

SOS3003

Anvendt statistisk dataanalyse i samfunnsvitenskap

Forelesingsnotat 02

Erling Berge
Institutt for sosiologi og statsvitenskap
NTNU

Haut 2004

© Erling Berge 2004

1

Forelesing II

- Bivariat regresjon II
 - Hamilton Kap 2 s51-59
- Multivariat regresjon I
 - Hamilton Kap 3 s65-72

Haut 2004

© Erling Berge 2004

2

Repetisjon: Bivariat Regresjon: Modell for utval

- $Y_i = b_0 + b_1 x_{1i} + e_i$
- $i=1, \dots, n$ $n = \# \text{ case i utvalet}$

Eksempel frå første forelesning med

- $Y = \text{IMPORTANCE OF PUBLIC CONTROL OF SALES OF AGRIC. ESTATES}$
- $X = \text{NUMBER OF DEKAR LAND OWNED}$

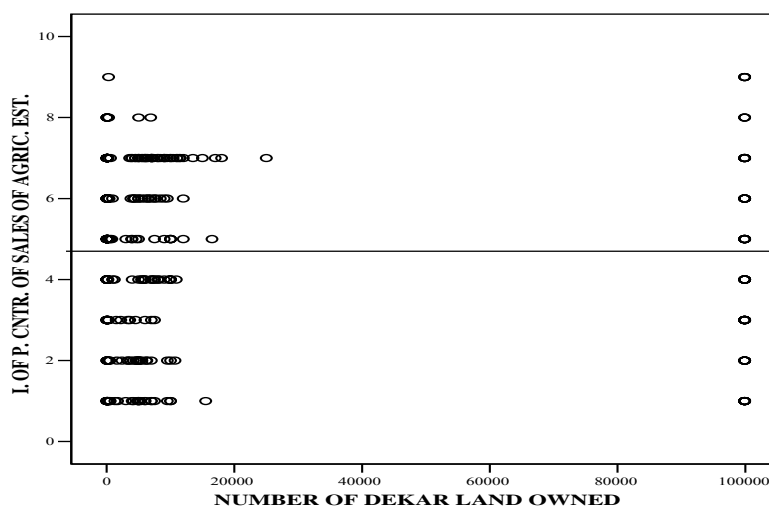
Kva var galt i eksempelet?

Haust 2004

© Erling Berge 2004

3

Kva er galt? Spreiingsdiagram med regresjonslinje



Haust 2004

© Erling Berge 2004

4

Generelt: Kva kan skape problem?

- Uteletne variablar
- Ikkje-lineære samanhengar
- Ikkje-konstant varians på feilen (heteroskedastisitet)
- Korrelasjon mellom feila (autokorrelasjon)
- Ikkje normale feil
- Case med uvanleg stor verknad

Haut 2004

© Erling Berge 2004

5

Ikkje-normale feil: Meir om variabelfordelingar

- Regresjon har **IKKJE føresetnader om fordelinga** til variablane
- MEN for å teste hypotesar må vi ha eit **normalfordelt feilledd**
- DERSOM **modellen er korrekt** og n er stor sikrar sentralgrenseteoremet at feilleddet er tilnærma normalfordelt
- MEN som regel er modellen feil eller ufullstendig. Derfor må vi **teste om residualen faktisk er normalfordelt.**

Haut 2004

© Erling Berge 2004

6

Residualanalyse:

Viktigaste innfallsporten til problema i regresjonsanalysen

Hjelpemiddel:

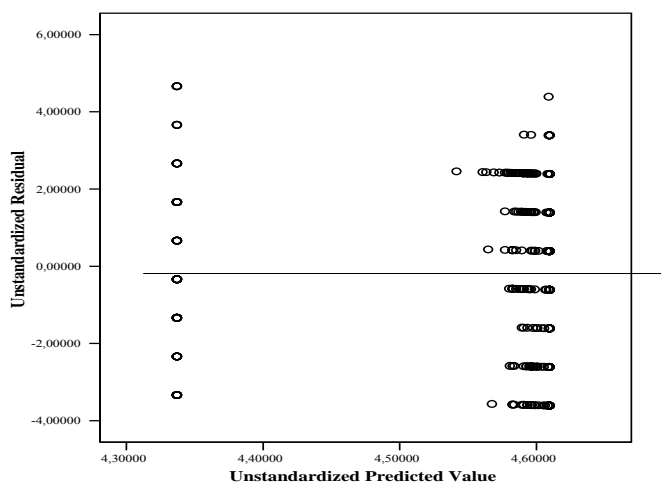
- Scatterplott – spreiingsplott
- Plott av residual mot predikert verdi
- Histogram
- Boksplott
- Symmetriplott
- Kvantil-normalplott

