

Methoden der Ökonometrie

Universität Duisburg-Essen

Christoph Hanck

Winter 2022

Überblick

- 1 Regressionsmodelle
- 2 Lineare Regression
- 3 Statistische Eigenschaften von Least Squares
- 4 Inferenz
- 5 Nichtlineare Regression
- 6 Generalized Least Squares
- 7 Instrumentvariablen

Was ist Ökonometrie und wofür brauchen wir sie?

Viele ökonomische Theorien spezifizieren Beziehungen zwischen Variablen. Interessant ist dabei oft eine quantitative Aussage. [Newton wäre auch nicht berühmt geworden für die Erkenntnis, *dass* ein Apfel *runterfällt*, wenn man ihn loslässt.]

Einige Beispiele:

- Um wie viel steigern Investitionen in das Humankapital (z.B. der Besuch der Universität) das Einkommen?
- Wie stark sinkt die Inflation, wenn die Zentralbank den Zinssatz anhebt?
- Um wie viel verringert sich die Kriminalität, wenn eine zusätzliche Million Euro für Polizisten auf den Straßen ausgegeben werden? `crime.R`
- Um wie viel erhöhen sich die Verkäufe, wenn Werbeausgaben erhöht werden?
- Wie viele zusätzliche Passagiere bringt eine Reduzierung der Zugticketpreise?
- ...

Allerdings sollen solche Aussagen typischerweise **ceteris paribus** sein, d.h., alle **anderen** wichtigen **Faktoren** (z.B. das Fähigkeitsniveau) **konstant** halten.

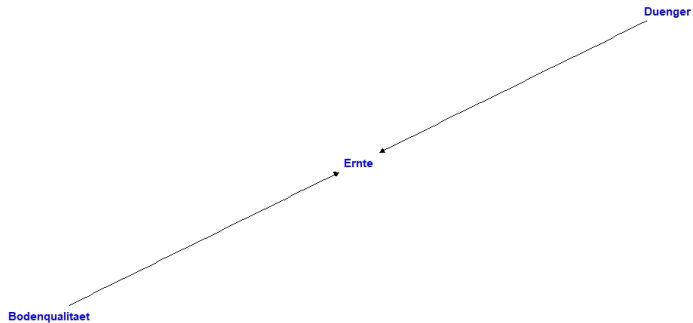
Anders formuliert handelt es sich um **kausale** Fragen, nicht nur Korrelationen. Erstere sind typischerweise interessanter. Siehe auch **hier**.

Ökonometrie ist die Kunst und Wissenschaft der Quantifizierung (Schätzung) solcher Beziehungen. Z.B., um *wie viel* erhöht sich das Einkommen *aufgrund* eines zusätzlichen Bildungsjahres?

Solche Schätzung ist recht einfach in **kontrollierten Experimenten**. Wir können z.B. die Einsatzmenge von Düngemitteln auf Äcker **zufällig zuteilen**, um den Effekt des Düngers auf den Ernteertrag zu messen.

In diesem Fall müssen andere wichtige Faktoren, wie Qualität des Bodens, *nicht* konstant gehalten werden, weil die Menge an Dünger den einzigen systematischen Unterschied zwischen den Äckern darstellt.

DAG Dünger



In den meisten interessanteren Szenarien (wie die o.g. Bildungsrenditen) sind solche Experimente aus unterschiedlichen Gründen nicht möglich.

Stattdessen müssen wir typischerweise mit **Beobachtungsdaten** arbeiten, die bereits bestehen werden und **nicht in einem Experiment gesammelt wurden**.

Dann müssen wir statistisch alle anderen wichtigen Faktoren konstant halten.

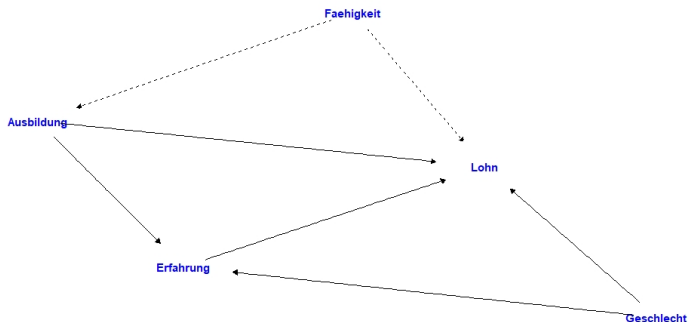
Beispielsweise werden fähigere Personen mehr verdienen, sich aber auch *entscheiden* mehr zu lernen. Höhergebildete haben zudem im Allgemeinen weniger Arbeitsmarkterfahrung, welche ebenfalls das Einkommen beeinflusst.

Die Arbeitsmarkterfahrung ist recht leicht zu beobachten, wohingegen „Fähigkeit“ schwer zu erfassen ist.

Wir werden sehen, dass das Berücksichtigen von Faktoren wie dem Ersteren leicht ist. Komplizierter wird es, wenn ein wichtiger Faktor nicht direkt beobachtet werden kann.

DAG Löhne

Die *angenommene* Problemstellung als DAG:



DAGs.R



Warum?

Siehe neben anderen:

<http://blog.revolutionanalytics.com/2014/05/companies-using-r-in-2014.html>

<http://r4stats.com/articles/popularity>

<http://www.r-bloggers.com/r-passes-sas-in-scholarly-use-finally>

<https://www.r-bloggers.com/oreilly-data-scientist-salary-and-tools-survey-november-2014/>

Das einfache Regressionsmodell

Der Einfachheit halber beginnen wir mit dem einfachen linearen Regressionsmodell.

Definition 1.1 (Einfaches **lineares** Regressionsmodell).

$$y_t = \beta_1 + \beta_2 X_t + u_t, \quad (1.1)$$

wobei

- y_t die **zu erklärende** oder **abhängige Variable**,
- X_t die **erklärende** oder **unabhängige Variable**,
- u_t der **Fehlerterm**, der alle unbeobachteten Einflüsse enthält,
- β_2 der **Steigungsparameter**, und β_1 der **Achsenabschnitt** ist.
- t ist der Index der Beobachtungen und läuft von 1 bis n .

y_t und X_t sind beobachtbar, der Rest nicht.

Das Modell ist linear, weil die Veränderung in y_t immer β_2 ist, wenn X_t sich um eine Einheit ändert, ungeachtet des Levels von X_t .

(1.1) ist interessant, wenn wir Annahmen bzgl. u_t treffen.

- u_t ist zufällig und unbeobachtbar—also modellieren wir unsere Unwissenheit über u_t .
- Besteht ein Zusammenhang zwischen u_t und X_t ?

Zentrales Ziel dieses Kurses:

Wir wollen mehr erfahren über β_2 (und manchmal über β_1).

Wiederholung von etwas Wahrscheinlichkeitstheorie

Zufallsvariablen formalisieren Zufälligkeit.

Einfachster Fall: Skalare Zufallsvariable (ZV) X

- Menge der möglichen Werte $x_i \in \mathbb{N}, \mathbb{R}, \mathbb{R}_+, \dots$
- Diskrete vs. stetige
- Wahrscheinlichkeiten
- Wahrscheinlichkeitsverteilung: weist x_i Wahrscheinlichkeit $p(x_i)$ zu
- Diskreter Fall:
 - ▶ ... $p(x_i)$ summiert sich zu eins, $\sum_{i=1}^{\infty} p(x_i) = 1$
 - ▶ ...Beispiel: binäre Daten, wie Zahlungsverzug/kein Zahlungsverzug

- Stetiger Fall:

- ▶ ...jedes mögliche Ereignis x hat die Wahrscheinlichkeit null,
- ▶ ...hat dennoch eine **Dichte** $f(x)$
- ▶ ...sodass $\int_{-\infty}^{\infty} f(x) dx = 1$
- ▶ ...Beispiel: Dinge, die in Zeit oder Euro gemessen werden
- ▶ ...können ebenfalls durch die **(kumulative) Verteilungsfunktion** $F(x)$ charakterisiert werden,
- ▶ ...definiert als

$$F(x) = P(X \leq x), \quad (1.2)$$

die Wahrscheinlichkeit, dass X gleich oder kleiner als ein Wert x ist.

- ▶ ...hat welche Verbindung zu $f(x)$?

Regeln für Wahrscheinlichkeitsverteilungen:

- Wahrscheinlichkeiten liegen zwischen 0 und 1;
- Die leere Menge hat die Wahrscheinlichkeit 0, die Grundgesamtheit hat die Wahrscheinlichkeit 1;
- Die Vereinigung zweier disjunkter Ereignisse hat als Wahrscheinlichkeit die Summe der Wahrscheinlichkeiten dieser disjunkten Ereignisse.

Daraus folgen folgende Eigenschaften...

- $F(x)$ ist monoton steigend
- $F(x)$ konvergiert gegen 0 wenn $x \rightarrow -\infty$ und $F(x) \rightarrow 1$ wenn $x \rightarrow \infty$
- $P(a \leq X \leq b) = F(b) - F(a)$
- Was erhält man für $a = b$?
- Dies motiviert die Dichte $f(x) = F'(x)$.
- $P(a \leq X \leq b) = F(b) - F(a) = \int_a^b f(x) dx$
- Diskrete ZV haben ebenfalls Verteilungsfunktionen.

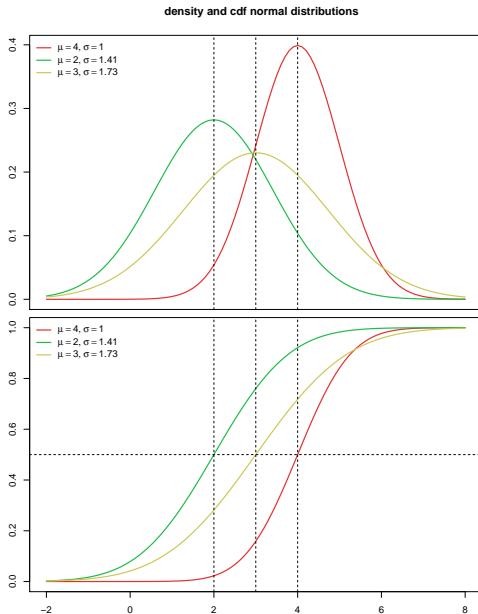
Beispiel 1.2 (Normalverteilung).

$$f(x) = \phi(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\} \quad \text{für } -\infty < x < \infty \quad (1.3)$$

Die Verteilungsfunktion (1.2) der Normalverteilung, Notation $\Phi(x)$, hat keine geschlossene Form. Alles, was wir sagen können, ist

$$\Phi(x) = \int_{-\infty}^x \phi(y) dy \quad (1.4)$$

Einige Beispiele (cf. pdfcdf.R):

 $\phi(x)$ $\Phi(x)$  x

Momente

- Am bekanntesten: Der Erwartungswert
- Für eine diskrete ZV, die m mögliche Werte x_1, x_2, \dots, x_m annehmen kann, ist der Erwartungswert

$$E(X) = \sum_{i=1}^m p(x_i)x_i,$$

die Summe jedes möglichen Wertes x_i multipliziert mit seiner Wahrscheinlichkeit.

- Für eine stetige ZV ist der Erwartungswert analog definiert als:

$$E(X) = \int_{-\infty}^{\infty} xf(x) dx.$$

Dennoch hat nicht jede ZV einen Erwartungswert!

- $E(X)$ (oft bezeichnet mit μ) bezieht sich auf die Grundgesamtheit und ist nicht zu verwechseln mit dem Stichprobenmittelwert.

- μ wird auch das erste Moment genannt...
- ...weil auch höhere Momente existieren (können). Im stetigen Fall:

$$m_k(X) = E(X^k) = \int_{-\infty}^{\infty} x^k f(x) dx$$

- Solche Momente werden unzentrierte Momente genannt...
- ...da auch **zentrale Momente** existieren (können):

$$\mu_k(X) = E(X - E(X))^k = \int_{-\infty}^{\infty} (x - \mu)^k f(x) dx$$

- Was ist μ_2 ?
- ...die Quadratwurzel heißt **Standardabweichung**
- ...und wird zur Konstruktion von **Standardfehlern** wichtig sein

Multivariate Verteilungen

- Mehrere ZV (X_1, X_2) haben **gemeinsame** Verteilungen.
- Formal

$$F(x_1, x_2) = P((X_1 \leq x_1) \cap (X_2 \leq x_2))$$

oder einfach

$$F(x_1, x_2) = P(X_1 \leq x_1, X_2 \leq x_2)$$

- Im stetigen Fall ergibt sich eine **gemeinsame Dichte**

$$f(x_1, x_2) = \frac{\partial F(x_1, x_2)}{\partial x_1 \partial x_2} \quad (1.5)$$

S. `BivariateNormalPlotly.R`

- Dann gilt

$$F(x_1, x_2) = \int_{-\infty}^{x_2} \int_{-\infty}^{x_1} f(y_1, y_2) dy_1 dy_2$$

Unabhängigkeit

Wir bezeichnen zwei ZV als **statistisch unabhängig**, wenn wir die gemeinsame Verteilungsfunktion schreiben können als:

$$F(x_1, x_2) = F(x_1, \infty)F(\infty, x_2) \equiv F_1(x_1)F_2(x_2)$$

Kenntnis von X_1 liefert dann keine Information über X_2 .

Unabhängigkeit definiert bezüglich der **Randdichten**

$$f_i(x_i) = \int_{-\infty}^{\infty} f(x_i, x_j) dx_j, \quad i \neq j \quad (1.6)$$

folgt analog als

$$f(x_1, x_2) = f_1(x_1)f_2(x_2)$$

Bedingte Wahrscheinlichkeiten

Für zwei beliebige Ereignisse A und B gilt,

$$P(A \cap B) = P(B)P(A|B) \quad \text{oder} \quad P(A|B) = \frac{P(A \cap B)}{P(B)},$$

gegeben $P(B) \neq 0$. Z.B.

- $A \cap B = \emptyset$ impliziert, dass $P(A|B) = 0$, genauso wie
- $B \subset A$ impliziert, dass $P(A \cap B) = P(B)$ und $P(A|B) = 1$.

Ein Venn-Diagramm hilft diese Notationen zu verstehen.

Beispiel: Was könnte die Wahrscheinlichkeit sein, dass eine Person über 2.000 Euro verdient, und wie hoch ist sie, wenn die Person älter als 30 ist?

Wir können ebenfalls mit **bedingten Dichten** arbeiten, welche auch funktionieren, wenn ein Ereignis ($X_2 = x_2$) die Wahrscheinlichkeit 0 hat (sofern $f_2(x_2) > 0$):

$$f(x_1|x_2) = \frac{f(x_1, x_2)}{f_2(x_2)}$$

Dies hilft uns den **bedingten Erwartungswert** von X_1 gegeben X_2 zu definieren

$$E(X_1|X_2) = \int_{-\infty}^{\infty} x_1 f(x_1|x_2) dx_1$$

Beachten Sie: Dies ist eine Funktion der ZV X_2 .

Daher ist der bedingte Erwartungswert selbst eine ZV und wir können seinen Erwartungswert betrachten.

$$\begin{aligned}\underbrace{\int_{-\infty}^{\infty} E(X_1|X_2)f_2(x_2) dx_2}_{=E[E(X_1|X_2)]} &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 f(x_1|x_2) dx_1 f_2(x_2) dx_2 \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 f(x_1|x_2)f_2(x_2) dx_1 dx_2 \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 f(x_1, x_2) dx_1 dx_2 \\ &= \int_{-\infty}^{\infty} x_1 \int_{-\infty}^{\infty} f(x_1, x_2) dx_2 dx_1 \\ &= \int_{-\infty}^{\infty} x_1 f_1(x_1) dx_1 \\ &= E(X_1)\end{aligned}\tag{1.7}$$

Dieses Resultat ist bekannt als das **Gesetz der iterierten Erwartungen**. Es ist sowohl intuitiv als auch wichtig.

Eine wichtige Regel ist, dass für jede Funktion h gilt,

$$E(X_1 h(X_2) | X_2) = h(X_2) E(X_1 | X_2). \quad (1.8)$$

Auch das macht intuitiv Sinn. Ein wichtiger Spezialfall ergibt sich, wenn $E(X_1 | X_2) = 0$, was impliziert, dass

$$\begin{aligned} E(X_1 h(X_2)) &= E(E(X_1 h(X_2) | X_2)) \\ &= E(h(X_2) E(X_1 | X_2)) \\ &= 0, \end{aligned} \quad (1.9)$$

bekannt als „(Erwartungswert-)Unabhängigkeit impliziert Unkorreliertheit.“

Zurück zur Ökonometrie, insbesondere (1.1). Nehme an, dass

$$E(u_t|X_t) = 0. \quad (1.10)$$

Dann ist, mit (1.1) und Nutzung von (1.8) und (1.10)

$$E(y_t|X_t) = \beta_1 + \beta_2 X_t + E(u_t|X_t) = \beta_1 + \beta_2 X_t$$

$E(u_t|X_t) = 0$ ist also letzten Endes die Annahme, dass unser Modell **korrekt spezifiziert** ist.

Betrachten wir

Beispiel 1.3 (Lineare Bildungsrenditen).

$$Lohn_t = \beta_1 + \beta_2 Ausbildung_t + u_t$$

Im Beispiel 1.3 enthält u_t angeborene Fähigkeiten. Für $X_t = \text{Ausbildung}_t$,

$$E(u_t|X_t) = 0$$

müssen die angeborenen Fähigkeiten immer gleich sein, unabhängig von den Bildungsjahren, z.B.

$$E(\text{Fähigkeit}_t|X_t = 8) = E(\text{Fähigkeit}_t|X_t = 16)$$

Wenn talentiertere Personen durchschnittlich mehr Bildungsjahre haben, ist die Annahme $E(u_t|X_t) = 0$ verletzt.

Wenn wir bei der Schätzung des Modells fälschlicherweise annehmen, dass $E(u_t|X_t) = 0$ ist, werden wir unser Ziel mehr über β_2 zu lernen nicht erreichen.

Für gewöhnlich wollen wir y_t **auf allgemeinere Informationssets** Ω_t bedingen, anstatt nur auf ein X_t . Z.B. $\Omega_t = \{X_t, X_t^2\}$ oder $\Omega_t = \{X_{t1}, X_{t2}\}$.

Wir wollen **exogene** Variablen in Ω_t einbinden, die außerhalb der Gleichung bestimmt werden. **Endogene** Variablen sind ebenfalls durch y_t bestimmt und dadurch auch durch u_t .

Neben der Spezifikation der Beziehung zu X_t müssen wir auch die Beziehung zwischen den u_t betrachten. Häufig treffen wir die Annahme, dass $\{y_t, X_t\}$ **unabhängig und identisch verteilt (u.i.v.)** sind.

Diese Annahme wäre bspw. verletzt, wenn **serielle Korrelation** vorläge, wenn z.B. auf große (kleine) Werte von u_t tendenziell große (kleine) Werte folgen. Mehr dazu in „Zeitreihenanalyse“.

Es liegt **Heteroskedastie** vor, wenn die bedingte Varianz verschiedener y_t unterschiedlich ist. $EVP.Rproj$

Lineare und nichtlineare Regressionsmodelle

Das lineare Modell erscheint sehr restriktiv, aber betrachten wir

$$y_t = \beta_1 + \beta_2 X_t + \beta_3 X_t^2 + u_t \quad (1.11)$$

$$y_t = \gamma_1 + \gamma_2 \log X_t + u_t$$

$$y_t = \delta_1 + \delta_2 \frac{1}{X_t} + u_t$$

Alle diese Modelle sind linear **in den Parametern**. Ersteres ist ein Beispiel für ein **multiple** Regressionsmodell.

Das Folgende ist nicht linear

$$y_t = e^{\beta_1} X_{t2}^{\beta_2} X_{t3}^{\beta_3} + u_t$$

Aber selbst

$$y_t = e^{\beta_1} X_{t2}^{\beta_2} X_{t3}^{\beta_3} e^{v_t}$$

ist linear (oder präziser **log-linear**) für unserer Zwecke: Logarithmieren liefert

$$\log y_t = \beta_1 + \beta_2 \log X_{t2} + \beta_3 \log X_{t3} + v_t$$

Wir werden uns bald mit multiplen Regressionsmodellen beschäftigen. Hierfür benötigen wir Matrixalgebra. Siehe hierzu den Matrixalgebra-Reader.

Das lineare Modell in Matrixnotation

$$\underset{(n \times 1)}{\mathbf{y}} = \underset{(n \times 2)}{\mathbf{X}} \underset{(2 \times 1)}{\boldsymbol{\beta}} + \underset{(n \times 1)}{\mathbf{u}} \quad (1.12)$$

wobei

$$\mathbf{X}_t = \begin{bmatrix} 1 & X_t \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix},$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix},$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_n \end{bmatrix} = \begin{bmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}.$$

Die Spalten von \mathbf{X} werden **Regressoren** genannt. Beachten Sie, dass Matrixalgebra nur eine kompaktere Notation darstellt — eine typische Zeile von (1.12) ist das Gleiche wie (1.1).

Das Tolle an Matrixalgebra ist, dass wir uns keine Gedanken machen müssen, wenn wir mehr als eine X_t -Variable haben, wie in (1.11). Wir schreiben einfach

$$\underset{(n \times 1)}{\mathbf{y}} = \underset{(n \times k)}{\mathbf{X}} \underset{(k \times 1)}{\boldsymbol{\beta}} + \underset{(n \times 1)}{\mathbf{u}} \quad (1.13)$$

wobei

$$\mathbf{X}_t = \begin{bmatrix} X_{t1} & \cdots & X_{tk} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_k \end{bmatrix},$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix},$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_n \end{bmatrix} = \begin{bmatrix} X_{11} & \cdots & X_{1k} \\ \vdots & \ddots & \vdots \\ X_{n1} & \cdots & X_{nk} \end{bmatrix}.$$

Normalerweise, aber nicht zwingend, ist $X_{t1} = 1$.

Eine typische Zeile ist dann

$$y_t = \mathbf{X}_t \boldsymbol{\beta} + u_t = \beta_1 + \sum_{i=2}^k \beta_i X_{ti} + u_t$$

Schätzung

Bisher haben wir uns nur mit der Frage des Modellaufbaus beschäftigt. In der Praxis kennen wir β nicht und müssen es aus unseren Daten \mathbf{y} , \mathbf{X} **schätzen**.

Stellen wir uns vor, dass die Daten aus einer **Population** gezogen, oder durch einen **Datengenerierenden Prozess (DGP)** erzeugt wurden, wie in (1.13).

Wir stellen nun selbst ein Modell auf und versuchen unser Modell so gut wie möglich an den DGP anzugleichen. Danach schätzen wir die Parameter des Modells.

Eine einfache aber sehr nützliche Schätzmethode ist die **Momentenmethode**. Die Idee ist, die Momentenbedingungen der Population ($\mu = g(\theta)$) durch die Momentenbedingungen der Stichprobe ($\bar{X} = g(\theta)$) zu ersetzen und dann nach den Parametern $\hat{\theta}$ aufzulösen, die diese Gleichungen lösen.

Beispiel 1.4 (Einfache lineare Regression). (1.1) impliziert, dass

$$u_t = y_t - \beta_1 - \beta_2 X_t$$

Die Annahme (1.10) ($E(u_t|X_t) = 0$) liefert mit dem Gesetz des iterierten Erwartungswertes (1.7), dass

$$E(u_t) = E(E(u_t|X_t)) = 0$$

Daraus erhalten wir die **Momentenbedingung**

$$E(u_t) = E(y_t - \beta_1 - \beta_2 X_t) = 0$$

Ersetzen wir den Erwartungswert durch den Stichprobenmittelwert, erhalten wir

$$\frac{1}{n} \sum_{t=1}^n u_t = \frac{1}{n} \sum_{t=1}^n (y_t - \tilde{\beta}_1 - \tilde{\beta}_2 X_t) = 0 \quad (1.14)$$

Wir nehmen für den Moment der Einfachheit halber an, dass $\beta_2 = 0$. Dann

$$\frac{1}{n} \sum_{t=1}^n (y_t - \tilde{\beta}_1) = 0$$

Der Wert, der diese Gleichung löst, folgt aus

$$\frac{1}{n} \sum_{t=1}^n y_t = \frac{n}{n} \tilde{\beta}_1$$

so dass

$$\hat{\beta}_1 = \frac{1}{n} \sum_{t=1}^n y_t.$$

Was, wenn $\beta_2 \neq 0$? Dann haben wir zwei Unbekannte und benötigen daher eine weitere Momentenbedingung.

Wegen (1.9) impliziert $E(u_t|X_t) = 0$ auch, dass

$$E(X_t u_t) = 0,$$

da

$$E(X_t u_t) \stackrel{(1.7)}{=} E[E(X_t u_t|X_t)] \stackrel{(1.8)}{=} E[X_t E(u_t|X_t)] \stackrel{(1.10)}{=} E[X_t 0] = E[0] = 0$$

Daher ist

$$E(X_t(y_t - \beta_1 - \beta_2 X_t)) = 0$$

Ersetzen wir die Populationsmomente durch die Stichprobenmomente, erhalten wir

$$\frac{1}{n} \sum_{t=1}^n (X_t(y_t - \tilde{\beta}_1 - \tilde{\beta}_2 X_t)) = 0. \quad (1.15)$$

Durch Umschreiben von (1.14) und (1.15) ergibt sich (nachrechnen!):

$$\begin{aligned}\tilde{\beta}_1 + \left(\frac{1}{n} \sum_{t=1}^n X_t \right) \tilde{\beta}_2 &= \frac{1}{n} \sum_{t=1}^n y_t \\ \left(\frac{1}{n} \sum_{t=1}^n X_t \right) \tilde{\beta}_1 + \left(\frac{1}{n} \sum_{t=1}^n X_t^2 \right) \tilde{\beta}_2 &= \frac{1}{n} \sum_{t=1}^n X_t y_t\end{aligned}$$

Multiplizieren mit n liefert

$$\begin{aligned}n\tilde{\beta}_1 + \left(\sum_{t=1}^n X_t \right) \tilde{\beta}_2 &= \sum_{t=1}^n y_t \\ \left(\sum_{t=1}^n X_t \right) \tilde{\beta}_1 + \left(\sum_{t=1}^n X_t^2 \right) \tilde{\beta}_2 &= \sum_{t=1}^n X_t y_t\end{aligned}$$

In Matrixschreibweise,

$$\begin{bmatrix} n & \sum_{t=1}^n X_t \\ \sum_{t=1}^n X_t & \sum_{t=1}^n X_t^2 \end{bmatrix} \begin{bmatrix} \tilde{\beta}_1 \\ \tilde{\beta}_2 \end{bmatrix} = \begin{bmatrix} \sum_{t=1}^n y_t \\ \sum_{t=1}^n X_t y_t \end{bmatrix}$$

oder noch kompakter

$$\mathbf{X}^\top \mathbf{X} \tilde{\boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{y} \quad (1.16)$$

Zum Auflösen von (1.16) nehmen wir die Inverse, um den **kleinsten Quadrate-Schätzer** zu erhalten.

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (1.17)$$

Formal ändert sich nicht viel in der multiplen Regression (1.13). Die Annahme, dass $E(u_t|X_t) = 0$ (siehe (1.10)) wird modifiziert zu

$$E(u_t|\mathbf{X}_t) = E(u_t|X_{t1}, \dots, X_{tk}) = 0. \quad (1.18)$$

Das führt wieder über (1.9) zu $E(X_{ti}u_t) = 0$, $i = 1, \dots, k$, und damit zu den k Momentenbedingungen

$$\frac{1}{n} \sum_{t=1}^n (X_{ti}(y_t - \mathbf{X}_t\tilde{\beta})) = 0 \quad i = 1, \dots, k \quad (1.19)$$

Zentral ist hierbei, dass (1.18) in Anwendungen wesentlich eher erfüllt sein wird als (1.10)!

