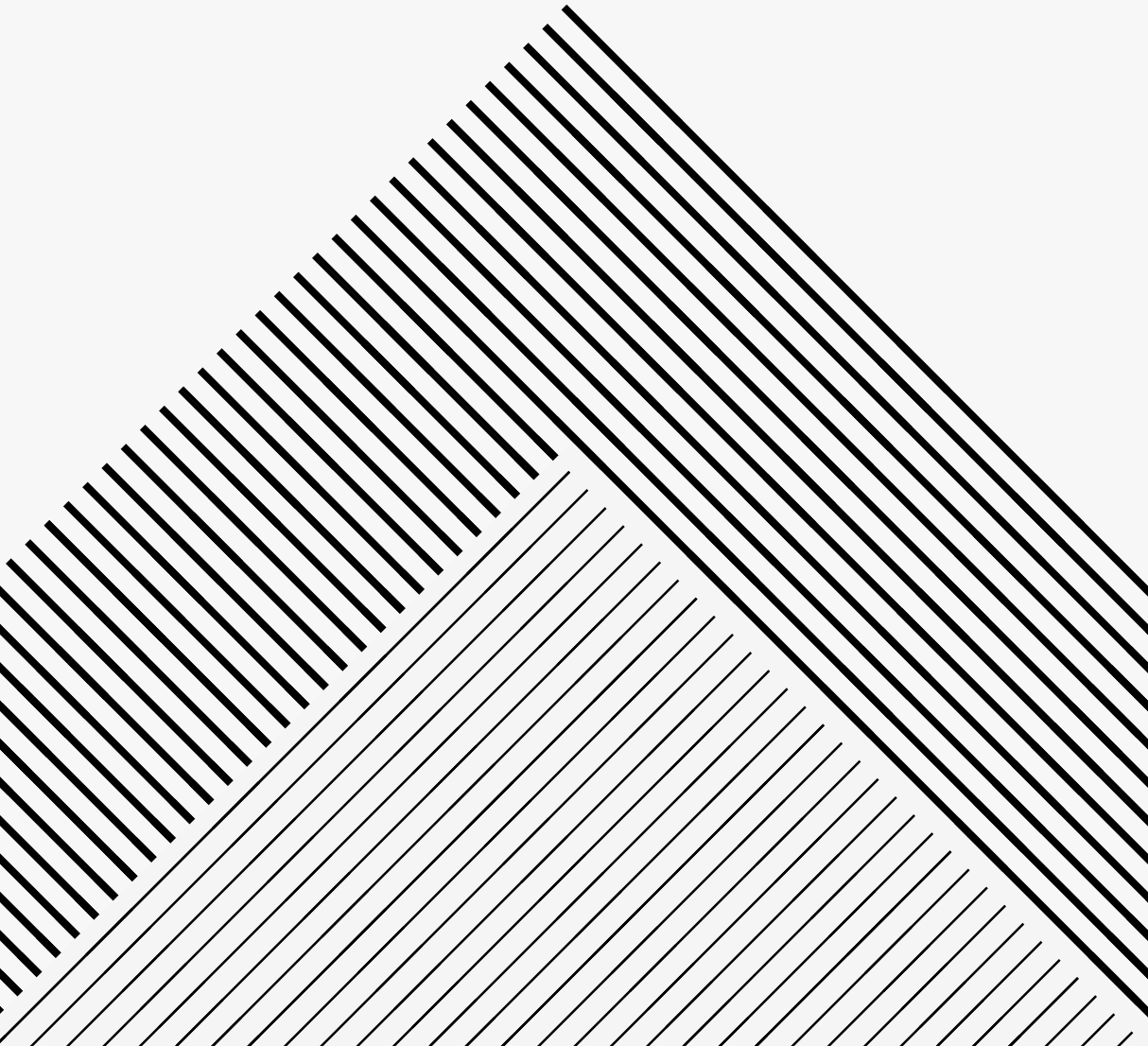


# Regression

CS229



# CHO5. Regression (화기)

## 5-1. 회귀소개

→ 여러개의 독립변수와 하나의 종속변수의 상관관계를 모델링하는 방법

→ 여러개 회귀 예측의 핵심 = 회귀의 regression hyperparameter를 찾는 것.

