



교육 일정

1일차 프로젝트 개요 및 데이터셋 이해

2일차 AI 적용을 위한 데이터셋 전처리

3일차 금속분말 생산공정 최적화를 위한 선형회귀 기법

4일차 금속분말 생산공정 최적화를 위한 비선형회귀 기법

5일차 금속분말 생산공정 최적화를 위한 딥러닝 기법 심화

6일차 AI 기법 성능 향상 방법론



머신러닝을 활용한 나노/탄소 소재 생산공정 최적화

03 금속분말 생산공정 최적화를 위한 머신러닝 기법





머신러닝을 활용한 나노/탄소 소재 생산공정 최적화

- 01 선형 회귀 분석
- 02 선형회귀모델
- 03 선형회귀모델의평가
- 04 선형회귀모델의구현

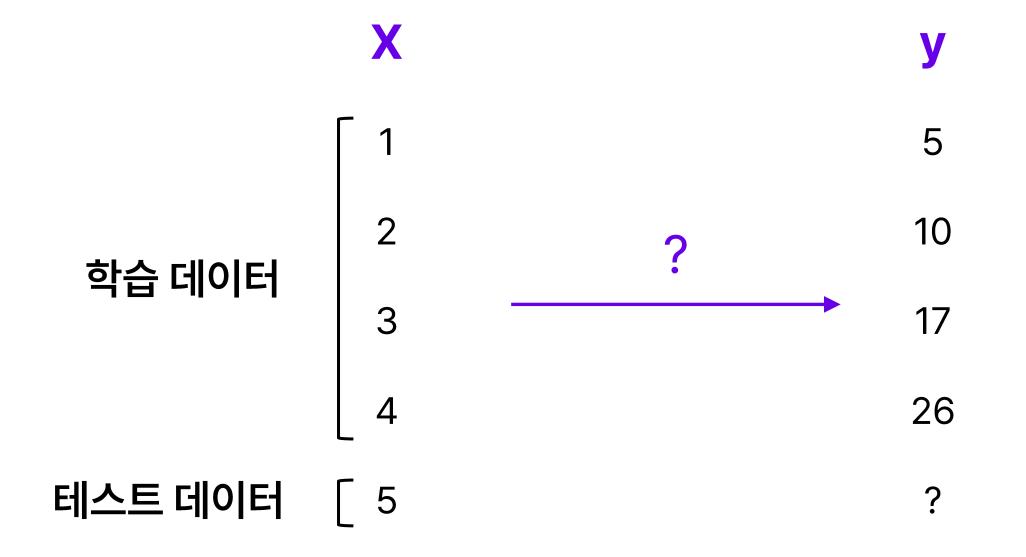


01 선형회귀분석



ਂ 회귀 분석

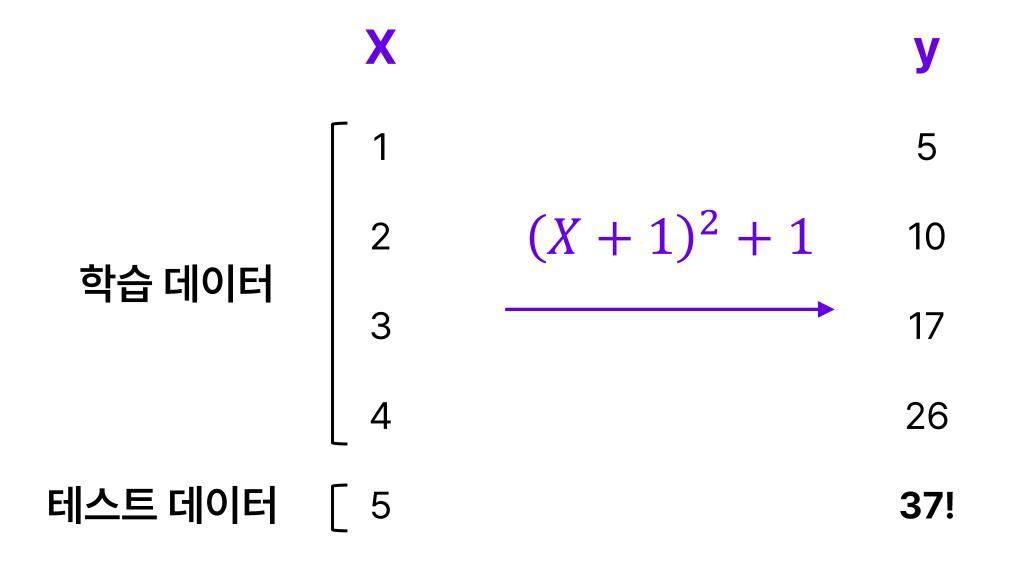
• 입력 변수와 출력 변수 사이의 관계를 모델링하여, 입력값에 대한 출력값을 예측하는 일.





