



ECONnect

NTNU

Faktor

- en eksamensavis utgitt av ECONnect



Pensumsammendrag: SØK3001 – Økonometri I

Forfatter: Drago Bergholt
E-post: bergholt@stud.ntnu.no
Skrevet: Våren 2010
Antall sider: 71



Om ECONnect:

ECONnect er en frivillig studentorganisasjon for studentene på samfunnsøkonomi- og finansøkonomistudiet ved NTNU. Vi arbeider for økt faglig kompetanse blant våre studenter samt tettere kontakt med næringslivet. Det gjør vi ved å arrangere fagdager, gjesteforelesninger, bedriftspresentasjoner m.m. I dag går det ca. 200 studenter på bachelornivå (1.-3. klasse) og ca. 70 studenter på masternivå (4.-5. klasse). Studentene på masternivå er fordelt på de to linjene samfunnsøkonomi (ca. 50 stk) og finansiell økonomi (ca. 20 stk). Mer om ECONnect og aktuelle arrangementer på www.econnect-ntnu.no.

ECONnect består av følgende personer ved utgivelsestidspunkt:

Bjørn Bergholt (Leder)	bjorn@econnect-ntnu.no
Sophie S. Strømman (Bedriftsansvarlig)	sophie@econnect-ntnu.no
Maiken Weidle (Fagdagsansvarlig)	maiken@econnect-ntnu.no
Joakim Bjørkhaug (Økonomi- og IT-ansvarlig)	joakim@econnect-ntnu.no
Elise Caspersen	elise@econnect-ntnu.no
Tiril Toftedahl	tiril@econnect-ntnu.no
Louis Dieffenthaler	louis@econnect-ntnu.no
Andreas H. Jung	andreas@econnect-ntnu.no
Mari Benedikte Ellingsen	mari@econnect-ntnu.no
Herman Westrum Thorsen	herman@econnect-ntnu.no

Post- og besøksadresse:

ECONnect, NTNU Dragvoll
 Institutt for samfunnsøkonomi
 Bygg 7, Nivå 5
 7491 Trondheim

Organisasjonsnummer:

NO 994 625 314

Hjemmeside:

www.econnect-ntnu.no

Merk: Alle pensumsammendrag og tekster som utgis av Faktor er skrevet av og for studenter. ECONnect står ikke ansvarlig for selve faginnholdet. Spørsmål om teksten kan rettes til tekstforfatteren.

Pensumsammendrag: SØK 3001 - ØKONOMETRI I

Enkel regression:

$$y_i = \beta_0 + \beta_1 x_i + u_i \quad (1)$$

Modell: $y = a + bx$ der $\bar{y} = \beta_0 + \beta_1 \bar{x} + \bar{u} \quad (2)$

Forudsætninger:

SLR1: $y = \beta_0 + \beta_1 x + u$ (linear, parametrene)

SLR2: Tilfældig trekning fra populationen

SLR3: Variation i $x \rightarrow \text{Var}(x) > 0$

SLR4: $E(u_i | x_i) = 0 \rightarrow E(u_i) = 0$
 $\rightarrow \text{COV}(x_i, u_i) = 0$

SLR5: $\text{Var}(u_i | x_i) = \sigma^2$

OLS - metoden: (MKM)

Samtidsmodell: $y_i = \beta_0 + \beta_1 x_i + u_i$
 Estimeret modell: $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$
 $\left. \begin{array}{l} y_i = \beta_0 + \beta_1 x_i + u_i \\ \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \end{array} \right\} y_i - \hat{y}_i = \hat{u}_i$

Vil minimere $\sum_{i=1}^n (\hat{u}_i)^2$:

$$U(\hat{\beta}_0, \hat{\beta}_1) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

$$\Rightarrow \frac{\partial MK}{\partial \hat{\beta}_0} = \sum_{i=1}^n 2 (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \cdot (-1) = 0 \quad (3)$$

$$\Rightarrow \frac{\partial MK}{\partial \hat{\beta}_1} = \sum_{i=1}^n 2 (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \cdot (-x_i) = 0 \quad (4)$$

(1) & (2) kaldes normal ligningerne

1. v. normalgleichung (1):

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i - n \hat{\beta}_0 - \hat{\beta}_1 \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} n \cdot \hat{\beta}_0 - \frac{1}{n} \hat{\beta}_1 \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x} \quad \Leftrightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (5)$$

Setze (5) in Minimierungsproblem:

$$\begin{aligned} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 &= \sum_{i=1}^n (y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) - \hat{\beta}_1 x_i)^2 \\ &= \sum_{i=1}^n [(y_i - \bar{y}) - \hat{\beta}_1 (x_i - \bar{x})]^2 \end{aligned}$$

=> Minimiere w.h.p. $\hat{\beta}_1$:

$$\sum_{i=1}^n 2 (y_i - \bar{y} - \hat{\beta}_1 (x_i - \bar{x})) \cdot (-1) (x_i - \bar{x}) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i - \bar{y}) (x_i - \bar{x}) - \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x}) (x_i - \bar{x}) = 0$$

$$\Rightarrow \hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y}) (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (6)$$

(5) & (6) kalles OLS-estimatoren:

$$(5) \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$(6) \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y}) (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Forventningsrett & konsistent estimator for β_1 :

- Forventningsrett:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (6)$$

Har at siden

$$y_i = \beta_0 + \beta_1 x_i + u_i \quad (7)$$

må også

$$\bar{y} = \beta_0 + \beta_1 \bar{x} + \bar{u} \quad (8)$$

Dette gir

$$(y_i - \bar{y}) = \beta_1 (x_i - \bar{x}) + (u_i - \bar{u}) \quad (9)$$

Setter (9) inn i (6):

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_{i=1}^n [\beta_1 (x_i - \bar{x}) + (u_i - \bar{u})] (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \frac{\sum_{i=1}^n \beta_1 (x_i - \bar{x})^2 + \sum_{i=1}^n (u_i - \bar{u})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \beta_1 \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} + \frac{\sum_{i=1}^n (u_i - \bar{u})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \beta_1 + \frac{\sum_{i=1}^n (u_i - \bar{u})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \beta_1 + \frac{\sum_{i=1}^n u_i (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (10) \end{aligned}$$

Tar forventningsverdien på begge sider og får:

$$E(\hat{\beta}_1) = \beta_1 + \frac{\sum_{i=1}^n E(u_i | x_i) (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \beta_1 \quad (11)$$

- Konsistens:

$$\hat{\beta}_1 = \beta_1 + \frac{\sum_{i=1}^n (u_i - \bar{u})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (10)$$

$$= \beta_1 + \frac{\frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})(x_i - \bar{x})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \quad (12)$$

Tar sannsynlighetsgrensen på begge sider og får:

$$\begin{aligned} \text{plim}_{n \rightarrow \infty} (\hat{\beta}_1) &= \text{plim}_{n \rightarrow \infty} \left[\beta_1 + \frac{\frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})(x_i - \bar{x})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \right] \\ &= \beta_1 + \frac{\text{plim}_{n \rightarrow \infty} \left[\frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})(x_i - \bar{x}) \right]}{\text{plim}_{n \rightarrow \infty} \left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right]} \\ &= \beta_1 + \frac{\text{Cov}(u, x)}{\text{Var}(x)} = \beta_1 \quad (13) \end{aligned}$$

der vi har brukt at $\text{Cov}(u, x) = 0$. Generelt antar vi at de empiriske momentene konvergerer i sannsynlighet mot sine respektive teoretiske momenter.

Alternativt kan vi ta utgangspunkt i (6):

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (6)$$

Sannsynlighetsgrensen gir:

$$\text{plim}_{n \rightarrow \infty} (\hat{\beta}_1) = \frac{\text{plim}_{n \rightarrow \infty} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) \right]}{\text{plim}_{n \rightarrow \infty} \left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right]} \quad (14)$$

$$= \frac{\text{Cov}(y, x)}{\text{Var}(x)} = \beta_1 \quad (15)$$

Dette alternativet viser til momentmetoden

Momentmetoden:

- Bruker at $E(u_i | x_i) = 0$ og $\text{Cov}(u_i | x_i) = 0$.
- Tar forventningsverdien på begge sider av (1):

$$E(y_i) = \beta_0 + \beta_1 E(x_i) + E(u_i) \quad (16)$$

- Trekker (16) fra (1) og får:

$$y_i - E(y_i) = \beta_1 (x_i - E(x_i)) + (u_i - E(u_i))$$

- Multipliserer gjennom med $x_i - E(x_i)$ og tar forventningen på begge sider:

$$(x_i - E(x_i))(y_i - E(y_i)) = \beta_1 (x_i - E(x_i))^2 + (u_i - E(u_i))(x_i - E(x_i))$$

$$\Rightarrow E[(x_i - E(x_i))(y_i - E(y_i))] = \beta_1 E[(x_i - E(x_i))^2] + E[(u_i - E(u_i))(x_i - E(x_i))]$$

$$\Rightarrow \text{Cov}(y, x) = \beta_1 \text{Var}(x) + \text{Cov}(u, x)$$

$$= \beta_1 \text{Var}(x)$$

$$\Rightarrow \beta_1 = \frac{\text{Cov}(y, x)}{\text{Var}(x)} \quad (15)$$

- Forholdet mellom teoretiske & empiriske momenter:

$\hat{\beta}_0$: (16) gir (5):

$$\hat{E}(y_i) = \hat{\beta}_0 + \hat{\beta}_1 \hat{E}(x_i) + \hat{E}(u_i) \Rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$\hat{\beta}_1$: (15) gir (14):

$$\hat{\beta}_1 = \frac{\text{Cov}(y, x)}{\text{Var}(x)} \Rightarrow \hat{\beta}_1 = \frac{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

Dekomponering av variasjonen i y - forklaringsmal:

- Normal ligningene (3) & (4) gir at:

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = \sum_{i=1}^n (\hat{u}_i) = 0 \quad (17)$$

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = \sum_{i=1}^n (\hat{u}_i x_i) = 0 \quad (18)$$

- Definerer sum of squares:

$$SST \equiv \sum_{i=1}^n (y_i - \bar{y})^2 \quad (19)$$

$$SSE \equiv \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad (20)$$

$$SSR \equiv \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (\hat{u}_i)^2 \quad (21)$$