## Linear Regression (선형 회귀) 94

→  $(\text{실제값} - \text{예측값})^2$ 을 최소화하는 최적의 회귀선을 회귀라

\* 대표적인 선형회귀 모델

귀제 X

- ① 일반 선형회귀  
→ 귀제 X  
→  $\sum (\text{예측} - \text{실제})^2$ 을 최소화하도록 회귀계수 회귀라

② 린지 (Ridge) → feature의 값의 제곱합  
→ ①에  $L_2$  regression을 추가  
→ 상대적으로 큰 회귀계수의 값에 예측 영향도를 감소시키기 위해서 회귀계수값을 더 작게 만들

귀제 0

- ③ 라쏘 (Lasso) → feature의 개수감소  
→ ①에  $L_1$  regression을 추가  
→ 예측 영향력이 적은 feature의 회귀계수는 0으로 만들어서 해당 feature이 선택되지 않게끔 한다.

- ④ 엘라스틱 (ElasticNet)  
→  $L_2 + L_1 + ①$   
→ feature이 많은 모델 사용

분류에 사용

- ⑤ 로지스틱회귀 (Logistic Regression)  
→ 사실 classification에 사용되 Linear model이다.  
→ 이진분류 + 희소 (sparse) 회귀 (ex. 특정도환기)에서 강하게

## 5-2. 다변량 회귀를 통한 회귀식

→ 다변량 회귀는 독립변수와 종속변수가 연속적인 선형 관계이다.

→  $\hat{y} = w_0 + w_1 x$   $x$  = 예측값 (=데이터) → 실제와 딱 가깝게 이리한 함수관계에서  
실제값만큼의 값을 빼거나 더함값이 된다.  
절편  
(=intercept)

$$RSS(w_0, w_1) = \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + w_1 x_i))^2$$

→ 이리한 함수, 즉 비용함수가 반환하는 값은 결국 적은 값이므로  
(=최소값)

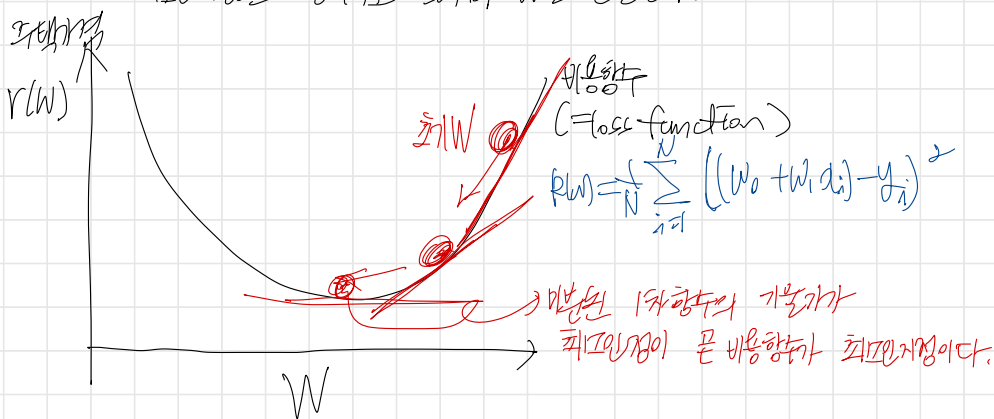
최소값으로 더이상 감소하지 않는 최선의 상태를 구하는것이 목적이다.

## 5-3. 비용최소화하기 - 경사하강법 (Gradient Descent)

- 비용함수가 최솟가 되는  $W$  parameter를 구해야 한다.
- 특히나  $W$  parameter의 개수가 많다면 고차방정식을 이용해서도 해결하기 힘들기 때문에 "점진적인 하강"을 이용하는 gradient descent를 사용하고자 한다.

→ 비용함수가 2차항의 형태라면 경사하강법, 2차원에서의 미분을 적용한 후에 이 미분값을 제곱 항의 상수 항으로 대체하고 일차식으로  $W$ 를 update한다.