ਂ 회귀 분석

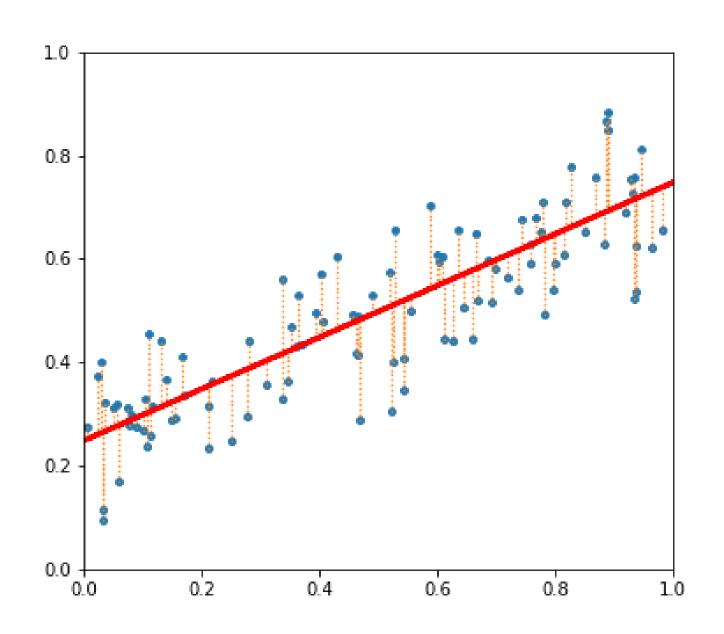
• 입력 변수와 출력 변수 사이의 관계를 모델링하여, 입력값에 대한 출력값을 예측하는 일.





❷ 회귀 분석

• 회귀의 기본 원칙은 잔차를 최소화 하는 것



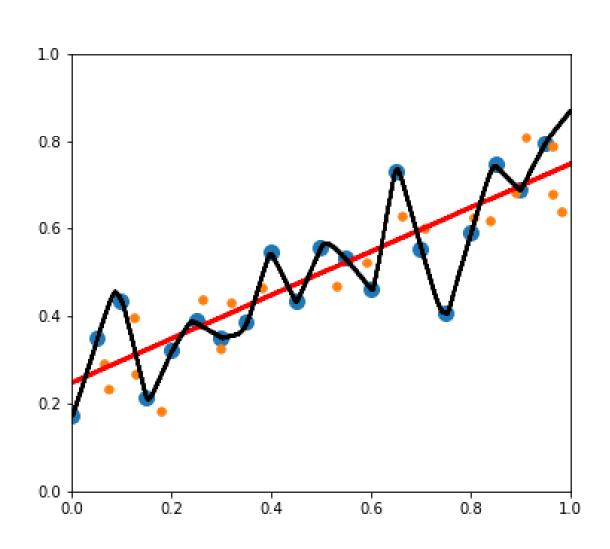
• 파란 점: 학습 데이터

• 빨간 선: 회귀 모델

• 주황 점선: 잔차



❷ 과적합



- 파란 점: 학습 데이터
- 주황 점: 테스트 데이터
- 검은 선: 과적합 된 모델
 - 학습 오차: 0
 - 테스트 오차: 0.19
- 빨간 선: 일반화가 잘 된 모델
 - 학습 오차: 0.07
 - 태스트 오차: 0.08



❷ 완벽한 모델을 만들 수는 없을까요?