- $SST = SSE + SSR$:

$$\begin{aligned} SST &= \sum_{i=1}^n (y_i - \bar{y})^2 \\ &= \sum_{i=1}^n [(y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})]^2 \\ &= \sum_{i=1}^n [(y_i - \hat{y}_i)^2 + 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) + (\hat{y}_i - \bar{y})^2] \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + 2 \sum_{i=1}^n \hat{u}_i (\hat{y}_i - \bar{y}) + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\ &= SSR + 2 \left(\sum_{i=1}^n \hat{u}_i \hat{y}_i - \bar{y} \sum_{i=1}^n \hat{u}_i \right) + SSE \\ &= SSE + SSR + 2 \left(\sum_{i=1}^n \hat{u}_i \hat{\beta}_0 + \sum_{i=1}^n \hat{u}_i \hat{\beta}_1 x_i - \bar{y} \sum_{i=1}^n \hat{u}_i \right) \\ &= SSE + SSR + 2 \left(\hat{\beta}_0 \sum_{i=1}^n \hat{u}_i + \hat{\beta}_1 \sum_{i=1}^n \hat{u}_i x_i - \bar{y} \sum_{i=1}^n \hat{u}_i \right) \end{aligned}$$

(17) & (18) gir at det siste leddet er null:

$$SST = SSE + SSR \quad (22)$$

Den multiple determinasjonskoeffisienten R^2 :

- Definisjon:

$$R^2 = \frac{SSE}{SST} \quad (23)$$

$$= \frac{SST - SSR}{SST}$$

$$= 1 - \frac{SSR}{SST} \quad (24)$$

der $0 \leq R^2 \leq 1$

- Justert R^2 : \bar{R}^2

$$\bar{R}^2 = 1 - \frac{\frac{SSR}{n-k-1}}{\frac{SST}{n-1}} \Rightarrow \text{"Straffer" parameteriske likninger}$$

Variansen til $\hat{\beta}_1$:

- Tar utgangspunkt i (10):

$$\hat{\beta}_1 = \beta_1 + \frac{\sum_{i=1}^n u_i (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \beta_1 + \frac{\sum_{i=1}^n u_i d_i}{SST_x} \quad (25)$$

der $d_i = (x_i - \bar{x})$ & $SST_x = \sum_{i=1}^n (x_i - \bar{x})^2$

Utnytter at $\text{Var}(ax) = a^2 \text{Var}(x)$ og får:

$$\text{Var}(\hat{\beta}_1 | x) = 0 + \frac{\sum_{i=1}^n \text{Var}(u_i | x) \cdot d_i^2}{SST_x^2} = \frac{\sum_{i=1}^n \sigma_i^2 d_i^2}{SST_x^2} \quad (26)$$

Hvis homoskedastisitet, dvs hvis $\sigma_i = \sigma_j = \sigma$:

$$\text{Var}(\hat{\beta}_1 | x) = \frac{\sigma^2 \sum_{i=1}^n d_i^2}{SST_x^2} = \frac{\sigma^2 SST_x}{SST_x^2} = \frac{\sigma^2}{SST_x} \quad (27)$$

Estimering av varians & standardavvikelse

- Variansen til u_i : $\text{Var}(u_i | x_i) = \sigma^2$

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (\hat{u}_i)^2 = \frac{\text{SSR}}{n-2} \quad (28)$$

- Standardavvikelse til u_i :

$$\hat{\sigma} = \sqrt{\hat{\sigma}^2} = \sqrt{\frac{\text{SSR}}{n-2}} = \frac{\sqrt{\text{SSR}}}{\sqrt{n-2}} \quad (29)$$

- Variansen til $\hat{\beta}_1$:

Erstatter den teoretiske variansen til u

i (27) med den empiriske:

$$\widehat{\text{Var}}(\hat{\beta}_1 | x) = \frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (30)$$

- Standardavvikelse til $\hat{\beta}_1$:

$$\text{Se}(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad (31)$$

Multipl regression:

- Modell: $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i$ (1)

der $\bar{y} = \beta_0 + \beta_1 \bar{x}_1 + \dots + \beta_k \bar{x}_k + \bar{u}$ (2)

- Forutsetninger:

MLR 1: $y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$
innebærer at $\text{Cov}(u_i, u_j | X) = 0$

MLR 2: Tilfeldig trekking fra populasjonen

MLR 3: Variasjon i X & ikke perfekt

kollinearitet mellom noen variabler

MLR 4: $E(u_i | x_{i1}, \dots, x_{ik}) = 0$
innebærer at $E(u_i) = 0$
 $\text{Cov}(u_i, x_{ij}) = 0$

MLR 5: $\text{Var}(u_i | x_{i1}, \dots, x_{ik}) = \sigma^2$

- forventning:

MLR 4 innebærer at $E(\hat{\beta}_j | X) = \beta_j$

- Multipl determinasjonskoeffisient:

har tidligere vist at

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

justert R^2 : \bar{R}^2

$$\bar{R}^2 = 1 - \frac{\frac{SSR}{n-k-1}}{\frac{SST}{n-1}}$$

(3) "straffer" parameterne likninger

- Variansen til $\hat{\beta}_j$:

$$\text{Var}(\hat{\beta}_j | X) = \frac{\hat{\sigma}^2}{\text{SST}_j (1 - R_j^2)} \quad (4)$$

der $\text{SST}_j = \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$

og R_j^2 er R^2 fra

$$x_{ij} = \alpha_0 + \alpha_1 x_{i1} + \dots + \alpha_{j-1} x_{i,j-1} + \alpha_{j+1} x_{i,j+1} + \dots + \alpha_k x_{ik} + u_i^j$$

- Estimering av varians & standardavvik:

Erstatter teoretisk moment med empirisk:

$$\widehat{\text{Var}}(\hat{\beta}_j | X) = \frac{\hat{\sigma}^2}{\text{SST}_j (1 - R_j^2)}$$

der $\hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{u}_i^2}{n - k - 1} = \frac{\text{SSR}}{n - k - 1}$

Standardavvik:

$$\text{Se}(\hat{\beta}_j) = \frac{\hat{\sigma}}{\sqrt{\text{SST}_j (1 - R_j^2)}}$$

Hypotese testing / inferens:

- Generell testing:

H_0 : Påstand A er sann

H_1 : Påstand A er ikke sann / H_0 er ikke sann

Hypotesene H_0 & H_1 er gjensidig utelukkende og fyller tilsammen hele utfallsrommet.

- t-testen: Testing av enkeltparameter

t-testens testobservator $t_{\hat{\beta}_1}$ er gitt ved:

$$t_{\hat{\beta}_1} = \frac{\hat{\beta}_1 - \beta_1^0}{\frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}} = \frac{\hat{\beta}_1 - \beta_1^0}{Se(\hat{\beta}_1)} \sim t_{n-k-1} \quad (1)$$

Behold H_0 hvis $|t_{\hat{\beta}_1}| < t_c$,

forkast H_0 hvis $|t_{\hat{\beta}_1}| \geq t_c$

- F-testen: Testing av flere parametre:

SSR vil alltid øke når en modell pålegges restriksjoner. Ønsker imidlertid å teste om

økingen i SSR er så stor at vi kan avvise hypotesen om at restriksjonen er sann.

Altså; er SSR_r tilstrekkelig større enn SSR_{ur} ? //

Testobservatoren F_{obs} er gitt ved:

$$F_{obs} \equiv \frac{\frac{SSR_r - SSR_{ur}}{q}}{\frac{SSR_{ur}}{n-k-1}} \sim F_{q, n-k-1} \quad (2)$$

Behold H_0 hvis $F_{obs} < F_c$,

forkast H_0 hvis $F_{obs} \geq F_c$

Alternativt:

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

$$\Rightarrow SSR = SST(1 - R^2)$$

$$\Rightarrow F_{obs} = \frac{\frac{SSR_r - SSR_{ur}}{q}}{\frac{SSR_{ur}}{n-k-1}}$$

$$= \frac{\frac{SST(1 - R_r^2) - SST(1 - R_{ur}^2)}{q}}{\frac{SST(1 - R_{ur}^2)}{n-k-1}}$$

$$= \frac{\frac{-R_r^2 + R_{ur}^2}{q}}{\frac{1 - R_{ur}^2}{n-k-1}}$$

$$= \frac{\frac{R_{ur}^2 - R_r^2}{q}}{\frac{1 - R_{ur}^2}{n-k-1}} \quad (3)$$

- Test av uvelketlig signifikans:

H_0 : Ingen variabler bidrar til forklaring.

$$\Rightarrow R_r^2 = 0 \quad \text{og} \quad q = k$$

$$\Rightarrow F_{obs} = \frac{\frac{R_r^2}{k}}{1 - R_r^2} \quad (4)$$

- Spesielt tilfelle: Bruk av t-test

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i$$

$$H_0: \beta_1 = \beta_2 \rightarrow \beta_1 - \beta_2 = 0$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_2 (x_{i1} - x_{i1}) + u_i$$

$$= \beta_0 + (\beta_1 - \beta_2) x_{i1} + \beta_2 (x_{i1} + x_{i2}) + u_i$$

$$= \beta_0 + \theta x_{i1} + \beta_2 z_i + u_i$$

$$\text{der } \theta = \beta_1 - \beta_2 \quad \& \quad z_i = x_{i1} + x_{i2}$$

$$H_0: \theta = 0$$

- Konfidensintervall:

1. Konfidensintervall for β_1 når σ er kjent:

Siden $\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{n}\right)$ vil den normaliserte

variabelen være fordelt slik at

$$\frac{\hat{\beta}_1 - \beta_1}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$$

Har da at et 95% konfidensintervall gis ved:

$$P\left(-1,96 \leq \frac{\hat{\beta}_1 - \beta_1}{\frac{\sigma}{\sqrt{n}}} \leq 1,96\right) = 0,95$$

$$\Rightarrow \hat{\beta}_1 \pm 1,96 \frac{\sigma}{\sqrt{n}}$$

2. Konfidensintervall for $\hat{\beta}_1$, når σ er ukjent:

$$\text{Estimerer } \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\text{Siden } \frac{\hat{\beta}_1 - \beta_1}{\frac{\hat{\sigma}}{\sqrt{n}}} \sim N(0,1)$$

er (når vi setter inn det empiriske mom.)

$$\frac{\hat{\beta}_1 - \beta_1}{\frac{\hat{\sigma}}{\sqrt{n}}} \sim t_{n-1}$$

=> 95% konfidensintervall gitt ved:

$$\hat{\beta}_1 \pm 1,96 \frac{\hat{\sigma}}{\sqrt{n}}$$

- Konfidensintervall i regresjonslikninger:

$$\text{Siden } \frac{\hat{\beta}_j - \beta_j}{\text{Se}(\hat{\beta}_j)} \sim t_{n-k-1}$$

er konfidensintervallet til den ukjente

β_j gitt ved

$$\hat{\beta}_j \pm t_c \cdot \text{Se}(\hat{\beta}_j) \quad (5)$$

der t_c er den $(1 - \frac{1}{2}p)$ persentilen i

en t_{n-k-1} -fordeling (p er p -verdi). Fls. $p=0,05$

gir $(1 - \frac{1}{2}p) = 1 - \frac{0,05}{2} = 0,975$ persentil.

$$\text{Eksempel: } \hat{y} = -4,38 + 1,084x_1 + 0,0217x_2 \quad n=32, R^2=0,918$$

$(0,47) \quad (0,060)$

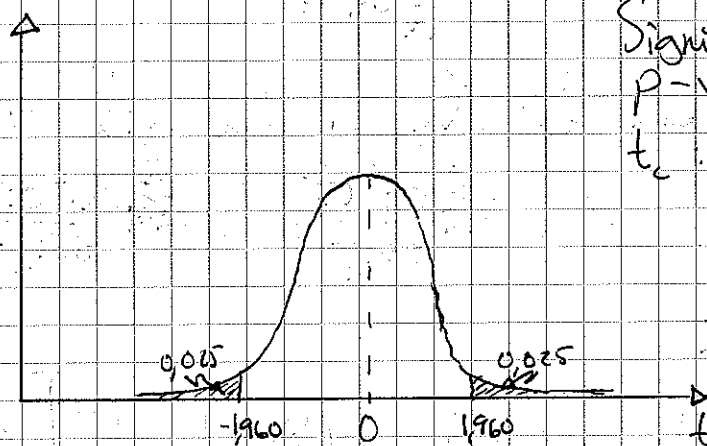
Konf.int.v.: $1,084 \pm 2,045 \cdot 0,060 \Rightarrow [0,961, 1,21]$

- p-verdi: Det minste signifikansnivået som en nullhypotese kan forkastes ved, gitt observert t-verdi.