Haust 2004

© Erling Berge 2004

7

Kva er galt? (1) residual-predikert verdi plott

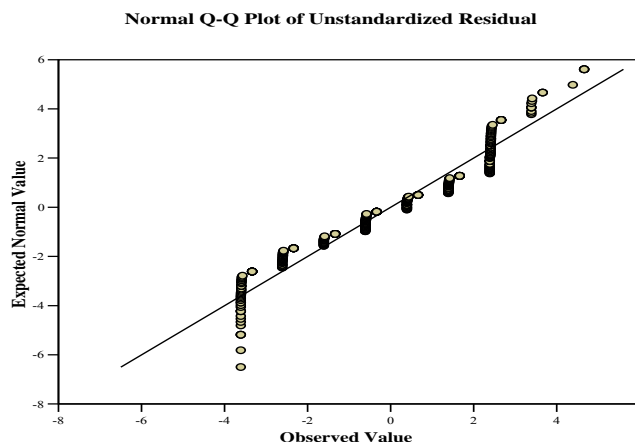


Haust 2004

© Erling Berge 2004

8

Kva er galt? (2) residual normal-kvantilplott



Haust 2004

© Erling Berge 2004

9

Potenstransformasjonar

kan løyse problem i samband med

- kurvelinearitet i modellen
- utliggarar
- case med stor innverknad
- ikkje-konstant varians hos feilleddet
- ikkje-normale feilledd

Transformasjon er m.a.o. eit generelt verkemiddel

Haust 2004

© Erling Berge 2004

10

Potenstransformasjonar (jfr H:17-22)

Y^* : les “transformert Y ” (transformasjon fra Y til Y^*)	Invers transformasjon (transformasjon fra Y^* til Y)
• $Y^* = Y^q$ $q > 0$	• $Y = [Y^*]^{1/q}$ $q > 0$
• $Y^* = \ln[Y]$ $q = 0$	• $Y = \exp[Y^*]$ $q = 0$
• $Y^* = - [Y^q]$ $q < 0$	• $Y = [- Y^*]^{1/q}$ $q < 0$