(1.19) lautet in Matrixform (nachdem mit n multipliziert wurde)

$$\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) = \mathbf{0},$$

wobei $\mathbf{0}$ ein Vektor mit k Nullen ist. Nun kann dieser Ausdruck umgeschrieben werden als (1.16), $\mathbf{X}^\top \mathbf{X}\tilde{\boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{y}$. Daher ist die Lösung ebenfalls (1.17).

Warum nennen wir

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

dann eigentlich „OLS“-Schätzer?

Die Methode der kleinsten Quadrate

Für eine Stichprobe (\mathbf{y}, \mathbf{X}) sei $\tilde{\beta}$ ein hypothetischer Schätzer von β . $\tilde{\beta}$ definiert dann die **Residuen**

$$y_t - \mathbf{X}_t \tilde{\beta} \quad (t = 1, \dots, n).$$

Wir definieren die **Summe der quadrierten Residuen** als:

$$\begin{aligned} \text{SSR}(\tilde{\beta}) &:= \sum_{t=1}^n (y_t - \mathbf{X}_t \tilde{\beta})^2 \\ &= (\mathbf{y} - \mathbf{X} \tilde{\beta})^\top (\mathbf{y} - \mathbf{X} \tilde{\beta}) \\ &\stackrel{\text{M20.4}}{=} \mathbf{y}^\top \mathbf{y} - 2 \tilde{\beta}^\top \mathbf{X}^\top \mathbf{y} + \tilde{\beta}^\top \mathbf{X}^\top \mathbf{X} \tilde{\beta}. \end{aligned} \tag{1.20}$$

Der **Kleinste-Quadrate (OLS) Schätzer** von β , \mathbf{b} , ist definiert als

$$\mathbf{b} := \arg \min_{\tilde{\beta}} \text{SSR}(\tilde{\beta}). \tag{1.21}$$

Die Matrixalgebra-Regeln 31, 1. und 2. sagen uns, dass

$$\frac{\partial(\tilde{\beta}^\top \mathbf{X}^\top \mathbf{y})}{\partial \tilde{\beta}} = \mathbf{X}^\top \mathbf{y}, \quad \frac{\partial(\tilde{\beta}^\top \mathbf{X}^\top \mathbf{X} \tilde{\beta})}{\partial \tilde{\beta}} = 2\mathbf{X}^\top \mathbf{X} \tilde{\beta}.$$

Um $SSR(\tilde{\beta})$ zu minimieren, muss $\tilde{\beta}$ die Bedingungen erster Ordnung erfüllen:

$$\frac{\partial SSR(\tilde{\beta})}{\partial \tilde{\beta}} = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \tilde{\beta} \stackrel{!}{=} \mathbf{0}$$

oder

$$\mathbf{X}^\top \mathbf{X} \tilde{\beta} = \mathbf{X}^\top \mathbf{y}.$$

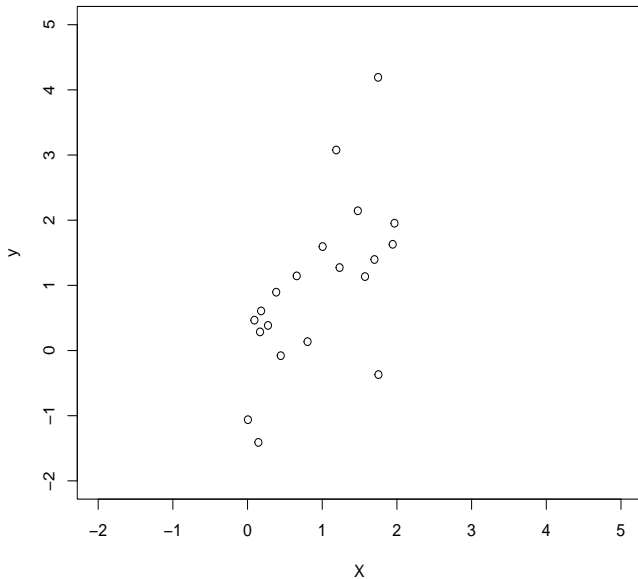
M6 liefert dann

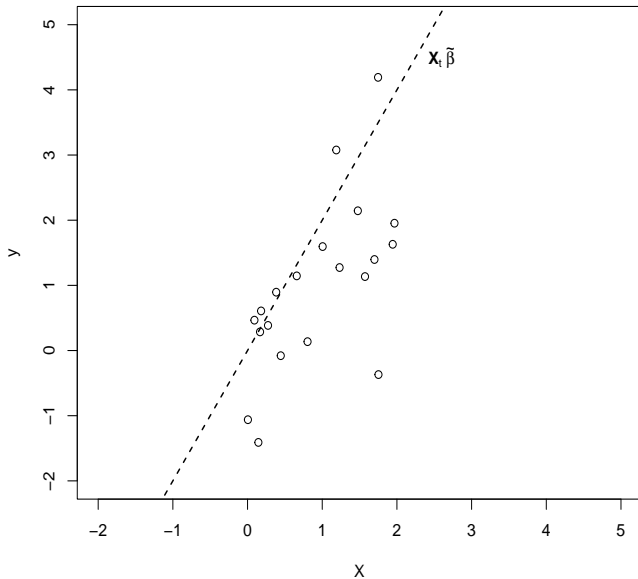
$$\mathbf{b} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

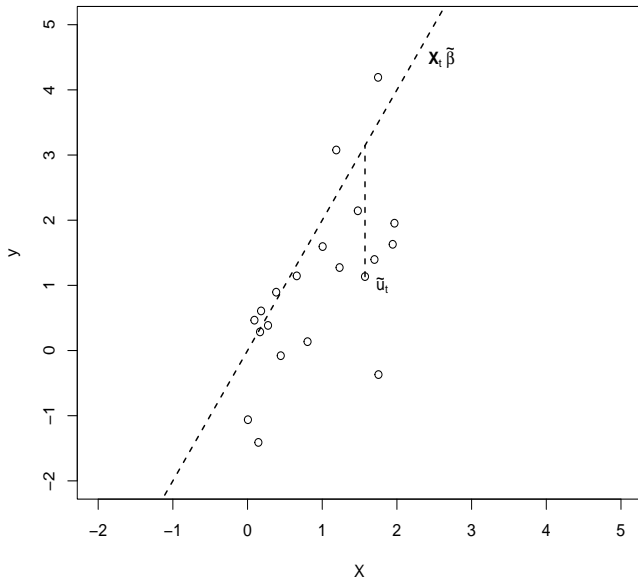
Das bedeutet, dass $\tilde{\beta}$ muss (1.16) erfüllen muss—in anderen Worten

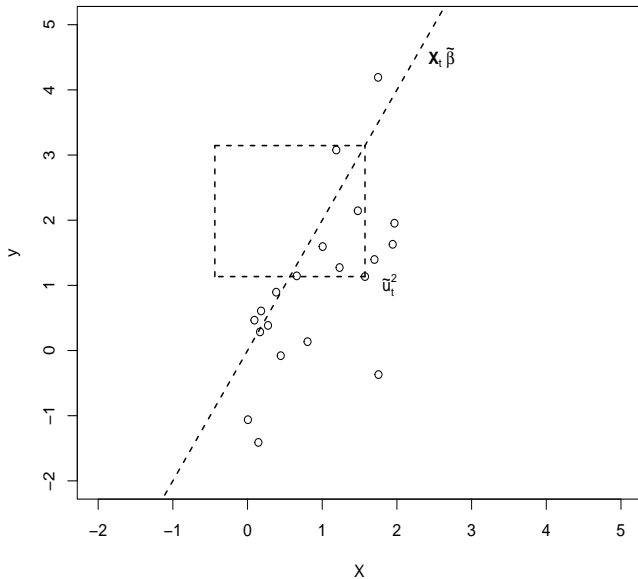
$$\mathbf{b} = \hat{\beta}_{OLS}!$$

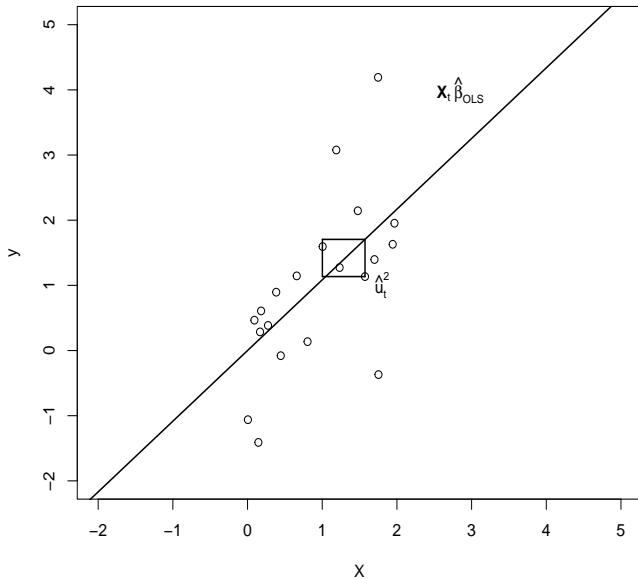
BechdelTest.R

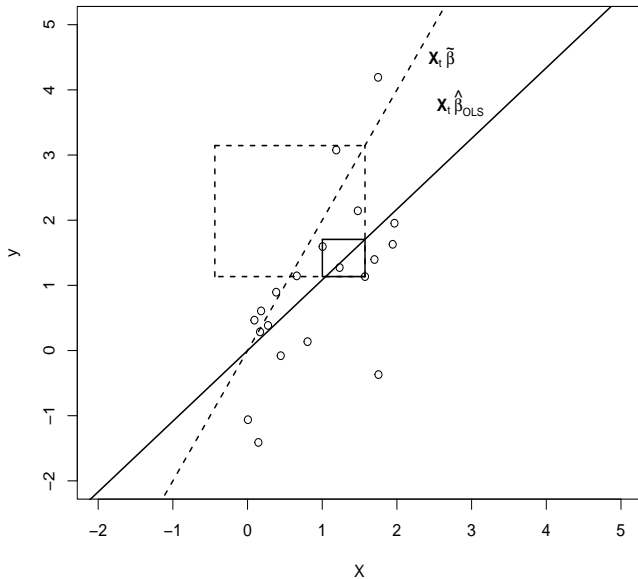












Überblick

- 1 Regressionsmodelle
- 2 **Lineare Regression**
- 3 Statistische Eigenschaften von Least Squares
- 4 Inferenz
- 5 Nichtlineare Regression
- 6 Generalized Least Squares
- 7 Instrumentvariablen

Bisher haben wir angenommen, dass die Inverse $(\mathbf{X}^\top \mathbf{X})^{-1}$ in (1.17) existiert. Dies ist jedoch nicht zwangsläufig der Fall. Die Inverse $(\mathbf{X}^\top \mathbf{X})^{-1}$ existiert, wenn $\mathbf{X}^\top \mathbf{X}$ den Rang k hat (Matrixalgebra-Reader Resultat 21.6).

Dies ist der Fall, wenn \mathbf{X} Rang k hat (siehe Resultat 20.8—ausgenommen den Fall, dass $n < k$). \mathbf{X} hat vollen Spaltenrang (= Rang k), wenn die Spalten $\mathbf{x}_1, \dots, \mathbf{x}_k$ linear unabhängig sind (siehe Definition 12). Dies ist wiederum der Fall, wenn der einzige Weg

$$\sum_{i=1}^k c_i \mathbf{x}_i = \mathbf{0}$$

zu erhalten ist, für alle i $c_i = 0$ zu setzen. Es darf demnach nicht möglich sein einen Regressor als eine lineare Kombination eines anderen zu schreiben, wie in

$$\mathbf{x}_j = \sum_{i \neq j} \frac{c_i}{c_j} \mathbf{x}_i.$$

Falls dies doch möglich sein sollte, liegt **Multikollinearität** vor.

Beispiel 2.1 Ein volles **Dummy-Variablen** Modell mit Konstante liefert Multikollinearität.

$$\log(wage)_t = \beta_1 + \beta_2 \cdot man_t + \beta_3 \cdot woman_t + \mathbf{X}_t\beta_4 + u_t$$

Wir haben $man_t + woman_t = 1$ für alle t (Bemerkung!). Eine Lösung ist es die Konstante oder einen der Dummies wegzulassen, z.B.

$$\log(wage)_t = \beta_1 \cdot man_t + \beta_2 \cdot woman_t + \mathbf{X}_t\beta_3 + u_t \quad (2.1)$$

Die **Interpretation der Koeffizienten** ändert sich. Zum Beispiel ist für (2.1)

$$E(\log(wage)_t | woman_t = 1) = \beta_2 + \mathbf{X}_t\beta_3$$

während $E(\log(wage)_t | woman_t = 1) = \beta_1 + \beta_2 + \mathbf{X}_t\beta_3$ in

$$\log(wage)_t = \beta_1 + \beta_2 \cdot woman_t + \mathbf{X}_t\beta_3 + u_t$$

BechdelTest.R

Definition 2.2 Die mit OLS **gefitteten Werte** sind

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{OLS}}.$$

Definition 2.3 Die **Projektionsmatrix** und die **Residuenmachermatrix** sind

$$\underset{(n \times n)}{\mathbf{P}_\mathbf{X}} := \mathbf{P} := \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \quad (2.2)$$

$$\underset{(n \times n)}{\mathbf{M}_\mathbf{X}} := \mathbf{M} := \mathbf{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top. \quad (2.3)$$

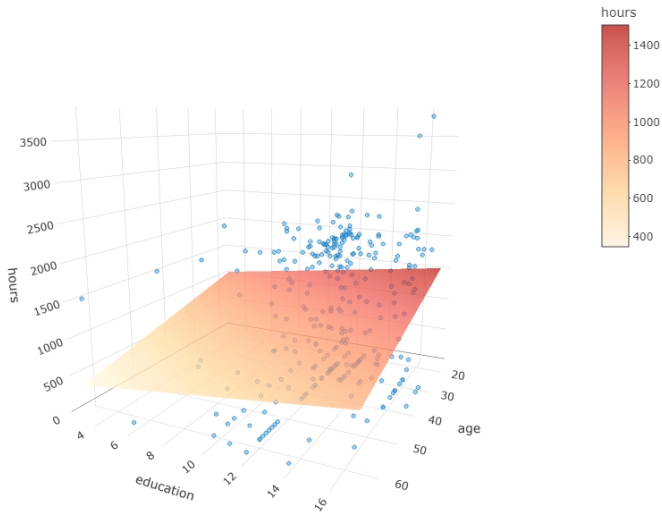
Diese Bezeichnungen kommen zustande, da

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{OLS}} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{P}\mathbf{y}$$

und

$$\hat{\mathbf{u}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{I} - \mathbf{P})\mathbf{y} = \mathbf{M}\mathbf{y}. \quad (2.4)$$

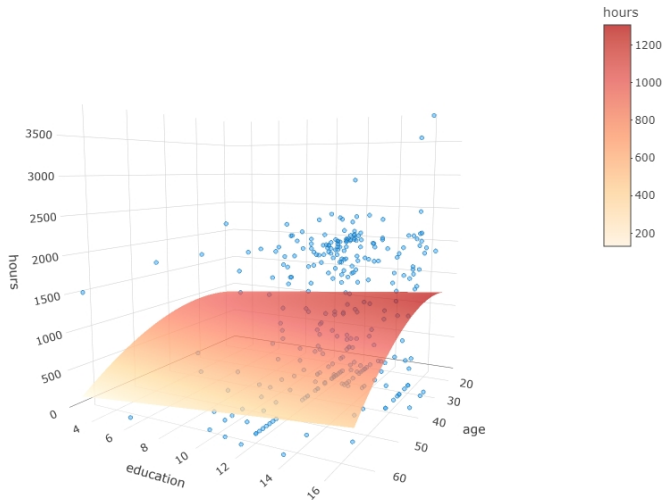
Mit zwei Regressoren kann Definition 2.2 noch visualisiert werden:



RegPlane.R. Erstellt in R mit dem Datensatz `Ecdat::workinghours`.



$$\text{hours}_t = \hat{\beta}_1 + \hat{\beta}_2 \text{age}_t + \hat{\beta}_3 \text{age}_t^2 + \hat{\beta}_4 \text{education}_t \quad (\text{s. auch (1.11)}):$$



Es ist nicht schwierig zu zeigen, dass

$$\mathbf{PX} = \mathbf{X} \quad (2.5)$$

und

$$\mathbf{MX} = \mathbf{0}. \quad (2.6)$$

Residuen und gefittete Werte sind **orthogonal** zueinander, da

$$\hat{\mathbf{u}}^\top \hat{\mathbf{y}} = 0 \quad (2.7)$$

BechdelTest.R

Dies ergibt sich aus einer Reihe von Eigenschaften. Zunächst ist \mathbf{P} symmetrisch

$$\begin{aligned}\mathbf{P}^\top &= (\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top \stackrel{\text{M20.4}}{=} (\mathbf{X}^\top)^\top ((\mathbf{X}^\top \mathbf{X})^{-1})^\top \mathbf{X}^\top \\ &\stackrel{\text{M21.3}}{=} \mathbf{X}((\mathbf{X}^\top \mathbf{X})^\top)^{-1} \mathbf{X}^\top \\ &= \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{P}\end{aligned}$$

Damit ist auch \mathbf{M} offensichtlich symmetrisch. Daher ist

$$\begin{aligned}\hat{\mathbf{u}}^\top \hat{\mathbf{y}} &= \mathbf{y}^\top \mathbf{M}^\top \mathbf{P} \mathbf{y} \\ &= \mathbf{y}^\top \mathbf{M} \mathbf{P} \mathbf{y}\end{aligned}$$

Nutze nun Idempotenz von \mathbf{P} ($\mathbf{P}\mathbf{P} = \mathbf{P}$, siehe Definition 9.2),

$$\begin{aligned}\mathbf{M}\mathbf{P} &= \mathbf{P} - \mathbf{P}\mathbf{P} \\ &= \mathbf{P} - \mathbf{X} \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}}_{\mathbf{I}} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \\ &= \mathbf{P} - \mathbf{P} = \mathbf{0}\end{aligned}\tag{2.8}$$

Dies impliziert demnach auch Orthogonalität von Residuen und Regressoren:

$$\mathbf{X}^\top \hat{\mathbf{u}} = \mathbf{0} \quad (2.9)$$

Übrigens ist auch \mathbf{M} idempotent:

$$\begin{aligned} \mathbf{M}\mathbf{M} &= (\mathbf{I} - \mathbf{P})(\mathbf{I} - \mathbf{P}) \\ &= \mathbf{I} - 2\mathbf{P} + \mathbf{P}\mathbf{P} \\ &= \mathbf{I} - 2\mathbf{P} + \mathbf{P} \\ &= \mathbf{I} - \mathbf{P} \\ &= \mathbf{M} \end{aligned} \quad (2.10)$$

Die Definition der Residuen sagt uns, dass

$$\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{OLS}} + \hat{\mathbf{u}} = \hat{\mathbf{y}} + \hat{\mathbf{u}}. \quad (2.11)$$

(2.7) erlaubt uns die **total sum of squares** $\text{TSS} := \mathbf{y}^\top \mathbf{y}$ wie folgt zu zerlegen:

$$\begin{aligned} \sum_{t=1}^n y_t^2 = \mathbf{y}^\top \mathbf{y} &= (\hat{\mathbf{y}} + \hat{\mathbf{u}})^\top (\hat{\mathbf{y}} + \hat{\mathbf{u}}) \\ &= \hat{\mathbf{y}}^\top \hat{\mathbf{y}} + 2\hat{\mathbf{u}}^\top \hat{\mathbf{y}} + \hat{\mathbf{u}}^\top \hat{\mathbf{u}} \\ &\stackrel{(2.7)}{=} \hat{\mathbf{y}}^\top \hat{\mathbf{y}} + \hat{\mathbf{u}}^\top \hat{\mathbf{u}} \\ &=: \text{ESS} + \text{SSR}, \end{aligned} \quad (2.12)$$

wobei ESS die **explained sum of squares** und SSR die **sum of squared residuals** ist.

Dies ist eine fundamentale OLS-Eigenschaft: Wir zerlegen die Varianz von \mathbf{y} in einen Teil, den wir erklären können und in einen, den wir nicht erklären können.

Solche Projektionsargumente sind nützlich, um einige Eigenschaften von OLS zu zeigen. Beispielsweise hat es keinen Einfluss auf die gefitteten Werte oder die Residuen, wenn wir die Maßeinheiten von \mathbf{X} ändern.

Betrachten wir eine Transformation von \mathbf{X} durch eine invertierbare $k \times k$ Matrix \mathbf{A} , \mathbf{XA} (ändern wir z.B. die Einheit von Monaten zu Jahren oder von Metern zu Zentimetern).

Dann:

$$\begin{aligned}
 P_{\mathbf{XA}} &:= \mathbf{XA}((\mathbf{XA})^\top \mathbf{XA})^{-1}(\mathbf{XA})^\top \\
 &\stackrel{\text{M20.4}}{=} \mathbf{XA}(\mathbf{A}^\top \mathbf{X}^\top \mathbf{XA})^{-1} \mathbf{A}^\top \mathbf{X}^\top \\
 &\stackrel{\text{M20.5}}{=} \mathbf{XAA}^{-1}(\mathbf{X}^\top \mathbf{X})^{-1}(\mathbf{A}^\top)^{-1} \mathbf{A}^\top \mathbf{X}^\top \\
 &= P_{\mathbf{X}}
 \end{aligned} \tag{2.13}$$

Wie sieht es bei $\hat{\beta}_{OLS}$ aus? Betrachten wir

$$\begin{aligned}\hat{\beta}_{OLS}^o &= \underbrace{(\mathbf{A}^\top \mathbf{X}^\top)}_{\text{„}\mathbf{X}^\top\text{“}} \underbrace{\mathbf{X} \mathbf{A}}_{\text{„}\mathbf{X}\text{“}}^{-1} \underbrace{\mathbf{A}^\top \mathbf{X}^\top}_{\text{„}\mathbf{X}^\top\text{“}} \mathbf{y} \\ &\stackrel{M20.5}{=} \mathbf{A}^{-1} (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{A}^\top)^{-1} \mathbf{A}^\top \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{A}^{-1} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{A}^{-1} \hat{\beta}_{OLS}\end{aligned}$$

Wenn also folgende Werte im obigen Beispiel verwendet werden,

$$\mathbf{A} = \begin{pmatrix} 1/12 & 0 \\ 0 & 100 \end{pmatrix} \quad \text{so dass} \quad \mathbf{A}^{-1} = \begin{pmatrix} 12 & 0 \\ 0 & 1/100 \end{pmatrix}$$

wird der Effekt einer Änderung in den Regressoren—intuitiverweise—angepasst.

Das Frisch-Waugh-Lovell Theorem ist (nicht nur) ein nützliches Stück Regressionsalgebra. Betrachten wir das Modell

$$\mathbf{y} = \underbrace{\mathbf{X}_1}_{n \times k_1} \beta_1 + \underbrace{\mathbf{X}_2}_{n \times k_2} \beta_2 + \mathbf{u} \quad (2.14)$$

Wir interessieren uns für die Schätzung von β_2 . Das FWL Theorem besagt, dass es zwei Möglichkeiten zur Berechnung von $\hat{\beta}_2$ gibt.

Theorem 2.4 *Folgende Berechnungen von $\hat{\beta}_2$ liefern numerisch identische Koeffizienten.*

1. *Regressiere \mathbf{y} auf \mathbf{X}_1 und \mathbf{X}_2 .*
2.
 - ▶ *Regressiere \mathbf{y} auf \mathbf{X}_1 und speichere die Residuen, nenne sie $\mathbf{M}_{\mathbf{X}_1} \mathbf{y}$, wobei $\mathbf{M}_{\mathbf{X}_1} := \mathbf{I} - \mathbf{X}_1(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top$.*
 - ▶ *Regressiere (jede Spalte in) \mathbf{X}_2 auf \mathbf{X}_1 und speichere die Residuen, nenne sie $\mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2$.*
 - ▶ *Regressiere die Residuen aufeinander, also $\mathbf{M}_{\mathbf{X}_1} \mathbf{y}$ auf $\mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2$, um den Schätzer $\hat{\beta}_2$ zu erhalten.*

Um dies zu zeigen, betrachten wir zunächst die Formel für $\hat{\beta}_2$ aus 2. Über die OLS-Formel (1.17) erhalten wir

$$\begin{aligned}\hat{\beta}_2 &= \underbrace{(\mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1}^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2)^{-1}}_{\text{„X“}} \underbrace{\mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1}^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{y}}_{\text{„y“}} \\ &\stackrel{(2.10)}{=} (\mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2)^{-1} \mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{y}\end{aligned}\quad (2.15)$$

Wenn wir die Regression wie in 1 durchführen, also direkt wie in (2.14), erhalten wir unter Verwendung von (2.4)

$$\begin{aligned}\mathbf{y} &= \mathbf{X} \hat{\beta}_{\text{OLS}} + \mathbf{My} \\ &= (\mathbf{X}_1 \quad \mathbf{X}_2) \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} + \mathbf{My} \\ &= \mathbf{X}_1 \hat{\beta}_1 + \mathbf{X}_2 \hat{\beta}_2 + \mathbf{My}\end{aligned}\quad (2.16)$$

Multiplizieren von (2.16) mit $\mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1}$ ergibt

$$\mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{y} = \mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2 \hat{\beta}_2 + \mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{My}$$

Hier fällt der erste Term in (2.16) weg, da $\mathbf{M}_{\mathbf{X}_1} \mathbf{X}_1 = \mathbf{0}$.

$\mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{M} \mathbf{y}$ ist ebenfalls gleich null. Denn:

$$\begin{aligned} \mathbf{M} \mathbf{M}_{\mathbf{X}_1} &= (\mathbf{I} - \mathbf{P})(\mathbf{I} - \mathbf{P}_{\mathbf{X}_1}) \\ &= \mathbf{I} - \mathbf{P} - \mathbf{P}_{\mathbf{X}_1} + \mathbf{P} \mathbf{P}_{\mathbf{X}_1} \\ &= \mathbf{I} - \mathbf{P} - \mathbf{P}_{\mathbf{X}_1} + \mathbf{P}_{\mathbf{X}_1} \\ &= \mathbf{I} - \mathbf{P} \\ &= \mathbf{M} \end{aligned}$$

Dass $\mathbf{P} \mathbf{P}_{\mathbf{X}_1} = \mathbf{P}_{\mathbf{X}_1}$ (dritte Gleichheit) gilt, geht auf die Tatsache zurück, dass die gefitteten Werte einer Regression einiger Spalten von \mathbf{X} (hier also \mathbf{X}_1), $\mathbf{P} \mathbf{X}_1$, \mathbf{X}_1 entsprechen. Um dies besser zu sehen, schreiben wir $\mathbf{P} \mathbf{X} = \mathbf{X}$ als

$$\mathbf{P}(\mathbf{X}_1 : \mathbf{X}_2) = (\mathbf{X}_1 : \mathbf{X}_2),$$

so dass die ersten k_1 Spalten von $\mathbf{P}(\mathbf{X}_1 : \mathbf{X}_2)$ der Matrix $\mathbf{P} \mathbf{X}_1 = \mathbf{X}_1$ entsprechen. Daraus folgt, dass

$$\mathbf{P} \mathbf{P}_{\mathbf{X}_1} = \mathbf{P} \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top = \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top = \mathbf{P}_{\mathbf{X}_1}$$

Daraus folgt ebenfalls, dass

$$\mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{M} = \mathbf{X}_2^\top \mathbf{M} = \mathbf{0}, \quad (2.17)$$

da $\mathbf{M}\mathbf{X}_2$ den Residuen einer Regression von \mathbf{X}_2 auf \mathbf{X} entspricht, was \mathbf{X}_2 enthält. Demzufolge hat die Regression perfekten Fit und die Residuen sind daher null.

