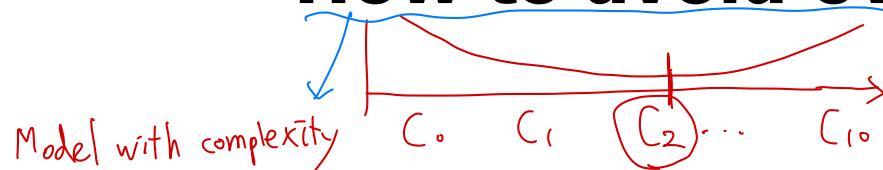


Regularization - How to avoid overfitting -

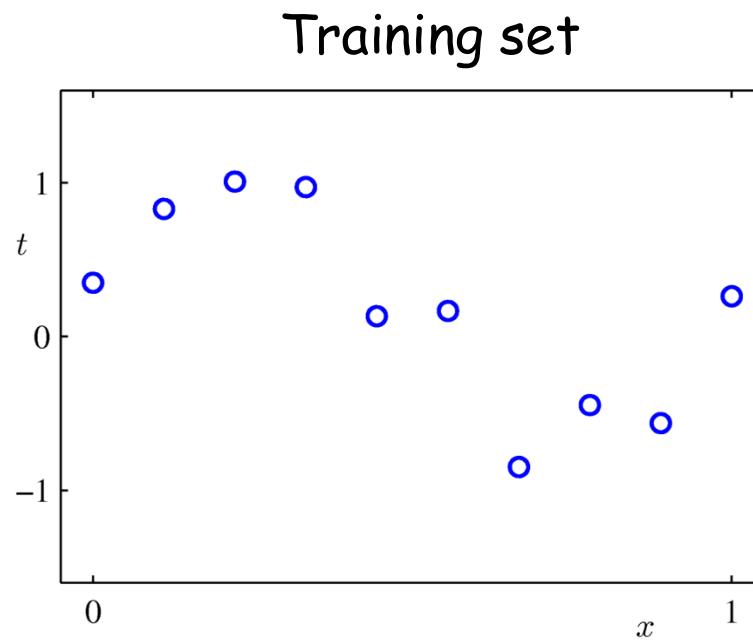


→ 문제점: 모든 모델을 빌드해야된다. (cost ↑)

Method

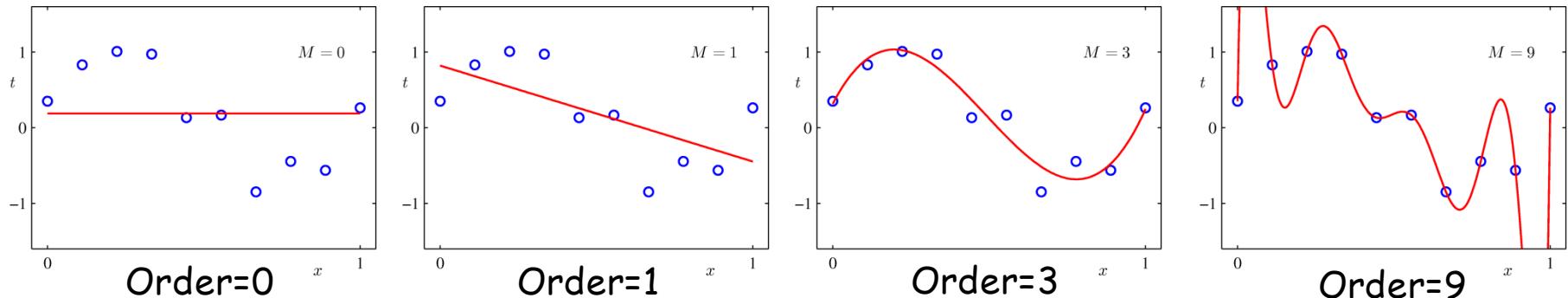
Introduction

- I want to build a good generalized model



Introduction

- **Of which model complexity??**
 - As high as possible but not too much to be overfitted
- **How can I choose the best model complexity?**
 - Try one by one and choose the best using Hold-out or CV
 - We usually use this approach



