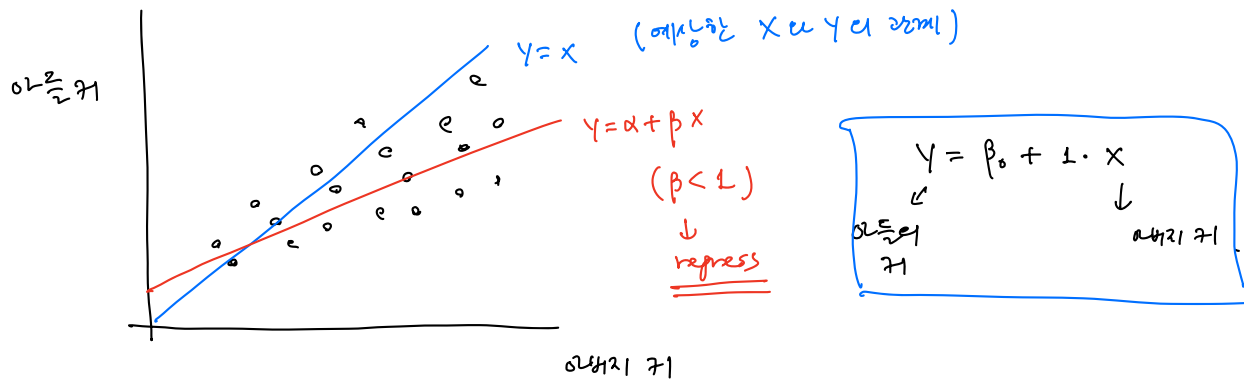


회귀분석 (regression analysis) 소개

Francis Galton : (아버지 키  $x_i$ , 아들 키  $y_i$ )  $i=1, 2, \dots, n$



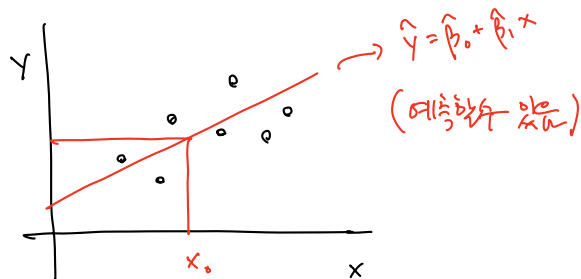
다항 선형 회귀

모형 :  $y = \beta_0 + \beta_1 x + \varepsilon$ ,  $E(\varepsilon) = 0$   
 $\text{Var}(\varepsilon) = \sigma^2$

$\beta_0, \beta_1 \Rightarrow$  회귀계수 (regression coefficient)

목적 ①  $\beta_0, \beta_1$  추정  $\Rightarrow$  예측 (prediction)

$\hat{\beta}_0, \hat{\beta}_1 \Rightarrow \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$



②  $H_0: \beta_1 = 0$  가정  $\Rightarrow$   $x$ 과  $y$ 의 관계를 가정

$\hookrightarrow$  under  $H_0$  ( $H_0$  is true)  $\Rightarrow$  모형  $y = \beta_0 + \varepsilon$   
 $\Rightarrow y$ 는  $x$ 의 관계가 없다.

\* 일반적으로 설명변수  $x$ 는 관측변수라 하지 않고 fixed & known

그러나  $E(y) = E(\beta_0 + \beta_1 x + \varepsilon) = \beta_0 + \beta_1 x + E(\varepsilon)$

$y = \beta_0 + \beta_1 x + \varepsilon$

$$= \beta_0 + \beta_1 x$$

= true regression line

o Back to model

$$\varepsilon_i \stackrel{iid}{\sim} \text{Normal}(0, \sigma^2)$$

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i=1, \dots, n$$

$$E(\varepsilon_i) = 0 \quad \& \quad \text{Var}(\varepsilon_i) = \sigma^2 \text{ for all } i.$$

define fitted values and residuals

$$\hat{\beta}_0, \hat{\beta}_1$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, \quad i=1, 2, \dots, n \quad (\text{fitted values})$$

$$e_i = y_i - \hat{y}_i, \quad i=1, 2, \dots, n \quad (\text{residuals})$$

Note

$$y_i = \underbrace{\beta_0 + \beta_1 x_i}_{\text{SS}} + \underbrace{\varepsilon_i}_{\text{SS}}$$

$$y_i = \underbrace{\hat{\beta}_0 + \hat{\beta}_1 x_i}_{\text{SS}} + \underbrace{e_i}_{\text{SS}}$$