Alternativ formulering:

Sannsynligheten for å observere t_{obs} gitt at nullhypotesen er sann.

- t-fordelingen:



Signifikansnivå: 95%
p-verdi: 0,05
 t_c : 1,96
for tosidig
test

Forkaster H_0 hvis $|t_{obs}| \geq 1,96$

Heteroskedastisitet:

Hartidligere antatt at $\text{Var}(u_i | X) = \sigma^2$, altså homoskedastisitet.

Generelt innebærer heteroskedastisitet at $\text{Var}(u_i | X) = \sigma_i^2$, altså ikke konstant varians i restleddet.

Tar utgangspunkt i

$$\hat{\beta}_1 = \beta_1 + \frac{\sum_{i=1}^n (u_i - \bar{u})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \beta_1 + \frac{\sum_{i=1}^n u_i d_i}{SST_x} \quad \text{der } d_i = x_i - \bar{x} \text{ og } SST_x = \sum_{i=1}^n (x_i - \bar{x})^2$$

Variansen til $\hat{\beta}_1$ er da gitt ved:

$$\text{Var}(\hat{\beta}_1) = 0 + \frac{\sum_{i=1}^n d_i^2 \sigma_i^2}{SST_x^2} = \frac{\sum_{i=1}^n d_i^2 \sigma_i^2}{SST_x^2}$$

fordi

$$\begin{aligned} \text{Var}(\hat{\beta}_1 | X) &= E\left[\left(\hat{\beta}_1 - E(\hat{\beta}_1)\right)^2\right] = E\left[\left(\beta_1 + \frac{\sum_{i=1}^n u_i d_i}{SST_x} - \beta_1\right)^2\right] \\ &= E\left[\left(\frac{\sum_{i=1}^n u_i d_i}{SST_x}\right)^2\right] = E\left[\frac{\left(\sum_{i=1}^n u_i d_i\right)^2}{SST_x^2}\right] \\ &= \frac{\sum_{i=1}^n \text{Var}(u_i | X) d_i^2}{SST_x^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{SST_x^2} \end{aligned}$$

Så lenge ikke $\sigma_i^2 = \sigma_j^2 = \sigma^2$ har vi heteroskedastisitet.

- Konsekvenser av heteroskedastisitet:

Ved heteroskedastisitet gjelder ikke de vanlige MKM-formulene for varians & standardavvik. Ei heller t-testen eller F-testen eller LM-testen.

Alle disse baseres nemlig på konstant varians i restleddet.

Heteroskedastisitet gir imidlertid ikke skjevhet eller inkonsistens i OLS-estimatene, ei heller i godness-of-fit-målene R^2 & \bar{R}^2 .

- Kan løse problemer knyttet til heteroskedastisitet vha korrigering eller hensyn til heteroskedastisiteten.

- Estimering av heteroskedastisitetshorrigerede eller robuste standardavvik:

Benytter metoden til Kenneth White:

$$\text{Var}(\hat{\beta}_1 | X) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{\text{SST}_x}$$

⇒ Setter inn for u_i^2 istedet for σ_i^2 :

$$\widehat{\text{Var}}(\hat{\beta}_1 | X) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 u_i^2}{\text{SST}_x} \quad (\text{i bivariate analysen})$$

⇒ Finner deretter standardavviket:

$$\begin{aligned} \widehat{\text{Se}}(\hat{\beta}_1 | X) &= \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2 u_i^2}{\text{SST}_x}} \\ &= \frac{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 u_i^2}}{\text{SST}_x} \end{aligned}$$

Fremgangsmåten er altså å først estimere parametrene vha OLS, & deretter bruke residualene til å estimere

heteroskedastisitetshorrigerede standardavvik.

Dermed kan man ta i bruk inferens.

Ved multippel regresjon:

$$\widehat{\text{Var}}(\hat{\beta}_j | X) = \frac{\sum_{i=1}^n \hat{r}_{ij}^2 u_i^2}{\text{SSR}_j}$$

$$\Rightarrow \widehat{\text{Se}}(\hat{\beta}_j | X) = \frac{\sqrt{\sum_{i=1}^n \hat{r}_{ij}^2 u_i^2}}{\text{SSR}_j}$$

der \hat{r}_{ij} er iende residual i regresjonen med x_j som avh. variabel

og alle andre x som uavhengige, og SSR_j er sum of squared residuals fra denne regresjonen.

- Testing for heteroskedastisitet:

Kan ta utgangspunkt i:

$$y_i = \beta_0 + \sum_{j=1}^k \beta_j x_{ij} + u_i$$

$$H_0: \text{Var}(u|X) = \sigma^2 \text{ eller } E(u^2|X) = \sigma^2$$

H_1 : H_0 er ikke sann.

H_1 kan spesifiseres på mange måter, f.eks

$$\text{slik: } u^2 = \delta_0 + \delta_1 x_1 + \delta_2 x_2 + \delta_3 x_1 x_2 + \text{støy}$$

Nullhypotesen $\delta_1 = \delta_2 = \delta_3 = 0$ kan testes

via F-testen eller LM-testen.

Slike testing for heteroskedastisitet

kalles Breusch-Pagan-testing for

heteroskedastisitet (test av helhetlig sign. der

Altså: $F_{obs} = \frac{\frac{R^2}{k}}{\frac{(1-R^2)}{n-k-1}}$)

Estimerer $\hat{u}_i^2 = \hat{\delta}_0 + \hat{\delta}_1 x_{i1} + \hat{\delta}_2 x_{i2} + \hat{\delta}_3 x_{i1} x_{i2}$

og tester $H_0: \delta_1 = \delta_2 = \delta_3 = 0$

Testing for heteroskedastisitet kan gjøres ved

å feste forskjellige funksjonsformer der

u^2 er på venstre siden & x_j er på høyresiden.

White: Test
 $u^2 = \delta_0 + \delta_1 x_1 + \delta_2 x_2$

- Vektet OLS (WLS):

Er en transformasjon slik at restleddet får konstant varians.

Antar at $y_i = \beta_0 + \sum_{j=1}^k \beta_j x_{ij} + u_i$ der variansen er gitt ved: $\text{Var}(u_i | X) = \sigma^2 \cdot h_i(X)$

Dividerer nå gjennom struktureligningen med $\frac{1}{\sqrt{h_i}}$ dette gir:

$$\frac{y_i}{\sqrt{h_i}} = \frac{\beta_0}{\sqrt{h_i}} + \sum_{j=1}^k \frac{\beta_j x_{ij}}{\sqrt{h_i}} + \frac{u_i}{\sqrt{h_i}}$$

Variansen blir da

$$\text{Var}\left(\frac{u_i}{\sqrt{h_i}} \mid X\right) = \left(\frac{1}{\sqrt{h_i}}\right)^2 \cdot \text{Var}(u_i | X) = \frac{1}{h_i} \cdot \sigma^2 \cdot h_i = \sigma^2$$

Her har altså restleddet konstant varians.

Ved transformasjon av struktureligningen får vi konstant varians i restleddet, og vi kan da benytte oss av vanlige inferensmetoder.

Metoden der vi estimerer h_i (\hat{h}_i) kalles Feasible GLS.

Funktionsform:

- Logaritmer:

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

$$\Rightarrow \frac{1}{y} \cdot dy = \beta_1 \cdot dx$$

\Rightarrow Marginal endring i x gir relativ endring i y (målt i prosent).

Ved betydelig effekt på dy av dx er følgende et mer presist estimat:

$$\% \Delta y = 100 \cdot (e^{\hat{\beta}_1 dx} - 1)$$

Eksempel:

$$\ln y = 9,23 + 0,306 x$$

$$\Rightarrow \% \Delta y = 100 \cdot (e^{0,306 \cdot 1} - 1) = 35,8\%$$

\Rightarrow hadde bare vært 30,6% hvis estimert direkte.

$$y = \beta_0 + \ln x$$

$$\Rightarrow dy = \beta_1 \cdot \frac{1}{x} \cdot dx$$

\Rightarrow Prosentvis endring i x gir endring i y .

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

$$\rightarrow \frac{1}{y} dy = \beta_1 \frac{1}{x} dx$$

=> Prosentvis endring i x gir prosentvis endring i y . Mål på elasticitetsform.

- Dummy:

$$y = \beta_0 + \beta_1 x_i \quad \text{der} \quad x_i = 1 \quad \text{hvis event inntreffer}$$

og $x_i = 0$ hvis event ikke inntreffer.

Forskjellige skjæringspunkt: & stigningsfall:

$$y = \beta_0 + \beta_1 D_{\text{MENN}} + \beta_2 x + \beta_3 D_{\text{MENN}} \cdot x$$

Skjæringspunkt:

Menn: $\beta_0 + \beta_1$

Kvinner: β_0

Stigningsfall:

Menn: $\beta_2 + \beta_3$

Kvinner: β_2

Dummyer kan også brukes til å skille/kategorisere variabler på nominalnivå, f.eks. kommuner.

- Sammenligning av to grupper: Chows test

Generell modell: $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i$

Ønsker å teste om f.eks. kvinner og menn

har forskjellige skjæringspunkt & stigningstall.