- 실제 데이터에는 **노이즈**가 있을 수 밖에 없음
- All models are wrong, but some are useful!
- 완벽한 모델이라는 것은 없음!
- •데이터 분석을 통해 제한된 시간과 자원으로 **적절한 모델**을 찾는 것이 중요



❷ 선형 회귀 모델

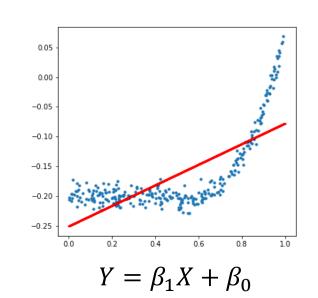
• 모델의 변수 β_i 와 Y가 선형 관계.

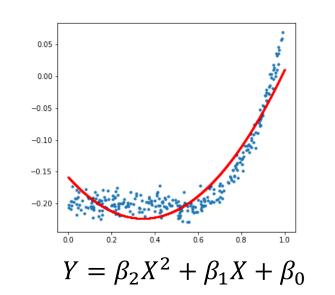
•
$$Y = \beta_1 X + \beta_0$$

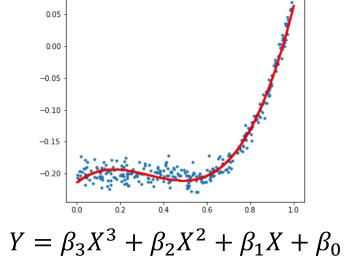
•
$$Y = \beta_2 X^2 + \beta_1 X + \beta_0$$

•
$$Y = \beta_3 X^3 + \beta_2 X^2 + \beta_1 X + \beta_0$$

•







- \bullet 목적: 잔차를 최소화하는 β 를 찾는 것.
- 장점: 빠르고 효율적이며, 최적화가 쉽다.
- 단점: 복잡한 관계를 모델링하기 어렵다.
- e.g., 최소 제곱 선형 회귀 모델, Ridge 회귀 모델



❷ 선형 회귀 모델

독립변수 종속변수	1개	2개 이상
1개	단변량 단일 선형 회귀	단변량 다중 선형 회귀
2개 이상	다변량 단일 선형 회귀	다변량 다중 선형 회귀



❷ 단변량 단일 선형 회귀

- •독립 변수: 1개
- 종속 변수: 1개
- e.g., 키로부터 몸무게 예측
- $\bullet y = \beta_1 x + \beta_0$



♥ 단변량 다중 선형 회귀

•독립 변수: 2개 이상

• 종속 변수: 1개

• e.g., 키와 체지방량으로부터 몸무게 예측

•
$$y = \beta_N x_N + \beta_{N-1} x_{N-1} + \dots + \beta_1 x_1 + \beta_0$$



○ 다변량 단일 선형 회귀

•독립 변수: 1개

• 종속 변수: 2개 이상

• e.g., 키로부터 체중과 체지방량 예측

•
$$[y_N, y_{N-1}, \dots, y_1] = \overrightarrow{\beta_1} x_1 + \overrightarrow{\beta_0}$$



ਂ 다변량 다중 선형 회귀

•독립 변수: 2개 이상

• 종속 변수: 2개 이상

• e.g., 키와 체중으로부터 체지방량과 골격근량 예측, 금속분말 데이터셋

$$\bullet \left[y_{N_{y}}, y_{N_{y}-1}, \dots, y_{1} \right] = \left[\overrightarrow{\beta}_{N_{y}}^{T} \left[x_{N_{x}}, x_{N_{x}-1}, \dots, x_{1}, 1 \right], \overrightarrow{\beta}_{N_{y}-1}^{T} \left[x_{N_{x}}, x_{N-1}, \dots, x_{1}, 1 \right], \dots, \overrightarrow{\beta}_{1}^{T} \left[x_{N_{x}}, x_{N_{x}-1}, \dots, x_{1}, 1 \right] \right]$$



❷ 선형 회귀에서 피쳐 선택

$$x_1$$
 x_2 지 x_2 지 x_2 x_1x_2 x_1^2 x_1^3 x_2 x_1^3 x_2 x_2 지 x_2 지 x_2 지 x_2 x_2 지 x_2 지