$\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$  라는 가정

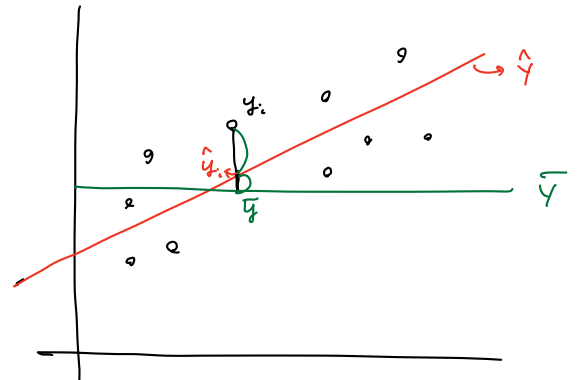
$e_i$  가 만족하는지 check 해야 함

// residual analysis (잔차분석)

o  $y_i$  분산 분해,  $R^2$  (결정 계수)

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$SS_T = SS_E + SS_R$$

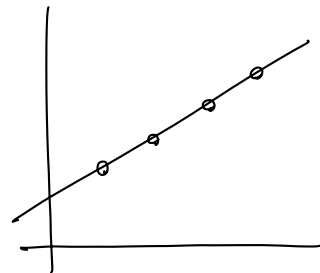


$$R^2 = \frac{SS_R}{SS_T}$$

\*  $R^2 = 1$  이면

$$SS_E = 0$$

\*  $R^2$  이 높을수록



$$F\text{-stat.} \quad F = \frac{MS_R}{MS_E} = \frac{SS_R / df_R}{SS_E / df_E} \underset{H_0}{\sim} F_{df_R, df_E}$$

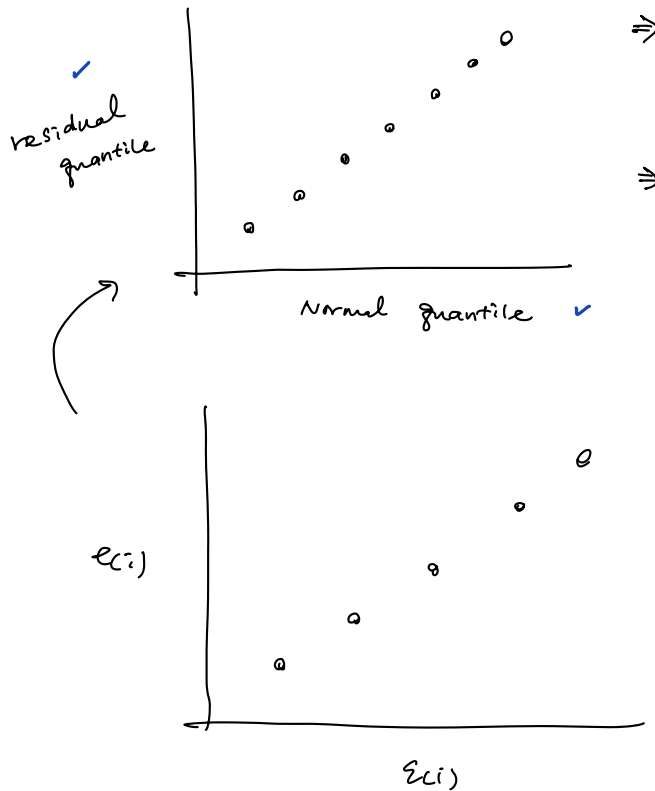
$F \uparrow \Rightarrow$  reject  $H_0: \beta_1 = 0 \Leftrightarrow t\text{-test}$

$\Rightarrow$  2214 Multiple regression model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon,$$

$$H_0: \beta_1 = \beta_2 \Rightarrow \text{use F-test.}$$

o Normal probability plot



$$\varepsilon_i \sim \text{Normal}$$

$\Rightarrow$  normal 분포를 따르는 데이터  $\varepsilon_1, \dots, \varepsilon_n$  추출하기

$$\varepsilon_{(1)} < \varepsilon_{(2)} < \dots < \varepsilon_{(n)} \quad \text{순서}$$

또한 각각  $e_1, \dots, e_n$  이 대해

$$e_{(1)} < e_{(2)} < \dots < e_{(n)} \quad \text{순서}$$

만약  $e_1, \dots, e_n$  이 normal 분포를 따르려면

o 이진형 설명 변수

만약  $x$  가 범주형 변수 (setosa, versicolor) 일때, 수치화한 것이라면 이렇듯

$$x = \begin{cases} 1 & \text{if setosa} \\ 0 & \text{if versicolor} \end{cases}$$

Back to iris data.