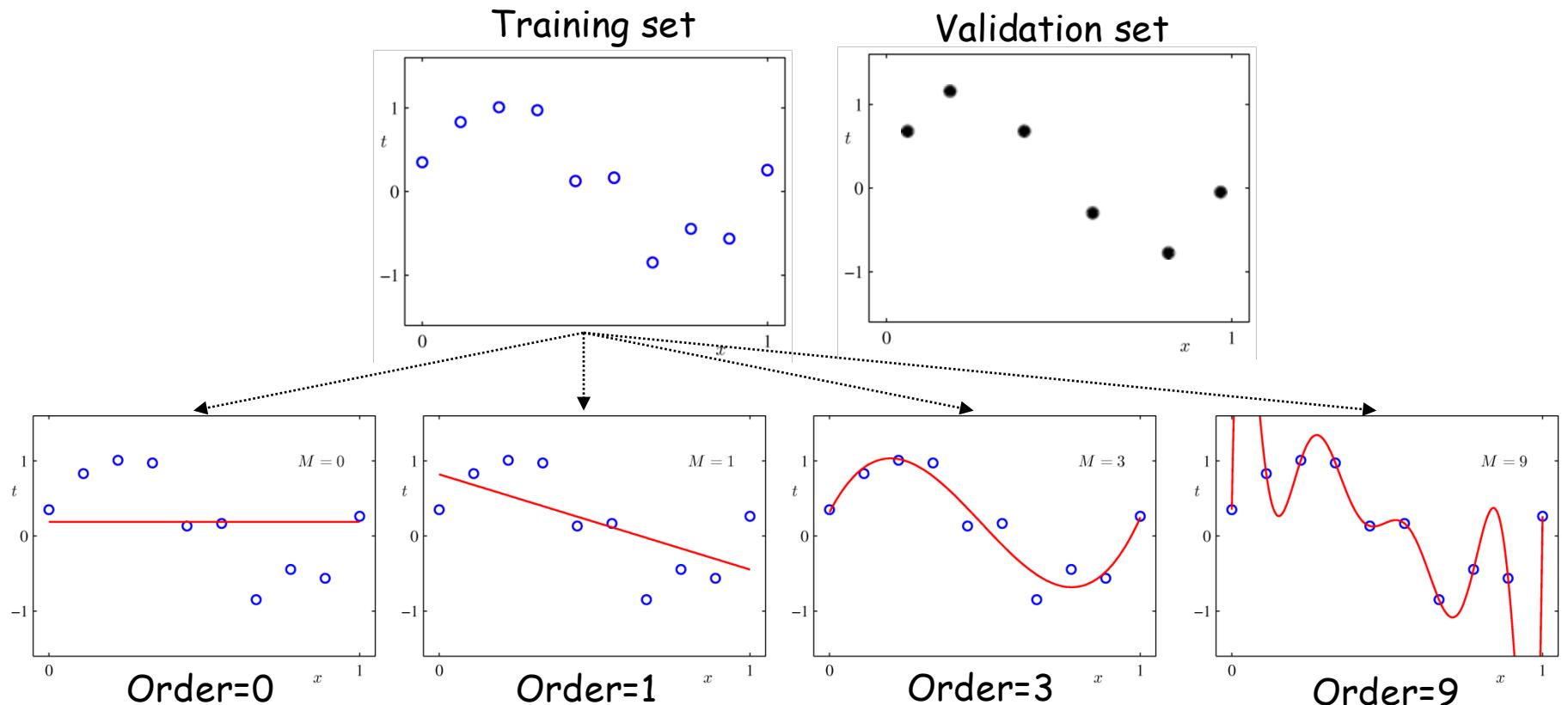
Introduction

- **Choosing the best model**
 - Train many models
 - Evaluate them with Cross-Validation or Hold-out method
 - Choose the best

- **Using Regularization Method**
 - Train one model with a regularization method

