

Multippel linær regresjon, 12.7 + 12.8  
Flere forklaringsvariabler per enhet (individ)

Responser  $y_i$ :

$$\text{Forklарingsvariabel } x_{i1} \quad i=1, \dots, n$$

$$\vdots$$

$$\sum_{i=1}^n k x_{ik}$$

$$\text{Modell: } Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i$$

der  $\varepsilon_i \sim N(0, \sigma^2)$  mørk ( $E\varepsilon_i = 0$ ,  $\text{Var}(\varepsilon_i) = \sigma^2$ )

Ex) Fødselsoverhett,  $Y_i =$  fødselsoverhett

$x_{i1} =$  Sveriges støpslengde

$x_{i2} =$  Mor alder

$x_{i3} =$  Antall tidlige fødsler

$x_{i4} =$  Barnets kjønn = {1 gutt

0 jente}

Ex) Kaffekonsum,  $Y_i =$  kaffekonsum

$x_{i1} =$  Aut. automaten =  $x_i$

$x_{i2} = x_{i1}^2$

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \varepsilon_i \quad \text{Polynomell regresjon}$$

okt 30-12:15

Formell • Multifaktorielle fenomener

• Mer præcise prediktjoner

• Lettre & påvisig signifikans.

• Kan løse opp i konflikten om hvilke variabler er hovedsaklig.

Fortolkning av  $\beta_j$ :

Endring i forventning  $\mu_j = EY_j$ :  
hvis  $x_{ij}$  endres til  $x_{ij+1}$ , når samtidig de øvrige forklaringsvariablene holdes konstant.  
Dette kan ikke alltid holde, f.eks. ved polynomell regresjon.

okt 30-12:27

Estimering: MKE = minste kvadrateters estimatorene.  
 $f(b_0, b_1, \dots, b_k) = \sum (Y_i - b_0 - b_1 x_{i1} - \dots - b_k x_{ik})^2$

og  $\hat{\beta}^T = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k)$  er den verdien av  $(b_0, b_1, \dots, b_k)$  som gjør  $f(b_0, b_1, \dots, b_k)$  minst mulig.

Normal ligningene

$$\frac{\partial f}{\partial b_0} = -2 \sum (Y_i - b_0 - b_1 x_{i1} - \dots - b_k x_{ik}) = 0$$

$$\frac{\partial f}{\partial b_1} = -2 \sum (Y_i - b_0 - b_1 x_{i1} - \dots - b_k x_{ik}) x_{i1} = 0$$

$$\vdots$$

$$\frac{\partial f}{\partial b_k} = -2 \sum (Y_i - b_0 - b_1 x_{i1} - \dots - b_k x_{ik}) x_{ik} = 0$$

Entydig løsning hvis  $k+1 \leq n$  og lign. linjeartig.

Ligningene om enkle utregninger med verdiene og matrisene:  
 $\underline{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \Rightarrow \underline{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} = E(\underline{Y}) = \begin{pmatrix} EY_1 \\ \vdots \\ EY_n \end{pmatrix} = \begin{pmatrix} \underline{E}x_1 \\ \vdots \\ \underline{E}x_n \end{pmatrix}$   
 der  $\underline{x}_i^T = (1, x_{i1}, x_{i2}, \dots, x_{ik}) \Rightarrow \underline{\beta}^T = (\beta_0, \beta_1, \dots, \beta_k)$   
 $\underline{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$   
 $\Rightarrow \text{design matrise} \quad \underline{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{nk} \end{pmatrix} = \begin{pmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{pmatrix}$   
 Ser da at  $\underline{\mu} = \underline{X} \underline{\beta}$   
 Vi kan da skrive modellen  
 $\underline{Y} = \underline{\mu} + \underline{\varepsilon} = \underline{X} \underline{\beta} + \underline{\varepsilon}$

okt 30-12:36

okt 30-12:55

På matriseform blir normal lign.  
 $-2 \underline{X}^T [\underline{Y} - \underline{\mu}] = 0 \quad (\text{Sjekk dette})$

og siden  $\underline{\mu} = \underline{X} \underline{\beta}$  kan vi forme dette til  
 $\underline{X}^T \underline{Y} = \underline{X}^T \underline{\mu} = \underline{X}^T \underline{X} \underline{\beta}$