꽃종 종류 (species)  $\rightarrow$  (setosa, versicolor, virginica) 범주형 자료

[수치화하는 방법]

$$x_1 = \begin{cases} 1 & \text{if versicolor} \\ 0 & \text{o.w.} \end{cases}$$

$$x_2 = \begin{cases} 1 & \text{if virginica} \\ 0 & \text{o.w.} \end{cases}$$

	$x_1$	$x_2$
Setosa	0	0

Versicolor	1	0
Virginica	0	1

꽃받침 너비

$$y (\text{Sepal. width}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

$$E(Y | \text{Setosa}) = \beta_0$$

$$E(Y | \text{versicolor}) = \beta_0 + \beta_1$$

$$E(Y | \text{virginica}) = \beta_0 + \beta_2$$

$H_0: \beta_1 = \beta_2 = 0 \Rightarrow$  under  $H_0$ , 3종류 꽃의 꽃받침 너비는 차이가 없다.

만약 이

$$x_1 = \begin{cases} 1 & \text{if Setosa} \\ 2 & \text{if versicolor} \\ 3 & \text{if virginica} \end{cases}$$

$$y = \beta_0 + \beta_1 x_1 + \varepsilon$$

$$E(Y | \text{Setosa}) = \beta_0 + \beta_1$$

$$E(Y | \text{versicolor}) = \beta_0 + 2\beta_1$$

$$E(Y | \text{virginica}) = \beta_0 + 3\beta_1$$

$\left. \begin{array}{l} \beta_1 \text{ 차이} \\ \beta_1 \text{ 차이} \end{array} \right\} ?$

## ○ 분산분석 (Analysis of Variance, ANOVA)

군집가공품 인장강도 (반응변수, 특성값):  $y$

$\Rightarrow$  여러개의 군집  $1, 2, \dots, k$

	1	2	...	...	k
$\left\{ \begin{array}{l} \text{반응} \\ \text{변수} \end{array} \right.$	$y_{11}$	$y_{12}$			$y_{1k}$
	$\vdots$	$\vdots$			$\vdots$
	$y_{1n}$	$y_{2n}$			$y_{kn}$
	$\overline{y_{1.}}$	$\overline{y_{2.}}$			$\overline{y_{k.}}$

$$\overline{y_{..}} = \frac{1}{nk} \sum_{i,j} y_{ij}$$

= grand average

전체평균

$$\sum_i \sum_j (y_{ij} - \bar{y}_{..})^2 = \sum_i \sum_j (y_{ij} - \bar{y}_{i.})^2 + n \sum_{i=1}^k (\bar{y}_{i.} - \bar{y}_{..})^2 \quad \text{분해}$$

$$SS_T = SS_W \text{ (공변량)} + SS_B \text{ (공법간)}$$

반응변수의 변동량  
(고정된 값이다.)

$SS_E$  1변동량  
i번째 공법 내에서의 변동량

$SS_{tot}$  2변동량

2변동량

공법에 의해 설명되는 변동량  
i번째 공법으로 인한.  
공법간 차이가 크다면 값이 커지고,  
공법간 차이가 없다면 값이 작아질 것이다.

$$\Rightarrow SS_B \uparrow \quad \text{공법간 차이가 크면}$$

$$SS_B = 0 \quad \text{"} \quad \text{없으면}$$

○ 일인식 배치 방법 모형

$$y_{ij} = \mu + \tau_i + \varepsilon_{ij} \quad \Leftrightarrow \quad y_{ij} = \mu_i + \varepsilon_{ij}$$

$\mu$ : 전체 평균

$$\mu_i \triangleq \mu + \tau_i$$

$\tau_i$ : i번째 처리 효과

$$\hat{\mu} + \hat{\tau}_i$$

$$\varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2) \quad \rightarrow \quad \text{잔차항이 필요함} \quad e_{ij} = y_{ij} - \hat{y}_{ij}$$

$$\text{모형: } H_0: \tau_i = 0 \quad \text{for all } i$$

$$H_0: \mu_1 = \dots = \mu_k = \mu$$

Under  $H_0 \Rightarrow$

$$y_{ij} = \mu + \varepsilon_{ij}$$

$$E(y_{ij}) = \mu \quad (\because E(\varepsilon_{ij}) = 0)$$

= 모든 처리 수준에서 특성값이 같음

= 공법간 차이가 없다.

$\Rightarrow$  F-test

ANOVA table

Source	SS	DF	MS	F <sub>0</sub>	p-value
Treat	$SS_{treat}$	$k-1$	$MS_{treat} = \frac{SS_{treat}}{k-1}$	$F_0 = \frac{MS_{treat}}{MS_E}$	
Error	$SS_E$	$N-k$	$MS_E = \frac{SS_E}{N-k}$		
Total	$SS_T$	$N-1$			

$H_0$