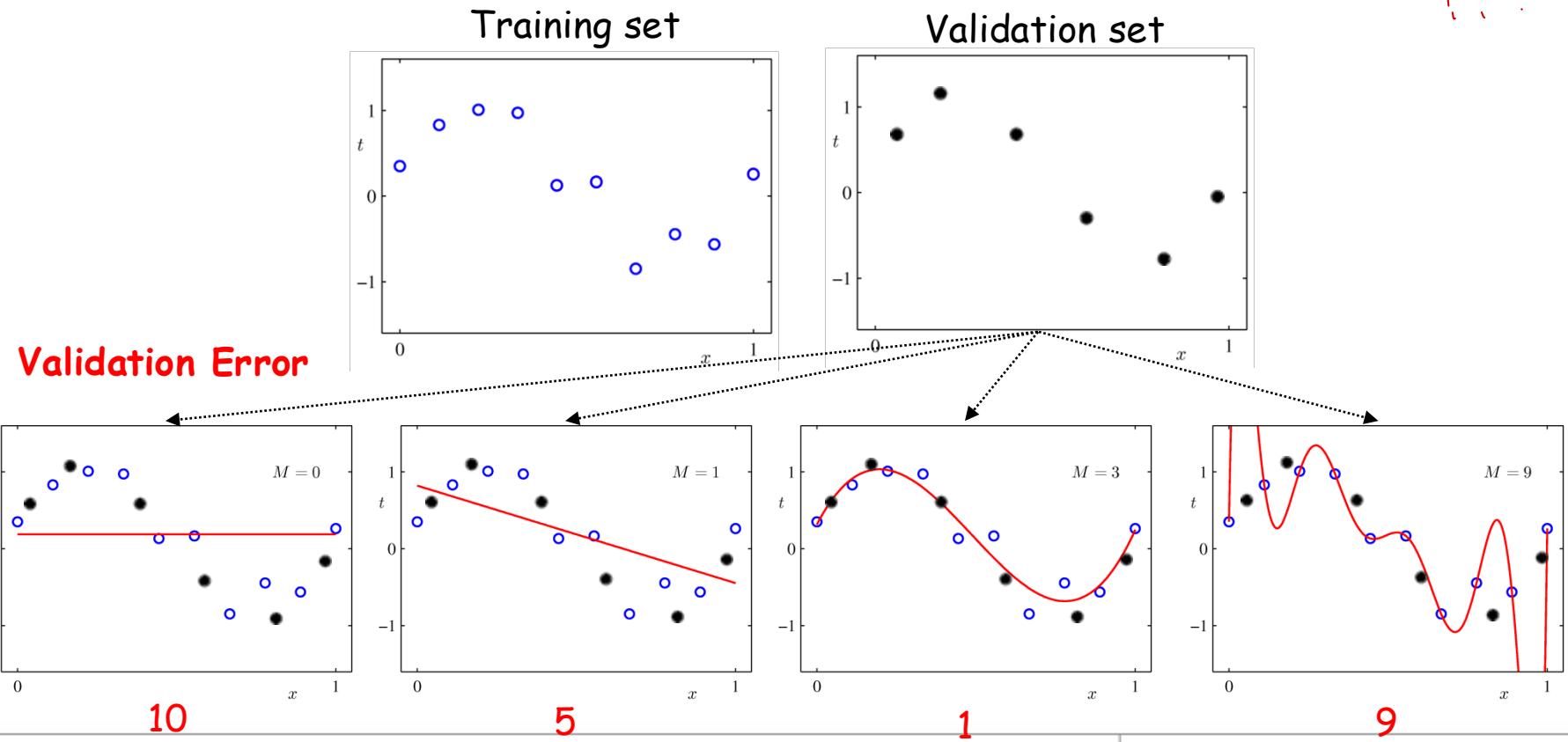
Choosing the Best Model

- I want to build a good generalized model
 - Step 1: build several models with training set



Choosing the Best Model

- I want to build a good generalized model
 - Step 2: evaluate the models and choose the best *Time consuming !!!*