→ 특히나 미분된 1차항의 기울기가 양도하지 않는 H점까지를 비용함수가 최솟점임을 간파하고 그때의  $W$ 를 반환한다.



$$\frac{\partial R(W)}{\partial w_0} = \frac{2}{N} \sum_{i=1}^N -(y_i - (w_0 + w_1 x_i)) = -\frac{2}{N} \sum_{i=1}^N (y_i - \text{예측값})$$

$$\frac{\partial R(W)}{\partial w_1} = \frac{2}{N} \sum_{i=1}^N -x_i \times (y_i - (w_0 + w_1 x_i)) = -\frac{2}{N} \sum_{i=1}^N (y_i - \text{예측값}) \times x_i$$

→ 그러나 일반적인 gradient descent의 경우에는 결과에 따라서는  
Stochastic Gradient Descent (확률적 경사하강법)도 사용

→ 전체 데이터에서 랜덤하게 batch-size만큼 데이터를 추출해서  
 이들을 기반으로  $w_1$ -update,  $w_0$ -update의 값을 계산하게 된다.

$$\begin{array}{c}
 \hat{Y} \\
 \downarrow \\
 Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \text{shape}(n, 1)
 \end{array}
 \quad
 \begin{array}{c}
 X_{\text{mat}} \\
 \begin{array}{c} \text{Feature 1} \quad \text{Feature 2} \quad \dots \quad \text{Feature M} \\
 \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{bmatrix} \\
 (M, n)
 \end{array}
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 * \\
 \downarrow \\
 (M, n) * (n, 1) \rightarrow (M, 1)
 \end{array}$$

$$[w_1 \ w_2 \ \dots \ w_m]^T + w_0 \quad (M, 1)$$

$$= (M, 1)$$

\* 표는  $w_0$ 를 어떤 constant로  $X_{\text{mat}}$ 의 Feature 0을 (1로 초기화) 하기  
 위해  $w_0$ 를  $w$ 배열에 앞에 추가해놓은 것이다.

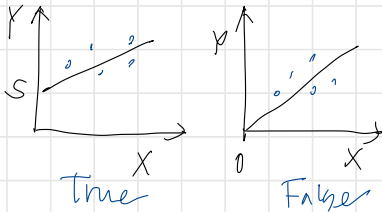
## 1-4. \*learn.linear-model.LinearRegression \*용하기 보편회귀분석 예측

→ 보편회귀 분석의 가장 기본적인 알고리즘을 위해 OLS (ordinary least squares) 라는 방법을 구현한 class 이다.

→ fit() 메서드는 x, y 배열을 입력받으면 리턴해주는 W를 coef\_ 변수에 저장한다.

### \* PARAMETERS

#### ① fit\_intercept



#### ② normalize

→ fit\_intercept = True 일 때  
데이터는 자동적으로  
평균제거이치를 갖게 된다

### \* 특성

#### ① coef\_

→ fit() 메서드를 수행했을 때  
리귀분석의 매개변수를 저장하는 특성.  
(Target 매, feature 매)

#### ② intercept\_

→ ordinary Least Squares 기반의 회귀계수 계산을 각 feature의 독립성에 영향을 많이 받기 때문에 feature들간의 상관관계가 높은 경우 회귀이커져서 값에 안정해진다.

≡ 다중공선성 (multi-collinearity) 문제 발생

→ 따라서 상관관계가 높은 feature이 많은 경우 독립적인 필요한 feature만 남기고 제거해서 문제를 직면한다.

→ 또는 PCA를 이용해 차원축소를 해준다.

### 회귀평가 지표

→ 회귀의 오차를 측정하는 지표는 여러가지가 있지만 대표적인 것만 소개한다.

$$MAE \quad \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i|$$

$$MSE \quad \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

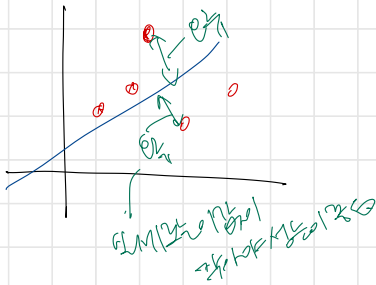
$$RMSE \quad \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2} = \sqrt{MSE}$$

$$R^2 \quad \frac{\text{회귀값 분산}}{\text{실제값 분산}}$$

⊕ 회귀 평가 지표를 적용할 때의 주의사항

MAE → Scoring 함수 적용 →  $\text{neg\_mean\_absolute\_error}$

MSE → Scoring 함수 적용 →  $\text{neg\_mean\_squared\_error}$



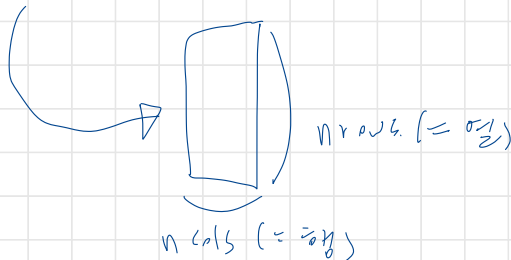
sklearn의 scoring 함수는 score이 높을수록  
좋은 모델이라고 가정하고 평가하는데,  
값이 낮을수록 좋은 모델이라고 가정하는 차이가  
있는데 sklearn의 경우 나쁜 모델이므로  
sklearn의 scoring 함수에 반전시켜서  
보정해주어야 한다. →  $-1$ 은 음수화 연산  
(최대값이 최선이)