Problem: Kan være vanskelig å oppnå signifikans når mange variabler er involvert. vha. dummyer

Løsning:

Beregne modellen kun for kvinner, finn SSR_1 .

Beregne modellen kun for menn, finn SSR_2 .

Beregne modellen for hele populasjonen, finn SSR_p .

Chows F-test:

$$F = \frac{SSR_p - (SSR_1 + SSR_2)}{\frac{k+1}{n-2(k+1)}(SSR_1 + SSR_2)}$$

H_0 : Det er ingen forskjell i parametrene mellom kvinner & menn.

- Dummy i venstre sidevariabelen: Linear probability model (LPM). Merk: Kan ta verdier >1 & <0 . 23

- Annengradsledd:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + u_i$$

$$\Rightarrow dy = \beta_1 dx + 2\beta_2 dx \cdot x$$

$$\Rightarrow \frac{dy}{dx} = \beta_1 + 2\beta_2 x$$

- Samspillsledd:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i1} x_{i2} + u_i$$

$$\Rightarrow dy = \beta_1 dx_1 + \beta_3 x_2 \cdot dx_1$$

$$\Rightarrow \frac{dy}{dx_1} = \beta_1 + \beta_3 x_2$$

Valg mellom like-nøstede modeller

- Valg mellom nøstede modeller, dvs. der én modell fremkommer som en annen modell med null restriksjoner, kan gjøres via en F-test eller LM-test

- like-nøstede modeller: Kan ikke skrives som spesialtilfeller av hverandre.

Eksempel: $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i$ (1)

$$y_i = \gamma_0 + \gamma_1 z_{i1} + \gamma_2 z_{i2} + v_i$$
 (2)

- Merk: Hvis begge modellene har samme venstresidevariabel, vil vi alt annet likt velge den med størst forklaringskraft (\bar{R}^2). Imidlertid mye annet enn \bar{R}^2 som spiller inn.

- Alternativ testprosedyre 1: Mizon & Richard

Slå sammen (1) & (2):

$$y_i = \alpha_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \gamma_1 z_{i1} + \gamma_2 z_{i2} + w_i$$

Tester deretter om: A: $\beta_1 = \beta_2 = 0$

B: $\gamma_1 = \gamma_2 = 0$

- Hvis kun én av hypotesene forkastes kan vi velge modell.
 - Hvis begge forkastes er ingen gyldige forenklinger av en mer generell modell.
 - Hvis ingen forkastes er testen inkonklusiv.
- Alternativ testprosedyre? Davidson & Mackinnon
intuisjon: Hvis (1) er sann er de estimerte verdiene fra (2) insignifikante når de puttes inn i (1).

Fremgangsmåte:

1. Estimer (2) for å finne \hat{y}_i .
2. Estimer varianten $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \delta_1 \hat{y}_i + u_i$ av (1).
3. Gjennomfør en t-test av $\hat{\delta}_1$. Forkast (1) dersom H_0 forkastes.
4. Gjenta samme prosedyre med utg. p. i (2).

Merk: Også her kan testen være inkonklusiv.

Endogen forklaringsvariabel

Kilder til endogenitet:

1. Utelatte variable
2. Målefeil
3. Simultanitet

Krav til eksogenitet i forklaringsvariablen:

$$\text{Cov}(x, u) = 0$$

1. Utelatt variabelskjevhet:

Sann modell: $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i$

Estimert modell: $y_i = \beta_0 + \beta_1 x_{i1} + v_i$ der $v_i = \beta_2 x_{i2} + u_i$

MKM-estimatoren for β_1 :

$$\begin{aligned}\tilde{\beta}_1 &= \frac{\sum_{i=1}^n (y_i - \bar{y})(x_{i1} - \bar{x}_1)}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2} \quad \text{der } y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i \\ &= \frac{\sum_{i=1}^n [\beta_1 (x_{i1} - \bar{x}_1) + \beta_2 (x_{i2} - \bar{x}_2) + (u_i - \bar{u})] (x_{i1} - \bar{x}_1)}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2} \quad \text{og } \bar{y} = \beta_0 + \beta_1 \bar{x}_1 + \beta_2 \bar{x}_2 + \bar{u} \\ &= \frac{\sum_{i=1}^n \beta_1 (x_{i1} - \bar{x}_1)^2 + \sum_{i=1}^n \beta_2 (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2) + \sum_{i=1}^n (u_i - \bar{u})(x_{i1} - \bar{x}_1)}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2} \\ &= \beta_1 \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2} + \beta_2 \frac{\sum_{i=1}^n (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2} + \frac{\sum_{i=1}^n (u_i - \bar{u})(x_{i1} - \bar{x}_1)}{\sum_{i=1}^n (x_{i1} - \bar{x}_1)^2} \\ &= \beta_1 + \beta_2 \frac{\frac{1}{n} \sum_{i=1}^n (x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{\frac{1}{n} \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2} + \frac{\frac{1}{n} \sum_{i=1}^n (u_i - \bar{u})(x_{i1} - \bar{x}_1)}{\frac{1}{n} \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2}\end{aligned}$$

$$\lim_{n \rightarrow \infty} (\tilde{\beta}_1) = \beta_1 + \beta_2 \cdot \frac{\text{Cov}(x_1, x_2)}{\text{Var}(x_1)} = \beta_1 + \beta_2 \delta_1$$

$$\text{der } \delta_1 = \frac{\text{Cov}(x_1, x_2)}{\text{Var}(x_1)}$$

Positiv variabelskjevhet: $\tilde{\beta}_1 > \beta_1$

1.1. $\beta_2 > 0$ & $\delta_1 > 0$

Utelatt variabel påvirker avhengig variabel positivt og korrelerer samtidig positivt med inkludert variabel.

1.2. $\beta_2 < 0$ & $\delta_1 < 0$

Utelatt variabel påvirker avhengig variabel negativt & korrelerer samtidig negativt med inkludert variabel.

Negativ variabelskjevhet: $\tilde{\beta}_1 < \beta_1$

1.3. $\beta_2 < 0$ & $\delta_1 > 0$

Utelatt variabel påvirker avhengig variabel negativt & korrelerer samtidig positivt med inkludert variabel.

1.4. $\beta_2 > 0$ & $\delta_1 < 0$

Utelatt variabel påvirker avhengig variabel positivt & korrelerer samtidig negativt med inkludert variabel.

Ingen variabelskjevhet: $\tilde{\beta}_1 = \beta_1$

1.5. $\beta_2 = 0$

Utelatt variabel har ingen innvirkning på avhengig variabel.

1.6. $\delta_1 = 0$

Utelatt variabel korrelerer ikke med inkludert variabel

Variansen til $\tilde{\beta}_1$:

$$\text{Har at } \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\text{SST}_1 (1 - R_1^2)}$$

$$\text{Var}(\tilde{\beta}_1) = \frac{\sigma^2}{\text{SST}_1}$$

Ser at så lenge $\text{Cov}(x_1, x_2) \neq 0$, er $\text{Var}(\tilde{\beta}_1) < \text{Var}(\hat{\beta}_1)$

Dette innebærer at dersom $\beta_2 = 0$, er det best å utelate x_2 fra modellen!

Intuisjon: Å inkludere irrelevante variable

fører bare til høyere varians i $\hat{\beta}_1$ gitt multi-kollinearitetsproblemet. α

2 Målefeil:

Antar at den riktige modellen er gitt ved:

$$y^* = \beta_0 + \beta_1 x^* + u$$

Målefeil når de sanne y^* & x^* ikke observeres, men hhv:

$$y = y^* + \varepsilon$$

$$x = x^* + e$$

Presisering av forutsetninger:

$$E(u) = E(\varepsilon) = E(e) = 0$$

$$\text{Var}(u) = \sigma_u^2, \quad \text{Var}(\varepsilon) = \sigma_\varepsilon^2, \quad \text{Var}(e) = \sigma_e^2$$

$$\text{Cov}(u, x) = \text{Cov}(u, \varepsilon) = \text{Cov}(u, e) = \text{Cov}(\varepsilon, e) = 0$$

Tilfelle 1: Kun målefeil i venstresidevariabelen.

Estimerer altså:

$$y_i = \beta_0 + \beta_1 x_i + (u_i + \varepsilon_i) = \beta_0 + \beta_1 x + v \quad \text{der } v = u + \varepsilon$$

hvorvidt $\hat{\beta}_1$ er forventningsrett eller ikke

avhenger av $\text{Cov}(v, x)$, siden $\tilde{\beta}_1 = \beta_1 + \frac{\text{Cov}(v, x)}{\text{Var}(x)}$

Har allerede antatt at $\text{Cov}(u, x) = 0$. Så lenge $\text{Cov}(e, x) = 0$ er da også $\text{Cov}(v, x) = 0$. Dette er normalt en rimelig antagelse.

Inndertid er variansen til restleddet større:

$$\text{Var}(v) = \text{Var}(u + e) = \sigma_u^2 + \sigma_e^2 > \sigma_u^2$$

Dette innebærer høyere varians & standardavvik for OLS-estimatoren.

Tilfelle 2: Kun målefeil i forklaringsvariabelen, målefeiten er ukorrelert med observert verdi ($\text{Cov}(e, x) = 0$).

Estimer altså:

$$y = \beta_0 + \beta_1(x - e) + u = \beta_0 + \beta_1 x + v \quad \text{der } v = u - \beta_1 e$$

Siden vi har antatt at $\text{Cov}(u, x) = 0$, & samtidig at $\text{Cov}(e, x) = 0$, må også

$$\text{Cov}(v, x) = 0. \quad \text{Dermed er } \tilde{\beta}_1 = \beta_1 + \frac{\text{Cov}(v, x)}{\text{Var}(x)} = \beta_1$$

Variansen til restleddet er inndertid større:

$$\text{Var}(v) = \text{Var}(u - \beta_1 e) = \sigma_u^2 + \beta_1^2 \sigma_e^2 > \sigma_u^2$$

Dermed høyere varians & st. avvik for OLS-estimatoren.

Tilfelle 3: Kun målefeil i forklaringsvariabelen,
målefeilen er ukorrelet med
observerbar variabel ($\text{Cov}(e, x^*) = 0$)

Estimerer altså:

$$y = \beta_0 + \beta_1(x - e) + u = \beta_0 + \beta_1 x + v \quad \text{der } v = u - \beta_1 e$$

Siden $\text{Cov}(e, x^*) = 0$ er

$$\text{Cov}(e, x) = \text{Cov}(e, x^* + e) = \sigma_e^2$$

Dermed er

$$\text{Cov}(v, x) = \text{Cov}(u - \beta_1 e, x^* + e) = -\beta_1 \sigma_e^2$$

Her er mas. restleddet i den estimerte
likningen korrelert med den inkluderte
forklaringsvariabelen.