Jetzt lösen wir

$$\mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{y} = \mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2 \hat{\beta}_2$$

nach $\hat{\beta}_2$ auf und sind fertig. Siehe (2.15).

fwl.R

Eine Anwendung

(Interessantere kommen später)

Beispiel 2.5 Nehmen wir $\mathbf{X}_1 = \boldsymbol{\iota} := (1, \dots, 1)^\top$. Das Folgende ergibt identische $\hat{\beta}_2$, die Koeffizienten für \mathbf{X}_2 .

- Regressiere \mathbf{y} auf $\boldsymbol{\iota}$ (d.h. eine Konstante) und \mathbf{X}_2 .
- Regressiere $\mathbf{y} - \bar{y}$ auf $\mathbf{X}_2 - \bar{\mathbf{X}}_2$, wobei letzteres die Spalten von \mathbf{X}_2 zentriert um deren jeweilige Spaltenmittelwerte bezeichnet.

Um dies zu sehen, betrachte

$$\mathbf{M}_{\boldsymbol{\iota}} = \mathbf{I} - \boldsymbol{\iota}(\boldsymbol{\iota}^\top \boldsymbol{\iota})^{-1} \boldsymbol{\iota}^\top = \mathbf{I} - \frac{\boldsymbol{\iota} \boldsymbol{\iota}^\top}{n}, \quad (2.18)$$

so dass (wobei \mathbf{x}_i der Spaltenvektor der n Beobachtungen einer Variable aus \mathbf{X}_2 sei)

$$\mathbf{M}_{\boldsymbol{\iota}} \mathbf{x}_i = \mathbf{x}_i - \boldsymbol{\iota} n^{-1} \boldsymbol{\iota}^\top \mathbf{x}_i = \mathbf{x}_i - \boldsymbol{\iota} \bar{x}_i =: \mathbf{x}_i - \bar{\mathbf{x}}_i.$$

Daher sind die Residuen einer Regression von Variablen auf eine Konstante $\mathbf{M}_{\boldsymbol{\iota}} \mathbf{x}_i$ lediglich die um den Mittelwert bereinigten Variablen.

Goodness of fit

Erinnern wir uns an die Zerlegung aus (2.12),

$$\text{TSS} = \mathbf{y}^\top \mathbf{y} = \hat{\mathbf{y}}^\top \hat{\mathbf{y}} + \hat{\mathbf{u}}^\top \hat{\mathbf{u}} = \text{ESS} + \text{SSR}$$

Dies liefert ein Maß für die Anpassungsgüte (goodness of fit), das **Bestimmtheitsmaß** oder auch R^2 genannt wird.

Definition 2.6 (unzentriertes R^2).

$$R_u^2 = \frac{\text{ESS}}{\text{TSS}} = \frac{\text{TSS} - \text{SSR}}{\text{TSS}} = 1 - \frac{\text{SSR}}{\text{TSS}} = 1 - \frac{\hat{\mathbf{u}}^\top \hat{\mathbf{u}}}{\mathbf{y}^\top \mathbf{y}}$$

Wir interpretieren das R^2 als den Teil der Variation in den Daten, der durch das Modell erklärt wird.

- $R_u^2 = 1$, wenn die Regressionslinie perfekt passt.
- $R_u^2 = 0$, wenn $\hat{\beta}_{\text{OLS}} = \mathbf{0}$ (s. (2.11) mit dann $\mathbf{y} = \hat{\mathbf{u}}$).

Diese Version ist jedoch nicht robust gegen scheinbar harmlosen Transformationen von \mathbf{y} . Die folgende Variante, die in Anwendungen beliebter ist, vermeidet diesen Nachteil.

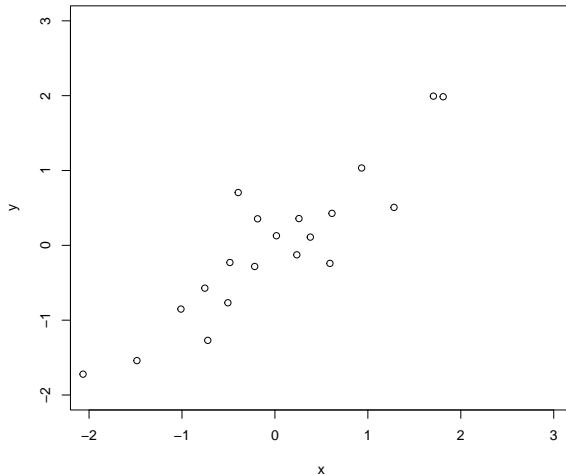
Definition 2.7 (zentriertes R^2).

$$\begin{aligned} R^2 &= 1 - \frac{\hat{\mathbf{u}}^\top \hat{\mathbf{u}}}{\mathbf{y}^\top \mathbf{M}_L \mathbf{y}} \\ &\stackrel{\text{Ex. 2.5}}{=} 1 - \frac{\hat{\mathbf{u}}^\top \hat{\mathbf{u}}}{(\mathbf{y} - \bar{\mathbf{y}})^\top (\mathbf{y} - \bar{\mathbf{y}})} \end{aligned} \quad (2.19)$$

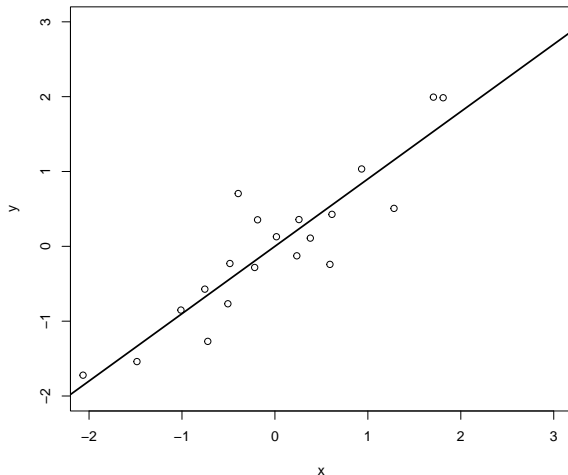
Dieses Maß macht jedoch nur Sinn, wenn eine Konstante in der Regression enthalten ist. Ansonsten ist ein $R^2 < 0$ möglich, was die Interpretation fragwürdig macht.

Wir müssen uns, unabhängig vom Wert für R^2 , immer bewusst sein, was das R^2 für uns tun kann und was nicht.

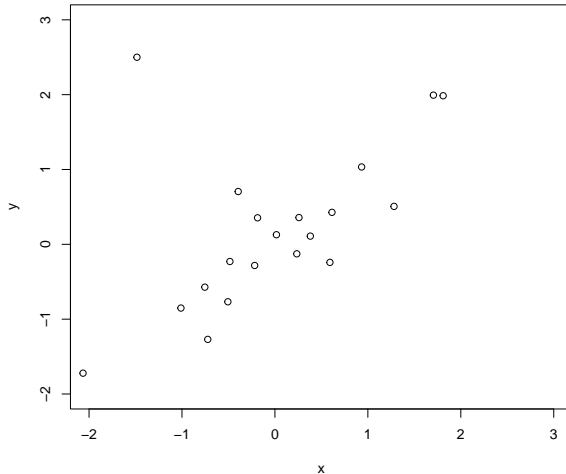
Die Minimierung der *quadrierten* Residuen in OLS impliziert, dass **Ausreißer (Outlier)** ein hohes Gewicht erhalten:



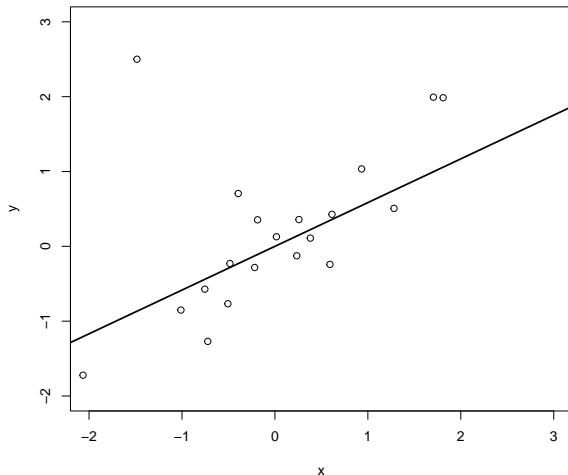
Die Minimierung der *quadrierten* Residuen in OLS impliziert, dass **Ausreißer (Outlier)** ein hohes Gewicht erhalten:



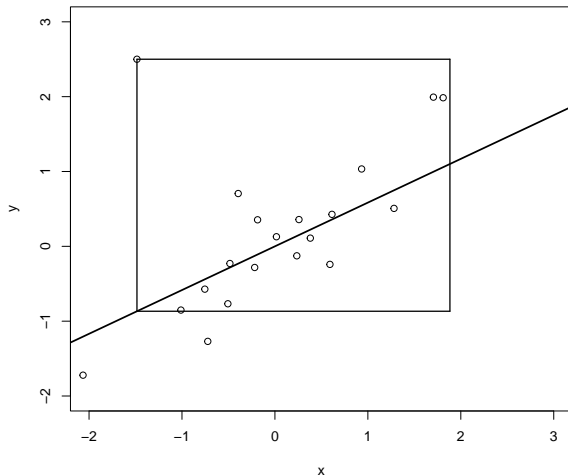
Die Minimierung der *quadrierten* Residuen in OLS impliziert, dass **Ausreißer (Outlier)** ein hohes Gewicht erhalten:



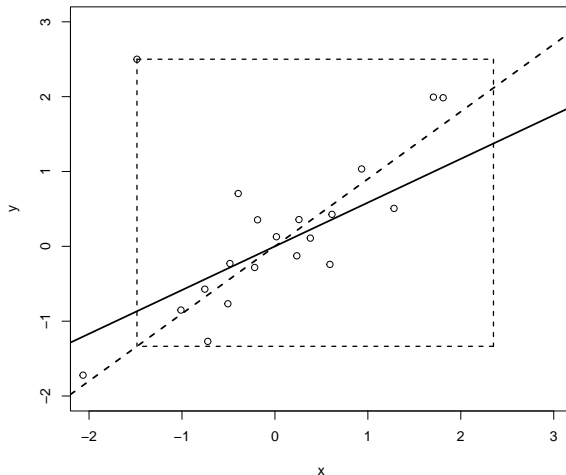
Die Minimierung der *quadrierten* Residuen in OLS impliziert, dass **Ausreißer (Outlier)** ein hohes Gewicht erhalten:



Die Minimierung der *quadrierten* Residuen in OLS impliziert, dass **Ausreißer (Outlier)** ein hohes Gewicht erhalten:



Die Minimierung der *quadrierten* Residuen in OLS impliziert, dass **Ausreißer (Outlier)** ein hohes Gewicht erhalten:



Daher ist es immer eine gute Idee die Daten erst einmal wie oben zu plotten.

Was wir hingegen mit Ausreißern machen sollen, ist weniger klar. Im einfachsten Fall wurde einfach falsch gemessen, oder fehlende Werte sind mit einem Zahlenwert codiert (oft -99), werden jedoch für eine echte Beobachtung gehalten. In diesem Fall können die Fehler korrigiert oder die Datenpunkte aus dem Datensatz entfernt werden.

Dennoch können auch ungewöhnlich große Realisationen vorkommen, wenn auch selten. Dann enthalten sie interessante Informationen über das Modell, gerade *weil* sie die Schätzungen so stark beeinflussen.

Meist ist es eine gute Strategie die Ergebnisse mit und ohne Ausreißer darzustellen und die Unterschiede zu interpretieren.

`MinimizeSumPowerResiduals.R`

Um Ausreißer etwas formaler zu analysieren, betrachten wir die Schätzung von β , wenn wir die t -te Beobachtung herauslassen. Wenn

$$\mathbf{y}_{(t)} = \begin{pmatrix} y_1 \\ \vdots \\ y_{t-1} \\ y_{t+1} \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{X}_{(t)} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_{t-1} \\ \mathbf{x}_{t+1} \\ \vdots \\ \mathbf{x}_n \end{pmatrix}, \quad \mathbf{u}_{(t)} = \begin{pmatrix} u_1 \\ \vdots \\ u_{t-1} \\ u_{t+1} \\ \vdots \\ u_n \end{pmatrix},$$

dann ist

$$\mathbf{y}_{(t)} = \mathbf{X}_{(t)}\beta + \mathbf{u}_{(t)}.$$

Eine OLS-Regression liefert

$$\hat{\beta}_{(t)} = (\mathbf{X}_{(t)}^\top \mathbf{X}_{(t)})^{-1} \mathbf{X}_{(t)}^\top \mathbf{y}_{(t)}.$$

Wenn die Differenz zwischen $\hat{\beta}$ und $\hat{\beta}_{(t)}$ „groß“ ist (die Differenz ist ein Vektor), dann bezeichnen wir die t -te Beobachtung als **einflussreich**.

Die Diagonalelemente der **Projektionsmatrix** \mathbf{P}_X sind gegeben durch $h_t = \mathbf{e}_t^\top \mathbf{P}_X \mathbf{e}_t$, $t = 1, \dots, n$, wobei die \mathbf{e}_t Einheitsvektoren sind, welche die Spalten der Einheitsmatrix bilden. Man kann zeigen, dass

$$0 \leq h_t \leq 1, \quad t = 1, \dots, n,$$

$$\sum_{t=1}^n h_t = k.$$

Um den Einfluss der t -ten Beobachtung zu ermitteln, betrachten wir das folgende Resultat. Erinnern wir uns an $\hat{\mathbf{u}} = \mathbf{y} - \hat{\mathbf{y}}$.

Resultat 2.8

$$\hat{\beta} - \hat{\beta}_{(t)} = \left(\frac{\hat{u}_t}{1 - h_t} \right) (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_t^\top \quad (2.20)$$

Wenn h_t groß ist, z.B. nah bei eins, sagen wir, dass \mathbf{X}_t starke **Hebelwirkung** (**Leverage**) hat. Wenn $|\hat{u}_t|$ ebenfalls groß ist, ist die t -te Beobachtung einflussreich.

Um das Resultat 2.8 auf eine andere Weise als im Lehrbuch zu zeigen, betrachten wir das folgende Matrixalgebra-Resultat.

Resultat 2.9 Gegeben sei eine nicht-singuläre Matrix \mathbf{A} , ein Vektor \mathbf{b} und ein Skalar λ . Wenn

$$\lambda \neq -\frac{1}{\mathbf{b}^\top \mathbf{A}^{-1} \mathbf{b}}, \quad (2.21)$$

dann ist

$$(\mathbf{A} + \lambda \mathbf{b} \mathbf{b}^\top)^{-1} = \mathbf{A}^{-1} - \left(\frac{\lambda}{1 + \lambda \mathbf{b}^\top \mathbf{A}^{-1} \mathbf{b}} \right) \mathbf{A}^{-1} \mathbf{b} \mathbf{b}^\top \mathbf{A}^{-1}. \quad (2.22)$$

Beachte, dass $\mathbf{A} + \lambda \mathbf{b} \mathbf{b}^\top$ singular ist, wenn (2.21) nicht erfüllt ist. Der **Beweis für Resultat 2.9** folgt direkt aus der Überprüfung von

$$\left\{ \mathbf{A}^{-1} - \left(\frac{\lambda}{1 + \lambda \mathbf{b}^\top \mathbf{A}^{-1} \mathbf{b}} \right) \mathbf{A}^{-1} \mathbf{b} \mathbf{b}^\top \mathbf{A}^{-1} \right\} (\mathbf{A} + \lambda \mathbf{b} \mathbf{b}^\top) = \mathbf{I}.$$

Das folgende Resultat ist hilfreich beim Beweis von 2.8.

Resultat 2.10

$$(\mathbf{X}_{(t)}^\top \mathbf{X}_{(t)})^{-1} \mathbf{X}_t^\top = \left(\frac{1}{1 - h_t} \right) (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_t^\top. \quad (2.23)$$

Beweis von Resultat 2.10: Durch Resultat 2.9 erhalten wir

$$\begin{aligned} (\mathbf{X}_{(t)}^\top \mathbf{X}_{(t)})^{-1} &= (\mathbf{X}^\top \mathbf{X} - \mathbf{X}_t^\top \mathbf{X}_t)^{-1} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} + \frac{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_t^\top \mathbf{X}_t (\mathbf{X}^\top \mathbf{X})^{-1}}{1 - \mathbf{X}_t^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_t}. \end{aligned}$$

Daher erhalten wir

$$\begin{aligned} (\mathbf{X}_{(t)}^\top \mathbf{X}_{(t)})^{-1} \mathbf{X}_t^\top &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_t^\top + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_t^\top \left(\frac{\mathbf{X}_t (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_t^\top}{1 - \mathbf{X}_t^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_t} \right) \\ &= \left(\frac{1}{1 - h_t} \right) (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_t^\top. \end{aligned}$$

q.e.d.

Der Beweis von Resultat 2.8 folgt nun aus Resultat 2.10.

Beweis von Resultat 2.8: Da

$$\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{y},$$

erhalten wir

$$(\mathbf{X}_{(t)}^\top \mathbf{X}_{(t)} + \mathbf{x}_t^\top \mathbf{x}_t) \hat{\boldsymbol{\beta}} = \mathbf{X}_{(t)}^\top \mathbf{y}_{(t)} + \mathbf{x}_t^\top y_t,$$

oder

$$\left\{ \mathbf{I}_k + (\mathbf{X}_{(t)}^\top \mathbf{X}_{(t)})^{-1} \mathbf{x}_t^\top \mathbf{x}_t \right\} \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_{(t)} + (\mathbf{X}_{(t)}^\top \mathbf{X}_{(t)})^{-1} \mathbf{x}_t^\top (\mathbf{x}_t \hat{\boldsymbol{\beta}} + \hat{u}_t).$$

Daher

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \hat{\boldsymbol{\beta}}_{(t)} + (\mathbf{X}_{(t)}^\top \mathbf{X}_{(t)})^{-1} \mathbf{x}_t^\top \hat{u}_t \\ &= \hat{\boldsymbol{\beta}}_{(t)} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_t^\top \frac{\hat{u}_t}{1 - h_t}, \end{aligned}$$

wobei die letzte Gleichung aus Resultat 2.10 folgt

q.e.d.



Überblick

- 1 Regressionsmodelle
- 2 Lineare Regression
- 3 Statistische Eigenschaften von Least Squares**
- 4 Inferenz
- 5 Nichtlineare Regression
- 6 Generalized Least Squares
- 7 Instrumentvariablen

Das klassische lineare Modell ist gegeben durch

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad (3.1)$$

wobei wir den n -Vektor \mathbf{y} , die abhängige Variable, beobachten und die $n \times k$ Matrix \mathbf{X} , die erklärenden Variablen enthält.

Der n -Vektor der Störgrößen \mathbf{u} wird nicht beobachtet. Wir treffen Annahmen über seine Verteilung (wie beispielsweise (1.18)).

Auch \mathbf{X} ist (außer z.B. in Experimenten) stochastisch (denken wir z.B. an Erfahrung). Wir nehmen an, dass der k -Vektor der Koeffizienten $\boldsymbol{\beta}$ nicht zufällig aber unbekannt ist.

Wir nehmen zudem **Exogenität** an:

Annahme 3.1

$$E(\mathbf{u}|\mathbf{X}) = \mathbf{0}. \quad (3.2)$$

Dies ist für eine für eine u.i.v. Stichprobe dasselbe wie $E(u_t|\mathbf{X}_t) = 0$ in (1.18).

Unter der Annahme 3.1 ist der OLS-Schätzer bedingt unverzerrt: Wir haben mit (3.1)

$$\hat{\beta}_{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \beta + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}, \quad (3.3)$$

und daher ist

$$E(\hat{\beta}_{\text{OLS}} | \mathbf{X}) \stackrel{(1.9)}{=} \beta + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top E(\mathbf{u} | \mathbf{X}) = \beta. \quad (3.4)$$

Aufgrund des Gesetzes der iterierten Erwartungen (1.7) ist $\hat{\beta}_{\text{OLS}}$ auch unbedingt erwartungstreu.

Siehe `unbiasedness.R` für einige graphische Darstellungen.

Der Begriff einer **Varianz-Kovarianzmatrix** eines $m \times 1$ Vektors \mathbf{x} ,

$$\mathbf{x} = (X_1, \dots, X_m)^\top,$$

ist gleich hilfreich:

$$\begin{aligned} \text{Var}(\mathbf{x}) &= E\left((\mathbf{x} - E(\mathbf{x}))(\mathbf{x} - E(\mathbf{x}))^\top\right) \\ &= \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_m) \\ \text{Cov}(X_1, X_2) & \text{Var}(X_2) & \dots & \text{Cov}(X_2, X_m) \\ \vdots & \ddots & \dots & \vdots \\ \vdots & & \ddots & \vdots \\ \text{Cov}(X_1, X_m) & \dots & \text{Cov}(X_{m-1}, X_m) & \text{Var}(X_m) \end{pmatrix} \end{aligned} \quad (3.5)$$

Für

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$

erhalten wir außerdem

$$\begin{aligned} \text{Var}(\mathbf{y}) &= E((\mathbf{y} - E(\mathbf{y}))(\mathbf{y} - E(\mathbf{y}))^\top) \\ &= E((\mathbf{A}\mathbf{x} - E(\mathbf{A}\mathbf{x}))(\mathbf{A}\mathbf{x} - E(\mathbf{A}\mathbf{x}))^\top) \\ &= E((\mathbf{A}\mathbf{x} - \mathbf{A}E(\mathbf{x}))(\mathbf{A}\mathbf{x} - \mathbf{A}E(\mathbf{x}))^\top) \\ &= E(\mathbf{A}(\mathbf{x} - E(\mathbf{x}))(\mathbf{A}(\mathbf{x} - E(\mathbf{x})))^\top) \\ &\stackrel{\text{M20.4}}{=} E(\mathbf{A}(\mathbf{x} - E(\mathbf{x}))(\mathbf{x} - E(\mathbf{x}))^\top \mathbf{A}^\top) \\ &= \mathbf{A}E((\mathbf{x} - E(\mathbf{x}))(\mathbf{x} - E(\mathbf{x}))^\top) \mathbf{A}^\top \\ &= \mathbf{A}\text{Var}(\mathbf{x})\mathbf{A}^\top \end{aligned} \tag{3.6}$$

Nehmen wir zusätzlich zu 3.1 ebenfalls an, dass

Annahme 3.2 $\text{Var}(\mathbf{u}|\mathbf{X}) = \sigma^2 \mathbf{I}_n$, wobei σ eine unbekannte Konstante ist, also ein zusätzlicher Parameter.

Aus 3.1 folgt $E(\mathbf{u}\mathbf{u}^\top|\mathbf{X}) = \text{Var}(\mathbf{u}|\mathbf{X})$, da

$$\text{Var}(\mathbf{u}|\mathbf{X}) = E\left(\underbrace{(\mathbf{u} - E(\mathbf{u}|\mathbf{X}))}_{=0}(\mathbf{u} - E(\mathbf{u}|\mathbf{X}))^\top|\mathbf{X}\right) = E(\mathbf{u}\mathbf{u}^\top|\mathbf{X})$$

Was bedeutet diese Annahme? `educearn.R`

Die Kovarianzmatrix stellt sich dann wie folgt dar:

$$\begin{aligned} \text{Var}(\hat{\beta}_{\text{OLS}}|\mathbf{X}) &= E\{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u} \mathbf{u}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} | \mathbf{X}\} \\ &\stackrel{(3.5)}{=} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Var}(\mathbf{u} | \mathbf{X}) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \{\sigma^2 \mathbf{I}_n\} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \underbrace{\mathbf{X}^\top \mathbf{I}_n \mathbf{X}}_{=\mathbf{I}} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}. \end{aligned} \tag{3.7}$$

$\text{Var}(\hat{\beta}_{\text{OLS}}|\mathbf{X})$ ist unbekannt, da σ^2 unbekannt ist.

σ^2 aus 3.2 kann unverzerrt geschätzt werden durch $\hat{\sigma}^2$. $\hat{\sigma} = \sqrt{\hat{\sigma}^2}$ wird der **Standardfehler der Regression** genannt.

$$\begin{aligned}
 \hat{\sigma}^2 &= \frac{(\mathbf{y} - \mathbf{X}\hat{\beta}_{\text{OLS}})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}_{\text{OLS}})}{n - k} & (3.8) \\
 &= \frac{\hat{\mathbf{u}}^\top \hat{\mathbf{u}}}{n - k} \stackrel{(2.4)}{=} \frac{\mathbf{y}^\top \mathbf{M}^\top \mathbf{M} \mathbf{y}}{n - k} \\
 &\stackrel{(2.10)}{=} \frac{\mathbf{y}^\top \mathbf{M} \mathbf{y}}{n - k} \\
 &= \frac{(\mathbf{X}\beta + \mathbf{u})^\top \mathbf{M}(\mathbf{X}\beta + \mathbf{u})}{n - k} \\
 &\stackrel{(2.6)}{=} \frac{\mathbf{u}^\top \mathbf{M} \mathbf{u}}{n - k} \\
 &= \frac{\text{tr}(\mathbf{u}^\top \mathbf{M} \mathbf{u})}{n - k} \\
 &\stackrel{\text{M20.6}}{=} \frac{\text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}^\top)}{n - k} \\
 &= \text{tr} \left\{ \frac{(\mathbf{I}_n - \mathbf{P}_X) \mathbf{u} \mathbf{u}^\top}{n - k} \right\}.
 \end{aligned}$$

Daher erhalten wir

$$\begin{aligned} E(\hat{\sigma}^2 | \mathbf{X}) &= E \left[\text{tr} \left\{ \frac{(\mathbf{I}_n - \mathbf{P}_X) \mathbf{u} \mathbf{u}^\top}{n - k} \right\} \middle| \mathbf{X} \right] \\ &= \text{tr} \left[E \left\{ \frac{(\mathbf{I}_n - \mathbf{P}_X) \mathbf{u} \mathbf{u}^\top}{n - k} \middle| \mathbf{X} \right\} \right] \\ &= \text{tr} \left[\frac{(\mathbf{I}_n - \mathbf{P}_X) E(\mathbf{u} \mathbf{u}^\top | \mathbf{X})}{n - k} \right] \\ &\stackrel{3.2}{=} \text{tr} \left[\frac{(\mathbf{I}_n - \mathbf{P}_X) \sigma^2 \mathbf{I}}{n - k} \right] \\ &= \sigma^2 \frac{\text{tr}(\mathbf{I}_n - \mathbf{P}_X)}{n - k} \end{aligned}$$

Außerdem ist

$$\begin{aligned}\text{tr}(\mathbf{I}_n - \mathbf{P}_X) &\stackrel{\text{M19.5}}{=} \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{P}_X) \\&= n - \text{tr}(\mathbf{P}_X) \\&= n - \text{tr}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \\&= n - \text{tr}(\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}) \\&= n - \text{tr}(\mathbf{I}_k) \\&= n - k\end{aligned}\tag{3.9}$$

Also ist $\hat{\sigma}^2$ ein unverzerrter Schätzer von σ^2 .

Frage 3.3 Was ist mit $E(\hat{\sigma})$? Ist es gleich σ , ist es kleiner oder ist es größer?

Gauss-Markov

Unter den Annahmen 3.1 und 3.2 können wir zeigen, dass der OLS-Schätzer $\hat{\beta}_{OLS}$ BLUE („Best Linear Unbiased Estimator“) ist.

Dies ist die Aussage des **Gauss-Markov-Theorems**.