Haust 2004

© Erling Berge 2004

11

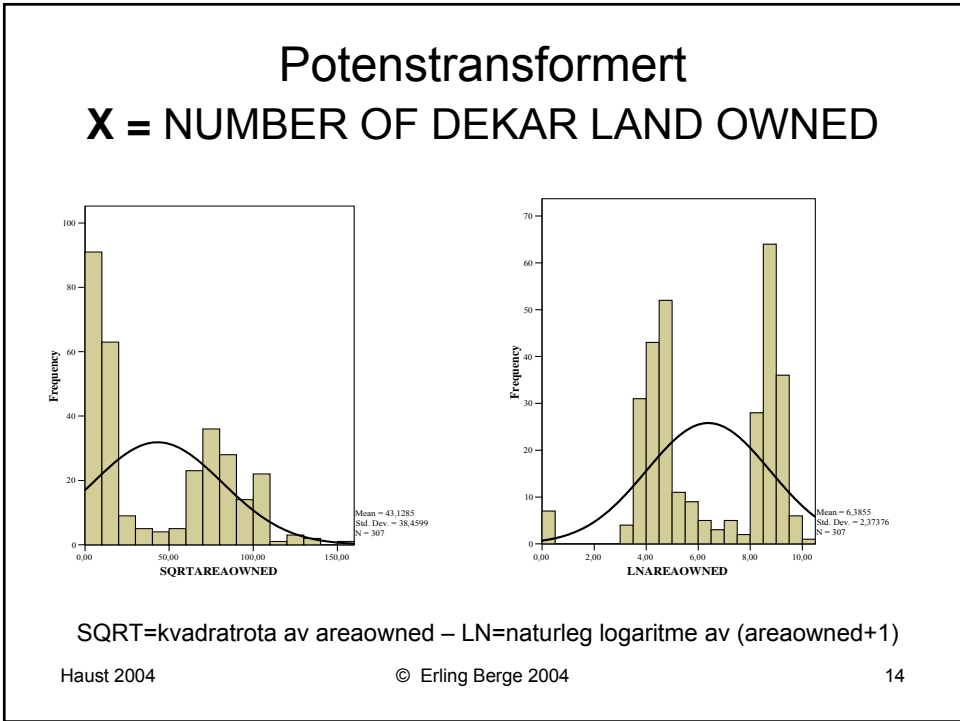
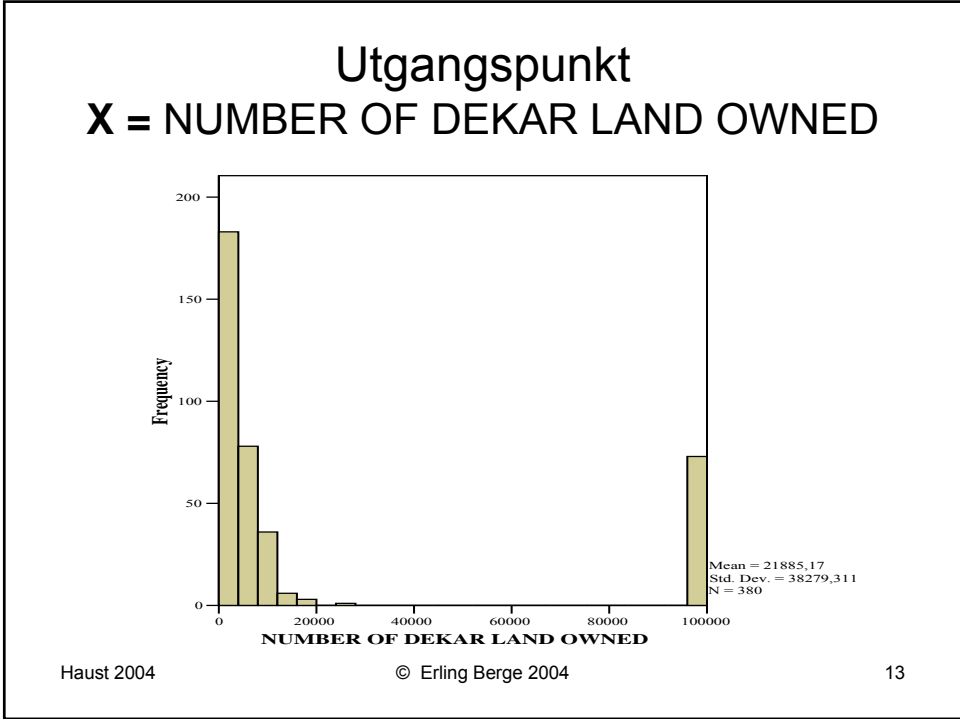
Potenstransformasjonar: konsekvensar

- $X^* = X^q$
 - $q > 1$ **aukar tyngda** til øvre hale relativt til nedre
 - $q = 1$ gir identitet
 - $q < 1$ **reduserer tyngda** til øvre hale relativt til nedre
- Dersom $Y^* = \ln(Y)$ vil regresjonskoeffesienten for ein intervallskala X variabel kunne tolkast som % endring i Y for ei einings endring i X
Dvs. dersom $\ln(Y) = b_0 + b_1 x + e$
Vil b_1 kunne tolkast som % endring i Y pr eining X