Videre er:

$$\text{Var}(x) = \text{Var}(x^* + e) = \sigma_{x^*}^2 + \sigma_e^2$$

Dette gir:

$$\tilde{\beta}_1 = \beta_1 + \frac{\text{Cov}(v, x)}{\text{Var}(x)} = \beta_1 + \frac{-\beta_1 \sigma_e^2}{\sigma_{x^*}^2 + \sigma_e^2} = \beta_1 \left(1 - \frac{\sigma_e^2}{\sigma_{x^*}^2 + \sigma_e^2} \right)$$

Sett at: $\beta_1 > 0$ gir $\tilde{\beta}_1 < \beta_1$

$\beta_1 < 0$ gir $\tilde{\beta}_1 > \beta_1$

Altså systematisk underestimering av absoluttverdien til effekten av x .

Videre er:

$$\begin{aligned}\tilde{\beta}_1 &= \beta_1 \left(1 - \frac{\sigma_e^2}{\sigma_{x^*}^2 + \sigma_e^2} \right) \\ &= \beta_1 \left(\frac{\sigma_{x^*}^2}{\sigma_{x^*}^2 + \sigma_e^2} \right) \quad \Leftrightarrow \quad \beta_1 \frac{\text{Var}(x^*)}{\text{Var}(x)} \\ &= \beta_1 \left(1 + \frac{\sigma_e^2}{\sigma_{x^*}^2} \right)\end{aligned}$$

Ser at desto større variansen til målefeilen er relativt til variansen til den samme forklaringsvariabelen, desto mer alvorlig er målefeilen.

Noise to signal ratio: $\frac{\sigma_e^2}{\sigma_{x^*}^2}$

3. Simultanitet:

Grunnleggende illustrasjon av
simultanitetsproblemet:

$$y = \beta_1 x + \alpha_1 z_1 + u_1 \quad (1)$$

men også

$$x = \beta_2 y + \alpha_2 z_2 + u_2 \quad (2)$$

Altså:

$$x \Rightarrow y$$

men også

$$y \Rightarrow x$$

Med determineres y & x simultant av
 z_1 & z_2 .

Estimering på redusert form ved å sette
(2) inn i (1):

$$y = \beta_1 (\beta_2 y + \alpha_2 z_2 + u_2) + \alpha_1 z_1 + u_1$$

$$\Rightarrow (1 - \beta_1 \beta_2) y = \beta_1 \alpha_2 z_2 + \alpha_1 z_1 + (u_1 + \beta_1 u_2)$$

$$\Rightarrow y = \frac{\alpha_1}{1 - \beta_1 \beta_2} z_1 + \frac{\beta_1 \alpha_2}{1 - \beta_1 \beta_2} z_2 + \frac{u_1 + \beta_1 u_2}{1 - \beta_1 \beta_2}$$

$$= \pi_1 z_1 + \pi_2 z_2 + v \quad (3)$$

der $\pi_1 = \frac{\alpha_1}{1-\beta_1\beta_2}$, $\pi_2 = \frac{\beta_1\alpha_2}{1-\beta_1\beta_2}$ & $v = \frac{u_1 + \beta_1 u_2}{1-\beta_1\beta_2}$

Estimering av (3) via OLS vil i prinsippet gi forventningsrette & konsistente estimatører for π_1 & π_2 .

Problemet er imidlertid å identifisere β_1 !

Generelt vil OLS-estimering av (1) gi skjeve estimatører av både β_1 & α_1 .

Hvis man ønsker å finne kausaleffekten av x på y :

1. Estimerer redusert form-regresjonen

$$x = \gamma_0 + \gamma_1 z_1 + \gamma_2 z_2 \quad (4)$$

der z_2 betraktes som mulig instrument.

2. Gitt at γ_2 er signifikant estimeres

(1) der x byttes ut med \hat{x} fra (4).

Kausaleffekten kan altså estimeres

via IV-2SLS så lenge (2) inneholder

minst én eksogen variabel som er

ekskudert i (1).

Mulige løsninger ved endogenitet i
forklaringsvariabelen:

1. Overse problemet.
2. Proxy & OLS
3. Paneldata & fixed effects eller first difference
4. Instrumentvariabelmetoden

Tidsseriedata:

Modell: $y_t = \beta_0 + \beta_1 x_t + u_t$

Forutsetninger:

TS 1: $y_t = \beta_0 + \beta_1 x_{t1} + \dots + \beta_k x_{tk} + u_t$

TS 2: Variasjon i X & ikke perfekt

kollinearitet mellom uavhengige variabler

TS 3: $E(u_t | X) = 0$

TS 4: $\text{Var}(u_t | X) = \sigma^2$

TS 5: $\text{Corr}(u_t, u_s | X) = 0$ for alle $t \neq s$

Altså ingen seriekorrelasjon

TS 6: Restleddene er normalfordelt;

$$N \sim (0, \sigma^2)$$

Under antagelsene TS1-TS3 er

OLS-estimatorene forventningsrette:

$$E(\hat{\beta}_j) = \beta_j$$

Finite distributed lag (FDL) model av orden q :

$$y_t = \alpha_0 + \beta_0 x_t + \beta_1 x_{t-1} + \dots + \beta_q x_{t-q} + u_t$$

Kortidseffekten (umiddelbar effekt): β_0

Langtidseffekten: $\beta_0 + \beta_1 + \dots + \beta_q$

(Merk at det her er snakk om varig endring i x).

Varians & standardavvik i $\hat{\beta}_j$:

Gitt TSS - TSS er:

$$\text{Var}(\hat{\beta}_j | X) = \frac{\sigma^2}{\text{SST}_j (1 - R_j^2)}$$

der $\text{SST}_j \equiv \sum_{i=1}^n (x_i - \bar{x}_j)^2$
og R_j^2 er R^2 fra

Estimatet er:

$$\widehat{\text{Var}}(\hat{\beta}_j | X) = \frac{\hat{\sigma}^2}{\text{SST}_j (1 - R_j^2)}$$

$$x_i = \alpha_0 + \alpha_1 x_{t-1} + \dots$$

der $\hat{\sigma}^2 = \frac{\text{SSR}}{n - k - 1}$

Standardavviket estimeres følgende som:

$$\text{Se}(\hat{\beta}_j) = \frac{\hat{\sigma}}{\sqrt{\text{SST}_j (1 - R_j^2)}}$$

Autoregressiv modell: Koyck-modellen

Utgångspunkt: en infinite distributed lag
(IDL) modell

$$y_t = \alpha + \beta_0 x_t + \beta_1 x_{t-1} + \dots + u_t \quad (1)$$

Koyck: Antar geometriska avtagande vikt
på parametrene:

$$\beta_j = \gamma \cdot p^j \text{ slikt att } \beta_0 = \gamma, \beta_1 = \gamma \cdot p, \dots \text{ osv.} \quad (2)$$

der $0 < p < 1$

Poeng: För nå bare to parametre å estimere:

$$y_t = \alpha + \gamma x_t + \gamma p x_{t-1} + \gamma p^2 x_{t-2} + \dots + u_t \quad (3)$$

Tilbakefører alle variabler i (3) med 1:

$$y_{t-1} = \alpha + \gamma x_{t-1} + \gamma p x_{t-2} + \gamma p^2 x_{t-3} + \dots + u_{t-1} \quad (4)$$

Multipliserer (4) med p :

$$p y_{t-1} = p \alpha + \gamma p x_{t-1} + \gamma p^2 x_{t-2} + \gamma p^3 x_{t-3} + \dots + p u_{t-1} \quad (5)$$

Treker nå (5) fra (3):

$$y_t - p y_{t-1} = (1-p)\alpha + \gamma x_t + (u_t - p u_{t-1})$$

$$\Rightarrow y_t = \alpha_0 + p y_{t-1} + \gamma x_t + v_t \quad (7)$$

$$\alpha_0 = (1-p)\alpha \text{ \& } v = u_t - p u_{t-1} \quad 39$$

(7) kalles en autoregressiv modell.

Denne kan i prinsippet estimeres via OLS, bortsett fra at v_t bryter med prinsippet om ingen seriekorrelasjon. Skal foreløpig se bort fra dette.

Kortids-effekt & langtids-effekt i en autoregressiv modell:

$$\text{Modellen: } y_t = \beta_0 + \beta_1 x_t + \rho y_{t-1} + u_t$$

$$\text{Kortids-effekt: } dy_t = \beta_1 dx_t \quad \Rightarrow \quad \frac{dy_t}{dx_t} = \beta_1 \quad (8)$$

$$\text{Langtids-effekt: } d\bar{y} = \frac{\beta_1}{1-\rho} d\bar{x} \quad \Rightarrow \quad \frac{d\bar{y}}{d\bar{x}} = \frac{\beta_1}{1-\rho} \quad (9)$$

$$\text{sidan } y_t = y_{t-1} = \bar{y} \text{ \& } x_t = \bar{x}$$

Merk: Når $0 < \rho < 1$ er prosessen stabil og langtids-effekten av en varig økning i

$$x \text{ er } \frac{\beta_1}{1-\rho}$$

Desto nærmere ρ er 1, desto tregere skjer tilpasningen.

Autoregressiv Distributed Lag (ADL) model:

Utvæder den autoregressive modellen med tilbage daterte verdier på forklaringsvariablene:

$$y_t = \alpha + \rho y_{t-1} + \beta_0 x_t + \beta_1 x_{t-1} + u_t \quad (10)$$

Kortidseffekten: β_0

Langtidseffekten: $\frac{\beta_0 + \beta_1}{1 - \rho}$

fordi:

$$y_t = y_{t-1} = \bar{y} \quad \& \quad x_t = x_{t-1} = \bar{x}$$

$$\Rightarrow \bar{y} = \alpha + \rho \bar{y} + (\beta_0 + \beta_1) \bar{x}$$

$$\Rightarrow \bar{y} = \frac{\alpha}{1 - \rho} + \frac{\beta_0 + \beta_1}{1 - \rho} \bar{x} = \mu + \beta \bar{x} \quad (11)$$

$$\Rightarrow \frac{d\bar{y}}{d\bar{x}} = \frac{\beta_0 + \beta_1}{1 - \rho} = \beta \quad (12)$$

Trend:

Trendproblem hvis følgende ligninger er signif.:

$$y_t = \beta_0 + \beta_1 t + u_t$$

$$x_t = \alpha_0 + \alpha_1 t + u_t$$

Får da spurious sammenheng i $y_t = \gamma_0 + \gamma_1 x_t + u_t$

Kan forstås som et utelatt variabelproblem.