cf. sklearn의 Scoring 함수

⇒ cross\_val\_score, GridSearchCV

cf. matplotlib을 이용한 여러개의 plot 그리기

→ matplotlib.subplots()는 여러개의 그래프를 한번에 그리기 위해  
사용되는데, 인자로 ncol과 nrow를 입력





# ⑧ 회귀분석 // 사례별 모형에서의 예측오차 측정방법

$$1. ME = \frac{\sum_1^n (Y_i - \hat{Y}_i)}{n}$$

$$2. RMSE = \sqrt{\frac{\sum_1^n (Y_i - \hat{Y}_i)^2}{n}} \quad (\text{구체적 평가 내용})$$

$$3. MAE = \frac{\sum_1^n |Y_i - \hat{Y}_i|}{n}$$

$$4. MPE = \frac{\sum_1^n \frac{Y_i - \hat{Y}_i}{Y_i} \times 100}{n}$$

$$5. MAPE = \frac{\sum_1^n \frac{|Y_i - \hat{Y}_i|}{Y_i} \times 100}{n}$$

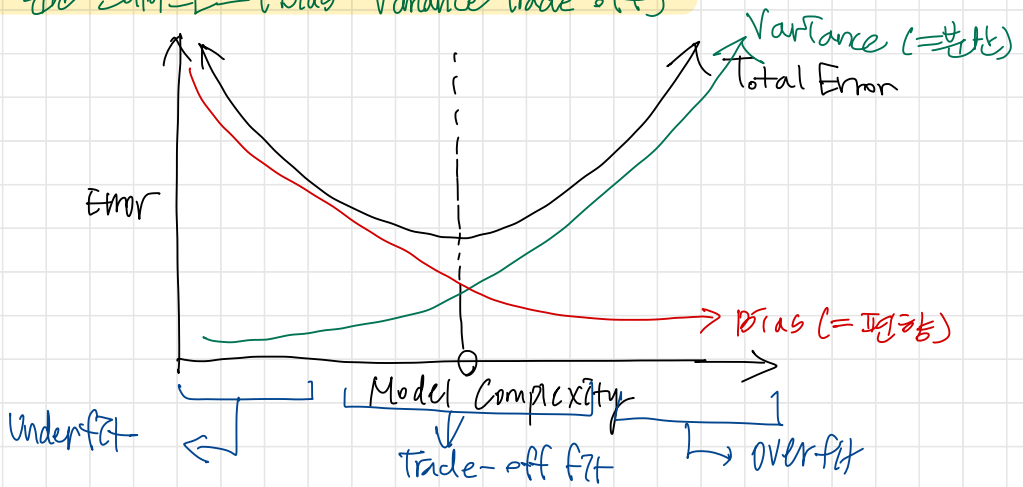
$$6. MASE = \frac{\sum_1^n \left( \frac{|Y_i - \hat{Y}_i|}{\frac{\sum_1^n |Y_i - Y_{i-1}|}{n-1}} \right)}{n} = \frac{\sum_1^n |Y_i - \hat{Y}_i|}{\frac{n}{n-1} \times \sum_2^n |Y_i - Y_{i-1}|}$$

⇒ RMSE, MAE와 같은 지표들은 모두 강함과 예측량의 적대적인 예측오차에 기반하여 예측오차의 크기를 나타내는데, 절대적인 예측오차 뿐만 아니라 상대적인 의미의 예측오차도 필요하다.  
→ 예를 들면 강함 10에 대한 예측오차 1과 강함 100에 대한 예측오차 1은 상대적대로 크기가 다르다.

이제는 이동평균이나 관계한도 등이 세상에 유행하는 것보다 이 이원화라는 MPE의 문제점을 개선한 방법이다.

