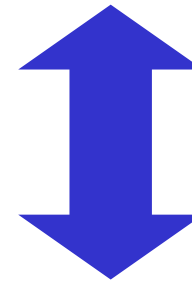
Lasso

$$\min_{w_1, \dots, w_p} \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \dots + w_p x_{pi}) \right\}^2 + \lambda \sum_{j=1}^p |w_j|$$

$\lambda \rightarrow 0$: 최소제곱법

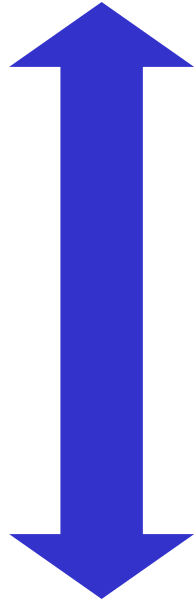
$\lambda \rightarrow \infty$: $\mathbf{w} \rightarrow \mathbf{0}$

λ 값을 어떻게 설정할 것인가?



몇 개의 변수를 선택할 것인가?

λ 값을 크게 잡으면?



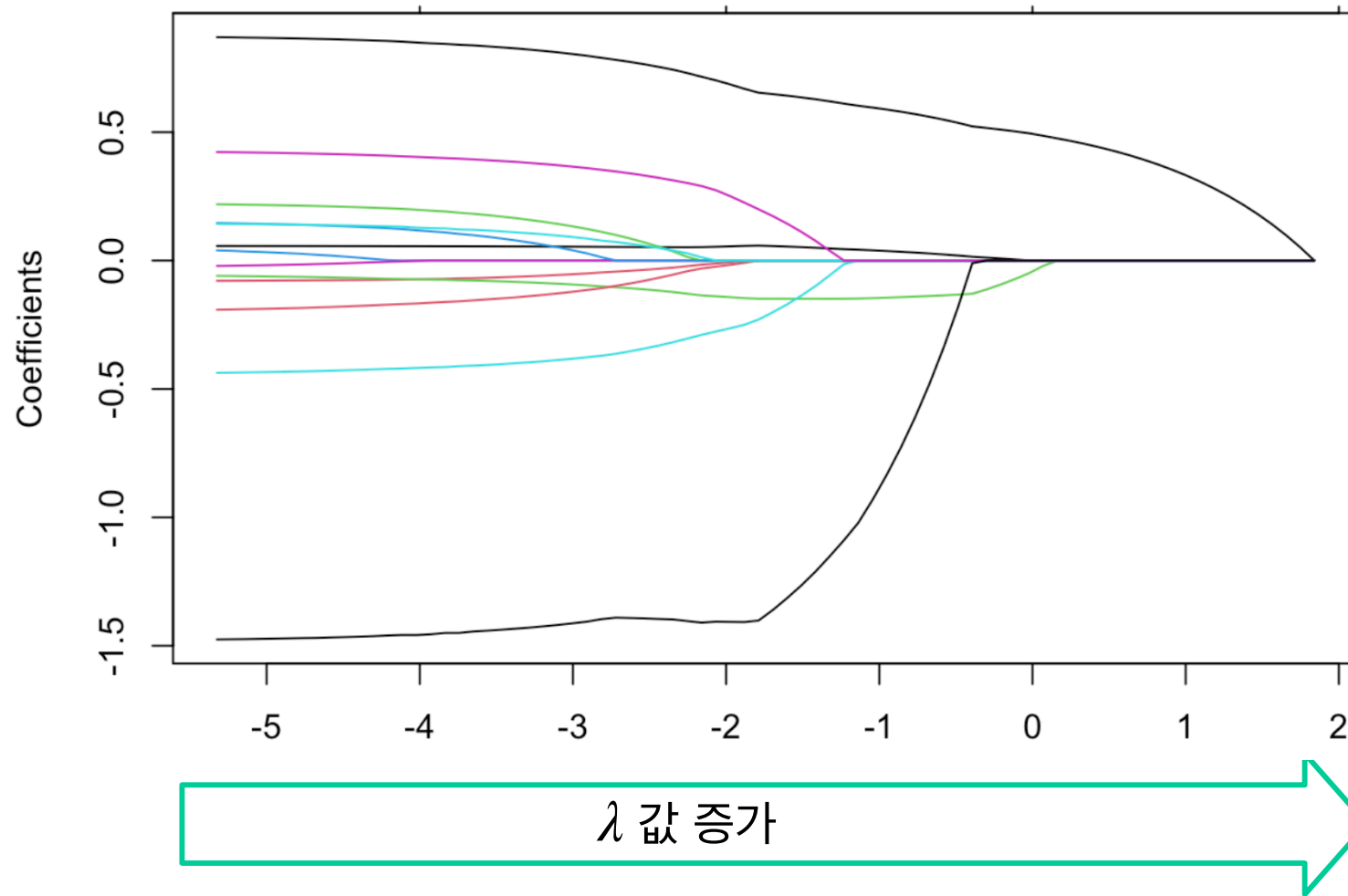
λ 값을 작게 잡으면?