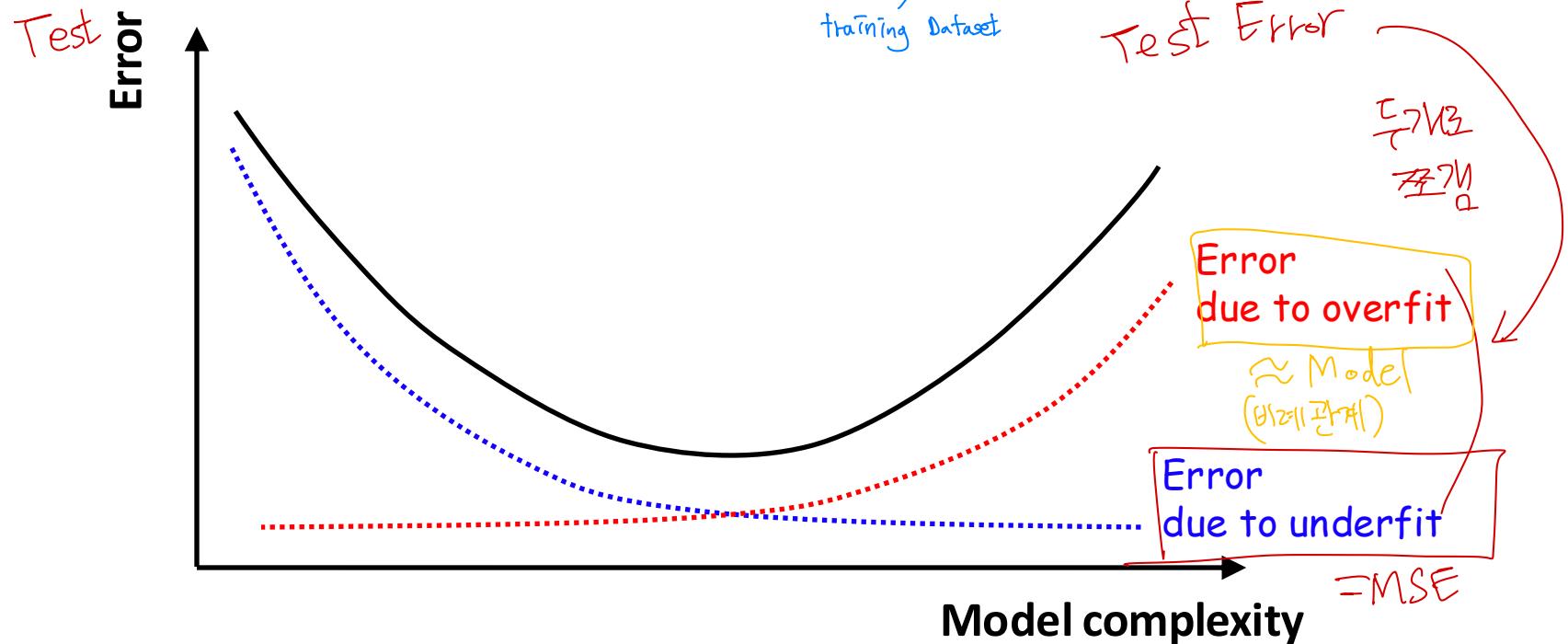


Regularization

- **Regularization** *만족의 best model 찾는법*
 - Techniques for automatic adjustment of model complexity
 - Most techniques used while model training.
Some are applied to the model after training
- **Each ML technique has its own regularization techniques**
 - Additive ~~regularizer~~ models: Weight decay
 - DNN: Weight decay, Drop-out ...
 - Decision Tree: Pruning, Max-Height Limitation ...
 - Naïve Bayes: Smoothing

Weight Decay

- Two kinds of errors
 - Error due to underfit
 - Error due to overfit



Weight Decay

- Two kinds of errors

- Error due to underfit \propto Training Error

- It may be proportional to the error of training dataset

$$E(\mathbf{w}) = \sum_{(\mathbf{x}, y) \in Data} (y - f(\mathbf{x}))^2$$

MSE

easy to be overfitted

- Error due to overfit \propto Model complexity

- It may be proportional to the complexity of models

$$C(\mathbf{w}) = \text{Complexity of the model}$$

- If we can minimize the total error, ...

