

Regression

: $X-Y$ 간 일차선형 관계 및 인과관계를 표상하는 식을 찾는 것

$$\hat{y} = \beta x + \beta_0$$

How to find? → 정주정

기분가게

① 오차들 e_i 는 확률변수이며 서로 독립 & 평균 0 분산 σ^2

② Y 는 서로 독립 & 평균 $\mu = \alpha + \beta X$, 분산 σ^2

$$L = \frac{1}{2} \sum (y_i - \hat{y}_i)^2 = \frac{1}{2} \sum (y_i - \beta x - \beta_0)^2$$

$$\textcircled{1} \frac{\partial L}{\partial \beta_0} = \sum (y_i - \beta x - \beta_0) (-1) = 0$$

$$= -\sum y_i + \beta_1 \sum x + \sum \beta_0 = 0$$

$$\rightarrow \hat{\beta}_0 = \frac{1}{N} (\sum y_i - \beta_1 \sum x) = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\textcircled{2} \frac{\partial L}{\partial \beta_1} = \sum (y_i - \beta_1 x - \beta_0) (-\sum x) = 0$$

$$= -\sum y_i \sum x_i + \sum \beta_1 x \sum x + \sum x \beta_0$$

$$= -\sum y_i \sum x_i + \sum \beta_1 x \sum x + \sum x_i (\bar{y} - \hat{\beta}_1 \bar{x})$$

$$= -\sum y_i \sum x_i + \hat{\beta}_1 \sum x_i \sum x_i + \bar{y} \sum x_i - \hat{\beta}_1 \bar{x} \sum x_i$$

$$\rightarrow \hat{\beta}_1 = \frac{\sum (y_i - \bar{y})}{\sum (x_i - \bar{x})} = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}$$

$\hat{\beta}_0, \hat{\beta}_1$ 로 식 구할 수 있다!

Ⓟ 귀환식이 통계적일 수 있는지 검증이 중요!

How to find? → 구간추정

기분가게

① 오차들 e_i 는 확률변수이며 $e_i \sim N(0, \sigma^2)$

② Y 는 서로 독립 & $Y \sim N(\alpha + \beta X_i, \sigma^2)$

① β 추정

$$\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{S_{xx}}\right) \quad \sigma^2 \approx s^2 \quad s^2 = \frac{SSE}{n-2} \quad \text{2 변인}$$

$$\left[\hat{\beta} - t_{\alpha/2, n-2} \sqrt{\frac{SSE}{n-2} \frac{1}{S_{xx}}}, \hat{\beta} + t_{\alpha/2, n-2} \sqrt{\frac{SSE}{n-2} \frac{1}{S_{xx}}} \right]$$

② α 추정

$$\hat{\alpha} \sim N\left(\alpha, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)\right)$$

$$\left[\hat{\alpha} - t_{\alpha/2, n-2} \sqrt{\frac{SSE}{n-2} \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}, \hat{\alpha} + t_{\alpha/2, n-2} \sqrt{\frac{SSE}{n-2} \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \right]$$

오차 제곱합의 추정량은 "불편추정량"

∴ 추정해야 하므로 $df = n-2$

σ^2 이 추정량인 $S^2 = \frac{SSE}{n-2} = MSE$

$$E(S^2) = E\left[\frac{SSE}{n-2}\right] = \sigma^2 \rightarrow \text{bias X.}$$

$$S_{xx} = \sum (x_i - \bar{x})^2$$

$$S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$S_{yy} = \sum (y_i - \bar{y})^2$$

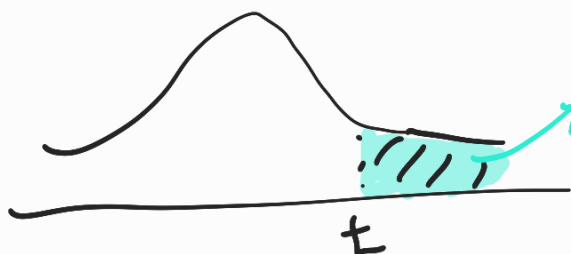
How to test?