입력값을 조합하여 피쳐를 늘림

늘린 피쳐 중 의미있는 피쳐 선택



02 선형회귀모델

02 선형 회귀 모델



❷ 선형 회귀 모델의 목표

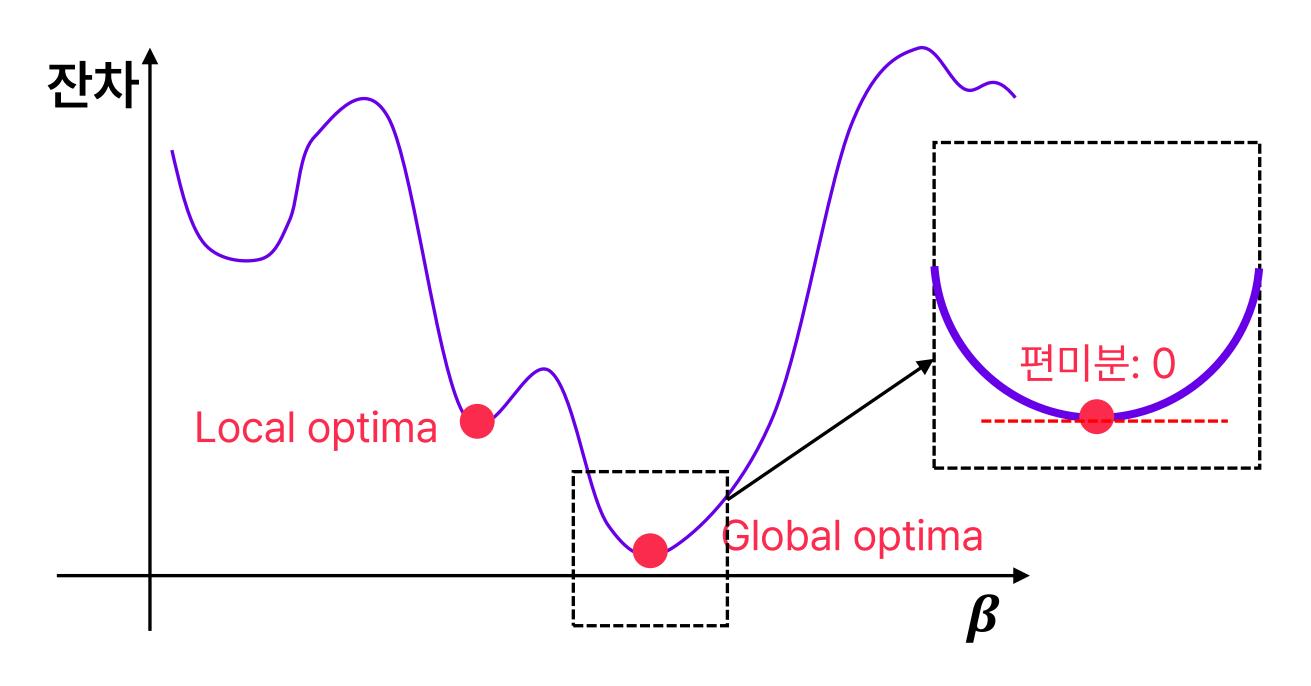
- 잔차를 최소화 하는 모델의 변수 β 를 찾는 것
- 과적합 되지 않고 일반화가 잘 되어야 한다.
- \bullet 다변량 다중 선형 회귀 모델에서의 eta

$$\bullet \ [y_N, y_{N-1}, \dots, y_1] = \left[\overrightarrow{\beta}_N^T[x_N, x_{N-1}, \dots, x_1], \overrightarrow{\beta}_{N-1}^T[x_N, x_{N-1}, \dots, x_1], \dots, \overrightarrow{\beta}_1^T[x_N, x_{N-1}, \dots, x_1]\right]$$

•
$$\beta = [\overrightarrow{\beta}_{N}^{T}, \overrightarrow{\beta}_{N-1}^{T}, ..., \overrightarrow{\beta}_{1}^{T}]$$



☑ 머신러닝의 최적화 방법론



- 학습 변수의 변화에 따라 잔차가 정해진다.
- Global optima 지점의 β 를 찾는 것이 최적화이다.
- Global optima 지점에서 잔차를
 β에 대해 편미분하면 0이 됨을 이용한다.
- 그러나, 그 지점이 local optima 일수도 있다.