Betrachten wir einen anderen linearen Schätzer

$$\hat{\hat{\beta}} = \mathbf{C}\mathbf{y} = \mathbf{C}\mathbf{X}\beta + \mathbf{C}\mathbf{u}.$$

Damit $\hat{\hat{\beta}}$ unverzerrt ist, muss gelten, dass

$$E(\hat{\hat{\beta}}|\mathbf{X}) = \mathbf{C}\mathbf{X}\beta = \beta.$$

Dies sollte ungeachtet des Wertes von β gelten, das heißt es sollte für alle $\beta \in \mathbb{R}^k$ gelten. Damit $\hat{\hat{\beta}}$ unverzerrt ist, benötigen wir muss also

$$\mathbf{C}\mathbf{X} = \mathbf{I}_k. \tag{3.10}$$

Die Kovarianzmatrix von $\hat{\beta}$ ist gegeben durch

$$\begin{aligned} \text{Var}(\hat{\beta}|\mathbf{X}) &= \text{Var}(\mathbf{C}\mathbf{u}|\mathbf{X}) \stackrel{3.2\&(3.6)}{=} \mathbf{C}\{\sigma^2\mathbf{I}_n\}\mathbf{C}^\top \\ &= \sigma^2\mathbf{C}\mathbf{C}^\top. \end{aligned}$$

Die kleinste Kovarianzmatrix (im Matrixsinn) wird über die Matrix \mathbf{C} gefunden, welche (3.10) erfüllt und $\mathbf{C}\mathbf{C}^\top$ minimiert.

Die Lösung ist einfach: Erinnern wir uns, dass $\mathbf{P}_\mathbf{X} = \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top$. Dann gilt

$$\mathbf{C}\mathbf{C}^\top \geq \mathbf{C}\mathbf{P}_\mathbf{X}\mathbf{C}^\top.$$

Um dies zu sehen, müssen wir zeigen, dass

$$\mathbf{C}\mathbf{C}^\top - \mathbf{C}\mathbf{P}_\mathbf{X}\mathbf{C}^\top = \mathbf{C}(\mathbf{I} - \mathbf{P}_\mathbf{X})\mathbf{C}^\top = \mathbf{C}\mathbf{M}_\mathbf{X}\mathbf{C}^\top \geq \mathbf{0}, \quad (3.11)$$

das heißt, dass $\mathbf{C}\mathbf{M}_\mathbf{X}\mathbf{C}^\top$ p.s.d. (positiv semi-definit) ist, also dass $\mathbf{a}^\top \mathbf{C}\mathbf{M}_\mathbf{X}\mathbf{C}^\top \mathbf{a} \geq 0$ für jeden Vektor \mathbf{a} gilt.

Da wir einfach $\mathbf{d} = \mathbf{C}^\top \mathbf{a}$ definieren können, ist die Bedingung erfüllt, wenn $\mathbf{d}^\top \mathbf{M}_X \mathbf{d} \geq 0$, das heißt, wenn \mathbf{M}_X p.s.d. ist. Wir wissen bereits, dass \mathbf{M}_X symmetrisch und idempotent ist (vgl. (2.10)).

Daher sagt uns Resultat 27.3, dass alle Eigenwerte 0 oder 1 sind, also nicht negativ. Resultat 27.7 sagt uns dann, dass \mathbf{M}_X p.s.d. ist.

Wir haben dann

$$\mathbf{C} \mathbf{P}_X \mathbf{C}^\top = \underbrace{\mathbf{C} \mathbf{X}}_I (\mathbf{X}^\top \mathbf{X})^{-1} \underbrace{\mathbf{X}^\top \mathbf{C}^\top}_I = (\mathbf{X}^\top \mathbf{X})^{-1}$$

Wenn (3.10) erfüllt ist, sehen wir, dass die Kovarianzmatrix von $\hat{\hat{\beta}}$ nicht kleiner sein kann als $\sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}$. Daher gilt

$$\text{Var}(\hat{\hat{\beta}} | \mathbf{X}) \geq \text{Var}(\hat{\beta}_{\text{OLS}} | \mathbf{X}).$$

Beachten Sie, dass die Untergrenze durch $\hat{\beta}_{OLS}$ erreicht wird, wobei

$$\mathbf{C} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top,$$

da

$$\begin{aligned} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_X ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \end{aligned}$$

Daher bezeichnen wir $\hat{\beta}_{OLS}$ als **Best Linear Unbiased Estimator**: BLUE.

Dies ist ein starkes Resultat, aber wir müssen uns vor Augen halten, dass es nur unter sehr starken Annahmen erfüllt ist.

Siehe [hier](#) für eine Illustration.

Ausschluss und Inklusion von Variablen

Nehmen wir an, dass das wahre Modell (mit $E(u|\mathbf{X}) = \mathbf{0}$; definiere $\mathbf{X} = (\mathbf{X}_1 : \mathbf{X}_2)$) gegeben ist durch

$$\mathbf{y} = \mathbf{X}_1\beta + \mathbf{X}_2\gamma + \mathbf{u}, \quad (3.12)$$

jedoch führen wir die Schätzung basierend auf

$$\mathbf{y} = \mathbf{X}_1\beta + \mathbf{v} \quad (3.13)$$

durch. Das bedeutet, dass wir β mit

$$\hat{\beta}_{\text{OLS},K} = (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y},$$

schätzen, während Modell (3.12) das wahre Modell ist. Das heißt, dass wir **relevante Variablen** auslassen (wenn $\gamma \neq \mathbf{0}$).

Die Konsequenz dieses Auslassens ist eine Verzerrung in der Schätzung von β .

$$\begin{aligned} E(\hat{\beta}_{\text{OLS},K}|\mathbf{X}) &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top E(\mathbf{y}|\mathbf{X}) \\ &\stackrel{(3.12)}{=} (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top E(\mathbf{X}_1\beta + \mathbf{X}_2\gamma + \mathbf{u}|\mathbf{X}) \\ &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top (\mathbf{X}_1\beta + \mathbf{X}_2\gamma) \\ &= \beta + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2\gamma, \end{aligned} \tag{3.14}$$

Dies ist ungleich β , außer wenn $\gamma = \mathbf{0}$ (also wenn die zusätzlichen Regressoren irrelevant sind) oder $\mathbf{X}_1^\top \mathbf{X}_2 = \mathbf{0}$ (also wenn \mathbf{X}_1 und \mathbf{X}_2 orthogonal sind).

Betrachten wir nochmal unbiasedness.R.

Nehme jetzt an, dass (3.13) das wahre Modell ist (mit dann $E(\mathbf{v}|\mathbf{X}) = \mathbf{0}$), wir jedoch auf Basis von Modell (3.12) schätzen.

Dies bedeutet **Inklusion irrelevanter Variablen**, da

$$\gamma = \mathbf{0}.$$

Beide Modelle (3.12) und (3.13) sind also korrekt, jedoch verwendet (3.13) die Information, dass $\gamma = \mathbf{0}$, (3.12) hingegen nicht.

Als Konsequenz folgt daraus, dass der Schätzer (siehe (2.15))

$$\hat{\beta}_{\text{OLS,L}} = \{\mathbf{X}_1^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}_2}) \mathbf{X}_1\}^{-1} \mathbf{X}_1^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}_2}) \mathbf{y},$$

unverzerrt ist, da

$$\begin{aligned} E(\hat{\beta}_{\text{OLS,L}} | \mathbf{X}) &= \{\mathbf{X}_1^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}_2}) \mathbf{X}_1\}^{-1} \mathbf{X}_1^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}_2}) E(\mathbf{X}_1 \beta + \mathbf{v} | \mathbf{X}) \\ &= \{\mathbf{X}_1^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}_2}) \mathbf{X}_1\}^{-1} \mathbf{X}_1^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}_2}) \mathbf{X}_1 \beta = \beta. \end{aligned}$$

Er hat jedoch eine größere Kovarianzmatrix (die man analog zu (3.7) herleitet) als die des OLS-Schätzers von (3.13) (vgl. (3.7)):

$$\text{Var}(\hat{\beta}_{\text{OLS,L}} | \mathbf{X}) = \sigma^2 \{\mathbf{X}_1^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}_2}) \mathbf{X}_1\}^{-1} \geq \sigma^2 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} = \text{Var}(\hat{\beta}_{\text{OLS,K}} | \mathbf{X}).$$

Dies gilt, da diese Bedingung äquivalent ist zu (siehe z.B. Magnus and Neudecker [1988, Thm. 1.24] oder [hier](#))

$$\frac{1}{\sigma^2}(\mathbf{X}_1^\top \mathbf{X}_1) \geq \frac{1}{\sigma^2} \{ \mathbf{X}_1^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}_2}) \mathbf{X}_1 \}$$

Dies wiederum ist der Fall, wenn

$$\mathbf{X}_1^\top \mathbf{X}_1 - \{ \mathbf{X}_1^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}_2}) \mathbf{X}_1 \} = \mathbf{X}_1^\top \mathbf{X}_1 - \mathbf{X}_1^\top \mathbf{X}_1 + \mathbf{X}_1^\top \mathbf{P}_{\mathbf{X}_2} \mathbf{X}_1 = \mathbf{X}_1^\top \mathbf{P}_{\mathbf{X}_2} \mathbf{X}_1 \geq 0$$

Dass dies gilt, folgt mit dem Argument, das (3.11) gezeigt hat.

Zusätzliche Regressoren reduzieren daher das Risiko des **omitted variable bias**, was (sehr) vorteilhaft ist.

Auf der anderen Seite verringern sie die Präzision des Schätzers (außer im Fall von Orthogonalität $\mathbf{P}_{\mathbf{X}_2} \mathbf{X}_1 = \mathbf{0} \Leftrightarrow \mathbf{X}_2^\top \mathbf{X}_1 = \mathbf{0}$), was weniger gut ist.

long short reg.R

Was machen wir also?

Wenn alle Stichproben extrem groß wären, erscheint Überspezifizierung naheliegend. Der Bias durch Unterspezifizierung verschwindet nicht mit zunehmender Stichprobengröße, jedoch konvergiert die Varianz „konsistenter“ (s.u., bspw. (4.32)) Schätzer gegen null.

Um dies zu sehen, betrachten wir ein einfaches Beispiel einer Regression auf eine Konstante. Nach Einsetzen von (1.17) und (3.7) erhalten wir:

$$\hat{\beta} = \bar{y}$$

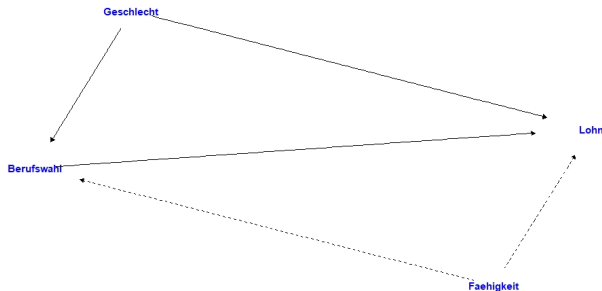
und für eine u.i.v.-Stichprobe mit der Varianz σ^2 ,

$$\text{Var}(\bar{y}) = \frac{\sigma^2}{n} \rightarrow 0$$

Auf der anderen Seite mag bei kleinem n der Effizienzverlust durch die Hinzunahme von zu vielen Variablen schwerer wiegen als die Konsequenzen des Omitted Variable Bias.

Bedeutet dies dann, dass man bei großem n alle auffindbaren Regressoren verwenden sollte? Nein!

Beispiel: Berufswahl und Lohndiskriminierung



Siehe [hier](#) bzw. [hier](#).

Überblick

- 1 Regressionsmodelle
- 2 Lineare Regression
- 3 Statistische Eigenschaften von Least Squares
- 4 Inferenz**
- 5 Nichtlineare Regression
- 6 Generalized Least Squares
- 7 Instrumentvariablen

Die Idee

Es wäre ein Glückstreffer, wenn $\hat{\beta}_{OLS}$ dem wahren Vektor β in einer endlichen Stichprobe entspräche. Daher müssen wir die Zufälligkeit von $\hat{\beta}_{OLS}$ miteinbeziehen, wenn wir Schlussfolgerungen über β treffen wollen.

Im einfachsten Fall: Der Erwartungswert einer Population, aus der eine Zufallsstichprobe gezogen wurde. Nehmen wir den folgenden DGP an

$$y_t = \beta + u_t, \quad u_t \stackrel{\text{u.i.v.}}{\sim} (0, \sigma^2)$$

wobei β der Erwartungswert ist. Der kleinste-Quadrate-Schätzer von β und seine Varianz sind gegeben durch (s.o.)

$$\hat{\beta} = \bar{y}$$

und

$$\text{Var}(\bar{y}) = \frac{\sigma^2}{n}$$

Wir wollen folgende **Nullhypothese** testen:

$$H_0 : \beta = \beta_{H_0}, \quad (4.1)$$

wobei β_{H_0} einen spezifischen Wert von β darstellt.

Wir brauchen eine **Teststatistik** z , die eine Zufallsvariable (ZV) ist, die (zumindest asymptotisch) eine bekannte Verteilung hat, wenn H_0 zutrifft.

Wenn der realisierte Wert von z nicht „ungewöhnlich“ ist, wenn H_0 zutrifft, dann liefert der Test keine Evidenz gegen H_0 .

Wenn der Wert jedoch nur selten unter H_0 (ökonometrische Abkürzung für „falls die Nullhypothese wahr ist“) erreicht würde, ist dies ein Anhaltspunkt gegen H_0 .

Wir entscheiden uns dann, H_0 **abzulehnen**.

Der Einfachheit halber treffen wir zunächst zwei sehr starke Annahmen:

1. Die u_t aus unserem Regressionsmodell sind normalverteilt.
2. Wir wissen, dass ihre Varianz σ^2 ist.

Unter diesen Annahmen können wir folgende Teststatistik benutzen

$$z = \frac{\hat{\beta} - \beta_{H0}}{(\text{Var}(\hat{\beta}))^{1/2}} = \frac{n^{1/2}}{\sigma}(\hat{\beta} - \beta_{H0}).$$

Unter der Nullhypothese (4.1), $\beta = \beta_{H0}$. Des Weiteren gilt $E(\hat{\beta}) = E(\bar{y}) = \beta$, so dass

$$E(z) = E\left(\frac{n^{1/2}}{\sigma}(\hat{\beta} - \beta_{H0})\right) = \frac{n^{1/2}}{\sigma}[E(\hat{\beta}) - \beta_{H0}] = \frac{n^{1/2}}{\sigma}[\beta - \beta] = 0$$

und

$$\text{Var}(z) = \text{Var}\left(\frac{n^{1/2}}{\sigma}\hat{\beta}\right) = \frac{n}{\sigma^2} \frac{\sigma^2}{n} = 1$$

Daraus folgt $z \sim \mathcal{N}(0, 1)$ — z hat einen Erwartungswert von 0, eine Varianz von 1 und ist eine lineare Kombination von normalverteilten Zufallsvariablen, die ebenfalls normalverteilt ist. Daher hat z unter H_0 eine bekannte Verteilung.

Es gibt außerdem (zumindest implizit) eine **Alternativhypothese**, meist als H_1 bezeichnet. Wir testen H_0 gegen H_1 . Nehmen Sie an, dass

$$H_1 : \beta \neq \beta_{H0}.$$

Unter H_1 ist z nicht $\mathcal{N}(0, 1)$ -verteilt. Wenn $\beta = \beta_1$, dann gilt $E(z) \neq 0$ und kann daher nicht $\mathcal{N}(0, 1)$ -verteilt sein.

Es ist nicht schwierig zu zeigen, dass

$$z \sim \mathcal{N}\left(\frac{n^{1/2}}{\sigma}(\beta_1 - \beta_{H0}), 1\right) \quad (4.2)$$

Wenn n groß oder die Differenz $(\beta_1 - \beta_{H0})$ absolut groß ist, ist auch $E(z)$ absolut groß. Wir lehnen H_0 ab, wenn z „ausreichend“ (siehe unten) weit entfernt von 0 ist—egal in welche Richtung: Da wir gegen

$$H_1 : \beta \neq \beta_{H0}$$

testen, führen wir einen **zweiseitigen Test** durch und verwerfen H_0 , wenn $|z|$ groß ist.

Wenn wir gegen $H_1 : \beta > \beta_{H0}$ testen wollen, führen wir einen **einseitigen Test** durch und verwerfen H_0 , wenn z ausreichend groß (im Sinne von positiv) ist.

Da z jeden Wert annehmen kann, auch unter H_0 , ist kein Wert z absolut inkompatibel mit H_0 . Daher können wir uns nie sicher sein, dass H_0 falsch ist. Wir brauchen also eine **Ablehnungsregel**, die ablehnt, wenn z in eine bestimmte Ablehnungsregion fällt.

Wenn wir eine wahre H_0 verwerfen, machen wir einen **Fehler erster Art**. Die Wahrscheinlichkeit eines solchen Fehlers ist demnach die Wahrscheinlichkeit unter H_0 , dass z in den Ablehnungsbereich fällt. Diese Wahrscheinlichkeit nennen wir **Signifikanzniveau**, bezeichnet als α . Wir *konstruieren* Ablehnungsregionen so, dass α kleine Werte wie .05 oder .01 annimmt.

Die Wahrscheinlichkeit, dass ein Test eine falsche H_0 ablehnt, heißt **Macht**. Die Macht ist abhängig von der Datengenerierung und von n . So steigt die Macht mit der Distanz $|\beta_1 - \beta_{H0}|$ in (4.2).

Eine falsche H_0 fälschlicherweise nicht abzulehnen wird als **Fehler zweiter Art** bezeichnet. Die Wahrscheinlichkeit dieses Fehlers ist $1 - \text{Macht}$.

Um eine Ablehnungsregion zum Niveau α zu konstruieren, berechnen wir zunächst den zu α gehörigen **kritischen Wert** c_α . Für eine $\mathcal{N}(0, 1)$ Teststatistik eines zweiseitigen Tests unter H_0 ist der kritische Wert implizit gegeben durch

$$1 - \alpha/2 = \Phi(c_\alpha) \quad (4.3)$$

wobei Φ die Verteilungsfunktion (engl. CDF) der $\mathcal{N}(0, 1)$ -Verteilung (1.4) bezeichnet. Daher ist

$$\Phi^{-1}(1 - \alpha/2) = c_\alpha,$$

wobei Φ^{-1} die **Quantilsfunktion** bezeichnet.

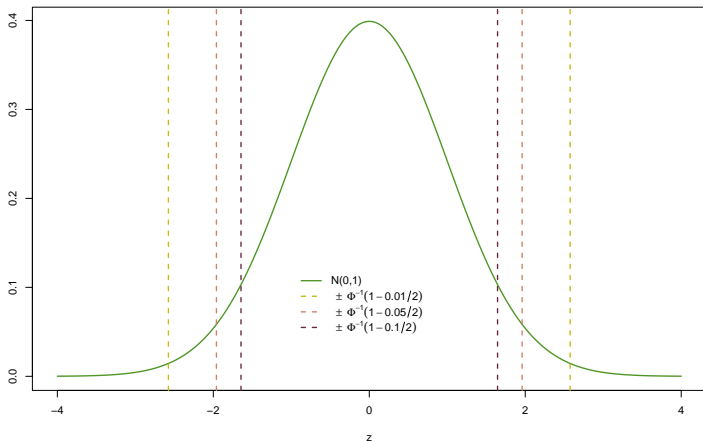
Die Wahrscheinlichkeit, dass $z > c_\alpha$, ist $1 - (1 - \alpha/2) = \alpha/2$ und durch Symmetrie auch $P(z < -c_\alpha) = \alpha/2$. Daher ist

$$P(|z| > c_\alpha) = \alpha,$$

wie gewünscht. Bei $\alpha = 0.05$ erhalten wir beispielsweise

$$\Phi^{-1}(1 - \alpha/2) = 1.96.$$

Graphisch können wir die jeweiligen Elemente auch von der Normalverteilungsdichte (1.3) ablesen.

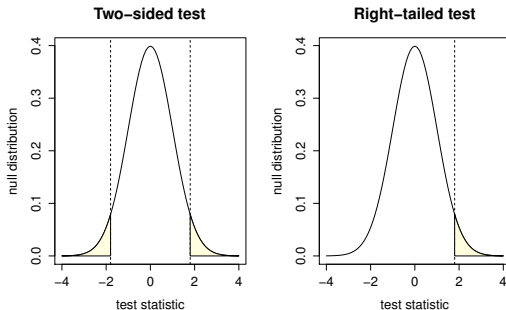


Das Ergebnis eines Tests ist Ja oder Nein: Ablehnen oder nicht ablehnen.

Ein ausgefeilterer Ansatz ist es den **p-Wert**, auch bezeichnet als das marginale Signifikanzniveau, der beobachteten Teststatistik \hat{z} zu berechnen.

Der p -Wert ist definiert als das kleinste Niveau, zu dem ein Test basierend auf \hat{z} gerade noch abgelehnt wird. Alternativ ist er die Wahrscheinlichkeit unter H_0 genau so oder extremere Statistiken als die beobachtete zu erhalten.

Siehe pValues.R.



Für einen zweiseitigen Test in unserem vorherigen Beispiel gilt

$$p(\hat{z}) = 2(1 - \Phi(|\hat{z}|))$$

Um dies zu sehen, erinnern wir uns, dass der Test basierend auf \hat{z} ablehnt, wenn

$$|\hat{z}| > c_\alpha$$

Dies ist äquivalent zu

$$\Phi(|\hat{z}|) > \Phi(c_\alpha),$$

da Φ strikt steigt.

Weiter gilt mit (4.3)

$$\Phi(c_\alpha) = 1 - \alpha/2$$

Der kleinste Wert α , für den obige Ungleichung erfüllt ist, wird über die Lösung der folgenden Gleichung nach α ermittelt:

$$\Phi(|\hat{z}|) = 1 - \alpha/2$$

Wir erhalten $2(1 - \Phi(|\hat{z}|))$.

Der p -Wert präsentiert die Information direkter. Beispielsweise führen Teststatistiken 2.02 und 5.77 beide zur Ablehnung bei $\alpha = 0.05$. 5.77 liefert jedoch eine „deutlichere“ Ablehnung der Nullhypothese als 2.02.

p -Werte verdeutlichen das Ausmaß der Differenz. Der p -Wert zu 2.02 ist .0434, während der p -Wert zu 5.77 gleich 7.93×10^{-9} und damit deutlich kleiner ist.

Dennoch sollte der p -value nicht überschätzt werden! Siehe hierfür [hier](#) und [hier](#) oder detaillierter [hier](#).

Oder wie andere gesagt haben:

„... surely, God loves the .06 nearly as much as the .05. Can there be any doubt that God views the strength of evidence for or against the null as a fairly continuous function of the magnitude of p ?“ —Rosnow, R. L., & Rosenthal, R. (1989). Statistical procedures and the justification of knowledge in psychological science. *American Psychologist*, 44(10), 1276-1284.

Ein verwandtes Problem:

<u>P-VALUE</u>	<u>INTERPRETATION</u>
0.001	} HIGHLY SIGNIFICANT
0.01	
0.02	
0.03	
0.04	} SIGNIFICANT
0.049	
0.050	} OH CRAP. REDO CALCULATIONS.
0.051	
0.06	
0.07	} ON THE EDGE OF SIGNIFICANCE
0.08	
0.09	
0.099	} HIGHLY SUGGESTIVE, SIGNIFICANT AT THE $p < 0.10$ LEVEL
≥ 0.1	
	HEY, LOOK AT THIS INTERESTING SUBGROUP ANALYSIS

Nun werden wir die starken Annahmen, dass Fehlerterme normalverteilt sind, und dass die Varianz der Fehlerterme bekannt ist, lockern. Zunächst beschäftigen wir uns mit der letzteren Annahme und kommen dann später auf erstere zurück.

Zusätzlich haben wir unsere Aufmerksamkeit auf eine einzige Restriktion bezüglich eines einzigen Parameters beschränkt.

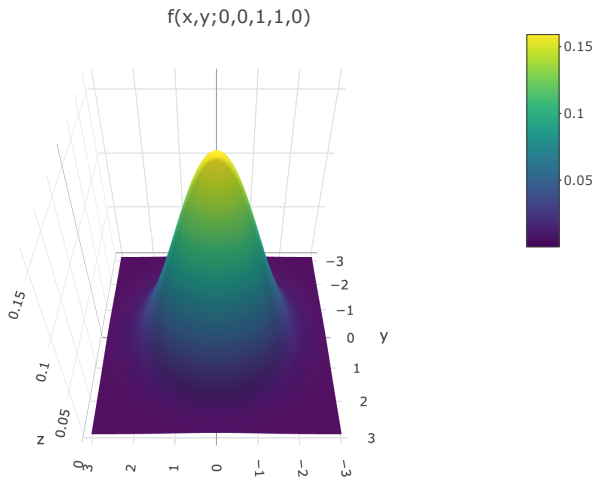
Wir werden uns einige Verteilungen ansehen müssen.

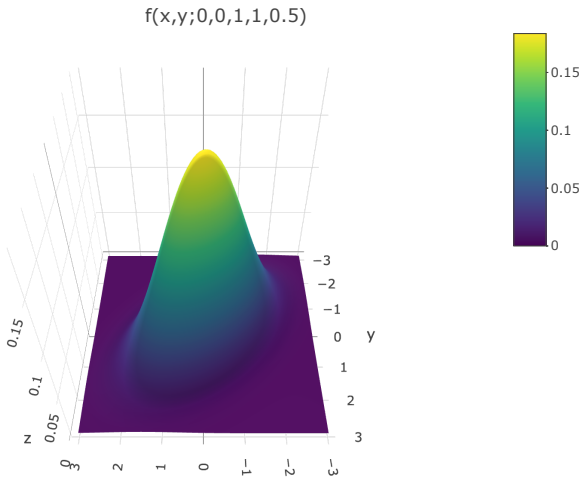
Die bivariate Normalverteilung

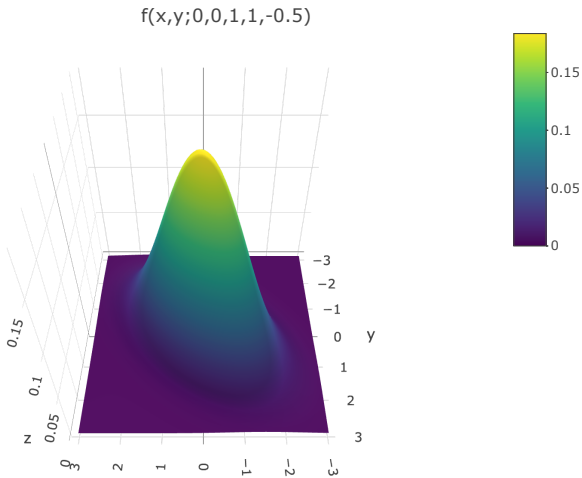
Definition 4.1 Ein Paar von Zufallsvariablen X und Y hat die **bivariate Normalverteilung**, wenn deren Dichte (1.5) gegeben ist durch

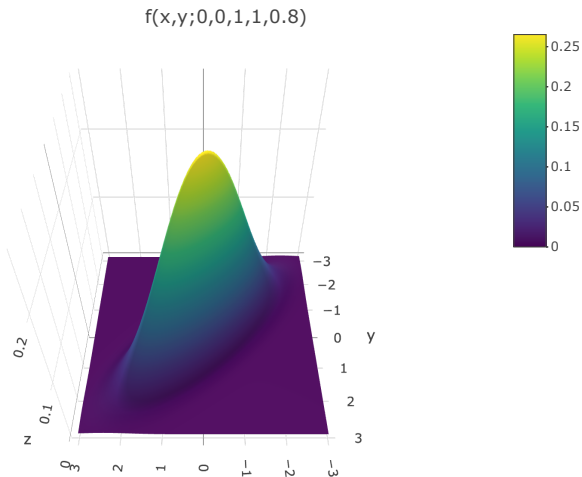
$$f(x, y; \mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left\{ -\frac{1}{2(1-\rho^2)} \times \right. \\ \left. \left[\left(\frac{x-\mu_1}{\sigma_1} \right)^2 + \left(\frac{y-\mu_2}{\sigma_2} \right)^2 - 2\rho \left(\frac{x-\mu_1}{\sigma_1} \right) \left(\frac{y-\mu_2}{\sigma_2} \right) \right] \right\}$$

mit $\sigma_1, \sigma_2 > 0$, $-1 < \rho < 1$ und $-\infty < x, y < \infty$.