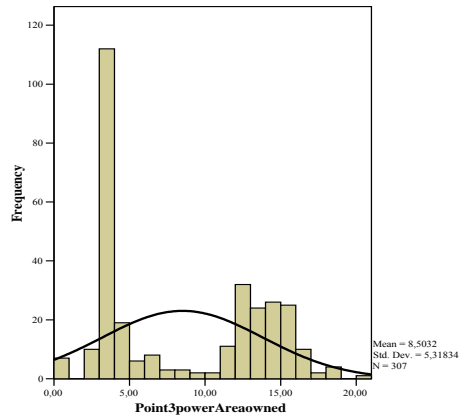
Haust 2004

© Erling Berge 2004

12



Potenstransformert X = NUMBER OF DEKAR LAND OWNED



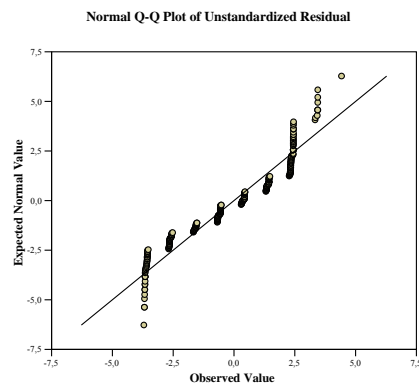
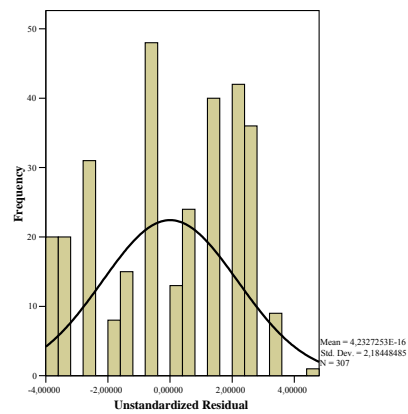
Poin3power = 0,3 potensen av areaowned

Haust 2004

© Erling Berge 2004

15

Hjelper det med potenstransformasjon?



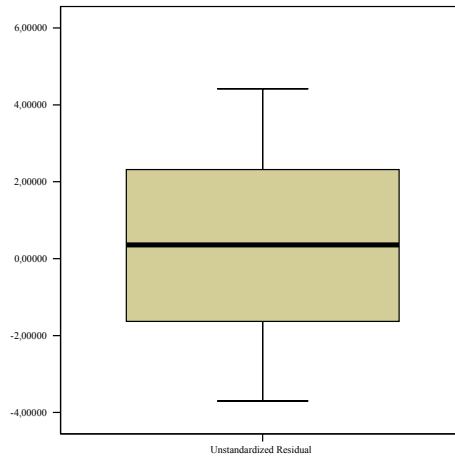
0.3potens-transformasjon gir lette halar og ingen utliggarar

Haust 2004

© Erling Berge 2004

16

Boksplott av residualen viser tilnærma symmetri utan utliggarar

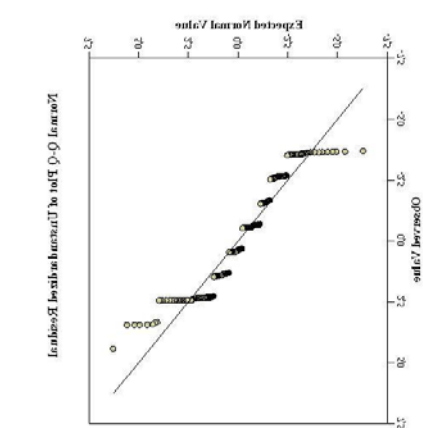
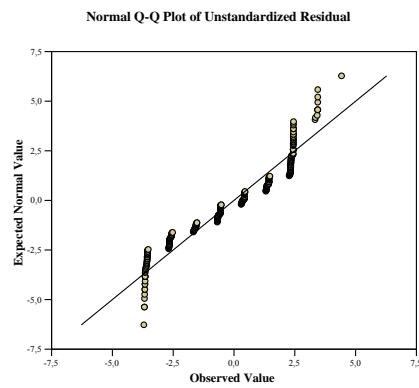


Haust 2004

© Erling Berge 2004

17

SPSSutskrift vs boka (jfr. s16)



Haust 2004

© Erling Berge 2004

18

Lesing av utskrifter fra SPSS (1)

Descriptive Statistics	Mean	Std. Deviation ¹	N ²
I. OF P. CNTR. OF SALES OF AGRIC. EST.	4.61	2.185	307
Point3powerAreaowned	8.5032	5.31834	307

Model	R	R Square ³	Adjusted R Square ⁴	Std. Error of the Estimate ⁵	Change Statistics				
					R Square Change	F Change	df1	df2	Sig. F Change
1	.024(a)	.001	-.003	2.188	.001	.182	1	305	.670

a Predictors: (Constant), Point3powerAreaowned

b Dependent Variable: I. OF P. CNTR. OF SALES OF AGRIC. EST.