최소 제곱 선형 회귀 (Linear least square regression)

- $X = \{X_1, ..., X_N\}, y = \{y_1, ..., y_N\}, X$ 는 입력데이터, y는 출력데이터, N은 데이터 개수.
- $X_i = \{x_{i1}, ..., x_{id}\}, d$ 는 입력데이터 차원 수.
- $\hat{y}_i = f(X_i; \beta) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + ... + \beta_d x_{id} = \beta_0 + \beta^T X_i$.

- $\hat{y} = X\beta$
- $e = y \hat{y} = y X\beta$ 에서 $e^T e = \sum_i^N (y_i - \hat{y}_i)^2 = (y - X\beta)^T (y - X\beta)$ 를 최소화하는 $\beta = \{\beta_0, \dots, \beta_d\}$ 최적화.
- $e^T e$ 를 잔차제곱합 (RSS: residual sum of squares)라고 한다.



최소 제곱 선형 회귀 (Linear least square regression)

- $e^T e$ 를 최소화하는 $\beta = \{\beta_0, ..., \beta_d\}$ 최적화 $\rightarrow e^T e = \beta$ 로 편미분하여 0이 되는 $\beta =$ 찾자!
- $e^T e = (y X\beta)^T (y X\beta) = y^T y 2X^T y\beta + \beta^T X^T X\beta$

$$\cdot \frac{d(e^T e)}{d\beta} = -2X^T y + 2X^T X \beta = 0$$

- $\bullet X^T X \beta = X^T y$
- 이 때, X^TX 가 역행렬을 가진다면 최적의 β 인 β^* 는 $\beta^* = (X^TX)^{-1}X^Ty$



☑ Ridge 선형 회귀 모델 (Ridge regression)

- 최소제곱 선형 회귀 모델은 **과적합**이 발생할 수 있음.
- Ridge 회귀 모델은 잔차제곱합 (RSS)을 최소화하는 대신, 아래의 식을 최소화하는 것을 목표로 한다.

$$RSS + \lambda \sum \beta^2$$
, λ 는 작은 상수.

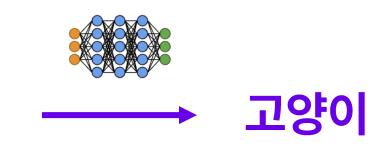
- 즉, 목적 함수에 **학습 계수의 제곱합**을 더하여 특정 계수가 너무 커지는 것을 방지하여 과적합을 예방한다.
- 위와 같이 학습 계수의 제곱합을 목적함수에 더하는 것을 **Ridge regularization** (Ridge 정규화)라고 한다.
- 학습 계수의 절대값의 합을 목적함수에 더하는 것은 **Lasso regularization** (Lasso 정규화)라고 한다.



❷ 비선형 회귀 모델

- •모델의 변수 β_i 와 Y가 비선형 관계.
 - $\bullet Y = \frac{\beta_0 X + \beta_1}{\beta_2 e^X + \beta_3}$
 - •
- 복잡한 관계를 모델링하기 적합하다.
- 선형 회귀 모델보다 복잡하며, 최적화가 어려울 수 있다.
- e.g., 딥러닝 (Multi layer perceptron, convolutional neural network)







03 선형회귀모델의평가



◎ 회귀 평가 지표

- 학습한 선형 회귀 모델의 성능 평가를 위한 지표
- 실제값과 모델이 예측한 값의 차이에 기반한 평가 방법 사용
- e.g., RSS, MSE, MAE, MAPE, R²



☑ RSS (잔차제곱합, Residual Sum of Squares)