Løsning (MKE) fra  $\underline{\beta} = \underline{X}^{-1} \underline{X}^T \underline{Y}$

nennt  $\underline{X}^T \underline{X}$  en invertibel  $\Leftrightarrow \underline{X}$  har full rang

Videre, så  $\underline{\beta} = \underline{X} \hat{\beta}$ ,  
 $f(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k) = \sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_k x_{ik})^2$   
 $= (\underline{Y} - \hat{\underline{\mu}})^T (\underline{Y} - \hat{\underline{\mu}}) = \| \underline{Y} - \hat{\underline{\mu}} \|^2$   
 $= SSE$

Form. av H estimasjon for  $\sigma^2$ :  $\sigma^2 = \frac{SSE}{n-k-1}$

Betyr at  $k+1$  frihetsgraden til  $\sigma^2$  er  $n-k-1$ .

La videre  $\hat{\underline{Y}} = (\hat{Y}_1, \hat{Y}_2, \dots, \hat{Y}_n)$   
 $\underbrace{\text{n elementer}}$

Da er  
 $SST = \sum (Y_i - \bar{Y})^2 = (\underline{Y} - \hat{\underline{Y}})^T (\underline{Y} - \hat{\underline{Y}}) = \| \underline{Y} - \hat{\underline{Y}} \|^2$   
 og vi kan uttrykke  
 $SSR = \text{kvariatsummen overende til regresjonen}$   
 $= SST - SSE = \dots = \| \hat{\underline{Y}} - \hat{\underline{Y}} \|^2$   
 Ut fra dette defineres  
 $R^2 = \text{Forhåndt andel av variansen} = 1 - \frac{SSE}{SST} = \frac{SSR}{SST}$   
 Fortolkning av multiplikativ  $R^2$   
 $= (\text{korrelasjon mellom } Y_i \text{ og } \hat{Y}_i)^2$

okt 30-13:16

okt 30-13:26

Forsøksgjennomgang  
Læringskurven  
Exempel 2.4.3.5

$$\begin{aligned} E[\hat{\beta}] &= E[(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{Y}] = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T E[\tilde{Y}] \\ &= (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \beta = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T X \beta = \beta \end{aligned}$$

Det finnes mange versjoner av kovariansene for  $\hat{\beta}$ .  
Husk at  $\hat{\beta}_j$  er linje som krysser  $y_i$ , og hvis disse  
er nøyaktig vil  $\hat{\beta}_j$  være standardfeil for  $\hat{\beta}_j$ . Da blir  
La  $s_{\hat{\beta}_j} = \text{standardfeil for } \hat{\beta}_j$   
 $t_j = \frac{\hat{\beta}_j - \beta_j}{s_{\hat{\beta}_j}} \sim t_{n-k-1}$

og vi kan teste  $H_0: \beta_j = 0$  vs  $\beta_j \neq 0$   
med  $|t_j|$ . For hvert kvis deles i større sum  
 $t_{adj} = (1-\alpha)100\%$  parallel:  $t_{n-k-1}$   
Et  $(1-\alpha)100\%$  KI for  $\beta_j$  gir oss

$$\hat{\beta}_j \pm t_{adj} s_{\hat{\beta}_j}$$

Dersom  $F = \frac{SSR/k}{SSE/(n-k-1)} \sim F_{n-k-1}$   
så er  $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$   
Alt.  $H_1: \text{Minst en } \beta_j \neq 0$   
Justert  $R^2$ :  $R_{adj}^2 = 1 - \frac{SSE/(n-k-1)}{SST/(n-1)}$

okt 30-13:33