Haust 2004

© Erling Berge 2004

19

Fotnotar til tabellen ovanfor (1)

1. Standard-avviket til gjennomsnittet (mean)
2. Talet av case brukt i analysen
3. Determinasjonskoeffisienten
4. Den justerte determinasjonskoeffisienten, sjå Hamilton side 41
5. Standard-avviket til residualen
 $s_e = \text{SQRT} (\text{RSS}/(n-K))$,
der SQRT = kvadratrota av (*)

Haust 2004

© Erling Berge 2004

20

Lesing av utskrifter fra SPSS (2)

Model		Sum of Squares ³	df	Mean Square	F ¹	Sig. ²
1	Regression	.870	1	.870	.182	.670(a)
	Residual	1460.224	305	4.788		
	Total	1461.094	306			

•Kvadratsummar: $TSS = ESS + RSS$

• $RSS = \sum_i (e_i)^2 = \sum_i (Y_i - \hat{Y}_i)^2$: sum kvadrert (avstand observert – estimert verdi)

•Mean Square = RSS / df For RSS har vi at $df = n - K$

K er lik talet på parametrar som modellen estimerer (dvs b_0 og b_1)

Her er $n=307$ og $K=2$, dvs. $Df = 305$

Haust 2004

© Erling Berge 2004

21

Fotnotar til tabellen ovanfor (2)

1. F-observatoren for nullhypotesa $\beta_1 = 0$ (sjå Hamilton side 45)
2. p-verdien for F-observatoren: dvs sannsynet for å finne ein så stor eller større F-verdi gitt at nullhypotesa er rett
3. Kvadratsummar
 1. $TSS = ESS + RSS$
 2. $RSS = \sum_i (e_i)^2 = \sum_i (Y_i - \hat{Y}_i)^2$ avstand observert – estimert verdi
 3. $ESS = \sum_i (\hat{Y}_i - \bar{Y})^2$ avstand estimert verdi – gjennomsnitt
 4. $TSS = \sum_i (Y_i - \bar{Y})^2$ avstand observert verdi – gjennomsnitt

Haust 2004

© Erling Berge 2004

22

Lesing av utskrifter fra SPSS (3)

M o d e l		Unstandardized Coefficients		Standardized Coefficients	t ⁴	Sig. ⁵	95% Confidence Interval for B	
		B ¹	Std. Error ²	Beta ³			Lower Bound	Upper Bound
1	(Constant)	4.524	.236		19.187	.000	4.060	4.988
	Point3-power Area-owned	.010	.024	.024	.426	.670	-.036	.056

Haust 2004 © Erling Berge 2004 23

Fotnotar til tabellen ovanfor (3)

1. Estimat av regresjonskoeffesientane b_0 og b_1
2. Standardavviket (standardfeilen) til parameterstimata b_0 og b_1
3. Dei standardiserte regresjonskoeffesientane: $b_1^{st} = b_1 \cdot (s_x/s_y)$ sjå Hamilton side 38-40
4. t-observatoren for nullhypotesa $\beta_1 = 0$ (sjå Hamilton side 44)
5. p-verdien for t-observatoren: dvs sannsynet for å finne ein så stor eller større t-verdi gitt at nullhypotesa er rett

Kurvelineær regresjon

- I eksempelet ovanfor brukte vi variabelen "Point3powerAreaowned", dvs 0.3 potensen av tal dekar areal ein eig. Dvs.:

- $\text{Point3powerAreaowned} = (\text{NUMBER OF DEKAR LAND OWNED})^{0.3}$

Modellen vi har estimert er altså

$$y_i = b_0 + b_1 x_i + e_i$$

$$y_i = b_0 + b_1 \text{Point3powerAreaowned}_i + e_i$$

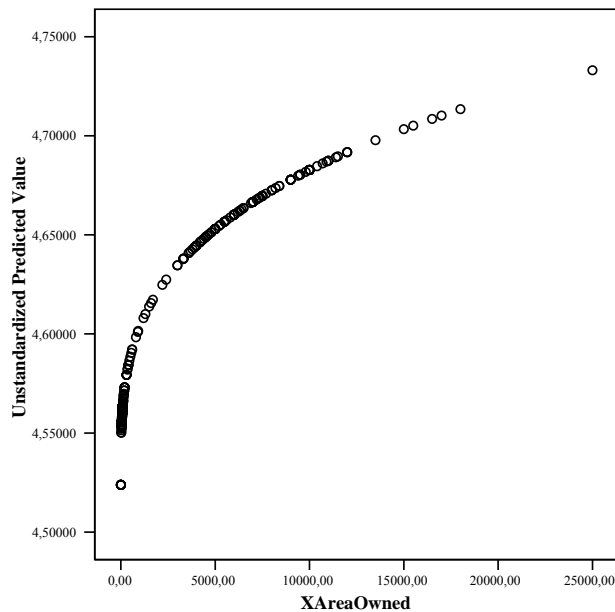
$$\hat{y}_i = 4.524 + 0.010 * (\text{NUMBER OF DEKAR LAND OWNED})^{0.3}$$

Haust 2004

© Erling Berge 2004

25

**Bruk av potens-
transformerte
variablar betyr
at regresjonen
blir kurvelineær**



Haust 2004

© Erling Berge 2004

26

Oppsummering

- I bivariat regresjon kan ein seie at OLS-metoden freistar finne den beste LINJA eller KURVA som passar til eit to-dimensjonalt spreingsmønster
- Scatter-plott og residualanalyse er hjelpemiddel for å diagnostisere problem i regresjonen
- Transformasjonar er eit generelt hjelpemiddel mot fleire typar problem, som t.d.:
 - Kurvelinearitet
 - Heteroskedastisitet
 - Ikkje-normalitet
 - Case med stor innverknad
- Regresjon med transformerte variablar er alltid kurvelineær. Vi tolkar resultatet letast ved hjelp av grafar

Haust 2004

© Erling Berge 2004

27

Multippel regresjon: modell (1)

- Målet med multippel regresjon er å finne nettoeffekten av ein variabel, kontrollert for variasjonen i alle dei andre
- Sett K = talet på parametrar i modellen (dvs. $K-1$ er talet på variablar).
Da kan (populasjons) modellen skrivast
- $$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_{K-1} x_{i,K-1} + \varepsilon_i$$

Haust 2004

© Erling Berge 2004

28

Multippel regresjon: modell (2)

- Dette kan skrivast

$$y_i = E[y_i] + \varepsilon_i ,$$

dette tyder at

- $E[y_i] = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_{K-1} x_{i,K-1}$
 $E[y_i]$ les vi som forventa verdi av y_i

Haust 2004

© Erling Berge 2004

29

Multippel regresjon: modell (3)

- Vi finn OLS estimata av modellen som dei b-verdiane i

$$\hat{y}_i = b_0 + b_1 x_{i1} + b_2 x_{i2} + b_3 x_{i3} + \dots + b_{K-1} x_{i,K-1}$$

(\hat{y}_i les vi som estimert eller "predikert" verdi av y_i)

som minimerer kvadratsummen av residualane

$$RSS = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n e_i^2$$

Haust 2004

© Erling Berge 2004

30

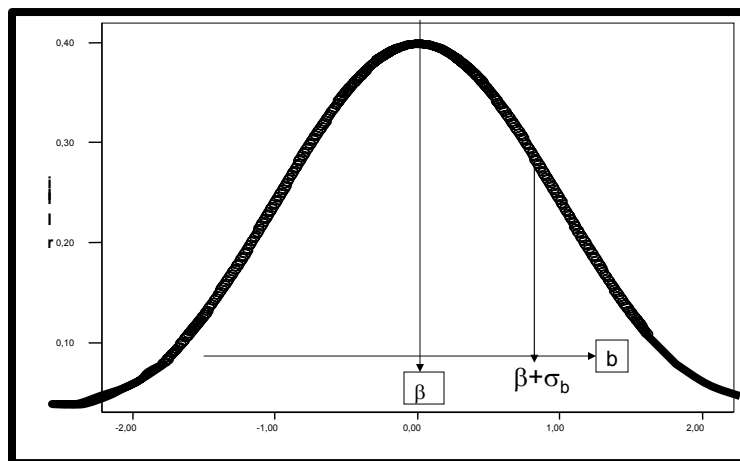
Meir om hypotesetesting

- Frå ein populasjon kan det trekkjast mange utval
- I kvart nytt utval vil vi kunne estimere nye regresjonsparametrar
- Lagar vi eit histogram over ulike estimerte verdiar av t.d. β_1 vil vi sjå at b_1 har ei fordeling (ei samplingfordeling)
- Ulike parametrar og observatorar har ulike samplingfordelingar
- Regresjonsparametrane (b 'ane) er t-fordelt