- •전체 데이터에 대한 실제 값과 예측 값의 오차 제곱 합
- RSS = $\sum (y_i \overline{y}_i)^2$
- 장점
 - 가장 간단한 평가 방법
 - 직관적 해석 가능
- 단점
 - 데이터의 개수에 의존적임
 - 데이터의 스케일에 의존적임



☑ MSE (평균제곱오차, Mean Squared Error)

- •전체 데이터에 대한 실제 값과 예측 값의 오차 제곱 평균
- MSE = $\frac{RSS}{N}$
- 이상치에 민감함
 - 실제값과 크게 다르게 예측한 값에서 오차가 커짐
- 장점
 - 간단한 방법으로 직관적 해석 가능
 - RSS와 달리, 데이터의 개수에 의존적이지 않음
- 단점
 - 여전히 데이터의 스케일에 의존적임



☑ MAE (평균 절대값 오차, Mean Absolute Error)

- •전체 데이터에 대한 실제 값과 예측 값의 오차 절대값 평균
- MAE = $\frac{\sum |(y_i \overline{y}_i)|}{N}$
- 이상치에 덜 민감함
 - 변동성이 큰 지표와 낮은 지표를 함께 예측할 시 유용
- 장점
 - 간단한 방법으로 직관적 해석 가능
 - RSS와 달리, 데이터의 개수에 의존적이지 않음
- 단점
 - MSE와 마찬가지로 여전히 데이터의 스케일에 의존적임