Der Parameter ρ ist der **Korrelationskoeffizient** zwischen X und Y , das heißt

$$\rho = \frac{\text{Cov}(X, Y)}{\sigma_1 \sigma_2}$$

Die Notation μ_1, μ_2 genauso wie σ_1^2, σ_2^2 deutet marginale Erwartungswerte und Varianzen von X und Y an. Dies ist auch der Fall. Herausintegrieren von Y wie in (1.6) (was ein bisschen umständlich ist), d.h.

$$\int_{-\infty}^{\infty} f(x, y; \mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho) dy,$$

liefert die marginale Dichte von X ,

$$g(x) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left\{-\frac{1}{2\sigma_1^2}(x - \mu_1)^2\right\}, \quad (4.4)$$

aus der wir ablesen können, dass $E(X) = \mu_1$ und $\text{Var}(X) = \sigma_1^2$. Durch Symmetrie gilt das Gleiche für Y .

Hier finden Sie eine Heatmap mit einer bivariaten Normalverteilung: `BivariateNormal.R`. Siehe auch `Regression to the mean.R` für eine Anwendung.

Diese Ergebnisse können zur **multivariaten Normalverteilung** generalisiert werden.

Erinnern wir uns an (3.5) und betrachten zunächst einen Vektor aus Zufallsvariablen $\mathbf{z} = (z_1, \dots, z_m)^\top$, wobei jedes $z_i \sim \mathcal{N}(0, 1)$ unabhängig von den anderen z_j sei. Dies bedeutet, dass $\text{Var}(\mathbf{z}) = \mathbf{I}$. Transformieren wir \mathbf{z} mit der $m \times m$ Matrix \mathbf{A} mit vollem Rang

$$\mathbf{x} = \mathbf{A}\mathbf{z},$$

erhalten wir

$$E(\mathbf{x}) = \mathbf{A}E(\mathbf{z}) = \mathbf{0}$$

Des Weiteren gilt

$$\begin{aligned} \text{Var}(\mathbf{x}) &\stackrel{(3.6)}{=} \mathbf{A}E(\mathbf{z}\mathbf{z}^\top)\mathbf{A}^\top \\ &= \mathbf{A}\text{Var}(\mathbf{z})\mathbf{A}^\top \\ &= \mathbf{A}\mathbf{A}^\top \\ &=: \mathbf{\Omega} \end{aligned}$$

Studentsche t -, F - and Chi-Quadrat-Verteilungen

Verwandt zu der Normalverteilung haben wir:

- Die Chi-Quadrat-Verteilung:

Wenn $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$, dann

$$\mathbf{z}^\top \mathbf{z} \sim \chi_n^2, \quad (4.5)$$

- Die Studentsche t -Verteilung:

Wenn $z \sim \mathcal{N}(0, 1)$ und davon unabhängig $y \sim \chi_n^2$, dann

$$\frac{z}{\sqrt{y/n}} \sim t_n, \quad (4.6)$$

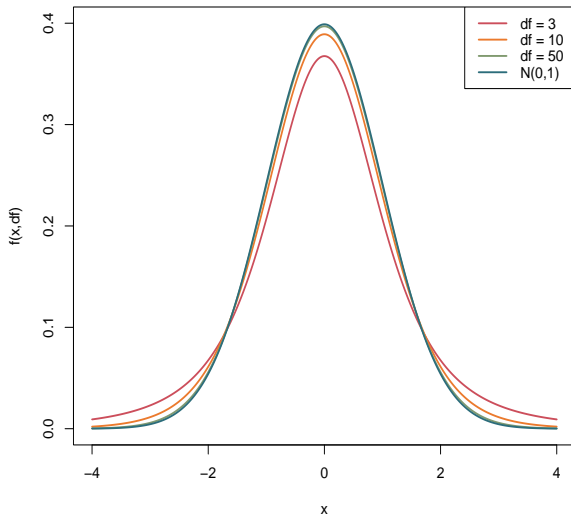
- (Fishers) F -Verteilung

Wenn $y_1 \sim \chi_{n_1}^2$ und unabhängig, $y_2 \sim \chi_{n_2}^2$, dann

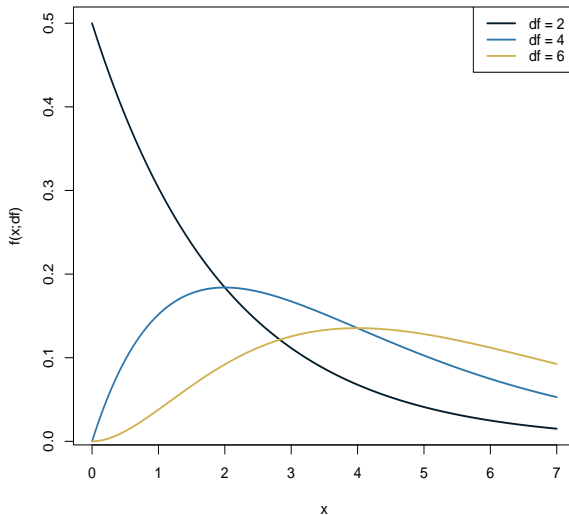
$$\frac{y_1/n_1}{y_2/n_2} \sim F_{n_1, n_2}. \quad (4.7)$$

Hier sind ein paar Bilder.

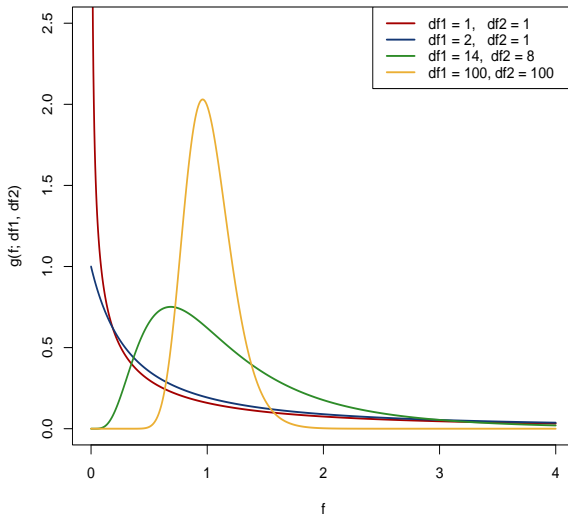
Die t -Verteilung:



Die χ^2 -Verteilung:



Die F -Verteilung:



Wir betrachten weiter das lineare Modell (3.1)

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}.$$

Wir erweitern nun Annahme 3.2 zu

$$\mathbf{u}|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n). \quad (4.8)$$

In diesem Fall überzeugen wir uns, dass

$$\begin{aligned} \mathbf{y}|\mathbf{X} &\sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n), \\ \hat{\boldsymbol{\beta}}_{\text{OLS}}|\mathbf{X} &\sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}). \end{aligned} \quad (4.9)$$

Um Tests ohne die Annahme, dass σ^2 bekannt ist, entwickeln zu können, brauchen wir

Resultat 4.2 Die Zufallsvektoren

$$(\mathbf{I}_n - \mathbf{P}_\mathbf{X})\mathbf{u} \text{ und } \mathbf{X}^\top \mathbf{u} \quad (4.10)$$

sind bedingt auf \mathbf{X} unabhängig und gemeinsam normalverteilt.

Beweis von Resultat 4.2

$$\begin{pmatrix} I_n - P_X \\ X^\top \end{pmatrix} u | X \quad (4.11)$$

hat eine multivariate Normalverteilung, da die Elemente lineare Funktionen von u sind. Der Erwartungswert ist null und die Kovarianzmatrix block-diagonal, da

$$E((I_n - P_X)uu^\top X | X) = \sigma^2(I_n - P_X)X \stackrel{(2.6)}{=} \mathbf{0}. \quad (4.12)$$

Die zwei Terme in (4.10) sind unabhängig, da sie eine multivariate Normalverteilung und eine Korrelation von null haben. q.e.d.

Erinnern wir uns, dass

$$y - X\hat{\beta}_{OLS} = (I_n - P_X)y \stackrel{(2.6)}{=} (I_n - P_X)u.$$

Daraus folgt, dass

$$\hat{\sigma}^2 = \frac{(y - X\hat{\beta}_{OLS})^\top (y - X\hat{\beta}_{OLS})}{n - k}$$

aus (3.8) eine Funktion von $(I_n - P_X)u$ ist.

Des Weiteren ist $\hat{\beta}_{\text{OLS}} = \beta + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}$ eine Funktion von $\mathbf{X}^\top \mathbf{u}$. Daher sind $\hat{\beta}_{\text{OLS}}$ und $\hat{\sigma}^2$ unabhängig. Außerdem ist

$$(n - k) \frac{\hat{\sigma}^2}{\sigma^2} \bigg| \mathbf{X} \sim \chi_{n-k}^2 \quad (4.13)$$

Warum? Beachten Sie, dass durch Symmetrie und Idempotenz gilt

$$(n - k) \hat{\sigma}^2 = \mathbf{u}^\top (\mathbf{I}_n - \mathbf{P}_\mathbf{X})^\top (\mathbf{I}_n - \mathbf{P}_\mathbf{X}) \mathbf{u} = \mathbf{u}^\top (\mathbf{I}_n - \mathbf{P}_\mathbf{X}) \mathbf{u}$$

Dividieren durch σ^2 liefert

$$(n - k) \frac{\hat{\sigma}^2}{\sigma^2} = \frac{\mathbf{u}^\top}{\sigma} (\mathbf{I}_n - \mathbf{P}_\mathbf{X}) \frac{\mathbf{u}}{\sigma},$$

wobei (4.8) impliziert, dass $\frac{\mathbf{u}}{\sigma} | \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Dies ist eine quadratische Form (siehe Definition 14 im Matrixalgebrareader) von Standardnormal-Zufallsvektoren in einer idempotenten Matrix des Rangs $n - k$, siehe (3.9) und Resultat 26.2.

Warum ist das interessant? Theorem 4.1 des Buches sagt uns, dass eine quadratische Form in einem standardnormalverteilten n -Vektor und einer idempotenten Matrix \mathbf{J} des Rangs m χ_m^2 ist. Daher gilt (4.13).

Ein alternativer Beweis nutzt die Spektralzerlegung (Resultat 29):

$$\mathbf{x}^\top \mathbf{J} \mathbf{x} = \mathbf{x}^\top \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top \mathbf{x}$$

Wir wissen, dass die Spur die Summe der Eigenwerte ist (Resultat 27.8), welche entweder eins oder null (Resultat 27.3) sind.

Daher hat $\mathbf{\Lambda}$ m Eigenwerte, die gleich 1 sind und $n - m$, die gleich 0 sind. Da \mathbf{U} orthogonal ist, gilt

$$E(\mathbf{U}^\top \mathbf{x} \mathbf{x}^\top \mathbf{U}) = \mathbf{U}^\top \mathbf{U} = \mathbf{I}$$

so dass auch

$$\mathbf{z} := \mathbf{U}^\top \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

Dann gilt

$$\mathbf{z}^\top \mathbf{\Lambda} \mathbf{z} = \sum_{i=1}^m z_i^2 \sim \chi_m^2$$

Wir erhalten insgesamt

(i) (4.9): $\hat{\beta}_{\text{OLS}} - \beta | \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}),$

(ii) (4.13): $(n - k) \frac{\hat{\sigma}^2}{\sigma^2} \Big| \mathbf{X} \sim \chi_{n-k}^2,$

(iii) Resultat 4.2: die Terme (i) und (ii) sind unabhängig.

Diese Bausteine ermöglichen es uns t -, F - and Chi-Quadrat-Teststatistiken zu konstruieren.

t -Tests

Betrachten wir eine einzige lineare Restriktion

$$H_0 : \mathbf{a}^\top \boldsymbol{\beta} = r \quad (4.14)$$

wobei $\mathbf{a} : k \times 1$ ein bekannter Vektor und r ein bekanntes Skalar ist.

Der häufigste Fall einer Nullhypothese ist $H_0 : \beta_i = 0$, was dem Standardoutput der meisten Regressionspakete entspricht.

Wir haben

$$\mathbf{a}^\top (\hat{\boldsymbol{\beta}}_{\text{OLS}} - \boldsymbol{\beta}) | \mathbf{X} \sim \mathcal{N}(0, \sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}).$$

Daher ist

$$\frac{\mathbf{a}^\top (\hat{\boldsymbol{\beta}}_{\text{OLS}} - \boldsymbol{\beta})}{\sqrt{\sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}} \sim \mathcal{N}(0, 1).$$

Dieser Ausdruck ist (Resultat 4.2) unabhängig von

$$(n - k) \frac{\hat{\sigma}^2}{\sigma^2} \Big| \mathbf{X} \sim \chi_{n-k}^2.$$

Daher gilt

$$\underbrace{\frac{\mathbf{a}^\top (\hat{\beta}_{\text{OLS}} - \beta)}{\sqrt{\sigma^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}}}_{\sim \mathcal{N}(0,1)} \Bigg/ \underbrace{\sqrt{\frac{\hat{\sigma}^2}{\sigma^2}}}_{\substack{\text{bed. auf } \mathbf{X} \\ \sim \sqrt{\chi_{n-k}^2/(n-k)}}} = \frac{\mathbf{a}^\top (\hat{\beta}_{\text{OLS}} - \beta)}{\sqrt{\hat{\sigma}^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}} \stackrel{(4.6)}{\sim} t_{n-k}. \quad (4.15)$$

Mit (4.14), d.h. $H_0 : \mathbf{a}^\top \beta = r$, haben wir eine Teststatistik (d.h. etwas, das nicht von etwas Unbekannten abhängt)

$$t = \frac{\mathbf{a}^\top \hat{\beta}_{\text{OLS}} - r}{\sqrt{\hat{\sigma}^2 \mathbf{a}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{a}}}. \quad (4.16)$$

t hat unter der Nullhypothese eine t_{n-k} -Verteilung.

t -Werte werden standardmäßig für $H_0 : \beta_i = 0$ ausgegeben. Diesen Fall erhalten wir aus (4.16) für $\mathbf{a} = \mathbf{e}_i$ (einen Vektor mit einer 1 an der Position i und sonst null) und $r = 0$. Daher ist der t -Wert für $H_0 : \beta_i = 0$ gegeben durch

$$t = \frac{\hat{\beta}_i - 0}{\text{se}(\hat{\beta}_i)} = \frac{\hat{\beta}_i}{\text{se}(\hat{\beta}_i)}, \quad (4.17)$$

wobei der Standardfehler von $\hat{\beta}_i$ gegeben ist durch

$$\text{se}(\hat{\beta}_i) = \sqrt{\widehat{\text{Var}}(\hat{\beta}_i)} = \sqrt{\hat{\sigma}^2 \mathbf{e}_i^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{e}_i} = \sqrt{\hat{\sigma}^2 (\mathbf{X}^\top \mathbf{X})_{ii}^{-1}}.$$

Siehe `Testing.R`.

F-Test

Oft wollen wir mehrere Restriktionen auf einmal testen.

Sei $H_0 : \mathbf{R}^\top \boldsymbol{\beta} = \mathbf{r}$ mit \mathbf{R} und \mathbf{r} bekannt, nicht zufällig und $\mathbf{R} : k \times q$ hat vollen Spaltenrang q . Hiermit werden q lineare Restriktionen abgebildet. Ein Beispiel findet sich unten in (4.22).

Aus (4.9) und (3.6) erhalten wir

$$\mathbf{R}^\top (\hat{\boldsymbol{\beta}}_{\text{OLS}} - \boldsymbol{\beta}) | \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{R}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R})$$

Definiere

$$\mathbf{B} = \mathbf{R}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}$$

und sei $\mathbf{B}^{-1/2}$ eine „Matrixwurzel“ von \mathbf{B}^{-1} (siehe z.B. Resultat 30). Für

$$\mathbf{g} := \frac{\mathbf{B}^{-1/2}}{\sigma} \mathbf{R}^\top (\hat{\beta}_{\text{OLS}} - \beta)$$

gilt dann, dass

$$\mathbf{g} | \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_q), \quad (4.18)$$

da

$$\begin{aligned} \text{Var}(\mathbf{g} | \mathbf{X}) &= \frac{\mathbf{B}^{-1/2}}{\sigma} \mathbf{R}^\top \text{Var}(\hat{\beta}_{\text{OLS}} | \mathbf{X}) \mathbf{R} \frac{\mathbf{B}^{-1/2}}{\sigma} \\ &\stackrel{(3.7)}{=} \frac{\mathbf{B}^{-1/2}}{\sigma} \sigma^2 \mathbf{B} \frac{\mathbf{B}^{-1/2}}{\sigma} = \mathbf{I} \end{aligned} \quad (4.19)$$

Dies ist unabhängig von

$$d := (n - k) \frac{\hat{\sigma}^2}{\sigma^2}, \quad d | \mathbf{X} \stackrel{(4.13)}{\sim} \chi_{n-k}^2.$$

Daher gilt

$$\stackrel{\sim \chi_q^2, \text{vgl. (4.5)}}{\underbrace{\mathbf{g}^\top \mathbf{g}}_{/q}} = \frac{(\hat{\beta}_{\text{OLS}} - \beta)^\top \mathbf{R} \{ \mathbf{R}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R} \}^{-1} \mathbf{R}^\top (\hat{\beta}_{\text{OLS}} - \beta) / q}{d / (n - k)} \stackrel{\sim}{=} \frac{(\hat{\beta}_{\text{OLS}} - \beta)^\top \mathbf{R} \{ \mathbf{R}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R} \}^{-1} \mathbf{R}^\top (\hat{\beta}_{\text{OLS}} - \beta) / q}{\hat{\sigma}^2} \sim F_{q, n-k}.$$

D.h., unter $H_0 : \mathbf{R}^\top \beta = \mathbf{r}$ erhalten wir

$$F = \frac{(\mathbf{R}^\top \hat{\beta}_{\text{OLS}} - \mathbf{r})^\top \{ \mathbf{R}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R} \}^{-1} (\mathbf{R}^\top \hat{\beta}_{\text{OLS}} - \mathbf{r}) / q}{\hat{\sigma}^2} \sim F_{q, n-k}. \quad (4.20)$$

Zur Illustration betrachten wir den Spezialfall $\mathbf{R}^\top = \mathbf{I}$, $\mathbf{r} = \mathbf{0}$, $q = 2$, $\hat{\sigma}^2 = 1$ und $\mathbf{X}^\top \mathbf{X} = \mathbf{I}$. Dann gilt

$$F = \hat{\beta}_{\text{OLS}}^\top \hat{\beta}_{\text{OLS}} / 2 = \frac{\hat{\beta}_{\text{OLS},1}^2 + \hat{\beta}_{\text{OLS},2}^2}{2}, \quad (4.21)$$

dies ist die quadrierte euklidische Distanz zwischen dem OLS-Schätzer und dem Ursprung, standardisiert durch die Anzahl der Elemente.

Der F -Test kann in speziellen Fällen eine aufschlussreichere Form annehmen.
Betrachten wir den Test

$$H_0 : \beta_2 = \mathbf{0} \quad (4.22)$$

in

$$\mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \mathbf{u} \quad (4.23)$$

Wir zeigen nun, dass die F -Statistik sich auch durch den intuitiven Ausdruck

$$F_{\beta_2} = \frac{(\text{RSSR} - \text{USSR})/r}{\text{USSR}/(n - k)} \quad (4.24)$$

darstellen lässt. Hierbei ist

- r die Spaltendimension von \mathbf{X}_2 ,
- $\text{USSR} = \mathbf{y}^\top \mathbf{M}_{\mathbf{X}} \mathbf{y}$ die Summe der quadrierten Residuen einer vollen Regression (4.23),
- $\text{RSSR} = \mathbf{y}^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{y}$ die Summe der quadrierten Residuen einer Regression von \mathbf{y} auf \mathbf{X}_1 , d.h. unter Auferlegung von H_0 .

Diese Nullhypothese entspricht den Restriktionen $\mathbf{R}^\top = [\mathbf{O} \ \mathbf{I}]$ und $\mathbf{r} = \mathbf{0}$.

Definiere den partitionierten KQ-Schätzer $\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\hat{\boldsymbol{\beta}}_{\text{OLS},1}^\top, \hat{\boldsymbol{\beta}}_{\text{OLS},2}^\top)^\top$.

Dann folgt

$$\mathbf{R}^\top \hat{\boldsymbol{\beta}}_{\text{OLS}} = \hat{\boldsymbol{\beta}}_{\text{OLS},2}$$

und

$$\mathbf{R}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R} =: \tilde{\mathbf{D}},$$

dem südöstlichen Block von

$$\begin{aligned} (\mathbf{X}^\top \mathbf{X})^{-1} &= \begin{pmatrix} \mathbf{X}_1^\top \mathbf{X}_1 & \mathbf{X}_1^\top \mathbf{X}_2 \\ \mathbf{X}_2^\top \mathbf{X}_1 & \mathbf{X}_2^\top \mathbf{X}_2 \end{pmatrix}^{-1} \\ &=: \begin{pmatrix} \tilde{\mathbf{A}} & \tilde{\mathbf{B}} \\ \tilde{\mathbf{C}} & \tilde{\mathbf{D}} \end{pmatrix} \end{aligned}$$

Ergebnisse für partitionierte Inversen (vgl. Matrixreader, Resultat 23) liefern

$$\tilde{\mathbf{D}} = (\mathbf{X}_2^\top \mathbf{X}_2 - \mathbf{X}_2^\top \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2)^{-1} = (\mathbf{X}_2^\top \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2)^{-1}$$

wobei $\mathbf{M}_{\mathbf{X}_1} = \mathbf{I} - \mathbf{X}_1 (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top$.

Der Zähler der F -Statistik ist also (ohne Teilen durch r)

$$F_{\text{num}} = \hat{\beta}_{\text{OLS},2}^{\top} (\mathbf{X}_2^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2) \hat{\beta}_{\text{OLS},2}$$

Das FWL-Theorem 2.4 liefert

$$\hat{\beta}_{\text{OLS},2} \stackrel{(2.15)}{=} (\mathbf{X}_2^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2)^{-1} \mathbf{X}_2^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{y},$$

so dass

$$\begin{aligned} F_{\text{num}} &= \mathbf{y}^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2 (\mathbf{X}_2^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2)^{-1} (\mathbf{X}_2^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2) (\mathbf{X}_2^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2)^{-1} \mathbf{X}_2^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{y} \\ &= \mathbf{y}^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2 (\mathbf{X}_2^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2)^{-1} \mathbf{X}_2^{\top} \mathbf{M}_{\mathbf{X}_1} \mathbf{y} \end{aligned} \quad (4.25)$$

Wir zeigen nun, dass $F_{\text{num}} = \text{RSSR} - \text{USSR}$. Mit FWL (das auch zeigt, dass die Residuen der beiden Vorgehensweisen identisch sind) können wir USSR als die SSR der Regression

$$M_{X_1}y \quad \text{auf} \quad M_{X_1}X_2$$

schreiben. D.h.,

$$\begin{aligned} \text{USSR} &= y^\top M_{X_1}^\top M_{M_{X_1}X_2} M_{X_1}y \\ &= y^\top M_{X_1}^\top (I - P_{M_{X_1}X_2}) M_{X_1}y \\ &= y^\top M_{X_1}y - y^\top M_{X_1} M_{X_1}X_2 ((M_{X_1}X_2)^\top M_{X_1}X_2)^{-1} (M_{X_1}X_2)^\top M_{X_1}y \\ &= y^\top M_{X_1}y - y^\top M_{X_1}X_2 (X_2^\top M_{X_1}X_2)^{-1} X_2^\top M_{X_1}y \end{aligned}$$

Also ist $\text{RSSR} - \text{USSR}$ gleich (4.25):

$$\begin{aligned} \text{RSSR} - \text{USSR} &= y^\top M_{X_1}y - (y^\top M_{X_1}y - y^\top M_{X_1}X_2 (X_2^\top M_{X_1}X_2)^{-1} X_2^\top M_{X_1}y) \\ &= y^\top M_{X_1}X_2 (X_2^\top M_{X_1}X_2)^{-1} X_2^\top M_{X_1}y \end{aligned}$$

Der Nenner von (4.20) ist offenbar

$$\text{USSR}/(n - k).$$

SimFstat.R

Wenn $\mathbf{X}_1 = \iota$, wir also in (4.22) testen, ob alle Steigungskoeffizienten gleich null sind, ist

$$\text{RSSR} = \mathbf{y}^\top \mathbf{M}_\iota \mathbf{y} \quad (4.26)$$

Damit gilt dann wegen (2.19), dass

$$R^2 = 1 - \frac{\text{USSR}}{\text{RSSR}}$$

oder

$$\text{USSR} = (1 - R^2)\text{RSSR}$$

und daher

$$\frac{\text{RSSR} - \text{USSR}}{\text{USSR}} = \frac{R^2}{1 - R^2}$$

Chow-Test

Der Chow-Test auf Strukturbruch testet

$$H_0 : \beta_1 = \beta_2 \quad (4.27)$$

für zwei unabhängige Stichproben $\mathbf{y}_1 = \mathbf{X}_1\beta_1 + \mathbf{u}_1$, and $\mathbf{y}_2 = \mathbf{X}_2\beta_2 + \mathbf{u}_2$, wobei $\mathbf{u}_1|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{n_1})$ und $\mathbf{u}_2|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{n_2})$.

Beachten Sie, dass wir die Störtermvarianz als gleich voraussetzen und die Zahl der Spalten von \mathbf{X}_1 und \mathbf{X}_2 gleich ist, da anderenfalls H_0 nicht wahr sein kann. Siehe Chow.R.

Wenn \mathbf{X}_1 und \mathbf{X}_2 vollen Spaltenrang k haben, können wir einen Test für H_0 formulieren, der auf dem klassischen F -Test basiert.

Wir reduzieren dazu ein System mehrerer Gleichungen auf ein Modell mit nur einer Gleichung.

Im gegebenen Fall haben wir als unrestringiertes Modell

$$\mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{X}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_2 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} = \mathbf{X}\beta + \mathbf{u}, \quad (4.28)$$

wobei \mathbf{X} , β und \mathbf{u} implizit definiert sind, mit $\mathbf{u}|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{n_1+n_2})$,
 \mathbf{X} : $(n_1 + n_2) \times 2k$ und β : $2k \times 1$.

Die daraus folgende Summe der quadrierten Residuen des unrestringierten Modells ist dabei

$$\hat{\mathbf{u}}_{\text{OLS}}^\top \hat{\mathbf{u}}_{\text{OLS}} = \hat{\mathbf{u}}_1^\top \hat{\mathbf{u}}_1 + \hat{\mathbf{u}}_2^\top \hat{\mathbf{u}}_2 \quad (4.29)$$

wobei $\hat{\mathbf{u}}_1$ und $\hat{\mathbf{u}}_2$ die OLS-Residuen der Regressionen $\mathbf{y}_i = \mathbf{X}_i \hat{\beta}_i + \hat{\mathbf{u}}_i$ sind, mit

$$\hat{\beta}_i = (\mathbf{X}_i^\top \mathbf{X}_i)^{-1} \mathbf{X}_i^\top \mathbf{y}_i$$

für $i = 1, 2$.