적은 변수
간단한 모델
해석 쉬움
높은 학습 오차(training error 증가)
Underfitting 위험 증가

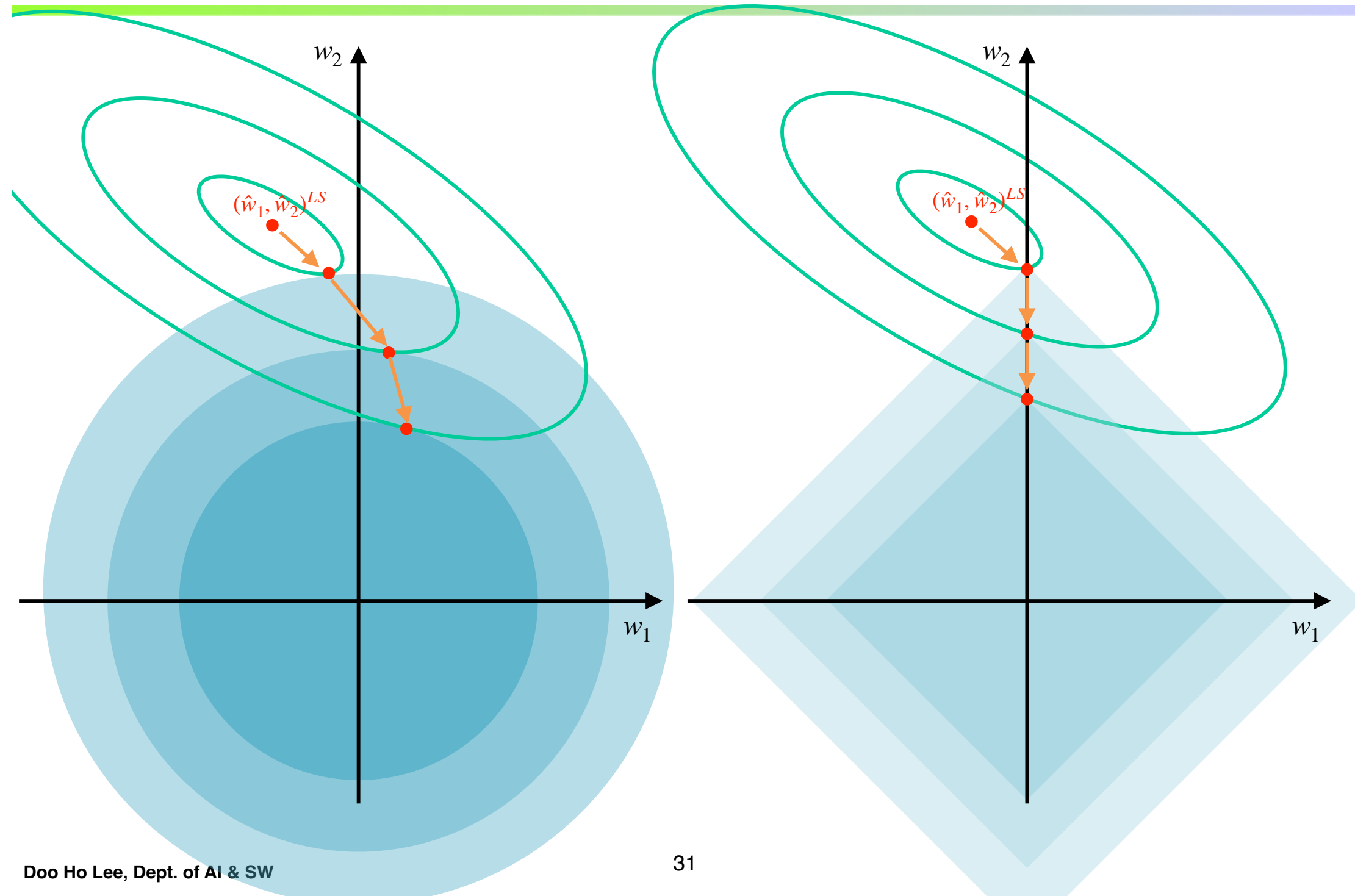
많은 변수
복잡한 모델
해석 어려움
낮은 학습 오차
Overfitting 위험 증가

Lasso parameter

$$\min_{w_1, \dots, w_p} \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \dots + w_p x_{pi}) \right\}^2 + \lambda \sum_{j=1}^p |w_j|$$



Solution Paths of Ridge and Lasso



Ridge vs Lasso

Ridge	Lasso
L_1 -norm regularization	L_2 -norm regularization
변수 선택 불가능	변수 선택 가능
Closed form solution 존재	Closed form solution 존재하지 않음 (수치해석적 방법 이용)
입력변수 간 상관관계가 높은 상황에서 좋은 예측 성능	입력변수 간 상관관계가 높은 상황에서 ridge에 비해 상대적으로 예측성능이 떨어짐

Elastic Net

- Elastic Net = Ridge + Lasso (L₁- and L₂- norm regularization)
- Elastic Net은 상관관계가 큰 입력변수를 동시에 선택/배제하는 특성