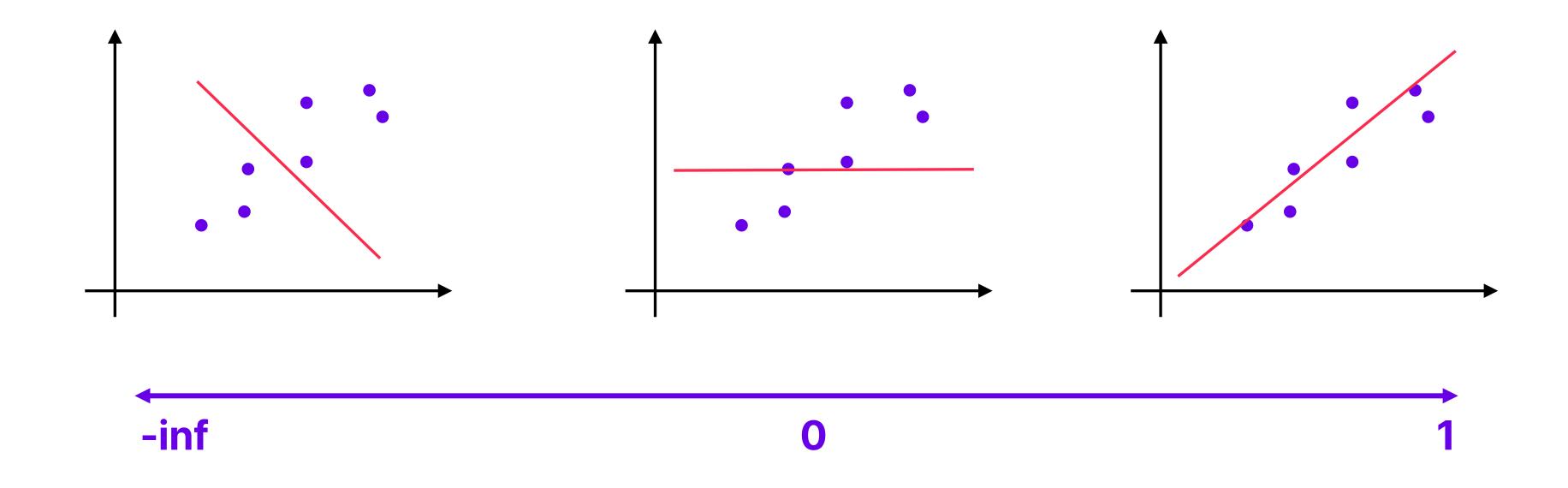
• 회귀 모델의 설명력을 표현하는 지표

•
$$R^2 = 1 - \frac{RSS}{TSS}$$

 \bullet TSS 는 데이터 평균값 ($\mu_{
m y}$) 과 실제값 차이의 제곱합

• TSS =
$$\sum (y_i - \mu_y)^2$$

- 해석
 - 1에 가까울수록 설명력이 좋음
 - 0이면 데이터의 평균값으로 예측한 모델
 - 음수이면 평균값 예측보다 성능이 좋지 못한 모델



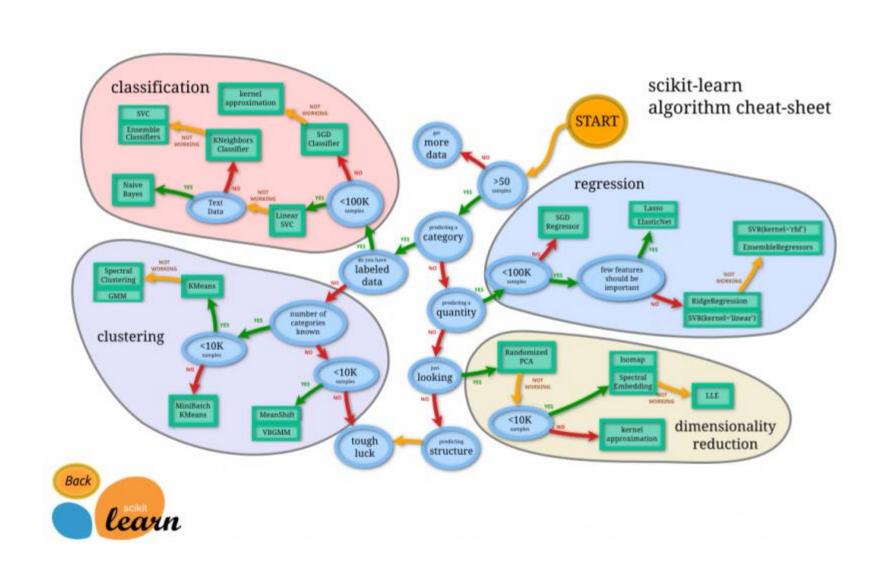


04 선형회귀모델의구현



❷ 사이킷런 (scikit-learn) 라이브러리

- Python 을 대표하는 머신러닝 라이브러리
- 오픈소스
- 다양한 머신러닝 알고리즘을 쉽게 구현 가능





❷ 최소 제곱 선형 회귀 모델 구현

```
from sklearn.linear_model import LinearRegression

fitter = LinearRegression() # 최소 제곱 선형 회귀 모델 생성

fitter.fit(X_train, y_train) # 모델 학습

pred = fitter.predict(X_test) # 테스트 데이터에 대해 예측
```



☑ Ridge 선형 회귀 모델 구현

```
from sklearn.linear_model import Ridge

fitter = Ridge(alpha = 1.0) # Ridge 선형 회귀 모델 생성, alpha는 시의 값
fitter.fit(X_train, y_train) # 모델 학습
pred = fitter.predict(X_test) # 테스트 데이터에 대해 예측
```



☑ RSS, MSE, MAE의 구현

```
def RSS(gt, pred):
   return ( (gt - pred) ** 2 ).sum()
def MSE(gt, pred):
   data_len = len(gt)
   return ( (gt - pred) ** 2 ).sum() / data_len
def MAE(gt, pred):
   data_len = len(gt)
   return ( np.abs(gt - pred) ).sum() / data_len
```




```
from sklearn.metrics import r2_score

y_pred = model.predict(X)

r2 = r2_score(y, y_pred)
```



• 독립변수에 의해 영향을 받는 값으로, 머신러닝 모델에서 출력값에 해당하는 변수를 무엇이라 부르

나요?

• 종속변수



• 머신러닝 모델 학습에서 찾고자 하는 지점으로, loss를 학습계수로 편미분 하였을 때 0이 되면서 Local optima 가 아닌 지점을 무엇이라 부르나요?

Global optima



• Minima 지점에서 가중치에 대한 목적함수의 편미분 값은 무엇인가요?

• 0



- 선형 회귀 분석의 평가 방법으로 올바르지 않은 것은 무엇인가요?
 - ① RSS
 - ② R² score
 - 3 MSE
 - 4 STD