Haust 2004

© Erling Berge 2004

31



Sampling fordeling for regresjonsparameteren b : $E[b] = \beta$

Haust 2004

© Erling Berge 2004

32

Om partielle effektar (1)

- Eit eksempel med 2 variablar
- Dersom vi estimerer ein modell

$$y_i = b_0 + b_1 x_{i1} + b_2 x_{i2} + e_i$$

er det i prinsippet 3 ulike korrelasjonar som er involvert:

- Mellom y og x_1
- Mellom y og x_2
- Mellom x_1 og x_2

Haust 2004

© Erling Berge 2004

33

Om partielle effektar (2)

- Dette kan teoretisk gi opphav til 3 ulike bivariate regresjonar der vi held den tredje variabelen konstant

$$(1) y = a_{y|x_1} + b_{y|x_1} x_1 + e_{y|x_1} \quad x_2 \text{ konstant}$$

$$(2) y = a_{y|x_2} + b_{y|x_2} x_2 + e_{y|x_2} \quad x_1 \text{ konstant}$$

$$(3) x_1 = a_{x_1|x_2} + b_{x_1|x_2} x_2 + e_{x_1|x_2} \quad y \text{ konstant}$$

indeksen "y|x₁" les vi "frå regresjonen av y på x₁"

- Likningane (2) og (3) kan vi skrive om

Haust 2004

© Erling Berge 2004

34

Om partielle effektar (3)

$$(2) e_{y|x_2} = y - (a_{y|x_2} + b_{y|x_2}x_2)$$

$$(3) e_{x_1|x_2} = x_1 - (a_{x_1|x_2} + b_{x_1|x_2}x_2)$$

Vi ser no at vi så å seie "fjerner" effekten av x_2 frå y og frå x_1

Vi ser også at $e_{y|x_2}$ og $e_{x_1|x_2}$ vert nye variablar der effekten av x_2 er fjerna

Haut 2004

© Erling Berge 2004

35

Om partielle effektar (4)

- Dersom vi no lagar ein ny regresjon

$$\hat{e}_{y|x_2} = a + b e_{x_1|x_2}$$

finn vi at

$$a = 0$$

$$b = b_1 \text{ frå regresjonen}$$

$$y_i = b_0 + b_1 x_{i1} + b_2 x_{i2} + e_i$$

- b_1 er altså effekten av x_1 på y etter at vi har "fjerna" effekten av x_2

Haut 2004

© Erling Berge 2004

36

Eksperiment og partielle effektar

- I eksperiment granskar ein kausalsamband mellom to variable med kontroll for alle mogelege andre kausale faktorar
- Multippel regresjon er ei form for etterlikning av eksperimentet – ei nest beste løysing - og ligg nært opp til det som heiter kvasi-eksperimentelt forskingsdesign

Haut 2004

© Erling Berge 2004

37

KAUSALANALYSE

- Eksperiment
 - randomisering av påverknad (“behandling”) gir presise kausale konklusjonar om verknader (“respons”) ved signifikant skilnad i gjennomsnitt
 - kan vere umogeleg på grunn av
 - praktisk tilhøve
 - økonomiske skrankar
 - etiske vurderingar
- Kvasi-eksperiment der eksperiment er umogeleg
 - t.d. regresjonsanalyse

Haut 2004

© Erling Berge 2004

38

Eksperiment plasserer "case" tilfeldig i ei av to grupper:

- **BEHANDLING (T)**
med observasjon
 - FØR behandling
 - ETTER behandling
- **KONTROLL (C)**
med observasjon
 - FØR "ikkje-behandling"
 - ETTER "ikkje-behandling"

Haut 2004

© Erling Berge 2004

39

Modell av kausaleffektar ^{Ref.:}

- Studiar av observasjonsdata brukar omgrep fra eksperimentell design
- "Påverknad/ Behandling", "Stimulus" (Treatment/ Stimulus)
- "Effekt", "Utfall" (Effect/ Outcome)

Ref.:

Winship, Christopher, and Stephen L. Morgan 1999 "The Estimation of Causal Effects from Observational Data", Annual Review of Sociology Vol 25: 659-707