Elastic Net

$$\begin{aligned} \min_{w_1, \dots, w_p} \quad & \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \dots + w_p x_{pi}) \right\}^2 \\ \text{s.t.} \quad & \alpha \sum_{j=1}^p |w_j| + (1 - \alpha) \sum_{j=1}^p w_j^2 \leq t \end{aligned}$$

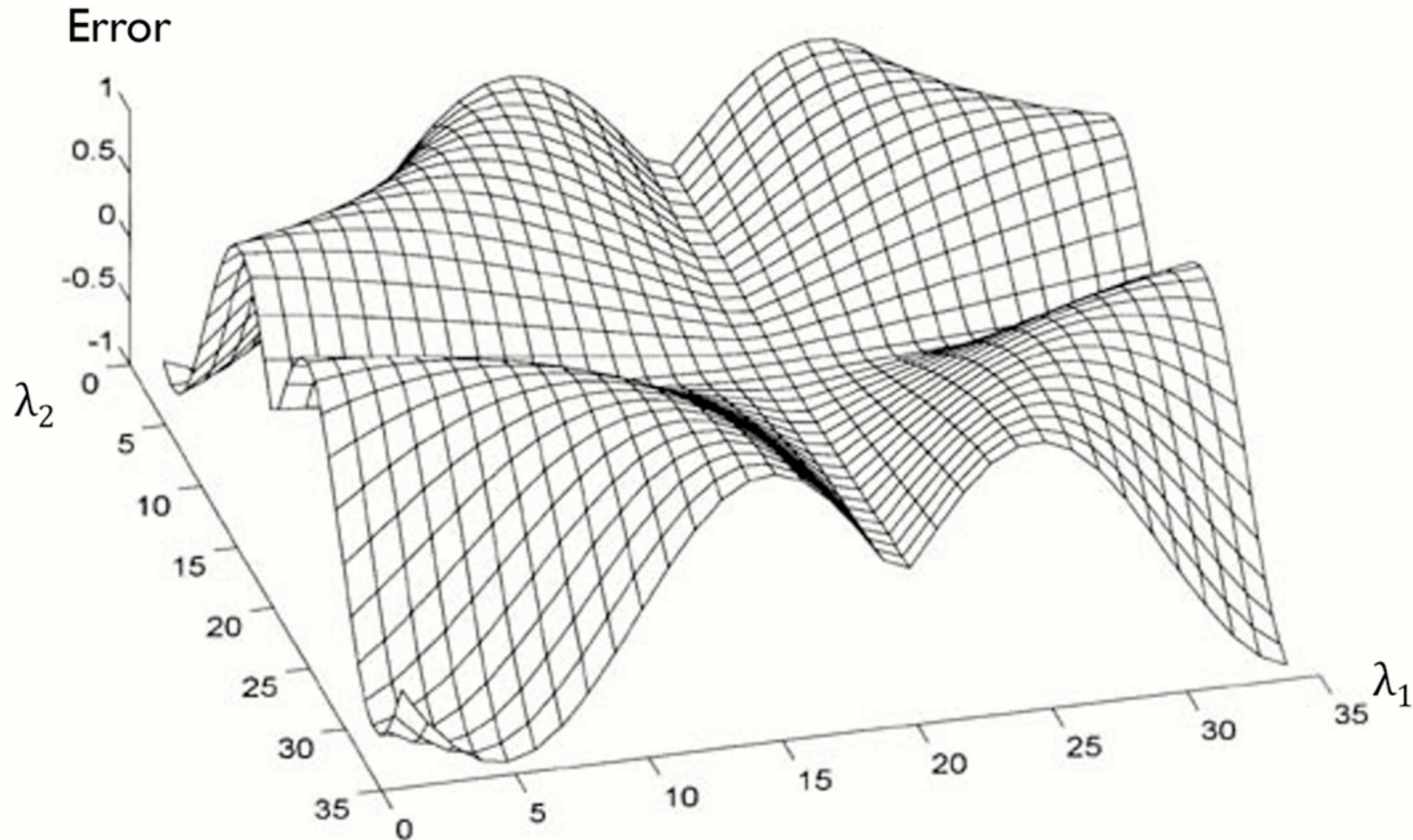
↕ Equivalent (Lagrangian multiplier)

$$\min_{w_1, \dots, w_p} \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \dots + w_p x_{pi}) \right\}^2 + \lambda \left(\alpha \sum_{j=1}^p |w_j| + (1 - \alpha) \sum_{j=1}^p w_j^2 \right)$$

Elastic Net Parameters

$$\min_{w_1, \dots, w_p} \sum_{i=1}^n \left\{ Y_i - (w_1 x_{1i} + \dots + w_p x_{pi}) \right\}^2 + \lambda_1 \sum_{j=1}^p |w_j| + \lambda_2 \sum_{j=1}^p w_j^2$$

- 일정 범위 내로 λ_1 과 λ_2 를 조정하여 오차가 가장 작은 결과를 보이는 λ_1 과 λ_2 값을 선정함



Ridge vs. Lasso vs. Elastic Net