Kan løses ved å inkludere tidstrender:

$$y_t = \gamma_0 + \gamma_1 x_t + \gamma_2 t + u_t \quad (3)$$

Alternative formuleringer av trend:

$$\ln y_t = \gamma_0 + \gamma_1 x_t + u_t$$

$$y_t = \gamma_0 + \gamma_1 x_t + \delta_0 D_t + \delta_1 D_{t-1} + \dots + u_t$$

Error correction model (likevekts- eller feiljusteringsmodell):

Antar følgende modell; (10):

$$y_t = \alpha + \rho y_{t-1} + \beta_0 x_t + \beta_1 x_{t-1} + u_t \quad (10)$$

$$\begin{aligned} \Rightarrow y_t - y_{t-1} &= (\rho - 1)y_{t-1} + \alpha + \beta_0 x_t + \beta_1 x_{t-1} + u_t \\ &= (\rho - 1)y_{t-1} + \alpha + \beta_0 x_t + \beta_1 x_{t-1} + \beta_0 (x_{t-1} - x_{t-1}) \\ &= (\rho - 1)y_{t-1} + \alpha + \beta_0 (x_t - x_{t-1}) + (\beta_0 + \beta_1)x_{t-1} + u_t \\ &= (\rho - 1) \left[y_{t-1} - \frac{\alpha}{1 - \rho} - \frac{\beta_0 + \beta_1}{1 - \rho} x_{t-1} \right] + \beta_0 (x_t - x_{t-1}) + u_t \\ \Rightarrow \Delta y_t &= -\lambda \left(y_{t-1} - \mu - \beta x_{t-1} \right) + \beta_0 \Delta x_t + u_t \quad (14) \end{aligned}$$

$$\begin{aligned} \text{der } \lambda &= 1 - \rho, & \Delta y_t &= y_t - y_{t-1}, \\ \mu &= \frac{\alpha}{1 - \rho}, & \Delta x_t &= x_t - x_{t-1}, \\ \beta &= \frac{\beta_0 + \beta_1}{1 - \rho}, \end{aligned}$$

i (14) kan $\mu + \beta x_{t-1}$ tolkes som likevektverdien på y på tidspunkt $t-1$ eller gitt x_{t-1} . Denne er også gitt av (11). Avviket i $(y_{t-1} - \mu - \beta x_{t-1})$ kan slik sett tolkes som avvik fra likevekten på tidspunkt $t-1$.

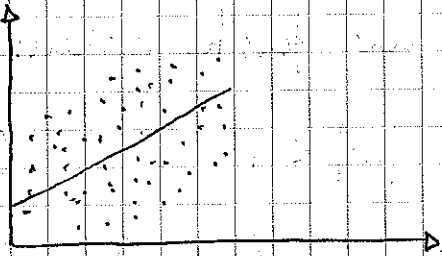
Dersom positivt avvik bidrar $0 < \lambda < 1$ til at $\Delta y_t < 0$.

Dersom negativt avvik bidrar $-1 < \lambda < 0$ til at $\Delta y_t > 0$.

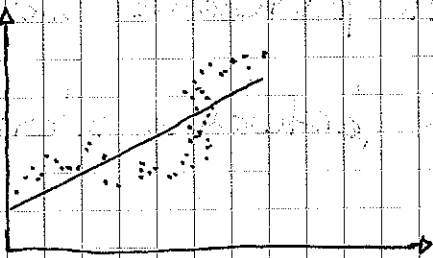
Generelt uttrykker $-\lambda (y_{t-1} - \mu - \beta_0 x_{t-1})$ en mekanisme i modellen som medfører at y utvikler seg i retning av den langsiktige likevektbanen over tid.

Seriekorrelasjon:

Fravær av seriekorrelasjon:



Tilstedeværelse av seriekorrelasjon:



Seriekorrelasjon dersom restleddene er avhengige av hverandre, dvs dersom $\text{Corr}(u_t, u_s) \neq 0$

- Dersom strikt eksogenitet ($E(u_t | X) = 0$) vil OLS gi forventningsrette estimatører.
- Dersom kontemporær eksogenitet ($E(u_t | x_t) = 0$) vil OLS gi konsistente, men ikke nødvendigvis forventningsrette estimatører.
- Seriekorrelasjon innebærer at de vanlige inferensmetodene ikke lenger er gyldige.

Hvordan variansen til OLS-estimatorene påvirkes av seriekorrelasjon:

Anta autoregressiv prosess av første orden:

$$u_t = \rho_1 u_{t-1} + e_t \quad \text{der } |\rho| < 1 \quad (1)$$

$$E(e_t | X) = 0, \quad \text{Var}(e_t) = \sigma_e^2, \quad \text{Cov}(e_t, e_s) = 0$$

Dersom $\rho > 0$ har vi positiv førsteordens s.korrelasjon.

Dersom $\rho < 0$ har vi negativ førsteordens s.korrelasjon.

Har at:

$$\begin{aligned} E(u_t u_{t-1}) &= E((\rho_1 u_{t-1} + e_t) u_{t-1}) = \rho_1 E(u_{t-1}^2) + E(e_t u_{t-1}) \\ &= \rho_1 \sigma_u^2 \end{aligned} \quad (2)$$

(1) innbærer at:

$$u_{t-1} = \rho_1 u_{t-2} + e_{t-1} \quad (3)$$

Setter (3) inn i (1):

$$\begin{aligned} u_t &= \rho_1 (\rho_1 u_{t-2} + e_{t-1}) + e_t \\ &= \rho_1^2 u_{t-2} + \rho_1 e_{t-1} + e_t \end{aligned}$$

Har da at:

$$\begin{aligned} E(u_t u_{t-2}) &= E[(\rho_1^2 u_{t-2} + \rho_1 e_{t-1} + e_t) u_{t-2}] \\ &= \rho_1^2 \sigma_u^2 + \rho_1 E(e_{t-1} u_{t-2}) + E(e_t u_{t-2}) = \rho_1^2 \sigma_u^2 \end{aligned}$$

Generelt:

$$E(u_t u_{t-s}) = \rho_1^s \sigma_u^2 \quad (4)$$

Korrelasjonskoeffisienten er da generelt gitt ved:

$$\begin{aligned} \text{Corr}(u_t, u_{t-s}) &= \frac{E(u_t u_{t-s})}{\sqrt{\text{Var}(u_t) \text{Var}(u_{t-s})}} \\ &= \frac{\rho_1^s \sigma_u^2}{\sigma_u^2} \\ &= \rho_1^s \end{aligned} \quad (5)$$

Hvis positiv seriekorrelasjon vil korrelasjonskoeffisienten avta i tid mellom restleddene.

Hvis negativ seriekorrelasjon vil også korrelasjonskoeffisienten avta i tid i absoluttverdi, men fortegnet vil skifte med verdien på s .

Har tidligere funnet at

$$\hat{\beta}_1 = \beta_1 + \frac{\sum_{t=1}^t (u_t - \bar{u})(x_t - \bar{x})}{\sum_{t=1}^t (x_t - \bar{x})^2}$$

Variansen er da:

$$\text{Var}(\hat{\beta}_1 | X) = \frac{\text{Var}\left(\sum_{t=1}^t (u_t - \bar{u})(x_t - \bar{x})\right)}{\left(\sum_{t=1}^t (x_t - \bar{x})^2\right)^2}$$

der $\text{Var}(u_t x_t | X)$ generelt blir ganske komplisert, og gjør vanlige inferensmetoder ugyldige.

Testing for seriekorrelasjon:

1. ordens seriekorrelasjon:

1. Estimer modellen

$$y_t = \beta_0 + \beta_1 x_{t1} + \dots + \beta_k x_{tk} + u_t$$

for å finne alle \hat{u}_t .

2. Estimer modellen

$$u_t = \rho_1 u_{t-1} + e_t \quad (\text{med eller uten konstantledd})$$

3. Test hypotesen $H_0: \rho_1 = 0$.

q. ordens seriekorrelasjon:

Tilsvarende, men estimer $u_t = \rho_1 u_{t-1} + \dots + \rho_q u_{t-q} + e_t$

i 2. Alternativt kan man også inkludere forkl.variablene.

$$H_0: \rho_1 = \rho_2 = \dots = \rho_q = 0$$

Durbin-Watson-testen for seriekorrelasjon:

$$\begin{aligned} DW &\equiv \frac{\sum_{t=1}^T (\hat{u}_t - \hat{u}_{t-1})^2}{\sum_{t=1}^T \hat{u}_t^2} = \frac{\sum_{t=1}^T (\hat{u}_t^2 - 2\hat{u}_t \hat{u}_{t-1} + \hat{u}_{t-1}^2)}{\sum_{t=1}^T \hat{u}_t^2} \\ &\approx \frac{\sum_{t=1}^T (2\hat{u}_t - 2\hat{u}_{t-1}) \hat{u}_t}{\sum_{t=1}^T \hat{u}_t^2} \\ &= 2 \cdot \frac{\sum_{t=1}^T \hat{u}_t - \sum_{t=1}^T \hat{u}_{t-1}}{\sum_{t=1}^T \hat{u}_t^2} \\ &= 2(1 - \hat{\rho}_1) \end{aligned} \quad (6)$$

Test hvorvidt DW er signifikant forskjellig fra 2.
Positiv s.korr. dersom $DW < 2$, negativ s.korr. dersom $DW > 2$.

Korrigerings for seriekorrelasjon av 1. orden:

Utgangspunkt i:

$$y_t = \beta_0 + \beta_1 x_t + u_t \quad (7)$$

$$u_t = \rho u_{t-1} + e_t \quad (8)$$

Ønsker nå å transformere (7) slik at seriekorrelasjonen forsvinner.

Dette kan gjøres ved først å tidsforsyve (7)

med 1 og deretter multiplisere med ρ :

$$\rho y_{t-1} = \rho \beta_0 + \rho \beta_1 x_{t-1} + \rho u_{t-1} \quad (9)$$

Trekker deretter (9) fra (8):

$$\begin{aligned} y_t - \rho y_{t-1} &= (1-\rho)\beta_0 + \beta_1 (x_t - \rho x_{t-1}) + (u_t - \rho u_{t-1}) \\ &= (1-\rho)\beta_0 + \beta_1 (x_t - \rho x_{t-1}) + (\rho u_{t-1} + e_t - \rho u_{t-1}) \\ &= (1-\rho)\beta_0 + \beta_1 (x_t - \rho x_{t-1}) + e_t \end{aligned} \quad (10)$$

$$\Rightarrow \tilde{y}_t = \gamma_0 + \beta_1 \tilde{x}_t + e_t \quad (11)$$

der $\tilde{y}_t = y_t - \rho y_{t-1}$, $\gamma_0 = (1-\rho)\beta_0$ & $\tilde{x}_t = x_t - \rho x_{t-1}$

Gitt antagelser om e_t er (11) ikke seriekorrelert, og kan derfor estimeres vha. OLS.