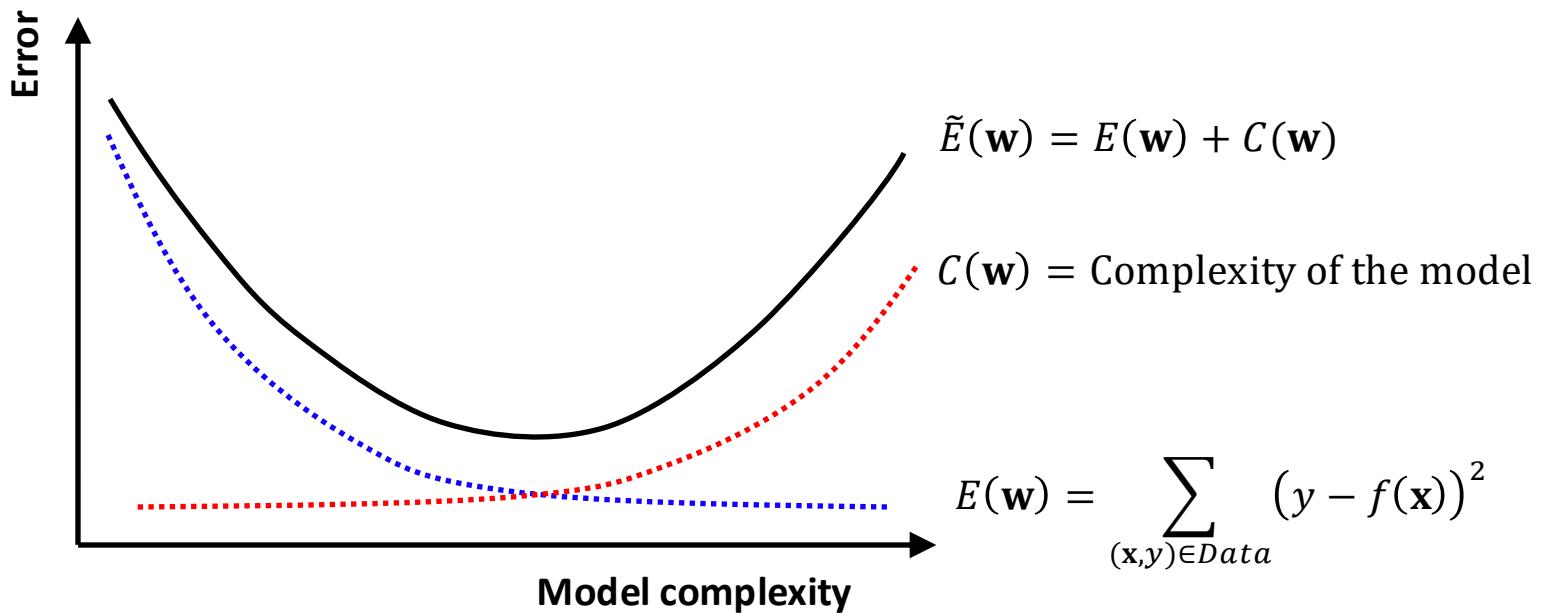
$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + C(\mathbf{w})$$

Weight Decay

- **New strategy**

- We tried to minimize $E(\mathbf{w})$, so we trouble with overfitting
- We can minimize $\tilde{E}(\mathbf{w})$, we may find the optimal model

But, how can we evaluate the model complexity, $C(\mathbf{w})$?



