

Medisinsk statistikk, del II, vår 2008 KLMED 8005

Eirik Skogvoll

Førsteamanuensis dr. med.
Enhets for Anvendt klinisk forskning
Det medisinske fakultet

1

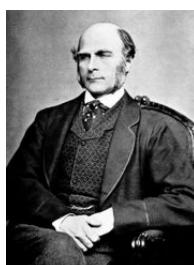
Lineær regresjon, Rosner 11.1 – 11.6

- Bakgrunn (11.1)
- Modell (11.2)
- Estimering av parametre i modellen (11.3)
- Statistisk inferens, testing (11.4)
- Intervallestimering (11.5)
- Modellens godhet

2

”Regresjon...?”

- Betyr ”tilbakegang”
- Sir Francis Galton (1822-1911) observerte at:
 - lange fedre fikk oftest lange barn, men de lengste fedrene fikk i gjennomsnitt kortere barn enn seg selv
 - kortvokste fedre fikk oftest korte barn, men de korteste fedrene fikk i gjennomsnitt lengre barn enn seg selv.



→ ”Regression toward the mean.” A principle stating that of related measurements, and selecting those where the first measurement is either higher or lower than the average, the expected value of the second is closer to the mean than the observed value of the first.

3

11.1 Innledning: østriol vs. fødselsvekt

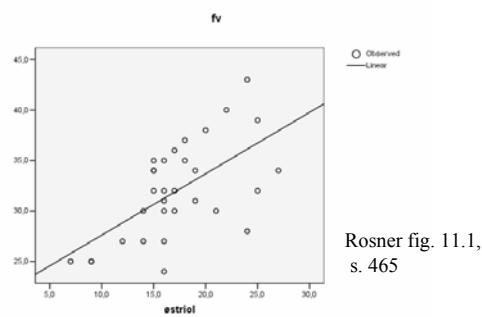
Observerer (x_i, y_i) for $i = 1, 2, \dots, n$
 x : prediktor, Y : respons

Rosner tabell 11.1, s. 466

4

i	østriol	fv
1	7	25
2	9	25
3	9	25
4	12	27
5	14	27
6	16	27
7	16	24
8	14	30
9	16	30
10	16	31
11	17	30
12	19	31
13	21	30
14	24	28
15	15	32
16	16	32
17	17	32
18	25	32
19	27	34
20	15	34
21	15	34
22	15	35
23	16	35
24	19	34
25	18	35
26	17	36
27	18	37
28	20	38
29	22	40
30	25	39
31	24	43

”én figur sier mer enn tusen t-tester”



5

11.2 Modell

$$Y_i = \alpha + \beta x_i + e_i \quad e_i \sim N(0, \sigma^2)$$

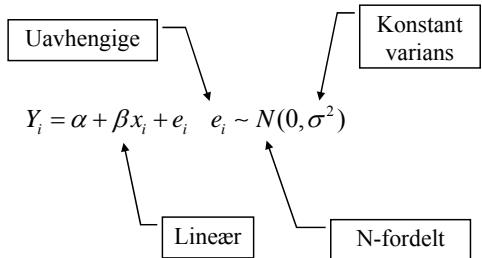
$$E(Y_i | x_i) = \alpha + \beta x_i$$

$$\text{Var}(Y_i | x_i) = \sigma^2$$

$$Y_i \sim N(\alpha + \beta x_i, \sigma^2)$$

6

Modellens forutsetninger



7

11.2 Modell (2)

observert: y_i

predikert: $\hat{y}_i = a + b x_i$

gjennomsnitt: \bar{y}

8

11.3 Estimering av parametere

- Minste kvadraters metode
- Maximum likelihood

$$\hat{\alpha} = a, \quad \hat{\beta} = b \qquad \text{Rosner fig. 11.4, s. 469}$$

$$d_i = y_i - \hat{y}_i = y_i - (a + b x_i)$$

$$S = \sum_{i=1}^n d_i^2 = \sum_{i=1}^n (y_i - a - b x_i)^2$$

9

11.3 Estimering av parametere (2)

observert: y_i

predikert: $\hat{y}_i = a + bx_i$

gjennomsnitt: \bar{y}

residual komponent: $y_i - \hat{y}_i$

regresjons-komponent: $\hat{y}_i - \bar{y}$ Rosner fig. 11.5, s. 473

10

Estimater

ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1	Regression 250,574	1	250,574	17,162	,000 ^b
	Residual 674,000	29	14,601		

a. Predictors: (Constant), estrol
b. Dependent Variable: fv

Coefficients ^a						
Model	Unstandardized Coefficients	Standardized Coefficients	Beta	Std. Error	t	95% Confidence Interval for B
1	(Constant) .623 2,620	.610 4,143	8,214 ,000	16,164 ,308	26,843 ,908	

a. Dependent Variable: fv

11

11.4 Inferens: hvor god er modellen?

Rosner fig. 11.6, s. 474:

- a) stor regresjonskomponent, liten residualkomponent
- b) stor regresjonskomponent, stor residualkomponent
- c) liten regresjonskomponent, liten residualkomponent
- d) liten regresjonskomponent, stor residualkomponent

best
↑
verst

12

11.4: Hvor god er modellen? (2)

$$\text{Total SS: } \sum_{i=1}^n (y_i - \bar{y})^2$$

$$\text{Regresjon SS: } \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$\text{Residual SS: } \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

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

eller... Total SS = Regresjon SS + Residual SS

11.4 Hvor god er modellen? (3)

- bedre enn ingen modell?

$$\text{Regresjon Mean Square: } \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{k} \quad k : \text{antall prediktorvariable}$$

(enkel lineær regresjon: $k = 1$)

$$\text{Residual Mean Square: } \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-k-1}$$

(enkel lineær regression: $k = 1$)

$$F = \frac{\text{Reg MS}}{\text{Res MS}} \sim F_{k, n-k-1} \quad (\text{enkel lineær regresjon: } F \sim F_{1, n-2})$$

14

ANOVA (ANalysis Of VAriance)

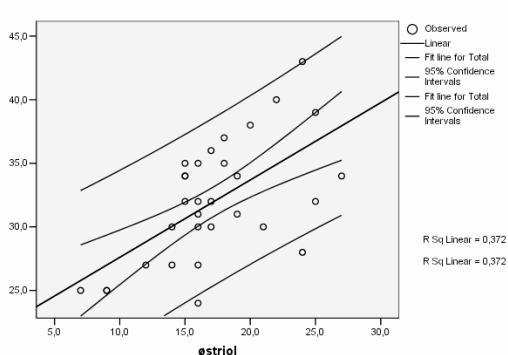
ANOVA ^a					
Model	Sum of Squares	df	Mean Square	F	Sig.
1	Regression	250,574	1	250,574	17,162
	Residual	423,426	29	14,601	
	Total	674,000	30		.000 ^a

a. Predictors: (Constant), østriol
 b. Dependent Variable: fy

Noen sammenhenger...

- I enkel lineær regresjon er hypotesen $\beta = 0$ ekvivalent med $\rho = 0$ som evalueres ved hjelpe Pearson's r i korrelasjonsanalysen
- R^2 ("coefficient of determination") tolkes best i regresjonsanalysen som andelen "forklart varians" via kvadratsummer. Ekvivalent med (Pearson's) r^2

Prediksjon



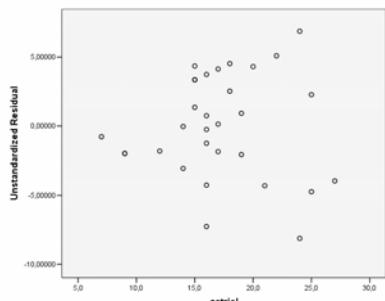
Modellens forutsetninger

$$Y_i = \alpha + \beta x_i + e_i \quad e_i \sim N(0, \sigma^2)$$

Uavhengige Konstant varians
Lineær N-fordelt

Sjekk av modellens forutsetninger

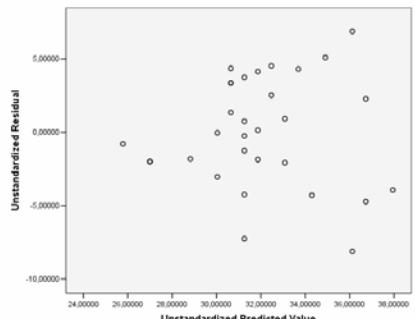
- Linearitet? residual vs. prediktor



22

Sjekk av modellens forutsetninger

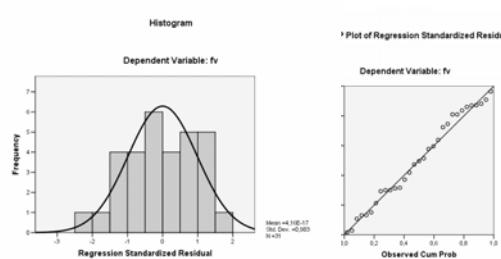
- Konstant varians? residual vs. predikert verdi



23

Sjekk av modellens forutsetninger

- Normalfordelte residualer? histogram og p-p plott



24

Sjekk av modellens forutsetninger

- Uavhengighet mellom feil-leddene

- Relevant for tidsserier (en prosess observert på regelmessige tidspunkter)
- Residual-korrelasjon? Durbin-Watson test