→ MASE 데이터의 평균 1에 가까울수록 좋다는 관점은 보단 나쁘게 해서 데이터를 적소라 하여 많은 비난을 받는다.

# \*) 편향-분산 트레이드오프 (Bias-Variance Trade-off)



지나치게 편향함성으로 치우친 경향이 있는 모델을 **고편향 (High Bias)** 라고 한다.

각종 하나 하나의 특징을 반영하면서 복잡함, 고차원 인코딩 되고 지나치게 높은 분산성을 가진 모델은 **고분산 (High Variance)** 라고 한다.

Low Variance

Low Bias

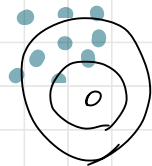


High Bias



분산이 커지면 예측의 분포가 넓어진다.

High Variance



편향이 높거나 분산이 높을 때 예측에서 벗어난다.

∴ 머신러닝의 목표

= 편향과 분산으로 trade-off을 이루면서 cost가 작아지도록 하는 모델을 구축하는 것. → ∴ 규제 (regularization)

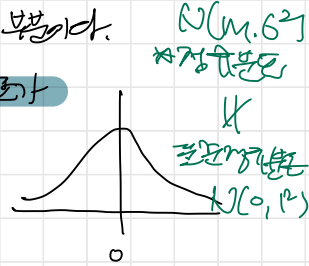
# \* 스미트위더데이터에 적용 - kagle Bike Sharing Demand

→ target값의 분포가 대략적인 형태를 이루고 있어서 적합해볼 수 있음

특히에서 나온 예측 결과가 범용성 때문에 가정인지 모르겠음 하는 부분이다.

이때 정규분포를 이루고 있다면 독립성에 차인 분포를 예측할 분포가

아기 때문에 보통 로그를 적용해 변환을 하더라.



→ sklearn을 이용한다면 이렇게 변경된 target값을  
기반으로 학습하고 예측한 값을 다시  $\exp(x)$  함수를 적용하여  
원래 scale로 원상복구하면 된다.



④ 스미트위더데이터는 이치에도

MinMaxScaler (최대값을 1, 최소값을 0으로  
변환해주는 방법)

StandardScaler (평균이 0, 분산이 1인 표준정규분포의  
형태로 변환함)

대부분의 데이터가  
크게만 있다.

feature의 개수가  
많은 경우에는

다양성으로 정상적인

feature의 개수가 너무 많아지면  
overfitting이 발생한다.

스케일링/정규화는 다양한 데이터셋에 대해  
이항 특성을 적용하여 범용성을 높인다

(대부분의 데이터는 0이므로 = 로그를 함)

따라서 log는 '종에서' scale을 바꿔준다

로그변 환은 이치에서 정규분포로 최대한 가깝게

형태로 바꿔준다.

(\*) 데이터의 가용을 하여 모델 학습을 추가적으로 해가되해 고려하는

① Feature dataset의 데이터 분포도 → 유사도 높은 data () 행을

이용해 같은의 대응된 행들을 나눠서,  
이때 반환값이 1 이상일때 대응행도가  
높다고 판단을 일반적으로 한다.  
(만, 원-핫인코딩된 데이터를 볼때는 안된다)

② 이상치 데이터 처리  
(outlier)

→ 특이치 예측을 하게 있어서 target data가  
상당히 커서 높은 데이터인 경우에는 이상치를 제거하는 방법으로도  
크로스로 풀었다.

(\*) Stratified K-Fold

— k-fold의 경우를 k개의 dataset을 만들고 k번 반복한다.

→ 이에 stacking 모델은 차이가 필요하다.

① 개별적인 기반 모델

② 이 개별기반 모델의 예측데이터를 다시 stacking의 형태로 결합하여  
크로스 meta model의 학습 feature set과 테스트 feature set을 만든다

# \*) CV set 기반의 stacking

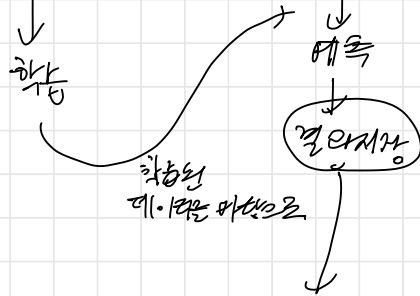
(stratified k-fold)을 여러 번 반복해서

(과제하는 개변수에 대해서 회귀계수 모형을 위한 데이터는 많은 회귀공제본으로 예측된 결과 데이터셋을 이용)