Das restringierte Modell ist

$$\mathbf{y} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix} \beta + \mathbf{u}. \quad (4.30)$$

Die Residuen des restringierten Modells $\hat{\mathbf{u}}_{\text{RLS}}$ erhalten wir durch eine Regression von \mathbf{y} auf $(\mathbf{X}_1^\top, \mathbf{X}_2^\top)^\top$, das k Regressoren hat. Die SSR der hieraus resultierenden OLS Residuen ist $\hat{\mathbf{u}}_{\text{RLS}}^\top \hat{\mathbf{u}}_{\text{RLS}}$.

Die Chow-Teststatistik ist nun (4.24), die Version des F -Tests, in welcher die SSR des restringierten und des unrestringierten Modells verglichen werden:

$$F_{\text{Chow}} = \frac{(\hat{\mathbf{u}}_{\text{RLS}}^\top \hat{\mathbf{u}}_{\text{RLS}} - \hat{\mathbf{u}}_1^\top \hat{\mathbf{u}}_1 - \hat{\mathbf{u}}_2^\top \hat{\mathbf{u}}_2)/k}{(\hat{\mathbf{u}}_1^\top \hat{\mathbf{u}}_1 + \hat{\mathbf{u}}_2^\top \hat{\mathbf{u}}_2)/(n_1 + n_2 - 2k)} \sim F_{k, n_1 + n_2 - 2k}. \quad (4.31)$$

okcupid.R

Konfidenzintervalle

Ein Konfidenzintervall zum Konfidenzniveau $(1 - \alpha)100\%$ ist die Menge von Parameterwerten θ^* für welche die Nullhypothese $H_0 : \theta = \theta^*$ zum Niveau α nicht abgelehnt würde. Nehmen wir an, dass H_0 abgelehnt wird, wenn

$$\tau(\mathbf{y}, \theta^*) \in \mathcal{A}(\mathbf{y}, \theta^*).$$

Das $(1 - \alpha)100\%$ Konfidenzintervall für θ ist dann

$$\{\theta^* | \tau(\mathbf{y}, \theta^*) \notin \mathcal{A}(\mathbf{y}, \theta^*)\}.$$

In der Praxis werden Konfidenzintervalle oft als eine Spanne „plausibler Werte“ für einen Parameter interpretiert.

Der einfachste Weg ein Konfidenzintervall zu konstruieren ist über ein **Pivot**, eine Funktion von Beobachtungen und Parametern, deren *Verteilung* nicht von den unbekannten Parametern abhängt.

Beispiel 4.3 Betrachte eine Stichprobe y_1, \dots, y_n aus einer $\mathcal{N}(\mu, 1)$ -Verteilung. Wir wissen, dass $\bar{Y} - \mu \sim \mathcal{N}(0, 1/n)$. Es sei wieder c_α das $\mathcal{N}(0, 1)$ -($1 - \alpha/2$)-Quantil. Daher gilt

$$\begin{aligned}
 1 - \alpha &= P\{(\bar{Y} - \mu)/(1/\sqrt{n}) \in (-c_\alpha, c_\alpha)\} \\
 &= P\{-c_\alpha \leq (\bar{Y} - \mu)\sqrt{n} \leq c_\alpha\} \\
 &= P\{c_\alpha \geq (\mu - \bar{Y})\sqrt{n} \geq -c_\alpha\} \\
 &= P\{-c_\alpha/\sqrt{n} \leq \mu - \bar{Y} \leq c_\alpha/\sqrt{n}\} \\
 &= P\{\bar{Y} - c_\alpha/\sqrt{n} \leq \mu \leq \bar{Y} + c_\alpha/\sqrt{n}\} \\
 &= P\{(\bar{Y} - c_\alpha/\sqrt{n}, \bar{Y} + c_\alpha/\sqrt{n}) \ni \mu\}
 \end{aligned}$$

ist ein Konfidenzintervall zum Niveau $1 - \alpha$. Die *Breite* des Intervalls hängt hier nicht von den Daten ab, sondern nur von α und n . Dies kommt daher, dass wir σ^2 als bekannt angenommen haben (und der Einfachheit halber gleich eins).

Setzen Sie einen Randwert des Konfidenzintervalls ein, um zu überprüfen, dass er tatsächlich ein μ^* repräsentiert, für welches gerade $H_0 : \mu = \mu^*$ nicht abgelehnt werden würde!

Es ist wichtig zu verstehen, was ein Konfidenzintervall aussagt und was nicht:

Das Intervall überdeckt den wahren Wert mit Wahrscheinlichkeit $1 - \alpha$. Was bedeutet diese Wahrscheinlichkeit?

Die Wahrscheinlichkeit ist eine Wahrscheinlichkeit über wiederholte Stichproben.

Sie besagt nicht, dass der wahre Parameter für ein konkret berechnetes Intervall mit Wahrscheinlichkeit 95% in dem Intervall liegt. Das tut er oder er tut es nicht. Welcher Fall eingetreten ist, wissen wir aber leider nicht, ähnlich wie wir auch nicht wissen, ob eine konkrete Verwerfung eines Tests eine korrekte Entscheidung oder aber ein Fehler erster Art ist.

Ein 95%-Konfidenzintervall überdeckt demnach den wahren Wert in etwa 95% der Fälle. Siehe `ConfidenceIntervals.R`.

Konfidenzintervalle

Betrachte etwas allgemeiner ein Pivot für einen Regressionskoeffizienten β_i , siehe (4.17)

$$\tau(\mathbf{y}, \beta_i) = \frac{\hat{\beta}_i - \beta_i}{\text{se}(\hat{\beta}_i)}.$$

Mit dem $1 - \alpha/2$ -Quantil der t_{n-k} -Verteilung $t_{n-k,\alpha}$ erhalten wir aus (4.15) das Konfidenzintervall

$$\{\beta_i \mid -t_{n-k,\alpha} \leq \tau(\mathbf{y}, \beta_i) \leq t_{n-k,\alpha}\} = \left[\hat{\beta}_i - \text{se}(\hat{\beta}_i)t_{n-k,\alpha}, \hat{\beta}_i + \text{se}(\hat{\beta}_i)t_{n-k,\alpha} \right].$$

Wir nennen diese Prozedur die **Invertierung von t -Tests**.

„Everybody believes in the [Normal] distribution: the experimenters, because they think it can be proved by mathematics; and the mathematicians, because they believe it has been established by observation.“

—Whittaker, E. T. and Robinson, G. „Normal Frequency Distribution.“
Ch. 8 in The Calculus of Observations: A Treatise on Numerical Mathematics, 4th ed. New York: Dover, pp. 164-208, 1967. p. 179.

Asymptotische Theorie

Wir beschäftigen uns nun mit der Lockerung der Annahme von normalverteilten Fehlern. Dafür wiederholen wir einige Konzepte der asymptotischen Theorie.

Für eine Folge von deterministischen Variablen gilt $(x_i)_{i=1}^{\infty} = x_1, x_2, x_3, \dots$, $x_n \rightarrow a$ oder $\lim_{n \rightarrow \infty} x_n = a$, wenn x_n beliebig nah an a für ausreichend große n liegt.

Bei Zufallsvariablen/-vektoren müssen wir hingegen verschiedene Arten von Konvergenz unterscheiden:

- (i) Konvergenz in quadratischem Mittel (q.m.)
- (ii) Konvergenz in Wahrscheinlichkeit (plim)
- (iii) Konvergenz in Verteilung.

Wir formulieren Gesetze der Großen Zahlen (Laws of Large Numbers/LLN) und zentrale Grenzwertsätze (Central Limit Theorems/CLT).

1. **(schwaches) Gesetz der großen Zahlen:** Für unabhängige und identisch verteilte (u.i.v.) Zufallsvariablen X_1, \dots, X_n gilt

$$\bar{X}_n = \frac{\sum_{i=1}^n X_i}{n} \xrightarrow{p} \mu = E(X_i), \quad (4.32)$$

vorausgesetzt dieser Erwartungswert existiert. ShinyApp:WLLN

Was bedeutet das? Das LLN ist ein Beispiel für **Konvergenz in Wahrscheinlichkeit**. Betrachten wir eine Folge von Zufallsvariablen Y_1, Y_2, \dots

2. **Konvergenz in Wahrscheinlichkeit:** Wir sagen, dass die Folge $\{Y_n\}$ in Wahrscheinlichkeit gegen c konvergiert, wenn für alle $\epsilon > 0$

$$\lim_{n \rightarrow \infty} P(|Y_n - c| < \epsilon) = 1. \quad (4.33)$$

Notation: $Y_n \xrightarrow{p} c$ oder $\text{plim}_{n \rightarrow \infty}(Y_n) = c$.

3. **Konvergenz im quadratischen Mittel:** Ein weiteres Konzept verwendet Momente (die nicht immer existieren): wenn $E(Y_n) \rightarrow c$ und $\text{Var}(Y_n) \rightarrow 0$, dann konvergiert $\{Y_n\}$ im quadratischen Mittel gegen c . Notation:

$$Y_n \xrightarrow{q.m.} c.$$

4. Wenn die Varianz existiert, gilt

$$Y_n \xrightarrow{q.m.} c \Rightarrow Y_n \xrightarrow{p} c \quad (4.34)$$

Dies ist oft der einfachste Weg, um 2. zu zeigen.

5. Für stetige Funktionen $g(\cdot)$, die nicht von n abhängen, gilt

$$\text{plim}_{n \rightarrow \infty} \{g(Y_n)\} = g(\text{plim}_{n \rightarrow \infty}(Y_n)) \quad (4.35)$$

6. **Konvergenz in Verteilung:** Wenn die Verteilungsfunktionen $F_{Y_n}(x)$ gegen die Verteilungsfunktion der Zufallsvariable X konvergiert, also $F_{Y_n}(x) \rightarrow F_X(x)$ für alle Stetigkeitsstellen x gilt, sagen wir

$$Y_n \xrightarrow{d} X, \quad \text{oder} \quad Y_n \overset{a}{\sim} F_X, \quad (4.36)$$

d.h. Y_n hat die asymptotische Verteilung F_X .

7. **Satz von Slutsky:** Nehmen wir an, es liegen Zufallsfolgen $X_n \xrightarrow{p} a$ und $Z_n \xrightarrow{d} Z$ vor, wobei a nicht zufällig ist. Dann gilt

$$X_n Z_n \xrightarrow{d} aZ, \quad (4.37)$$

$$X_n + Z_n \xrightarrow{d} a + Z. \quad (4.38)$$

Diese Resultate erlauben es asymptotische Verteilungen abzuleiten, die wiederum verwendet werden können, um Verteilungen in endlichen Stichproben zu approximieren.

Hierbei sind zentrale Grenzwertsätze zentraler Baustein.

8. **Zentraler Grenzwertsatz:** Für u.i.v. Beobachtungen X_1, \dots, X_n mit Erwartungswert μ und Varianz σ^2 gilt

$$n^{1/2}(\bar{X}_n - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2). \quad (4.39)$$

Anders ausgedrückt, gilt für $Y_n = n^{1/2}(\bar{X}_n - \mu)/\sigma$, dass

$$F_{Y_n}(x) \rightarrow \Phi(x)$$

für alle x , mit Φ der CDF der $\mathcal{N}(0, 1)$ -Verteilung (1.4).

Siehe clt.R.

Wir sagen oft etwas salopp, dass \bar{X}_n eine asymptotische Normalverteilung hat, obwohl eigentlich gilt $\bar{X}_n \xrightarrow{p} \mu$. Wir meinen, dass die Verteilung von

$$n^{1/2}(\bar{X}_n - \mu)$$

gegen eine Normalverteilung konvergiert.

9. Wenn $Y_n \xrightarrow{d} 0$, dann $Y_n \xrightarrow{p} 0$.

Wenn also, basierend auf Resultat 7 oben, $Y_n \xrightarrow{p} 0$ und $Z_n \xrightarrow{d} Z$, dann

$$Y_n Z_n \xrightarrow{p} 0. \quad (4.40)$$

Folglich gilt für

$$n^{1/2}(\bar{X}_n - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2),$$

dass $\text{plim}_{n \rightarrow \infty} \bar{X}_n = \mu$.

Dies ist der Fall, da $n^{-1/2} \rightarrow 0$ (was gelesen werden kann als $n^{-1/2} \xrightarrow{p} 0$) und $n^{1/2}(\bar{X}_n - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ daher gilt

$$\bar{X}_n - \mu = n^{-1/2} \left\{ n^{1/2}(\bar{X}_n - \mu) \right\} \xrightarrow{p} 0, \quad (4.41)$$

10. Die obigen Resultate wurden für skalarwertige Zufallsvariablen formuliert, es gelten jedoch sehr ähnliche Resultate für Zufallsvektoren und -Matrizen.

Asymptotik für Regressionen

Erinnern wir uns an die Darstellung

$$\hat{\beta}_{\text{OLS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \stackrel{(3.3)}{=} \beta + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}.$$

Beachten Sie

$$\mathbf{X}^\top \mathbf{X} = \sum_{t=1}^n \mathbf{x}_t^\top \mathbf{x}_t,$$

wobei \mathbf{x}_t die t -te Zeile von \mathbf{X} ist. Also ist

$$\frac{\mathbf{X}^\top \mathbf{X}}{n} = \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t^\top \mathbf{x}_t,$$

ein Mittelwert.

Wir treffen die folgenden Annahmen.

Annahme 4.4 Die Folge $\{y_t, \mathbf{X}_t\}$ ist gemeinsam u.i.v.

Annahme 4.4 impliziert mit dem Gesetz der großen Zahlen (vgl. (4.32)), dass

$$\frac{\mathbf{X}^\top \mathbf{X}}{n} \xrightarrow{p} E(\mathbf{X}_t^\top \mathbf{X}_t) =: \mathbf{S}_{\mathbf{X}^\top \mathbf{X}} > \mathbf{O}, \quad (4.42)$$

sofern die Momente existieren.

Analog ist

$$\frac{\mathbf{X}^\top \mathbf{u}}{n} = \frac{1}{n} \sum_{t=1}^n \mathbf{X}_t^\top u_t \quad (4.43)$$

ein Mittelwert. Wenn wir bei der Aufstellung unseres Modells einen guten Job gemacht haben, wird

Annahme 4.5

$$E(\mathbf{X}_t^\top u_t) = \mathbf{0} \quad (4.44)$$

erfüllt sein. Wie motiviert sich diese Annahme?

Unter Annahme 4.4 wird (3.2) zu $E(u_t | \mathbf{X}_t) = 0$ und wir erhalten durch das Gesetz der iterierten Erwartungen (LIE) (1.7)

$$\begin{aligned} E(\mathbf{X}_t^\top u_t) &= E(E(\mathbf{X}_t^\top u_t | \mathbf{X}_t)) \\ &= E(\mathbf{X}_t^\top E(u_t | \mathbf{X}_t)) \\ &= \mathbf{0} \end{aligned}$$

Die Annahme, dass $E(\mathbf{X}_t^\top u_t) = \mathbf{0}$ ist jedoch schwächer als $E(u_t | \mathbf{X}_t) = 0$!

Dann liefert Annahme 4.4 über das Gesetz der großen Zahlen (4.32), dass

$$\frac{\mathbf{X}^\top \mathbf{u}}{n} \xrightarrow{p} E(\mathbf{X}_t^\top u_t) = \mathbf{0}. \quad (4.45)$$

Mittels (4.37) und (4.38) erhalten wir

$$\begin{aligned} \text{plim}_{n \rightarrow \infty} \hat{\beta}_{\text{OLS}} &= \beta + \text{plim}_{n \rightarrow \infty} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u} \\ &= \beta + \text{plim}_{n \rightarrow \infty} \left(\frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \text{plim}_{n \rightarrow \infty} \left(\frac{\mathbf{X}^\top \mathbf{u}}{n} \right) \quad (4.46) \\ &= \beta + \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1} \mathbf{0} \\ &= \beta. \end{aligned}$$

Siehe ConsistencyOLS.R.

Um die asymptotische Verteilung herzuleiten, treffen wir zusätzlich

Annahme 4.6

$$\Omega = \text{Var}(\mathbf{X}_t^\top u_t) \quad (4.47)$$

existiert und ist endlich.

Unter (4.44) ist

$$\Omega = \text{Var}(\mathbf{X}_t^\top u_t) = E(u_t^2 \mathbf{X}_t^\top \mathbf{X}_t) \quad (4.48)$$

Der ZGWS (4.39) impliziert mit Annahme 4.4 und (4.44), dass

$$n^{1/2} \left(\frac{\mathbf{X}^\top \mathbf{u}}{n} \right) = n^{1/2} \left(\frac{\sum_{t=1}^n \mathbf{X}_t^\top u_t}{n} \right) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Omega) \quad (4.49)$$

Wir müssen nicht mehr annehmen, dass \mathbf{u} selbst normalverteilt ist!

Wir haben bisher keine zusätzliche Annahme getroffen, die besagen würde, dass

$$\Omega = \sigma^2 \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}. \quad (4.50)$$

(4.50) ergibt sich aus Annahmen 3.2, 4.4 und 4.5, da

$$\text{Var}(\mathbf{u}|\mathbf{X}) = \sigma^2 \mathbf{I}_n \stackrel{\text{u.i.v.}}{\Rightarrow} E(u_t^2|\mathbf{X}_t) = \sigma^2$$

und

$$\text{Var}(\mathbf{X}_t^\top u_t) = E(u_t^2 \mathbf{X}_t^\top \mathbf{X}_t) = E(E(u_t^2|\mathbf{X}_t) \mathbf{X}_t^\top \mathbf{X}_t) = \sigma^2 E(\mathbf{X}_t^\top \mathbf{X}_t) \quad (4.51)$$

Aus (4.37), (4.42) und (4.49) folgt

$$\begin{aligned} n^{1/2}(\hat{\beta}_{\text{OLS}} - \beta) &= \left(\frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \left(\frac{\mathbf{X}^\top \mathbf{u}}{n^{1/2}} \right) \\ &\xrightarrow{d} \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1} \mathcal{N}(\mathbf{0}, \Omega) \\ &= \mathcal{N}(\mathbf{0}, \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1} \Omega \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}), \end{aligned} \quad (4.52)$$

Wenn wir die Annahme bedingter Homoskedastie (4.50) treffen, erhalten wir

$$\begin{aligned} n^{1/2}(\hat{\beta}_{\text{OLS}} - \beta) &\xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1} \sigma^2 \mathbf{S}_{\mathbf{X}^\top \mathbf{X}} \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}) \\ &= \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}) \end{aligned} \quad (4.53)$$

Wegen (4.42) wird $\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}$ konsistent geschätzt durch $\mathbf{X}^\top \mathbf{X}/n$,

$$\frac{\mathbf{X}^\top \mathbf{X}}{n} \xrightarrow{p} \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}.$$

Daher ist

$$n^{1/2}(\hat{\beta}_{\text{OLS}} - \beta) \stackrel{\text{appr}}{\sim} \mathcal{N}\left(\mathbf{0}, \sigma^2 \left(\frac{\mathbf{X}^\top \mathbf{X}}{n}\right)^{-1}\right), \quad (4.54)$$

oder

$$\hat{\beta}_{\text{OLS}} \stackrel{\text{appr}}{\sim} \mathcal{N}\left(\beta, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}\right) \quad (4.55)$$

Die Verwendung einer konsistenten Schätzung von σ^2 komplettiert das Vorgehen: Betrachten wir das asymptotische Verhalten von (3.8),

$$\hat{\sigma}^2 = \frac{\mathbf{y}^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}}) \mathbf{y}}{n - k} = \frac{\mathbf{u}^\top (\mathbf{I} - \mathbf{P}_{\mathbf{X}}) \mathbf{u}}{n - k}. \quad (4.56)$$

Daher gilt

$$\hat{\sigma}^2 = \frac{\frac{\mathbf{u}^\top \mathbf{u}}{n} - \frac{\mathbf{u}^\top \mathbf{X}}{n} \left(\frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \frac{\mathbf{X}^\top \mathbf{u}}{n}}{\frac{n-k}{n}}. \quad (4.57)$$

Dies ist eine stetige Funktion in drei Argumenten

$$\begin{aligned} \frac{\mathbf{X}^\top \mathbf{X}}{n} &\xrightarrow{p} \mathbf{S}_{\mathbf{X}^\top \mathbf{X}} \\ \frac{\mathbf{u}^\top \mathbf{u}}{n} &= \frac{1}{n} \sum_{t=1}^n u_t^2 \xrightarrow{p} \sigma^2 \quad (\text{warum?}) \\ \frac{\mathbf{X}^\top \mathbf{u}}{n} &\xrightarrow{p} \mathbf{0} \quad (\text{vgl. (4.45)}). \end{aligned}$$

Daher gilt

$$\begin{aligned} \text{plim}_{n \rightarrow \infty} \hat{\sigma}^2 &= \frac{\text{plim} \left(\frac{\mathbf{u}^\top \mathbf{u}}{n} \right) - \text{plim} \left(\frac{\mathbf{u}^\top \mathbf{X}}{n} \right) \text{plim} \left(\frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \text{plim} \left(\frac{\mathbf{X}^\top \mathbf{u}}{n} \right)}{\lim_{n \rightarrow \infty} \frac{n-k}{n}} \\ &= \sigma^2. \end{aligned} \quad (4.58)$$

Damit kann die asymptotische Verteilung (4.53) konsistent geschätzt und gemäß (4.55) als Approximation für die Verteilung von $\hat{\beta}_{\text{OLS}}$ in endlichen Stichproben genutzt werden.

Ferner ist

$$\frac{(\mathbf{X}^\top \mathbf{X})^{1/2}(\hat{\beta}_{\text{OLS}} - \beta)}{\hat{\sigma}} \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{I}_k), \quad (4.59)$$

da aus (4.37), (4.42), (4.53) folgt, dass

$$\begin{aligned} \frac{(\mathbf{X}^\top \mathbf{X})^{1/2}(\hat{\beta}_{\text{OLS}} - \beta)}{\hat{\sigma}} &= \frac{(\mathbf{X}^\top \mathbf{X}/n)^{1/2}}{\hat{\sigma}} n^{1/2}(\hat{\beta}_{\text{OLS}} - \beta) \\ &\xrightarrow{d} \frac{\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{1/2}}{\sigma} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}) =: \mathcal{B} \end{aligned}$$

und, wie in (4.19),

$$\text{Var}(\mathcal{B}) = \frac{\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{1/2}}{\sigma} \sigma^2 \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1} \frac{\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{1/2}}{\sigma} = \mathbf{I}_k.$$

Damit sind t -Statistiken wie in (4.17) unter H_0 asymptotisch standardnormalverteilt.

Sie entsprechen Elementen des Vektors (4.59), wenn wir den unbekannten wahren Parameter β mit dem unter H_0 angenommenen Wert ersetzen und im Nenner das entsprechende Diagonalelement wie in (4.17) betrachten.

BechdelTest.R

Asymptotische Tests haben unter $H_0 : \theta = \theta_0$ die Eigenschaft

$$P(\tau(\mathbf{y}, \theta_0) \in \mathcal{A}(\mathbf{y}, \theta_0)) \rightarrow \alpha$$

wenn $n \rightarrow \infty$.

Asymptotische Konfidenzregionen und Tests

Mit (4.53) $(n^{1/2}(\hat{\beta}_{\text{OLS}} - \beta) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}))$ gilt, dass

$$n^{1/2} \mathbf{R}^\top (\hat{\beta}_{\text{OLS}} - \beta) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{R}^\top \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1} \mathbf{R}),$$

sowie, mit einer ähnlichen Überlegung zu (4.18) und (4.5),

$$\frac{n(\mathbf{R}^\top (\hat{\beta}_{\text{OLS}} - \beta))^\top (\mathbf{R}^\top \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1} \mathbf{R})^{-1} \mathbf{R}^\top (\hat{\beta}_{\text{OLS}} - \beta)}{\sigma^2} \xrightarrow{d} \chi_q^2, \quad (4.60)$$

wobei $\mathbf{R} : k \times q$ vollen Spaltenrang hat.

Da

$$\hat{\sigma}^2 \left(\frac{\mathbf{X}^\top \mathbf{X}}{n} \right)^{-1} \xrightarrow{p} \sigma^2 \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}, \quad (4.61)$$

erhalten wir ebenfalls (analog zu (4.59)) und mit (4.37) und (4.39).

$$\frac{(\mathbf{R}^\top (\hat{\beta}_{\text{OLS}} - \beta))^\top (\mathbf{R}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R})^{-1} \mathbf{R}^\top (\hat{\beta}_{\text{OLS}} - \beta)}{\hat{\sigma}^2} \xrightarrow{d} \chi_q^2.$$

Die Waldstatistik, um $H_0 : \mathbf{R}^\top \beta = \mathbf{r}$ zu testen, ist gegeben durch:

$$W = \frac{(\mathbf{R}^\top \hat{\beta}_{\text{OLS}} - \mathbf{r})^\top (\mathbf{R}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R})^{-1} (\mathbf{R}^\top \hat{\beta}_{\text{OLS}} - \mathbf{r})}{\hat{\sigma}^2}, \quad (4.62)$$

und unter H_0 gilt

$$W \xrightarrow{d} \chi_q^2. \quad (4.63)$$

Beachte, dass

$$W = q \cdot F \quad (4.64)$$

mit F aus (4.20).

BechdelTest.R.

Sei

$$F_{\chi_q^2}(c_\alpha^*) = 1 - \alpha,$$

also ist c_α^* das $1 - \alpha$ -Quantil der χ_q^2 -Verteilung.

Dann ist eine asymptotische $(1 - \alpha)100\%$ Konfidenzregion für \mathbf{r} gegeben durch

$$\left\{ \mathbf{r} \mid \frac{(\mathbf{R}^\top \hat{\boldsymbol{\beta}}_{\text{OLS}} - \mathbf{r})^\top \{ \mathbf{R}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R} \}^{-1} (\mathbf{R}^\top \hat{\boldsymbol{\beta}}_{\text{OLS}} - \mathbf{r})}{\hat{\sigma}^2} \leq c_\alpha^* \right\}.$$

Heteroskedastie

Was sollten wir tun, wenn die Annahme bedingter Homoskedastie (4.50),

$$\frac{\mathbf{X}^\top \mathbf{u}}{n^{1/2}} \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}),$$

zweifelhaft ist? Standardfehler gemäß (4.53) sind dann „falsch“, d.h. (4.59) trifft nicht mehr zu.

Wir benötigen dann für (4.52) noch einen Schätzer für Ω aus (4.48).

Wir können Ω schätzen mit

$$\hat{\Omega} := \frac{1}{n} \sum_{t=1}^n \hat{u}_t^2 \mathbf{x}_t^\top \mathbf{x}_t = \frac{\mathbf{X}^\top \hat{\mathbf{U}} \mathbf{X}}{n}, \quad (4.65)$$

wobei \hat{u}_t , $t = 1, \dots, n$, OLS-Residuen sind.

Hierbei ist

$$\hat{\mathbf{U}} = \begin{pmatrix} \hat{u}_1^2 & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \hat{u}_n^2 \end{pmatrix}.$$

Unter geeigneten Regularitätsbedingungen gilt $\hat{\mathbf{\Omega}} \xrightarrow{p} \mathbf{\Omega}$.

Der **Heteroskedastie-konsistente-Kovarianz-Matrix-Schätzer** (HCCME) (oder Eicker-White-Schätzer) von $\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1} \mathbf{\Omega} \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}$ ist dann gegeben durch

$$n(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{U}} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (4.66)$$

Die hieraus resultierenden Standardfehler sind extrem nützlich. In der t -Statistik (4.17) ersetzen wir $\text{se}(\hat{\beta}_i)$ durch einen robusten Standardfehler basierend auf (4.66).