Haut 2004

© Erling Berge 2004

40

Modell av kausaleffektar:

Den "Kontrafaktiske" hypotesen for studiet av kausalitet

- Individet "i" kan i utgangspunktet tenkjast "selektert" til ei av to grupper
 - behandlingsgruppa, T, eller kontrollgruppa, C.
- Behandlinga, t (treatment), så vel som ikkje-behandling, c, kan i utgangspunktet tenkjast gitt til individ både i T- og C-gruppa
- Faktisk vil vi berre kunne observere t i T-gruppa og c i C-gruppa

Haut 2004

© Erling Berge 2004

41

Modell av kausaleffektar:

Den "Kontrafaktiske" hypotesen

- For kvart individ "i" kan ein tenkje seg fire moglege utfall
 - $Y_i(c, C)$ eller $Y_i(t, C)$; ved plassering i kontrollgruppe
 - $Y_i(c, T)$ eller $Y_i(t, T)$; ved plassering i behandlingsgruppa
- Berre $Y_i(c, \text{gitt "i" er med i C})$ eller $Y_i(t, \text{gitt "i" er med i T})$ kan observerast for eit gitt individ

Haut 2004

© Erling Berge 2004

42

Modell av kausaleffektar: Den "Kontrafaktiske" hypotesen

Litt meir formelt kan ein skrive dei mogelege utfalla for person i:

	Behandling: t	Ikkje beh.: c
T-gruppa	$Y_i^t \in T$	$Y_i^c \in T$
C-gruppa	$Y_i^t \in C$	$Y_i^c \in C$

Haut 2004

© Erling Berge 2004

43

Modell av kausaleffektar: Den "Kontrafaktiske" hypotesen

- Kausaleffekten for individ i er da

- $\delta_i = Y_i(t) - Y_i(c)$

- Berre ein av desse to storleikane kan observerast for eit gitt individ

Derfor "Den kontrafaktiske hypotesen"

Haut 2004

© Erling Berge 2004

44

Modell av kausaleffektar: Den "Kontrafaktiske" hypotesen

- Vi kan til dømes observere $Y_i(c | i \in C)$, men ikkje $Y_i(t | i \in C)$
- Problemet kan seiest å vere manglande data
- I staden for individeffektar vil ein estimere gjennomsnittseffektar i heile populasjonen

Haut 2004

© Erling Berge 2004

45

Modell av kausaleffektar:

- Gjennomsnittseffektar lar seg estimere, men som regel berre med store vanskar
- Ein føresetnad er at effekten av påverknad vil vere den same for eit gitt individ uavhengig av kva gruppe individet er plassert i
- Dette er likevel ikkje sjølvsgagt

Haut 2004

© Erling Berge 2004

46

Modell av kausaleffektar:

Den “Kontrafaktiske” hypotesen føreset:

- at endring av behandlingsgruppe for eitt individ ikkje verkar inn på utfallet for andre individ (fråvær av interaksjon)
- at behandlinga, “påverknaden”, faktisk er manipulerbar (t.d.: kjønn er ikkje manipulerbar)

Modell av kausaleffektar:

- Ein av vanskane er at i eit utval vil den prosessen som plasserer personen ”i” i kontroll- eller behandlings-gruppa kunne verke inn på det estimerte gjennomsnittsutfallet (seleksjonsproblemet)
- I somme høve er likevel den interessante storleiken gjennomsnitteffekten for dei som faktisk får påverknaden

Modell av kausaleffektar:

- Det kan visast at det er to kjelder til feil (bias) i estimata av gjennomsnitteffekten
 - ein eksisterande skilnad mellom C- og T-gruppene
 - behandlinga verkar i prinsippet ulikt for dei som er i T-gruppa samanlikna med dei som er i C-gruppa
- For å handtere dette må vi utvikle modellar for korleis folk hamnar i C- og T-gruppene

Haut 2004

© Erling Berge 2004

49

Modell av kausaleffektar:

- Ein generell klasse metodar som kan nyttast til å estimere kausaleffektar er regresjonsmodellane
- Dei vil kunne “kontrollere” for observerbare skilnader mellom T- og C-gruppene, men ikkje for ulik respons på behandling

Haut 2004

© Erling Berge 2004

50