Weight Decay



■ Model Complexity \rightsquigarrow Complexity 를 적절하게 설정

Model Complexity increases

$$\begin{array}{ll} f(x) = w_0 & \xrightarrow{\text{간주}} f(x) = 0 \cdot x^3 + 0 \cdot x^2 + 0 \cdot x + w_0 \\ f(x) = w_1 x + w_0 & \xrightarrow{} f(x) = 0 \cdot x^3 + 0 \cdot x^2 + w_1 x + w_0 \\ f(x) = w_2 x^2 + w_1 x + w_0 & \xrightarrow{} f(x) = 0 \cdot x^3 + w_2 x^2 + w_1 x + w_0 \\ f(x) = w_3 x^3 + w_2 x^2 + w_1 x + w_0 & \xrightarrow{} f(x) = w_3 x^3 + w_2 x^2 + w_1 x + w_0 \end{array}$$

↑
small
C 항 수
일정하게
증가함
Large

$$\text{Model Complexity} \propto \# \text{ of Non-zero w's}$$

approximately

$$\sum_i |w_i| \quad \text{or} \quad \sum_i w_i^2$$

■ We may say

$$C(\mathbf{w}) = \sum_i |w_i| \quad \text{or} \quad \sum_i w_i^2$$

approximately

$\hookrightarrow W_3 = 0,001 / 0.12 \approx 0.0083$ \neq Exactly
 $f(x)$ 가 3차지면 0차랑 같아짐 (그러나 approximately)

Weight Decay

어떤 것인가?

■ Regularization: Weight Decay

- Penalizing the Model Complexity

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} [E(\mathbf{w}) + \lambda \cdot C(\mathbf{w})]$$

E, C 둘 중 영향을 더 주고 싶은 차이 따위

$E(\mathbf{w})$: Training error function, such as MSE

$C(\mathbf{w})$: A function reflecting model complexity

λ : A balancing factor given by experts (hyper parameter)

- In the case of Regression

• Ridge Regression: $C(\mathbf{w}) = \sum_{i=1}^n w_i^2$

• Lasso Regression: $C(\mathbf{w}) = \sum_{i=1}^n |w_i|$

기존 regression은 $E(\mathbf{w})$ 만

Weight Decay

- **Ridge Regression**

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n (y_i - f(\mathbf{x}_i, \mathbf{w}))^2 + \lambda \cdot \sum_{i=1}^m w_i^2$$

- **Lasso Regression**

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n (y_i - f(\mathbf{x}_i, \mathbf{w}))^2 + \lambda \cdot \sum_{i=1}^m |w_i|$$

How they work ???

Weight Decay

- How they works?

- For example, Ridge Regression

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n (y_i - f(\mathbf{x}_i, \mathbf{w}))^2 + \lambda \cdot \sum_{i=1}^m w_i^2$$

$$f(\mathbf{x}, \mathbf{w}) = w_9x_2^3 + w_8x_1^3 + w_7x_2^2x_1 + w_6x_2x_1^2 \\ + w_5x_2^2 + w_4x_1^2 + w_3x_1x_2 + w_2x_2 + w_1x_1 + w_0$$

- What happens if x_2^3 is less affective to model accuracy?

$$f(\mathbf{x}, \mathbf{w}) = 0 \cdot x_2^3 + 0 \cdot x_1^3 + w_7x_2^2x_1 + 0 \cdot x_2x_1^2 \\ + 0 \cdot x_2^2 + w_4x_1^2 + 0 \cdot x_1x_2 + w_2x_2 + w_1x_1 + w_0$$

Model is simplified, with less hurt of performance !!