Siehe `TestingWithRobustStandardErrors.R`.

Die Waldstatistik (4.62) für $H_0 : \mathbf{R}^\top \boldsymbol{\beta} = \mathbf{r}$ wird zu

$$(\mathbf{R}^\top \hat{\boldsymbol{\beta}}_{\text{OLS}} - \mathbf{r})^\top \left(\mathbf{R}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{U}} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R} \right)^{-1} (\mathbf{R}^\top \hat{\boldsymbol{\beta}}_{\text{OLS}} - \mathbf{r}) \xrightarrow{d} \chi_q^2$$

Verbesserungen können resultieren, wenn Alternativen für \hat{u}_t^2 genutzt werden:

- $\hat{u}_t^2 n / (n - k)$,
- $\hat{u}_t^2 / (1 - h_t)$, $h_t = \mathbf{e}_t^\top \mathbf{P}_X \mathbf{e}_t$,
- $\hat{u}_t^2 / (1 - h_t)^2$.

Diese kompensieren, dass Beobachtungen mit großen Varianzen die Schätzungen stark beeinflussen.

Überblick

- 1 Regressionsmodelle
- 2 Lineare Regression
- 3 Statistische Eigenschaften von Least Squares
- 4 Inferenz
- 5 Nichtlineare Regression**
- 6 Generalized Least Squares
- 7 Instrumentvariablen

Die Annahme der Linearität (vgl. (1.12)) kann abgeschwächt werden, so dass $E(y_t|X_t)$ eine nichtlineare Funktion von X_t und, was noch entscheidender ist, der Parameter sein kann.

Dieses Thema wird ausführlich im Kurs „Microeconometrics“ besprochen.

Überblick

- 1 Regressionsmodelle
- 2 Lineare Regression
- 3 Statistische Eigenschaften von Least Squares
- 4 Inferenz
- 5 Nichtlineare Regression
- 6 Generalized Least Squares**
- 7 Instrumentvariablen

Einleitung

In diesem Kapitel widmen wir uns Modellen, in denen die Störterme \mathbf{u} eine Kovarianzmatrix \mathbf{V} haben, die sich von einem skalaren Vielfachen der Einheitsmatrix unterscheiden:

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{u} \\ \mathbf{u}|\mathbf{X} &\sim (\mathbf{0}, \mathbf{V}). \end{aligned} \tag{6.1}$$

Wenn

- $\mathbf{V} = \sigma^2 \mathbf{I}$ gilt, dann sind wir wieder im klassischen linearen Modell
- $\mathbf{V} \neq \sigma^2 \mathbf{I}$ und \mathbf{V} diagonal ist, dann liegt Heteroskedastie vor.
- \mathbf{V} nicht diagonal ist, sind die Störterme autokorreliert. (Siehe z.B. ARMA-Modelle weiter unten.)

Es gibt zwei Strategien im Umgang mit $\mathbf{V} \neq \sigma^2 \mathbf{I}$:

1. Wir benutzen weiter OLS und nutzen Standardfehler, die gegen $\mathbf{V} \neq \sigma^2 \mathbf{I}$ robust sind. Für diagonale \mathbf{V} haben wir das im vorletzten Kapitel (HCCME) gemacht. (Es gibt auch gegen Autokorrelation robuste Kovarianzmatrix „HAC“-Schätzer, die wir in diesem Kurs jedoch nicht betrachten.)
2. Wir benutzen neue, möglicherweise effizientere, Schätzer, die die Struktur von \mathbf{V} berücksichtigen.

Wir beschäftigen uns jetzt mit der zweiten Alternative: Als Alternative zum HCCME versuchen wir die Struktur von \mathbf{V} miteinzubeziehen, um effizientere Schätzer zu erhalten.

Wenn \mathbf{V} in

$$\mathbf{V} = \sigma^2 \boldsymbol{\Sigma} \quad (6.2)$$

bis auf ein skalares Vielfaches bekannt ist— $\boldsymbol{\Sigma}$ wird also als bekannt vorausgesetzt—können wir den **generalized least squares (GLS)**-Schätzer formulieren und zeigen, dass er BLUE ist.

Betrachten wir eine „Matrixwurzel“ Ψ von Σ^{-1} , so dass (mehr Details wie man Ψ findet unten)

$$\Sigma^{-1} = \Psi \Psi^{\top}. \quad (6.3)$$

Wir transformieren (6.1) via

$$\tilde{\mathbf{y}} = \Psi^{\top} \mathbf{y}, \quad \tilde{\mathbf{X}} = \Psi^{\top} \mathbf{X}, \quad \tilde{\mathbf{u}} = \Psi^{\top} \mathbf{u}. \quad (6.4)$$

Wir erhalten dann

$$\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\beta + \tilde{\mathbf{u}} \quad \text{mit} \quad \tilde{\mathbf{u}}|\tilde{\mathbf{X}} \sim (\mathbf{0}, \sigma^2 \mathbf{I}),$$

da

$$\text{Var}(\tilde{\mathbf{u}}|\tilde{\mathbf{X}}) = \Psi^{\top} \text{Var}(\mathbf{u}|\mathbf{X}) \Psi \stackrel{(6.2)}{=} \Psi^{\top} \sigma^2 \Sigma \Psi = \sigma^2 \Psi^{\top} (\Psi^{\top})^{-1} \Psi^{-1} \Psi,$$

weil (6.3) impliziert, dass

$$\Sigma = (\Psi \Psi^{\top})^{-1} \stackrel{\text{M20.5}}{=} (\Psi^{\top})^{-1} \Psi^{-1}$$

Wir befinden uns nun wieder im klassischen Modell.

Der OLS-Schätzer basierend auf den transformierten Daten $\tilde{\mathbf{y}}$ und $\tilde{\mathbf{X}}$,

$$\tilde{\beta} = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{y}},$$

ist der **GLS-Schätzer**

$$\begin{aligned} \hat{\beta}_{\text{GLS}} &\stackrel{(6.4)}{=} (\mathbf{X}^\top \Psi \Psi^\top \mathbf{X})^{-1} \mathbf{X}^\top \Psi \Psi^\top \mathbf{y} \\ &\stackrel{(6.3)}{=} (\mathbf{X}^\top \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Sigma^{-1} \mathbf{y}. \end{aligned} \quad (6.5)$$

Er ist BLUE, da $\tilde{\beta}$ der KQ-Schätzer für das klassische Modell ist.

Basierend auf ähnlichen Annahmen wie denen in Kapiteln 4 erhalten wir

$$n^{1/2}(\hat{\beta}_{\text{GLS}} - \beta) \xrightarrow{d} \mathcal{N}\left(\mathbf{0}, \sigma^2 \text{plim}_{n \rightarrow \infty} \left(\frac{\mathbf{X}^\top \Sigma^{-1} \mathbf{X}}{n}\right)^{-1}\right),$$

was wie zuvor GLS-basierte t-Tests etc. liefert.

Für den Spezialfall, bei dem Σ eine Diagonalmatrix ist mit Elementen σ_{tt} , $t = 1, \dots, n$, ist der GLS-Schätzer gegeben durch

$$\hat{\beta}_{\text{GLS}} = \hat{\beta}_{\text{WLS}} = \left(\sum_{t=1}^n \sigma_{tt}^{-1} \mathbf{x}_t^\top \mathbf{x}_t \right)^{-1} \sum_{t=1}^n \sigma_{tt}^{-1} \mathbf{x}_t^\top \mathbf{y}_t$$

und wird normalerweise als **weighted least squares** (WLS) bezeichnet.

Je größer die Varianz σ_{tt} ist, desto weniger stark wird die Beobachtung gewichtet.

Alternative Motivation für GLS: $\hat{\beta}_{\text{GLS}}$ minimiert die folgende Funktion in Bezug auf β :

$$\tilde{\mathbf{u}}^\top \tilde{\mathbf{u}} = \mathbf{u}^\top \Sigma^{-1} \mathbf{u} = (\mathbf{y} - \mathbf{X}\beta)^\top \Sigma^{-1} (\mathbf{y} - \mathbf{X}\beta).$$

Feasible GLS

In der Praxis ist Σ normalerweise unbekannt, so dass GLS nicht direkt implementierbar ist. Wir betrachten daher nun Fälle, in denen \mathbf{V} eine bekannte Funktion einer fixen Anzahl an Parametern $\boldsymbol{\theta} : m \times 1$ ist, die nicht größer wird, wenn n steigt.

(Die Anzahl nicht doppelter Elemente von Σ entspräche ansonsten $n(n+1)/2$. Dies steigt mit n . Eine wichtige Bedingung für die Konsistenz des Schätzers ist jedoch, dass die Anzahl der Parameter fix ist. Wir haben daher ein Problem, wenn die Anzahl der Parameter von n abhängt.)

Dann können wir einen konsistenten Schätzer für $\boldsymbol{\theta}$ suchen, also einen Schätzer $\hat{\boldsymbol{\theta}}$, für den gilt

$$\hat{\boldsymbol{\theta}} \xrightarrow{P} \boldsymbol{\theta}.$$

Wir arbeiten dann mit $\mathbf{V}(\hat{\boldsymbol{\theta}})$, was wir berechnen können. Dies bringt uns zum **feasible GLS (FGLS)**-Schätzer.

FGLS hat typischerweise dieselbe asymptotische Verteilung wie GLS:

Nehmen wir an, dass (mit Verwendung der Notation $\hat{\mathbf{V}} = \mathbf{V}(\hat{\boldsymbol{\theta}})$) gilt

$$\frac{\mathbf{X}^\top \hat{\mathbf{V}}^{-1} \mathbf{X}}{n} - \frac{\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X}}{n} \xrightarrow{p} \mathbf{0} \quad (6.6)$$

Wir müssen dies explizit annehmen, da—obwohl $\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X}$ eine $k \times k$ Matrix ist, wobei k nicht von n abhängt und für ein einziges n eine stetige Funktion von $\boldsymbol{\theta}$ ist, da \mathbf{V}^{-1} eine stetige Funktion von $\boldsymbol{\theta}$ ist—wir die 5. asymptotische Regel nicht verwenden können. Der Grund hierfür ist, dass die Funktion sich in Abhängigkeit von n verändert.

Wenn jedoch Regularitätsbedingungen erfüllt sind (was für die folgenden Beispiele der Fall ist), gilt (6.6) sowie auch das folgende stärkere Resultat (wir überspringen technische Details)

$$\frac{\mathbf{X}^\top \hat{\mathbf{V}}^{-1} \mathbf{u}}{n^{1/2}} - \frac{\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{u}}{n^{1/2}} \xrightarrow{p} \mathbf{0}.$$

Dann hat

$$\hat{\beta}_{\text{FGLS}} = (\mathbf{X}^\top \hat{\mathbf{V}}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{V}}^{-1} \mathbf{y}$$

dieselbe asymptotische Verteilung wie $\hat{\beta}_{\text{GLS}}$, d.h.

$$n^{1/2}(\hat{\beta}_{\text{FGLS}} - \hat{\beta}_{\text{GLS}}) \xrightarrow{p} \mathbf{0}.$$

ARMA Modelle

Die Kovarianzmatrix der Störterme wird oft mittels eines ARMA-Modells konstruiert.

Häufig wird ein AR(1)-Prozess verwendet:

$$u_t = u_{t-1}\rho + \epsilon_t, \quad \epsilon_t \stackrel{\text{u.i.v.}}{\sim} (0, \sigma_\epsilon^2),$$

mit ϵ_t unabhängig von u_{t-j} , $j > 0$.

Stationarität gilt, wenn die unbedingte Verteilung von u_t nicht von t abhängt, oder schwächer, wenn die ersten zwei Momente (wenn sie existieren) nicht von t abhängen.

Wenn also $E(u_t) = 0$ und $Var(u_t) = \sigma^2$ für alle t , dann gilt, bei Unabhängigkeit von ϵ_t und u_{t-1}

$$Var(u_t) = \rho^2 Var(u_{t-1}) + \sigma_\epsilon^2,$$

und Stationarität impliziert

$$\sigma^2 = \rho^2 \sigma^2 + \sigma_\epsilon^2.$$

Also gilt

$$\sigma^2 = \frac{\sigma_\epsilon^2}{1 - \rho^2}.$$

Wenn $|\rho| < 1$ liegt Stationarität vor, wenn der Prozess bei $t \rightarrow -\infty$ beginnt oder wenn

$$u_0 \sim \left(0, \frac{\sigma_\epsilon^2}{1 - \rho^2}\right).$$

Die Kovarianzmatrix ist gegeben durch (siehe Übung)

$$\text{Var}(\mathbf{u}|\mathbf{X}) = \mathbf{V}(\rho, \sigma_\epsilon^2) = \frac{\sigma_\epsilon^2}{1 - \rho^2} \begin{pmatrix} 1 & \rho & \cdots & \rho^{n-2} & \rho^{n-1} \\ \rho & 1 & \cdots & \rho^{n-3} & \rho^{n-2} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \rho^{n-2} & \rho^{n-3} & \cdots & 1 & \rho \\ \rho^{n-1} & \rho^{n-2} & \cdots & \rho & 1 \end{pmatrix}.$$

Die Matrixzerlegung gemäß (6.3) ist gegeben durch

$$\Psi(\rho) = \begin{pmatrix} (1 - \rho^2)^{1/2} & -\rho & 0 & \cdots & 0 & 0 \\ 0 & 1 & -\rho & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & -\rho \\ 0 & 0 & 0 & \cdots & 0 & 1 \end{pmatrix}, \quad (6.7)$$

was eine sogenannte **Bandmatrix** mit zwei Diagonalen ist.

Denn: Direkte Multiplikation zeigt, dass

$$\boldsymbol{\Psi}^\top(\rho) \mathbf{V}(\rho, \sigma_\epsilon^2) \boldsymbol{\Psi}(\rho) = \sigma_\epsilon^2 \mathbf{I}_n.$$

Daraus folgt mit M20.5

$$(\boldsymbol{\Psi}^\top(\rho) \mathbf{V}(\rho, \sigma_\epsilon^2) \boldsymbol{\Psi}(\rho))^{-1} = \boldsymbol{\Psi}(\rho)^{-1} \mathbf{V}(\rho, \sigma_\epsilon^2)^{-1} \boldsymbol{\Psi}^\top(\rho)^{-1} = \sigma_\epsilon^{-2} \mathbf{I}_n.$$

Wir sehen, dass

$$\mathbf{V}^{-1}(\rho, \sigma_\epsilon^2) = \sigma_\epsilon^{-2} \boldsymbol{\Psi}(\rho) \boldsymbol{\Psi}^\top(\rho),$$

und der GLS-Schätzer (6.5) ist gegeben durch (σ_ϵ^{-2} kürzt sich raus)

$$\hat{\beta}_{\text{GLS}} = (\mathbf{X}^\top \boldsymbol{\Psi}(\rho) \boldsymbol{\Psi}^\top(\rho) \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Psi}(\rho) \boldsymbol{\Psi}^\top(\rho) \mathbf{y}.$$

Dies wäre der BLUE-Schätzer, wenn ρ bekannt wäre.

Für den FGLS-Schätzer ersetze ρ durch eine konsistente Schätzung $\hat{\rho}$.

Z.B. können wir die Residuen \hat{u}_t einer KQ-Regression von \mathbf{y} auf \mathbf{X} auf deren verzögerte Werte \hat{u}_{t-1} regressieren.

Wir können andere ARMA-Prozesse in ähnlicher Weise verwenden. Z.B. erhalten wir für einen MA(1)-Prozess

$$u_t = \epsilon_t + \alpha\epsilon_{t-1}, \quad \epsilon_t \stackrel{\text{u.i.v.}}{\sim} (0, \sigma_\epsilon^2)$$

dass

$$\mathbf{V}(\alpha, \sigma_\epsilon^2) = \sigma_\epsilon^2 \begin{pmatrix} 1 + \alpha^2 & \alpha & 0 & \cdots & 0 & 0 \\ \alpha & 1 + \alpha^2 & \alpha & \cdots & 0 & \vdots \\ \vdots & \vdots & \vdots & & \vdots & \alpha \\ 0 & 0 & 0 & \cdots & \alpha & 1 + \alpha^2 \end{pmatrix}.$$

Faktorisierung

Es war ein bisschen ad hoc die Matrix Ψ zu finden, so dass $\Sigma^{-1} = \Psi\Psi^\top$, wie wir es in den ARMA-Fällen vorher getan haben. Wir besprechen nun kurz Möglichkeiten für die allgemeine Wahl von Ψ .

Die untere Dreiecks-Cholesky-Zerlegung ist eine Möglichkeit (siehe Resultat 30).

Eine weitere Möglichkeit für Ψ ist gegeben durch eine symmetrische Matrix $\Sigma^{-1/2}$ über die Singulärwert- (SVD) oder die Eigenwertzerlegung von Σ :

$$\begin{aligned}\Sigma &= \mathbf{P}\mathbf{D}\mathbf{P}^\top \\ \mathbf{P}^\top \mathbf{P} &= \mathbf{P}\mathbf{P}^\top = \mathbf{I}, \\ \mathbf{D} &> \mathbf{0}, \quad \text{diagonal}\end{aligned}$$

Dann gilt

$$\Sigma^{1/2} = \mathbf{P}\mathbf{D}^{1/2}\mathbf{P}^\top,$$

und $\Sigma^{-1/2}$ ist dessen Inverse, gegeben durch $\mathbf{P}\mathbf{D}^{-1/2}\mathbf{P}^\top$.

Test auf Heteroskedastie

Oft wollen wir testen, ob Homoskedastie vorliegt, d.h. $H_0 : E(u_t^2) = \sigma^2$, $t = 1, \dots, n$ gegen die Alternative $E(u_t^2) \neq E(u_s^2)$, $s \neq t$, $s, t = 1, \dots, n$.

Betrachte die restriktivere Alternative

$$E(u_t^2) = h(\delta + \mathbf{Z}_t\gamma),$$

wobei h die **skedastische Funktion** ist. h sei stetig differenzierbar und $h \geq 0$, muss ansonsten jedoch nicht weiter spezifiziert werden. \mathbf{Z}_t ist ein beobachteter $1 \times r$ Vektor (der Elemente von \mathbf{X}_t enthalten kann).

H_0 wird zu $\gamma = \mathbf{0}$, während H_1 durch $\gamma \neq \mathbf{0}$ gegeben ist.

Um einen Test durchführen zu können, benötigen wir quadrierte Residuen, die für u_t^2 konsistent sind. Eine offensichtliche Wahl sind die quadrierten OLS-Residuen \hat{u}_t^2 .

White-Test

Bezeichne den Vektor von r unterschiedlichen (nicht konstanten) Elementen von $\mathbf{X}_t^\top \mathbf{X}_t$ als \mathbf{Z}_t . Unter H_0 sind die u_t^2 keine Funktion von \mathbf{Z}_t .

Berechne das R^2 einer Regression von \hat{u}_t^2 auf eine Konstante und \mathbf{Z}_t . Wir können dann zeigen, dass (dies ist ein recht allgemeines Resultat!)

$$nR^2 \xrightarrow{d} \chi_r^2$$

Schreibe dazu (4.24) im „ χ^2 -Format“ (4.64):

$$r \cdot F_{\beta_2} = \frac{\text{RSSR} - \text{USSR}}{\text{USSR}/(n - k)}$$

Ersetze, was asymptotisch keinen Unterschied macht, $n - k$ durch n und bemerke, dass unter H_0 $\text{USSR}/(n - k)$ und RSSR/n äquivalent sind. Nun ist mit (2.19) für diese H_0 (vgl. (4.26))

$$\frac{\text{RSSR} - \text{USSR}}{\text{RSSR}} = R^2$$

Zum Beispiel erhalten wir für $\mathbf{X}_t = (1, q_t, q_t^2, p_t)$

$$\mathbf{X}_t^\top \mathbf{X}_t = \begin{pmatrix} 1 \\ q_t \\ q_t^2 \\ p_t \end{pmatrix} (1, q_t, q_t^2, p_t) = \begin{pmatrix} 1 & q_t & q_t^2 & p_t \\ & q_t^2 & q_t^3 & q_t p_t \\ & & q_t^4 & q_t^2 p_t \\ \bullet & & & p_t^2 \end{pmatrix},$$

so dass

$$\mathbf{Z}_t = (q_t, q_t^2, p_t, q_t^3, q_t p_t, q_t^4, q_t^2 p_t, p_t^2)$$

und $r = 8$.

white test.R

Test auf Autokorrelation

Wir verwenden den unten beschriebenen t -Test. (Diese Prozedur ist asymptotisch valide, auch wenn erklärende Variablen nur „vorbestimmt“ statt exogen sind—wir verfolgen das aber nicht im Detail weiter.)

Die Prozedur ist wie folgt. Seien $\hat{\mathbf{u}}$ die OLS-Residuen aus

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{u},$$

Wir bezeichnen den Vektor der verzögerten Residuen mit $\hat{\mathbf{u}}_{-1} = (\hat{u}_1, \dots, \hat{u}_{n-1})^\top$. Dann können wir folgende Regression verwenden

$$\mathbf{y} = \mathbf{X}\beta + b_\rho \hat{\mathbf{u}}_{-1} + \text{Fehler},$$

um $H_0 : b_\rho = 0$ mit einer üblichen t -Statistik (4.17) zu testen (y_1 und \mathbf{X}_1 müssen aus dem Datensatz gestrichen werden oder $\hat{\mathbf{u}}_{-1}$ am Anfang um eine Null ergänzt werden). `breusch godfrey.R`

Was tun wir, wenn wir Autokorrelation finden?

Es ist verlockend das Problem Autokorrelation zu „lösen“, indem ein FGLS-Schätzer angewendet wird, der Autokorrelation berücksichtigt.

Allerdings können autokorrelierte Residuen oft einfach ein Hinweis auf eine andere Art von Misspezifikation des Modells sein, z.B. das Fehlen von quadratischen Termen. Ein lineares Modell wird dann zunächst viele positive (negative) und dann viele negative (positive) Residuen liefern, wodurch Autokorrelation in den Residuen entsteht.

Einfach einen FGLS-Schätzer zu verwenden löst das Problem fälschlicherweise ausgelassener quadratischer Terme also nicht.

Fazit: Nachdenken!

Paneldaten

Wir beschäftigen uns nun mit Daten, die über zwei Dimensionen erhoben wurden, auch **Paneldaten** genannt.

Wir könnten z.B. jährliche Beobachtungen zu verschiedenen europäischen Ländern haben („Querschnittseinheiten“) oder Beobachtungen von mehreren Kindern aus einer Familie.

Wir betrachten das folgende lineare Modell

$$y_{it} = \mathbf{X}_{it}\beta + u_{it} \quad i = 1, \dots, m, \quad t = 1, \dots, T \quad (6.8)$$

Hierbei entspricht nun m der Anzahl der Einheiten und T der Anzahl der Beobachtungen pro Einheit.

Wenn $E(u_{it}|\mathbf{X}_{it}) = 0$ gilt, können wir das Modell einfach mit **pooled OLS** schätzen, was die Paneldaten-Struktur ignoriert und einfach alle $n = mT$ Beobachtungen zusammenwirft.

Allerdings ist diese Annahme in Paneldaten oft unrealistisch. Häufig gibt es Faktoren, die ein Land über mehrere Zeitperioden oder alle Kinder einer Familie beeinflussen. Wenn diese Faktoren mit \mathbf{X}_{it} in Beziehung stehen, produziert OLS inkonsistente Schätzer (vgl. (3.14)).

Wir modellieren u_{it} mit dem **Fehlerkomponentenmodell**

$$u_{it} = \eta_i + \epsilon_{it} \quad (6.9)$$

Beispiel 6.1 (Geronimus and Korenman, Quarterly Journal of Economics 1992).

Wir interessieren uns für den Effekt von Teenagerschwangerschaften auf zukünftige ökonomische Ergebnisse, z.B. das Einkommen relativ zum Bedarf (wobei Letzteres von der Anzahl der Kinder abhängt). Betrachten wir

$$\begin{aligned} \log(\text{income}/\text{needs}_{fs}) = & \beta_0 + \delta_0 \text{sister2}_s + \beta_1 \text{teenbirth}_{fs} \\ & + \beta_2 \text{age}_{fs} + \text{other factors} + a_f + \epsilon_{fs}, \end{aligned}$$

wobei f der Index für die Familie und s der Index für eine Schwester innerhalb der Familie ist. D.h. t wird durch s und i durch f ersetzt. Zum Beispiel ist a_f hier der familienspezifische Effekt.

Unser Fokus liegt auf β_1 , d.h. der ökonomische Effekt einer Teenagergeburt ($\text{teenbirth}_{fs} = 1$). Eine Problem bei gepoolten oder Random Effects-Regressionen (siehe unten) ist, dass teenbirth_{fs} und a_f wahrscheinlich korreliert sind.

Die Bildung der ersten Differenz zwischen Schwestern bietet eine Lösung:

$$\Delta \log(\text{income/needs}_f) = \delta_0 + \beta_1 \Delta \text{teenbirth}_f + \beta_2 \Delta \text{age}_f + \dots + \Delta \epsilon_f \quad (6.10)$$

Dafür ist es offensichtlich notwendig, dass für mindestens ein Schwesternpaar $\Delta \text{teenbirth}_f \neq 0$ ist, um β_1 zu identifizieren. D.h. es muss Familien geben, in denen genau eine der Schwestern eine Teenagergeburt hatte.

Mittels Pooled OLS (d.h. die Familienstruktur der Daten wird ignoriert) fanden Geronimus und Korenman $\hat{\beta}_1 \approx -0.3$ (in Abhängigkeit von den genauen anderen Faktoren), wohingegen die Regression mit der ersten Differenz einen Wert von $\hat{\beta}_1 = -0.08$ ergibt.

Dies suggeriert, dass der Familienhintergrund der Mutter einen stärkeren negativen Effekt auf ökonomische Outcomes hat als die Teenagerschwangerschaft selbst.

sisters.R

Wir beschäftigen uns jetzt mit einer allgemeineren Formulierung eines Paneldaten-Schätzers, der für $T = 2$ äquivalent zum obigen First Difference-Schätzer (d.h. eine OLS-Regression von (6.10)) ist.

Im Folgenden nehmen wir an, dass \mathbf{X}_{it} exogen ist—verzögerte abhängige Variablen sorgen für einige zusätzliche Komplikationen mit Paneldaten.

In Matrixnotation können (6.8) und (6.9) kompakt als

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{D}\boldsymbol{\eta} + \boldsymbol{\varepsilon} \quad (6.11)$$

geschrieben werden, wobei \mathbf{y} und $\boldsymbol{\varepsilon}$ n -Vektoren mit typischen Elementen y_{it} und ε_{it} sind und \mathbf{D} eine $n \times m$ (eine Spalte pro Einheit) Matrix aus Dummyvariablen ist.

Wenn eine Zeile eine zu Einheit i gehörende Beobachtung enthält, dann hat \mathbf{D} eine eins in Spalte i und ansonsten eine null, für alle $i = 1, \dots, m$.

Die **individuenspezifischen Effekte** $\boldsymbol{\eta}$ müssen unabhängig von ε_{it} sein. Im **Fixed Effects** Paneldatenmodell dürfen sie allerdings mit den Variablen in \mathbf{X} korrelieren.

Der **Fixed Effects Schätzer** ist der OLS-Schätzer von β in der Regression von \mathbf{y} auf \mathbf{X} und \mathbf{D} in (6.11). Wegen seiner Berechnungsweise wird er auch **Least Squares Dummy Variablen-** oder **LSDV-Schätzer**, genannt.