Fremgangsmåte for korrigering av seriekorrelasjon

i praksis:

1. Estimer (I) $y_t = \beta_0 + \beta_1 x_t + u_t$ & finn residualene \hat{u}_t .

2. Estimer ρ ved å bruke OLS på hjelperegresjonen

$$(8) \hat{u}_t = \rho \hat{u}_{t-1} + \text{støy}.$$

3. Bruk $\hat{\rho}$ til å finne de transformerte

$$\tilde{y}_t = y_t - \hat{\rho} y_{t-1} \quad \& \quad \tilde{x}_t = x_t - \hat{\rho} x_{t-1}$$

4. Estimer (II) $\tilde{y}_t = \gamma_0 + \beta_1 \tilde{x}_t + e_t^*$

der e_t^* symboliserer støy pga $\hat{\rho}$.

Spesialtilfelle: $\rho = 1$

$$\Rightarrow \tilde{y}_t = y_t - y_{t-1} = \Delta y_t$$

$$\Rightarrow \tilde{x}_t = x_t - x_{t-1} = \Delta x_t$$

$$\Rightarrow \gamma_0 = (1-1)\beta_0 = 0$$

Da reduseres (II) til:

$$\Delta y_t = \beta_1 \Delta x_t + e_t \quad (12)$$

Alternativ strategi:

Antar at sann modell er gitt ved:

$$y_t = \beta_0 + \alpha_1 y_{t-1} + \gamma_0 x_t + \gamma_1 x_{t-1} + u_t \quad (13)$$

Anta at vi underestimerer modellen slik:

$$y_t = \beta_0 + \gamma_0 x_t + v_t \quad (14)$$

$$\text{der } v_t = \alpha_1 y_{t-1} + \gamma_1 x_{t-1} + u_t \quad (15)$$

Her er (15) korrelert med

$$v_{t-1} = \alpha_1 y_{t-2} + \gamma_1 x_{t-2} + u_{t-1} \quad (16)$$

siden y_{t-1} ifl. (13) avhenger av y_{t-2} & x_{t-2} .

Dette kan åpenbart korrigeres ved å inkludere utelatte variable y_{t-1} & x_{t-1} i (14) slik at vi får (13).

Dermed blir strategien å starte med en svært generell dynamisk modell, og stegvis utelukke variablene med lavest t -verdi.

Denne prosessen kan fortsette så lenge man ikke får seriekorrelerte restledd.

Sammenkoblede tverrsnittsdata (pooled cross sections) & paneldata:

Sammenkoblede tverrsnittsdata:

Forskjellige tverrsnittstrøg fra forskjellige tidsperioder kobles sammen (pooled).

Paneldata:

Skiller seg fra sammenkoblede tverrsnittsdata ved at hvert enkelt tverrsnitt inneholder informasjon om de samme enhetene.

Chow-test for strukturelle endringer over tid:

Alternativ 1:

Anta at vi ønsker å teste hvorvidt

$H_0: \delta_0 = \delta_1 = \delta_2 = 0$ i modellen

$$y_{it} = \beta_0 + \delta_0 D_i + \beta_1 x_{i1} + \delta_1 D_i x_{i1} + \beta_2 x_{i2} + \delta_2 D_i x_{i2} + u_{it} \quad (1)$$

der $D_i = 1$ hvis periode 2, og $= 0$ hvis periode 1.

Test H_0 vha. F_{3, n_1+n_2-6}

Alternativ 2:

1. Estimer

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i \quad (2)$$

kun basert på enhetene fra periode 1. Finn SSR_1 .

2. Estimer (2) kun basert på enhetene fra periode 2. Finn SSR_2 .

3. Estimer (2) med alle observasjoner.

Finn SSR_p .

4. Beregn

$$F = \frac{\frac{SSR_p - SSR_1 - SSR_2}{3}}{\frac{SSR_1 + SSR_2}{n_1 + n_2 - 6}} \sim F_{3, n_1 + n_2 - 6}$$

Alternativ 1 & 2 av Chow-testen er ekvivalente og gir samme resultat.

Politikanalyse med sammenkoblede tværstidsdata:

Difference-in-difference-estimatoren

Estimat for behandling:

$$\hat{y}_{t+1} = \hat{\beta}_0 + \hat{\beta}_1 D \quad \text{der } D=1 \text{ for de som skal behandles på tidspunkt } t.$$

Estimat etter behandling:

$$y_t = \hat{\gamma}_0 + \hat{\gamma}_1 D$$

Forskjell i forskjell for d etter behandling:

$$\hat{\delta} = \hat{\beta}_1 - \hat{\gamma}_1$$

Komplett modellgenerelt:

$$y = \beta_0 + \gamma_0 D_2 + \beta_1 D_T + \gamma_1 D_2 \cdot D_T + \text{kontrollvariable} + u$$

der $D_2 = 1$ hvis periode 2 (etter behandling)

og $D_T = 1$ hvis behandlet gruppe

β_0 - ikke-behandlede for behandling

γ_0 - Endring for alle over tid

β_1 - Tidsuavhengig forskjell mellom behandlede & ikke-behandlede

γ_1 - Effekten av behandling

Illustrasjon av diff-in-diff-estimatoren:

	Før	Etter	Etter - Før
Kontroll	β_0	$\beta_0 + \gamma_0$	γ_0
Behandling	$\beta_0 + \beta_1$	$\beta_0 + \gamma_0 + \beta_1 + \gamma_1$	$\gamma_0 + \gamma_1$
Beh. - kontr.	β_1	$\beta_1 + \gamma_1$	γ_1

Paneldata:

- Balanserte paneldata: Har opplysninger om alle enheter på alle tidspunkter.
- Ubalanserte paneldata har i prinsipp de samme egenskapene som balanserte paneldata.
- Antall observasjoner i balansert datasett: $N \cdot T$
- Paneldatamodell:

$$y_{it} = \beta_0 + \beta_1 x_{it} + u_{it} \quad (t=1,2) \quad (1)$$

Et problem med (1) er at denne antageligvis vil feilestimeres pga. utelatt variabel problem. Kan skille mellom variabler som er konstante og variabler som varierer over tid:

$$y_{it} = \beta_0 + \delta_0 D_2 + \beta_1 x_{it} + \alpha_i + u_{it} \quad (2)$$

Ved estimering av (2) får vi:

$$\hat{y}_{it} = \hat{\beta}_0 + \hat{\delta}_0 D_2 + \hat{\beta}_1 x_{it} + \hat{v}_{it} \quad \text{der } \hat{v}_{it} = \hat{\alpha}_i + \hat{u}_{it} \quad (3)$$

Denne vil gi sløyerhet så lenge α_i korrelerer med x_{it} .

First Difference (FD):

Kan transformere bort individualspecifikke & tidsuafhængige faktorer α_i ved at differensiere (3):

$$\text{År 1: } y_{i1} = \beta_0 + \beta_1 x_{i1} + \alpha_i + u_{i1} \quad (4)$$

$$\text{År 2: } y_{i2} = (\beta_0 + \delta_0) + \beta_1 x_{i2} + \alpha_i + u_{i2} \quad (5)$$

Trækker (5) fra (4):

$$y_{i2} - y_{i1} = \delta_0 + \beta_1 (x_{i2} - x_{i1}) + (u_{i2} - u_{i1})$$

$$\Rightarrow \Delta y_{it} = \delta_0 + \beta_1 \Delta x_{it} + \Delta u_{it} \quad (6)$$

Poeng: Har transformeret bort α_i i (6).

Globalt gjennomsnitt:

$$\text{Modell: } y_{it} = \beta_0 + \beta_1 x_{it} + \beta_2 z_i + \beta_3 q_t + u_{it} \quad (7)$$

$$\bar{y} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T y_{it} \quad (8)$$

$$\bar{x} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T x_{it} \quad (9)$$

$$\bar{z} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T z_i = \frac{1}{NT} \sum_{i=1}^N T z_i = \frac{T}{NT} \sum_{i=1}^N z_i = \frac{1}{N} \sum_{i=1}^N z_i \quad (10)$$

$$\bar{q} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T q_t = \frac{N}{NT} \sum_{t=1}^T q_t = \frac{1}{T} \sum_{t=1}^T q_t \quad (11)$$

Forutsetninger i paneldatamodell:

Modellen:

$$y_{it} = \alpha + \beta_1 x_{it1} + \dots + \beta_k x_{itk} + \gamma_1 z_i + u_{it} \quad (12)$$

$$u_{it} = \eta_i + \varepsilon_{it} \quad (13)$$

1. $E(\varepsilon_{it} | x_{it}, z_i) = 0$

2. $E(\varepsilon_{it} \cdot \varepsilon_{js}) = \begin{cases} \sigma_\varepsilon^2 & \text{for } i=j \text{ \& } t=s \\ 0 & \text{ellers} \end{cases}$

3. $E(\eta_i \cdot \eta_j) = \begin{cases} \sigma_\eta^2 & \text{for } i=j \\ 0 & \text{ellers} \end{cases}$

4. $E(\varepsilon_{it} \eta_j) = 0$ for alle i, j & t .

Dersom også

5. $E(\eta_i | x_{itk}, z_i) = 0$

er u_{it} uavhengig av både x_{itk} & z_i . Da kan OLS brukes direkte på (12).

Fordeler ved å anvende OLS direkte på (12):

1. Utnytter all variasjon i data; både i tverrsnitt & tid.
2. Gir mulighet til å estimere effekten av tverrsnittsspesifikke variable som ikke varierer over tid.

Ulempe ved å anvende OLS direkte på (12):

Vi får skjevete estimasjoner dersom

$$E(n_i | x_{itk}, z_{it}) \neq 0$$

Dette vil være tilfellet dersom vi har utelulket individspesifikke faktorer som er korrelert med de inkluderte

(2) kan skrives:

$$y_{it} = \alpha + \beta_1 x_{it1} + \dots + \beta_k x_{itk} + \gamma z_{it} + n_i + \varepsilon_{it}$$

$$= \alpha_i + \beta_1 x_{it1} + \dots + \beta_k x_{itk} + \gamma z_{it} + \varepsilon_{it} \quad (14)$$

$$\text{der } \alpha_i = \alpha + n_i$$

Estimeringsprosedyrer som kan anvendes på

(14) er:

1. Estimering av OLS der vi inkluderer dummyvariable for $N-1$ av tverrsnittsenhetene for å identifisere alle η_i (ekvivalent med å identifisere alle α_i).
2. Transformere bort η_i via Fixed Effects.
3. Transformere bort η_i via First Difference.