① 변수 테스트 : β_1 에 대한 가정 t-test

$H_0: \beta_1 = 0 \rightarrow$ 변수의 설명력이 없다

$H_1: \beta_1 \neq 0$

$$\hat{\beta}_1 \sim N\left(0, \frac{s}{\sqrt{SSR}}\right) \quad , \quad t = \frac{\hat{\beta}_1}{\frac{s}{\sqrt{SSR}}}$$



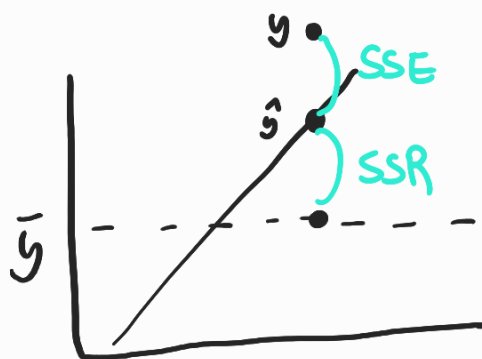
유의수준보다 작으면

p-value < $\frac{\alpha}{2}$ 면 기각!

② 전체 회귀식 test : F-test

SSE : Sum of square Error
 \rightarrow 회귀식이 설명 못하는 값

SSR : Sum of square regression.
 \rightarrow 회귀식이 설명하는 값



$H_0: \beta_1 = \beta_2 = \dots = 0$

$H_1: \text{적어도 하나의 } \beta \text{는 0이 아님}$

$$F = \frac{MSR}{MSE} \sim \chi^2$$

	df	
SSR	P	$MSR = SSR/P$
SSE	$n-p-1$	$MSE = SSE/(n-p-1)$
TSS	$n-1$	

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

$$(0 \leq R^2 \leq 1)$$

$$R^2_{Adj} = 1 - \frac{MSE}{MST}$$

변수 개수가 많아질수록 R^2 값 높아지는데
 이를 방지하기 위해 R^2_{Adj}

R^2 0.67면
 회귀식 설명

가정

① 선형성

: $X-Y$ 관계가 선형적이지

② 독립성 (\leftrightarrow 다중공선성)

: 다중회귀 모형에서 X 변수들간 서로 연관 있는 것

: 오차항들은 서로 독립.

① 예측 X ? \rightarrow 상관성 높은 변수 제거 / 단계적 레지레이션 (전체 데이터를 몇개 구간으로 나눠 레지레이션)

* 자기상관 (auto correlation)

: 시계열 자료인 경우 시간에 대해 상관관계 발생

: 시계열 데이터의 독립성 판단

- 잔차 분포도 확인
 - 독립성 X 면 t 에 따라 패턴 존재
- 테번왓슨 (DW) 통계량 ($0 \leq DW \leq 4$)
 - 2일 때가 상관관계 0이라 독립성 만족
 - 0일 때 $cor = 1$, 4일 때 $cor = -1$

$$DW = \frac{\sum_{i=2}^k (e_i - e_{i-1})^2}{\sum_{i=1}^k e_i^2}$$

다중공선성 판정

① 모현 경쟁 (F) 야 $\textcircled{\text{bot}}$ 개별 변수 경쟁 (t) 유익 X

② 개별 변수 분가 상자와 반대일 때

③ 상관분석 시 독립변수간 높은 상관성

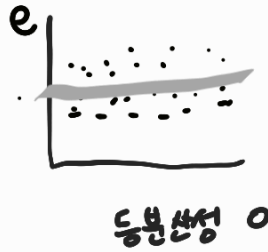
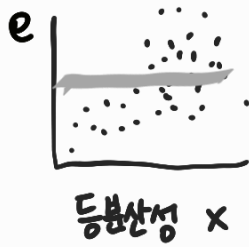
④ VIF 계속

③ 등분산성

잔차의 분산이 일정 (독립성 검정 상대없이 분산이 일정하다)
→ 오차항의 분산 σ^2 이 균등하게 분포해야한다

이분산성 판별

- 몇개 소그룹으로 나누어 그룹간 분산 차이가 있는지 비교
- 잔차 그림 이용.



만족 X? → 변수변환

④ 정규성

: 잔차가 정규분포를 따른다

만족 X? → 변수 변환 (제곱근, 로그, 역수, 지수, 이차식..)