

# Supervised learning I.

Linear Regression : 회귀의 기본 모델.

$\beta_0$ 은 유리한가?

선형 회귀 모델의 파라미터인 계수에 대해서

$\beta_0$ 은  $X$ 의 값으로 0인 경우  $Y$ 와 무관해짐.

통제리스트를 진행할 수 있음.

이를 T-test로 확인할 수 있음.

유용한 변수 선택에 대해 알아볼 수 있음.

$H_0: \beta_0 = 0, H_A: \beta_0 \neq 0$ , t값을 구하고.

비선형적인 데이터를 어떻게 포착할지 알 수 있음.

T분포 내에서 t값을 통해 p-value를 찾을 수 있음.

회귀는 데이터 포인트를 가장 잘 설명하는 회귀선을

모델과 데이터의 적합도를 측정하는 지표

찾아 주는 것임. 모델에서는  $Y$ 를 제로로 설명하는

$RSS = \sqrt{\frac{RSS}{n-2}}$ . (MS의 좋은 측정치임)

회귀계수와 관련이 있음. 이를 확인할 수 있음.

$R^2 = 1 - \frac{RSS}{TSS}$

고로, 모델의 MS를 통해 회귀계수를 찾아야 함.

TSS는 총제곱합, RSS는 오차합계의 제곱

RSS는 잔차제곱합으로, MS에서  $\frac{1}{n}$ 만을 제거한

$R^2$ 이 1에 가까워 수록 좋음.

$$RSS = \sum (y_i - \hat{y}_i - \beta_0 - \beta_1 x_i)^2 = \sum e_i^2$$

RSS가 2차 함수이므로 이를 미분하여 각 계수에 대한 최적의 선형회귀

대신 최적의 값을 찾아낼 수 있음.

$$Y \approx \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

이렇게 얻어진 회귀 계수는 모델로 부터 얻은

여기에서 RSS를 최소화하는  $\beta$ 를 찾음.

것이기 때문에 모델의 회귀계수의 측정값이

$$y = X\beta + e$$

측정된 회귀계수가 모델의 회귀계수와 일치

$$\rightarrow y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}, e = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}$$

동일한지는 통계적 문제임.

$$RSS = (y - \hat{y})^T (y - \hat{y})$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

모델로 찾은  $\hat{\beta}$ 과 모델의  $\beta$

$$\hat{y} = X (X^T X)^{-1} X^T y$$

$\mu$ 를 기준으로 신뢰구간을 찾을 때는 아래와 같음.

$$RSE = \sqrt{\frac{RSS}{n-p-1}}$$

$$p = Pr[\hat{\mu} - k SE(\hat{\mu}) < \mu < \hat{\mu} + k SE(\hat{\mu})]$$

SE matrix를 통해 각  $\beta$ 의 SE를 얻을 수 있음.

T-test를 통해서 실제  $\mu$ 가  $\hat{\mu}$ 인지를 확인함.

$\beta_0$ 은  $\hat{\beta}_0 \pm 1.96 SE(\hat{\beta}_0)$  사이에 95% 있음.

T-test는  $(\hat{\mu} - \mu) / SE(\hat{\mu})$  값을 관찰할 확률을

이후 T-test를 통해  $\beta$ 가 0인지 확인함.

T분포에서 확인함으로써 p-value를 구함.

이런 신뢰구간을 통해서 모델의 회귀계수와 변수 선택의 문제.

변수가 많을 수록 유망성이 높아짐.

$\hat{\beta} \pm k SE(\hat{\beta})$  범위 내에 몇 %로 존재한다고

훈련세트에서 RSS를 줄일 수 있음. 과대적합되기

결론 내릴 수 있음

수는 문제정답이 있음.

이때 변수가 유용한지 파악하기 어렵다

사용하는 것이 장래됨.