Fixed Effects

1. Beregn individspesifikke gjennomsnitt:

$$\bar{y}_i = \frac{1}{T} \sum_{t=1}^T y_{it}$$

$$\bar{\varepsilon}_i = \frac{1}{T} \sum_{t=1}^T \varepsilon_{it}$$

$$\bar{x}_{ij} = \frac{1}{T} \sum_{t=1}^T x_{itj}$$

$$\bar{z}_i = \frac{1}{T} \sum_{t=1}^T z_{it} = \frac{1}{T} T z_i = z_i$$

Tar utgangspunkt i (14), summerer over alle t & dividerer med T :

$$\bar{y}_i = \alpha_i + \beta_1 \bar{x}_{i1} + \dots + \beta_k \bar{x}_{ik} + \gamma_0 z_i + \bar{\varepsilon}_i \quad (15)$$

2. Trekker (15) fra (14):

$$\begin{aligned} y_{it} - \bar{y}_i &= \alpha_i - \alpha_i + \beta_1 (x_{it1} - \bar{x}_{i1}) + \dots + \beta_k (x_{itk} - \bar{x}_{ik}) \\ &\quad + \gamma_0 (z_i - z_i) + (\varepsilon_{it} - \bar{\varepsilon}_i) \\ &= \beta_1 (x_{it1} - \bar{x}_{i1}) + \dots + \beta_k (x_{itk} - \bar{x}_{ik}) + (\varepsilon_{it} - \bar{\varepsilon}_i) \end{aligned}$$

$$\Rightarrow \hat{y}_{it} = \beta_1 \hat{x}_{it1} + \dots + \beta_k \hat{x}_{itk} + \hat{\varepsilon}_{it} \quad (16)$$

Estimering av (16) gir de såkalte fixed effects-estimatorenene (within estimator).

First Difference:

1. Tar igjen utgangspunktet i (14). Tilbake-daterer denne med én periode og får:

$$y_{it-1} = \alpha_i + \beta_1 x_{it-1} + \dots + \beta_k x_{it-k} + \gamma_1 z_i + \varepsilon_{it-1} \quad (17)$$

2. Trekker (17) fra (14):

$$y_{it} - y_{it-1} = \alpha_i - \alpha_i + \beta_1 (x_{it} - x_{it-1}) + \dots + \beta_k (x_{itk} - x_{it-k}) + \gamma_1 (z_i - z_i) + (\varepsilon_{it} - \varepsilon_{it-1})$$

$$= \beta_1 (x_{it} - x_{it-1}) + \dots + \beta_k (x_{itk} - x_{it-k}) + (\varepsilon_{it} - \varepsilon_{it-1})$$

$$\Rightarrow \Delta y_{it} = \beta_1 \Delta x_{it} + \dots + \beta_k \Delta x_{it} + \Delta \varepsilon_{it} \quad (18)$$

Estimering av (18) gir de såkalte first difference-estimatorene.

Valg mellom FE & FD:

Hvis stor T & liten N foretrekkes gjerne FD.

Ved seriekorrelasjon i FE foretrekkes gjerne FE.

Ofte lunte å rapportere begge.

Fordeleer & ulemper ved FE & FD:

- Fjerner individspesifikke komponenter & eliminerer dermed et endogenitets- (utelat variabel-) problem. Både FE & FD gir konsistente & forventningsrette estimatorer.
- Utnytter bare variasjon i tid, siden individspesifikke egenskaper er transformert bort.
- Andre begrensninger: Vi har ikke alltid paneldata, og vi kan være interessert i individspesifikke variable som ikke varierer over tid. FE & FD løser heller ikke problemer med utelatte variable som varierer i tid & er korrelert med inkluderte forklaringsvariable.

Instrumentvariabelmetoden & 2-SLS:

Den afhængige variabelen er endogen når den er korreleret med restleddet:

$$y_{it} = \beta_0 + \beta_1 x_{it} + u_{it} \quad \text{der} \quad \text{Cov}(u_{it}, x_{it}) \neq 0$$

Med udgangspunkt i denne skrivt:

$$\text{Cov}(y, z) = \beta_1 \text{Cov}(x, z) + \text{Cov}(u, z) \quad (1)$$

$$\Rightarrow \beta_1 = \frac{\text{Cov}(y, z)}{\text{Cov}(x, z)} - \frac{\text{Cov}(u, z)}{\text{Cov}(x, z)}$$
$$= \frac{\text{Cov}(y, z)}{\text{Cov}(x, z)} \quad (2)$$

gilt at:

$$\text{Cov}(u, z) = 0 \quad \text{dvs. eksogen} \quad (3)$$

$$\text{Cov}(x, z) \neq 0 \quad \text{dvs. relevant} \quad (4)$$

Ser at $\text{plim}_{n \rightarrow \infty}(\hat{\beta}_1) = \beta_1$

Bytter ut teoretiske momenter med empiriske:

$$\hat{\beta}_1^{IV} = \frac{\sum_{i=1}^n (y_i - \bar{y})(z_i - \bar{z})}{\sum_{i=1}^n (x_i - \bar{x})(z_i - \bar{z})} \quad (5)$$

Homoskedastisitetssantagelsen:

$$E(u^2 | z) = \sigma^2 = \text{Var}(u) \quad (6)$$

Kan vises at gift (3), (4) & (6) er den asymptotiske variansen til $\hat{\beta}_1^{iv}$:

$$\text{Var}(\hat{\beta}_1^{iv}) = \frac{\hat{\sigma}^2}{n \sigma_x^2 \rho_{x,z}} \quad (7)$$

Den empiriske motparten er:

$$\text{Var}(\hat{\beta}_1^{iv}) = \frac{\hat{\sigma}^2}{\text{SST}_x \cdot R_{x,y}^2} \quad (8)$$

og

$$\text{SE}(\hat{\beta}_1^{iv}) = \frac{\hat{\sigma}}{\sqrt{\text{SST}_x \cdot R_{x,z}^2}} \quad (9)$$

der $\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^2$

$$\text{SST}_x = \sum_{i=1}^n (x_i - \bar{x})^2$$

$R_{x,z}^2$ er R^2 fra $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 z_i + \text{støy}$

Sammenlignet med $\text{Var}(\hat{\beta}_1^{OLS})$:

$$\text{Var}(\hat{\beta}_1^{iv}) = \frac{\hat{\sigma}^2}{\text{SST}_x \cdot R_{x,z}^2} > \frac{\hat{\sigma}^2}{\text{SST}_x} = \text{Var}(\hat{\beta}_1^{OLS})$$

\Rightarrow iv-estimering går alltid ut over variansen til estimatoren

Verdt å merke seg at (4) innebærer at:

$$x = \beta_0 + \beta_1 z + u \quad \text{må gi signifikant } \beta_1 \neq 0$$

siden $\beta_1 = \frac{\text{Cov}(x, z)}{\text{Var}(z)}$

Merk: R^2 gir ikke mye mening i instrumenterte regresjoner; kan ta negative verdier og > 1 . Gjelder også SSR, slik at F-tester ikke lenger kan brukes. Tilpassede F-tester kan imidlertid settes opp.

IV med multipl regression:

Anta:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + u_i \quad (10)$$

Kaller (10) for strukturlikningen der x_i mistenkes å være korrelert med restleddet u_i .

Normallikningene i IV-metoden er da gitt ved:

$$\begin{aligned} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 z_{i1}) &= 0 \\ \sum_{i=1}^n z_{i1} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 z_{i1}) &= 0 \\ \sum_{i=1}^n z_{i2} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 z_{i1}) &= 0 \end{aligned} \quad (11)$$

der z_{i2} er instrumentvariabelen. Likningssettet

(11) gir 3 likninger til å estimere de tre

ukjente $\hat{\beta}_0$, $\hat{\beta}_1$ & $\hat{\beta}_2$.

For å teste (4) kan man sette opp følgende redusert form-ligning:

$$X_{i1} = \beta_0 + \beta_1 Z_{i1} + \beta_2 Z_{i2} + u_{i2} \quad (12)$$

og teste $H_0: \beta_2 = 0$.

Redusert-form-ligningen skal inneholde alle de eksogene variablene Z .

2 Stage Least Squares (2-SLS):

Ved brude av mer enn én instrumentvariabel for x kan vi benytte oss av 2-SLS:

Tar igjen utgangspunkt i (10), men inkluder oss nå to instrumenter z_2 & z_3 .

Den beste instrumentvariabelen er da estimatet \hat{x}_1 i redusert form-regresjonen:

$$x_1 = \beta_0 + \beta_1 z_1 + \beta_2 z_2 + \beta_3 z_3 \quad (13)$$

Kravet (4) kan da testes der $H_0: \beta_2 = \beta_3 = 0$ må forkastes for at instrumentene skal være relevante.

Så lenge IV-variabelen er gyldig kan \hat{x}_1 erstatte z_1 i (11).

Dette er identisk med å beregne

$$y_1 = \beta_0 + \beta_1 \hat{x}_1 + z_1 + u \quad (14)$$

via OLS. Derfor kalles ofte 2-SLS-estimatoren for $\hat{\beta}_1^{2SLS}$.

2-SLS kan også brukes ved instrumentering av flere endogene forklaringsvariable. Metoden er ekvivalent med den beskrevet frem til nå, men man trenger minst én eksogen variabel (instrument) som ikke er inkludert i strukturlikningen per endogene forklaringsvariable.

Testing for endogenitet:

Motivasjon: OLS-estimatorene er mer effisiente når de er eksogene, derfor ikke et mål i seg selv med instrumentvariable.

Fremgangsmåte:

1. Estimer redusert form-regresjonen.

Finn residualene \hat{u}_2 .

2. Estimer strukturlikningen inkludert

leddet $\rho \hat{u}_2$. Test $H_0: \rho = 0$. Hvis H_0 forkastes kan vi konkludere at x_1 er endogen (Corr(u_1, u_2) $\neq 0$)

Testing for overidentifikasjon:

Når man bare har ett instrument for den endogene forklaringsvariabelen er denne eksakt identifisert. Da kan man bare resonnerer rundt gyldigheten av (3):

$$\text{Cov}(z, u) = 0$$

Når man har flere instrumenter kan man teste om man har flere enn nødvendig via overidentifikasjonstester.

Gitt homoskedastisitet kan man gjøre følgende:

1. Estimer strukturtilikningen via 2SLS for å finne 2SLS-residualene \hat{u}_i .
2. Estimer \hat{u}_i som funksjon av alle eksogene variable z . Finn R^2 kalt R_1^2 .
3. H_0 : Alle instrumenter er ukorrelert med u_i .

Test H_0 via at $n \cdot R_1^2 \sim \chi^2_g$ der g er

antall instrumenter minus antall endogene forklaringsvariable. Hvis H_0 forkastes er idet minste noen instrumenter eksogene.