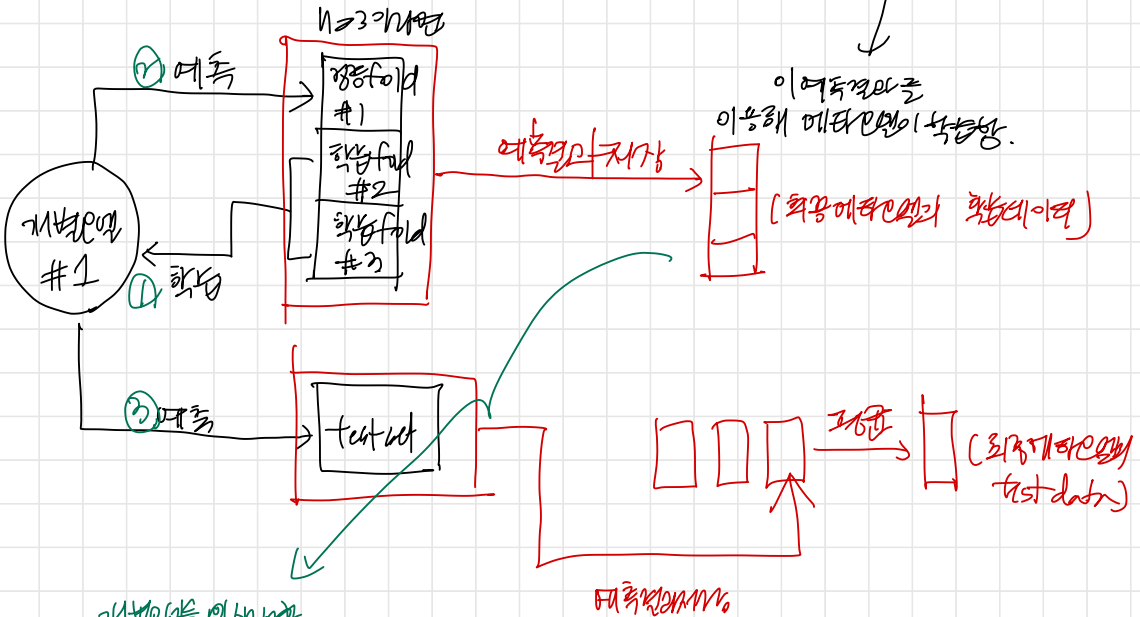
① 회귀 데이터를  $m$ 개의 fold로 나눔

② 아래의 반복에서 개변수의 예측값으로 학습데이터와 test data 생성

① 학습데이터를  $m$ 개 중  $n-1$ 개 학습을 위해, 나머지 1개는 검증용 데이터로



이 예측값을 이용해 메타모델이 학습됨.



개변수를 여러 번

데이터 = 회귀계수 모형을 위한 데이터 사용

개변수들의 test data와 = 회귀계수 모형을 위한 test data 사용

# ★정리

① 비용함수 (손실함수 또는 비용함수)를 최소화하는 매개변수의 값을 찾는 것이 목적

→ 이를 위해 평가방법을 사용

→ 비용함수의 최적화 기법

But 학습데이터에 overfitting의 문제가 발생 가능

해결책

→ Regression

$$L_{\text{Linear}}(L) \Rightarrow \text{RSS}(W) + \alpha \times \|W\|_1$$

$$L_{\text{Ridge}}(L) \Rightarrow \text{RSS}(W) + \alpha \times \|W\|_2^2$$

$$L_{\text{Elastic}}(L) \Rightarrow \text{RSS}(W) + \alpha_1 \|W\|_1 + \alpha_2 \times \|W\|_2^2$$



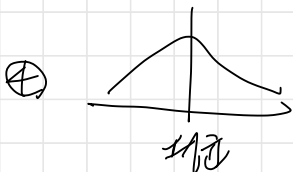
이제 원소라 하든지!!

② Logistic 회귀

→ 선형회귀와 마찬가지로 회귀의 Sigmoid 함수를 도출해 특징(feature)이 binomial 분포를 갖는 확률값을 하여 여러 label 표현



③ 히스토그램 — 2D 또는 3D data를 표현하기 위해 사용



회귀분석에 대한 분포

→ 이차원 lognormal 분포

⊕ object 데이터는 2D 또는 3D (특히 2D)