Weight Decay

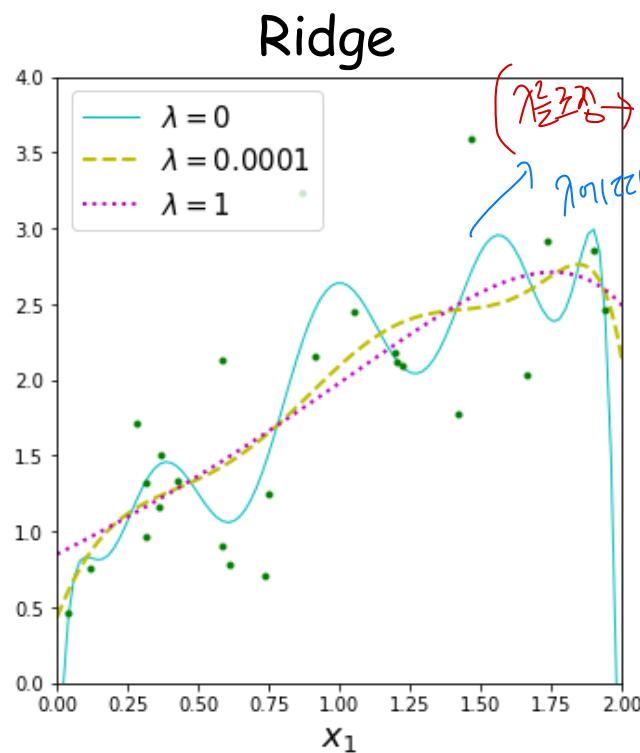
보통은 $\lambda_{\max} = \frac{1}{(데이터개수/2) * 디자인)$

Example

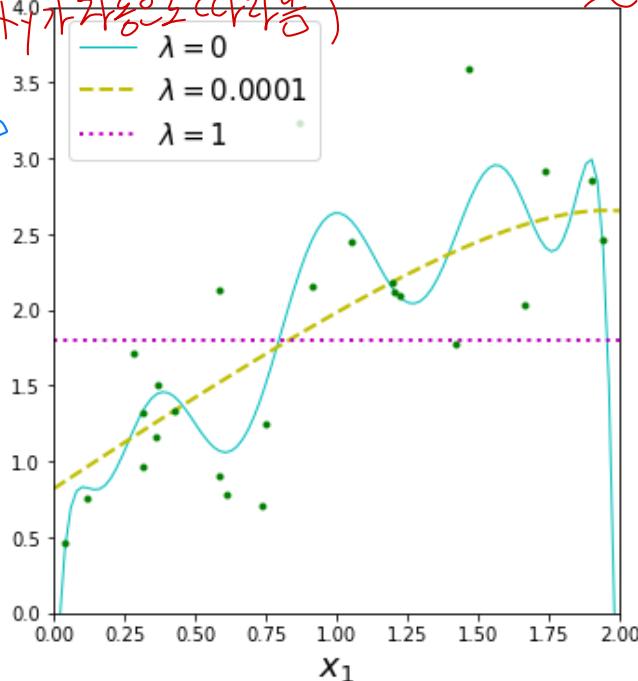
$$E(\mathbf{w}) + \lambda \cdot C(\mathbf{w})$$

$\lambda=0$, MSE만 고려

$\lambda=1$, Ridge/Lasso에 따라 달라짐



E by overfitting $\xrightarrow{\text{측정하기 어려움}}$
Lasso \sim model complexity
 $\xrightarrow{\text{로 측정함}}$



Weight Decay

L1-Regularization

$$C(\mathbf{w}) = \sum_{i=1}^m |w_i|$$

- Encourage sparsity by setting weight = 0.
- Used to select the most informative features.

(robust한 이유)

L_1	$w_3 \propto^3$	$w_2 \propto^2$	$w_1 \propto^1$	w_0	$\sum_{i=1}^m w_i $ 둘다 40
40	0	0	0	10	40
10	10	10	10	10	40

L_2	40	0	0	0	$\sum_{i=1}^m w_i^2 < 160$ 400
10	10	10	10	10	400

$$f(g) = w_1 g + w_0 \rightarrow f(g) \text{가 심플하면 좋은 모델}$$

$$f(g) = w_3 g^3 + w_2 g^2 + w_1 g + w_0 \text{ 찾기 어렵다}$$

↳ 복잡한 경우에 사용하자!

↳ 보다 더 좋은 방법이다.

방법이다.

L2-Regularization

$$C(\mathbf{w}) = \sum_{i=1}^m w_i^2$$

- Does not encourage sparsity → small but non-zero weights.
- Distributes weight across related features (robust).

안정적인

$$MSE + \sum_{i=1}^m w_i^2 \Rightarrow L_2 \text{ 정규화}$$

둘은 구별할 수 있다.

그리고 더 stable한 W를 찾을 수 있다.