Sei

$$\mathbf{M}_D = \mathbf{I} - \mathbf{D}(\mathbf{D}^\top \mathbf{D})^{-1} \mathbf{D}^\top.$$

Dann erhalten wir mittels FWL (siehe 2.4) den OLS-Schätzer durch eine Regression von $\mathbf{M}_D \mathbf{y}$, die Residuen von der Regression von \mathbf{y} auf \mathbf{D} , auf $\mathbf{M}_D \mathbf{X}$, die Matrix der Residuen aus der Regression jeder Spalte von \mathbf{X} auf \mathbf{D} .

Der Fixed Effects Schätzer ist daher

$$\hat{\beta}_{\text{FE}} = (\mathbf{X}^\top \mathbf{M}_D \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{M}_D \mathbf{y} \quad (6.12)$$

Es ist einfach zu überprüfen (vgl. Beispiel 2.5), dass das it -te Element des Vektors $\mathbf{M}_D \mathbf{z}$, für jeden Vektor \mathbf{z} , gleich $z_{it} - \bar{z}_i$ ist, der Abweichung vom Gruppenmittelwert.

Da alle Variablen in (6.12) mit \mathbf{M}_D multipliziert werden, verwendet dieser Schätzer lediglich die Variation um den Mittelwert für jede Einheit.

Damit beseitigen wir das einheitenweise konstante η_i !

Aus diesem Grund wird der Schätzer auch **Within-Schätzer** genannt. Siehe `panel data.R` und `panel data FD and FE.R`.

Dies impliziert unglücklicherweise, dass der Schätzer nicht für erklärende Variablen, die für jede Beobachtung einer Gruppe den selben Wert annehmen, verwendet werden kann, da er nach Transformation mit \mathbf{M}_D gleich null ist.

Können wir Paneldaten verwenden, um die Renditen von Schulbildung zu schätzen? Angeborene Fähigkeiten sind per Definition über die Zeit unveränderlich. Ein mögliches Modell ist:

$$\log(wage)_{it} = \beta_0 + \delta_0 d_{2t} + \beta_1 educ_{it} + \beta_2 gender_{it} + \eta_i + \epsilon_{it}, \quad t = 1, 2$$

Mit den ersten Differenzen erhalten wir

$$\Delta \log(wage)_i = \delta_0 + \beta_1 \Delta educ_i + \beta_2 \Delta gender_i + \epsilon_i \quad (6.13)$$

Das Problem ist, dass wir, um Einkommen zu beobachten, arbeitende Erwachsene betrachten, bei denen sich die Bildung (wie meist auch Geschlecht) über t nicht mehr verändert. Dann können wir (6.13) jedoch nicht schätzen.

Wenn sich Bildung nur bei einigen über t verändert, ist der Schätzer für β_1 unpräzise, da der Regressor wenig Varianz hat.

Dies illustriert die generelle Lektion, dass der FE-Schätzer nur den Effekt von **zeitvariierenden Variablen** schätzen kann, da erste Differenzen/Dummies zeitunveränderliche Effekten entfernen.

Dieses Manko macht es manchmal attraktiv, einen anderen Schätzer, den **Random-Effects-Schätzer**, zu verwenden. Wir kehren zurück zu Modell (6.8) und (6.9), benötigen jedoch für diesen Schätzer die striktere Annahme, dass

$$E(\eta_i | \mathbf{X}) = 0 \quad (6.14)$$

Unter dieser Annahme wäre pooled OLS unverzerrt, aber wir können einen effizienteren GLS-Schätzer (6.5) finden. GLS ist effizienter, da die Fehler-Kovarianz-Matrix in diesem Model nicht homoskedastisch ist.

Nehmen wir an, die η_i sind u.i.v. mit $E(\eta_i) = 0$ und $Var(\eta_i) = \sigma_\eta^2$. Diese Annahmen motivieren die Bezeichnung „zufällige Effekte“ (random effects). Wir nehmen des Weiteren an, dass

$$E(\varepsilon \varepsilon^\top | \mathbf{X}) = \sigma_\varepsilon^2 \mathbf{I}$$

Wenn die zwei Fehlerkomponenten ferner unabhängig sind, gilt in (6.9)

$$Var(u_{it}) = \sigma_\eta^2 + \sigma_\varepsilon^2$$

$$Cov(u_{it}, u_{is}) = \sigma_\eta^2$$

$$Cov(u_{it}, u_{js}) = 0 \quad \text{für alle } i \neq j.$$

Wir erhalten dann die folgende Varianz-Kovarianz-Matrix \mathbf{V} :

$$\mathbf{V} = \begin{pmatrix} \Sigma & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \Sigma & \cdots & \mathbf{0} \\ \vdots & \vdots & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \Sigma \end{pmatrix}$$

Hier gilt

$$\Sigma = \sigma_{\eta}^2 \boldsymbol{\iota} \boldsymbol{\iota}^{\top} + \sigma_{\epsilon}^2 \mathbf{I}_T$$

mit $\boldsymbol{\iota}$ einem T -Vektor aus Einsen.

Wie gewöhnlich benötigt FGLS einen konsistenten Schätzer für \mathbf{V} , und daher hier für σ_{η}^2 und σ_{ϵ}^2 . Solche existieren. Wir verfolgen dies jedoch nicht im Detail, da die Schlüsselannahme für Random Effects $E(\eta_i | \mathbf{X}) = 0$ (siehe (6.14)), in Anwendungen oft nicht glaubwürdig ist.

Man kann die Modelle mit einem **Hausman-Test** gegeneinander testen, siehe das folgende Kapitel und panel data FE and RE.R.

Überblick

- 1 Regressionsmodelle
- 2 Lineare Regression
- 3 Statistische Eigenschaften von Least Squares
- 4 Inferenz
- 5 Nichtlineare Regression
- 6 Generalized Least Squares
- 7 Instrumentvariablen**

Einführung

Wir betrachten nun

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad (7.1)$$

$$\mathbf{u} \sim (\mathbf{0}, \sigma^2 \mathbf{I}_n), \quad (7.2)$$

wobei wir nun erlauben, dass (möglicherweise alle) Elemente von \mathbf{X}_t mit u_t korreliert sind.

Die Randverteilung von u_t erfüllt somit (7.2). D.h. es kann sein, dass (4.44) nicht erfüllt ist, sondern

$$E(\mathbf{X}_t^\top u_t) \neq \mathbf{0}. \quad (7.3)$$

Unter (7.3) ist der OLS-Schätzer inkonsistent (vgl. auch (4.46)).

Zunächst betrachten wir Fälle, in denen $E(u_t|\mathbf{X}_t) \neq 0$ sein kann und untersuchen die Folgen für die Schätzer aus den vorherigen Kapiteln.

Danach führen wir neue Informationen in Form von **Instrumentvariablen** ein und beschreiben die klassische **Instrumentvariablen-Schätzung** (IV) ähnlich zur MM-Schätzung.

Danach befassen wir uns mit verschiedenen Tests:

- Heteroskedastie-robuste Tests,
- Tests auf überidentifizierende Restriktionen,
- „Tests auf Exogenität“,
- Hausmantests.

Korrelation zwischen Regressoren und Störtermen

Wir betrachten die folgenden Fälle, bei denen \mathbf{X}_t mit u_t korreliert sein kann.

Erster Fall: Omitted Variable Bias (OVB)

Es gibt Variablen, die \mathbf{y} erklären und die mit den bereits berücksichtigten \mathbf{X} korreliert sind, die wir jedoch nicht beobachten (können). Vgl. (3.14).

Betrachte z.B. Bildungsrenditen: Sei y_t das Einkommen, b_t eine geeignete Funktion des Geburtsjahres und s_t die in der Schule verbrachten Jahre.

Durchführen einer OLS-Regression für

$$y_t = \alpha + s_t\beta + b_t\gamma + u_t,$$

wobei $t = 1, \dots, n$, kann zur inkonsistenten Schätzung des „wahren“ Effekts von s_t führen, da sowohl s_t als auch u_t durch dieselbe unbeobachtbare Variable („Fähigkeiten“) beeinflusst sein kann. Also ist

$$E(u_t | s_t) \neq 0.$$

Fehler in Variablen

Zweiter Fall: Sei

$$\mathbf{y} = \mathbf{x}^* \beta + \mathbf{u},$$

wobei \mathbf{x}^* der Einfachheit halber $n \times 1$ sei.

Wir beobachten jedoch nicht \mathbf{x}^* , sondern $\mathbf{x} = \mathbf{x}^* + \mathbf{v}$, wobei die \mathbf{v} **Messfehler** darstellen.

Es seien $\{u_i, v_i\}$ u.i.v., $E(u_i) = E(v_i) = 0$ und

$$\text{Var}(u_i, v_i) = \begin{bmatrix} \sigma_u^2 & \sigma_{uv} \\ \sigma_{uv} & \sigma_v^2 \end{bmatrix}$$

sowie

$$\frac{\mathbf{x}^{*\top} \mathbf{x}^*}{n} \xrightarrow{p} M > 0$$

und

$$\text{Cov}(x_i^*, v_i) = E(x_i^* v_i) = 0 \quad (7.4)$$

Das Gesetz der großen Zahlen ($\{u_i v_i\}$ und $\{v_i^2\}$ sind u.i.v.) liefert

$$\begin{aligned}\frac{\mathbf{u}^\top \mathbf{v}}{n} &= \frac{1}{n} \sum_{i=1}^n u_i v_i \xrightarrow{p} E(u_i v_i) = \sigma_{uv}, \\ \frac{\mathbf{v}^\top \mathbf{v}}{n} &= \frac{1}{n} \sum_{i=1}^n v_i^2 \xrightarrow{p} E(v_i^2) = \sigma_v^2.\end{aligned}$$

Des Weiteren gilt mit dem GGZ (4.32) und (7.4), dass

$$\frac{\mathbf{x}^{*\top} \mathbf{v}}{n} \xrightarrow{p} E(x_i^* v_i) = 0$$

Analog zeigt man, dass unter geeigneten Exogenitätsannahmen wie üblich

$$\frac{\mathbf{x}^{*\top} \mathbf{u}}{n} \xrightarrow{p} 0. \quad (7.5)$$

Also ist

$$\begin{aligned}\text{plim}_{n \rightarrow \infty} \hat{\beta} &= \text{plim} \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}} \\&= \frac{\text{plim} \left(\frac{(\mathbf{x}^* + \mathbf{v})^\top (\mathbf{x}^* \beta + \mathbf{u})}{n} \right)}{\text{plim} \left(\frac{(\mathbf{x}^* + \mathbf{v})^\top (\mathbf{x}^* + \mathbf{v})}{n} \right)} \\&= \frac{\text{plim} \left(\frac{\mathbf{x}^{*\top} \mathbf{x}^*}{n} \right) \beta + \text{plim} \left(\frac{\mathbf{v}^\top \mathbf{x}^*}{n} \right) + \text{plim} \left(\frac{\mathbf{x}^{*\top} \mathbf{u}}{n} \right) \beta + \text{plim} \left(\frac{\mathbf{v}^\top \mathbf{u}}{n} \right)}{\text{plim} \left(\frac{\mathbf{x}^{*\top} \mathbf{x}^*}{n} \right) + 2 \text{plim} \left(\frac{\mathbf{x}^{*\top} \mathbf{v}}{n} \right) + \text{plim} \left(\frac{\mathbf{v}^\top \mathbf{v}}{n} \right)} \\&= \frac{M\beta + \sigma_{uv}}{M + \sigma_v^2} \\&= \beta + \frac{\sigma_{uv} - \sigma_v^2 \beta}{M + \sigma_v^2}.\end{aligned}$$

Demzufolge ist der OLS-Schätzer nicht konsistent, wenn die erklärenden Variablen fehlerhaft gemessen werden.

Beachte zudem, dass der $\text{plim} \hat{\beta}$ näher an null liegt als β , wenn $\sigma_{uv} = 0$:
Messfehler führen zum so genannten **Attenuation Bias**.

Simultane Gleichungssysteme

Ein weiterer Fall, bei dem die Annahme $E(u_t|\mathbf{X}_t) = 0$ verletzt sein könnte, tritt bei **simultanen Gleichungssystemen** ein:

$$\begin{aligned}y_{1t} &= \alpha_1 + x_t\beta_1 + y_{2t}\gamma + u_{1t}, \\y_{2t} &= \alpha_2 + x_t\beta_2 + u_{2t},\end{aligned}$$

wobei $t = 1, \dots, n$. Beide y_{1t} und y_{2t} sind endogen, jedoch ist y_{2t} ebenfalls ein Regressor für y_{1t} .

Wenn $\text{Cov}(u_{1t}, u_{2t}) \neq 0$, dann $\text{Cov}(u_{1t}, y_{2t}) \neq 0$ und eine Regression von \mathbf{y}_1 auf \mathbf{x} und \mathbf{y}_2 würde zu inkonsistenten OLS-Schätzern führen.

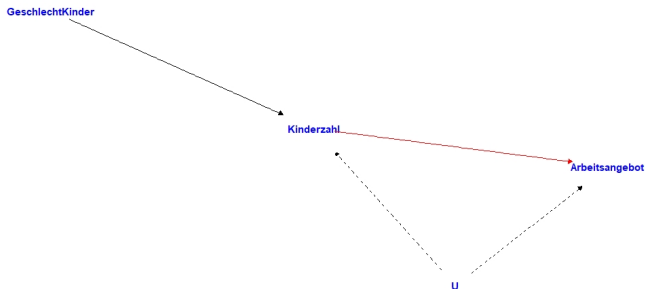
Wir befassen uns hier nicht weiter mit simultanen Gleichungssystemen.

Stichprobenverzerrungen führen ebenfalls zu Inkonsistenz.

Wenn die Stichprobe zufällig in eine **Behandlungs-** und eine **Kontrollgruppe** geteilt wird, kann der Effekt einer **Behandlung** valide geschätzt werden. Wird diese Zuteilung nicht randomisiert vorgenommen, sondern z.B. durch die Wahl der Probanden, kommt es zu Stichprobenverzerrungen.

Manchmal machen jedoch **natürliche Experimente** Beobachtungen verfügbar, die Beobachtungen aus einem Experiment mit zufällig zugewiesenen Behandlungen ähneln.

In unserem Beispiel der Bildungsrenditen könnte man möglicherweise gesetzliche Unterschiede zwischen verschiedenen Bundesstaaten ausnutzen, wenn Schüler in einigen Staaten länger in der Schule bleiben müssen als in anderen. Diese Information kann zur Konstruktion eines Instrumentes verwendet werden, um β konsistent zu schätzen.



Angrist, J. and W. Evans, "Children and Their Parents' Labor Supply: Evidence from Exogenous Variation in Family Size," *American Economic Review* 88(3), 1998, 450-477

Instrumentvariablenschätzung

Es gebe zusätzliche Informationen in Form von Instrumenten \mathbf{W}_t , die mit \mathbf{X}_t , jedoch nicht mit u_t korreliert sind.

Sei \mathbf{W} eine $n \times l$ Matrix, welche die Beobachtungen der Instrumente enthält.

Annahme 7.1

$$E(\mathbf{W}_t^\top \mathbf{u}_t) = \mathbf{0}. \quad (7.6)$$

Erweitere Annahme 4.4 auf

Annahme 7.2 Die Folge $\{y_t, \mathbf{X}_t, \mathbf{W}_t\}$ ist gemeinsam u.i.v.

Dann gilt mit dem GGZ (4.32), dass

$$\frac{\mathbf{W}^\top \mathbf{X}}{n} \xrightarrow{p} E(\mathbf{W}_t^\top \mathbf{X}_t) =: \mathbf{S}_{\mathbf{W}^\top \mathbf{X}}$$

Annahme 7.3 $\mathbf{S}_{\mathbf{W}^\top \mathbf{X}}$ hat vollen Spaltenrang.

Die Annahme impliziert $l \geq k$ (vgl. Matrixreader, Definition 13).

Analog gilt

$$\frac{\mathbf{W}^\top \mathbf{W}}{n} \xrightarrow{p} E(\mathbf{W}_t^\top \mathbf{W}_t) =: \mathbf{S}_{\mathbf{W}^\top \mathbf{W}} > 0. \quad (7.7)$$

Natürlich können (und sollten) vorbestimmte Regressoren (also solche Regressoren $\mathbf{X}_1 \subset \mathbf{X}$, für die wie in (4.44) $E(\mathbf{X}_{1t}^\top u_t) = \mathbf{0}$ gilt) als Instrumente verwendet werden.

Wenn $l = k$, ist das Modell exakt identifiziert (**just identified**). Wenn $l > k$ gilt, ist das Modell überidentifiziert (**overidentified**).

Der **IV-Schätzer** ist gegeben durch

$$\begin{aligned}\hat{\beta}_{IV} &= \arg \min_{\beta} \left\{ (\mathbf{y} - \mathbf{X}\beta)^\top \mathbf{P}_W (\mathbf{y} - \mathbf{X}\beta) \right\} \\ &= (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{y}.\end{aligned}\tag{7.8}$$

(7.8) zeigt man analog zu (1.21).

Wenn entweder $\mathbf{W} = \mathbf{X}$ oder $l = n$ gilt, ist $\hat{\beta}_{IV} = \hat{\beta}_{OLS}$.

Wie in (3.3) gilt $\hat{\beta}_{IV} = \beta + (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{u}$.

Wir sehen

$$\begin{aligned} \text{plim}_{n \rightarrow \infty} \left(\frac{\mathbf{X}^\top \mathbf{P}_W \mathbf{X}}{n} \right) &= \text{plim}_{n \rightarrow \infty} \left\{ \left(\frac{\mathbf{X}^\top \mathbf{W}}{n} \right) \left(\frac{\mathbf{W}^\top \mathbf{W}}{n} \right)^{-1} \left(\frac{\mathbf{W}^\top \mathbf{X}}{n} \right) \right\} \\ &= \mathbf{S}_{\mathbf{X}^\top \mathbf{W}} \mathbf{S}_{\mathbf{W}^\top \mathbf{W}}^{-1} \mathbf{S}_{\mathbf{W}^\top \mathbf{X}} =: \tilde{\mathbf{S}}. \end{aligned} \quad (7.9)$$

Annahme 7.2 und (7.6) liefern mit dem GGZ (4.32), dass

$$\text{plim}_{n \rightarrow \infty} \frac{\mathbf{W}^\top \mathbf{u}}{n} = \mathbf{0}.$$

Daraus folgt

$$\begin{aligned} \text{plim}_{n \rightarrow \infty} \left(\frac{\mathbf{X}^\top \mathbf{P}_W \mathbf{u}}{n} \right) &= \text{plim}_{n \rightarrow \infty} \left(\frac{\mathbf{X}^\top \mathbf{W} (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{u}}{n} \right) \\ &= \mathbf{S}_{\mathbf{X}^\top \mathbf{W}} \mathbf{S}_{\mathbf{W}^\top \mathbf{W}}^{-1} \mathbf{0} = \mathbf{0}. \end{aligned} \quad (7.10)$$

Damit gilt

$$\begin{aligned}\text{plim}_{n \rightarrow \infty} \hat{\beta}_{\text{IV}} &= \beta + \text{plim}_{n \rightarrow \infty} \left(\frac{\mathbf{X}^\top \mathbf{P}_W \mathbf{X}}{n} \right) \text{plim}_{n \rightarrow \infty} \left(\frac{\mathbf{X}^\top \mathbf{P}_W \mathbf{u}}{n} \right) \\ &= \beta,\end{aligned}$$

d.h. der IV-Schätzer ist unter Annahme (7.6) konsistent.

Die folgende Annahme 7.4 für den IV-Fall entspricht z.B. Annahme 4.6 mit (4.50).

Annahme 7.4

$$\text{Var}(\mathbf{W}_t^\top u_t) = \sigma^2 \mathbf{S}_{\mathbf{W}^\top \mathbf{W}}$$

Dann liefert der ZGWS (4.39) mit Annahme 7.2, dass

$$n^{1/2} \left(\frac{\mathbf{W}^\top \mathbf{u}}{n} \right) = n^{1/2} \left(\frac{\sum_{t=1}^n \mathbf{W}_t^\top u_t}{n} \right) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{S}_{\mathbf{W}^\top \mathbf{W}}) \quad (7.11)$$

Die asymptotische Verteilung folgt dann wie in (4.53) mit (7.11) und (7.9) als

$$\begin{aligned}
 n^{1/2}(\hat{\beta}_{IV} - \beta) &= \left(\frac{\mathbf{X}^\top \mathbf{P}_W \mathbf{X}}{n} \right)^{-1} \frac{\mathbf{X}^\top \mathbf{P}_W \mathbf{u}}{n^{1/2}} \\
 &\xrightarrow{d} \tilde{\mathbf{S}}^{-1} \mathbf{S}_{\mathbf{X}^\top \mathbf{W}} \mathbf{S}_{\mathbf{W}^\top \mathbf{W}}^{-1} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{S}_{\mathbf{W}^\top \mathbf{W}}) \\
 &\sim \mathcal{N}(\mathbf{0}, \mathbf{V}_{IV}),
 \end{aligned} \tag{7.12}$$

wobei mit (3.6)

$$\begin{aligned}
 \mathbf{V}_{IV} &= \tilde{\mathbf{S}}^{-1} \mathbf{S}_{\mathbf{X}^\top \mathbf{W}} \mathbf{S}_{\mathbf{W}^\top \mathbf{W}}^{-1} \sigma^2 \mathbf{S}_{\mathbf{W}^\top \mathbf{W}} (\tilde{\mathbf{S}}^{-1} \mathbf{S}_{\mathbf{X}^\top \mathbf{W}} \mathbf{S}_{\mathbf{W}^\top \mathbf{W}}^{-1})^\top \\
 &= \sigma^2 \tilde{\mathbf{S}}^{-1} \mathbf{S}_{\mathbf{X}^\top \mathbf{W}} \mathbf{S}_{\mathbf{W}^\top \mathbf{W}}^{-1} \mathbf{S}_{\mathbf{W}^\top \mathbf{X}} \tilde{\mathbf{S}}^{-1} = \sigma^2 \tilde{\mathbf{S}}^{-1} \tilde{\mathbf{S}} \tilde{\mathbf{S}}^{-1} \\
 &= \sigma^2 \tilde{\mathbf{S}}^{-1},
 \end{aligned}$$

was man auch schreiben kann als $\sigma^2 \text{plim}_{n \rightarrow \infty} \left(\frac{\mathbf{X}^\top \mathbf{P}_W \mathbf{X}}{n} \right)^{-1}$.

In der Praxis nennt man (7.8) vor allem IV, wenn $l = k$ und sonst **two-stage least squares**: 2SLS-Schätzer.

Dieser Name kann wie folgt erklärt werden. Betrachte

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \\ \mathbf{X} &= \mathbf{W}\boldsymbol{\Pi} + \mathbf{E}. \end{aligned}$$

In der **ersten Stufe** (first-stage) regressieren wir jede Spalte von \mathbf{X} auf die Spalten von \mathbf{W} :

$$\widehat{\boldsymbol{\Pi}} = (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X}.$$

In der **zweiten Stufe** (second-stage) regressieren wir \mathbf{y} auf

$$\widehat{\mathbf{X}} = \mathbf{W}\widehat{\boldsymbol{\Pi}}.$$

Damit erhalten wir

$$\begin{aligned}\hat{\beta}_{2\text{SLS}} &= (\widehat{\mathbf{X}}^\top \widehat{\mathbf{X}})^{-1} \widehat{\mathbf{X}}^\top \mathbf{y} \\ &= (\widehat{\boldsymbol{\Pi}}^\top \mathbf{W}^\top \mathbf{W} \widehat{\boldsymbol{\Pi}})^{-1} \widehat{\boldsymbol{\Pi}}^\top \mathbf{W}^\top \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{y} \\ &= \hat{\beta}_{\text{IV}}\end{aligned}$$

Hypothesentests

Tests können anhand der asymptotischen Verteilung (7.12) konstruiert werden:

$$n^{1/2}(\hat{\beta}_{IV} - \beta) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{V}_{IV})$$

Eine Teststatistik für $H_0 : \mathbf{R}^\top \beta = \mathbf{r}$, wobei die $k \times q$ Matrix \mathbf{R} vollen Spaltenrang hat, ist analog zu (4.62) gegeben durch

$$\begin{aligned} \tau &= n(\mathbf{R}^\top \hat{\beta}_{IV} - \mathbf{r})^\top (\mathbf{R}^\top \hat{\mathbf{V}}_{IV} \mathbf{R})^{-1} (\mathbf{R}^\top \hat{\beta}_{IV} - \mathbf{r}) \\ &= \frac{(\mathbf{R}^\top \hat{\beta}_{IV} - \mathbf{r})^\top \{ \mathbf{R}^\top (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{R} \}^{-1} (\mathbf{R}^\top \hat{\beta}_{IV} - \mathbf{r})}{\hat{\sigma}^2} \end{aligned} \quad (7.13)$$

wobei

$$\hat{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\beta}_{IV})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}_{IV})}{n}$$

Unter H_0 ist $\tau \xrightarrow{d} \chi_q^2$. Wie gewöhnlich verwerfen wir H_0 , wenn τ den kritischen Wert, das $1 - \alpha$ -Quantil der χ_q^2 -Verteilung, übersteigt.

Testen von überidentifizierenden Restriktionen

Nehmen wir an, dass $\mathbf{W} = (\mathbf{W}_1, \mathbf{W}_2)$ wobei \mathbf{W}_1 k Spalten und \mathbf{W}_2 $l - k$ Spalten hat. Der Grad der Überidentifizierung ist gegeben durch $l - k$.

Um β konsistent zu schätzen, benötigen wir lediglich \mathbf{W}_1 . Wir können das Modell formulieren als

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{W}_2\gamma + \mathbf{u}, \quad (7.14)$$

$$\mathbf{X} = \mathbf{W}\mathbf{\Pi} + \mathbf{E}, \quad (7.15)$$

und die überidentifizierenden Restriktionen $H_0: \gamma = \mathbf{0}$ testen.

Definiere $\hat{\mathbf{u}}_{IV} = \mathbf{y} - \mathbf{X}\hat{\beta}_{IV}$. Eine Teststatistik (von der man zeigen kann, dass sie ein Spezialfall von (7.13) ist) für H_0 ist

$$\tau_{IV} = \frac{n\hat{\mathbf{u}}_{IV}^{\top} \mathbf{P}_{\{\mathbf{P}_W(\mathbf{X}, \mathbf{W}_2)\}} \hat{\mathbf{u}}_{IV}}{\hat{\mathbf{u}}_{IV}^{\top} \hat{\mathbf{u}}_{IV}} = \frac{n\hat{\mathbf{u}}_{IV}^{\top} \mathbf{P}_W \hat{\mathbf{u}}_{IV}}{\hat{\mathbf{u}}_{IV}^{\top} \hat{\mathbf{u}}_{IV}}.$$

Die zweite Gleichung folgt aus

$$P_{\{P_W(X, W_2)\}} = P_W,$$

da

$$P_W(X, W_2) = W(W^\top W)^{-1}W^\top(X, W_2) = WC,$$

wobei C , welches implizit definiert ist, nicht singulär ist, ähnlich zu Annahme 7.3. Daher gilt wegen, wie in (2.13),

$$P_{WC} = WC((WC)^\top WC)^{-1}(WC)^\top = WCC^{-1}(W^\top W)^{-1}(C^\top)^{-1}C^\top W^\top,$$

dass

$$P_{\{P_W(X, W_2)\}} = P_{WC} = P_W.$$

Unter $H_0 : \gamma = \mathbf{0}$ hat die **Teststatistik für überidentifizierende Restriktionen** eine asymptotische Chi-Quadrat-Verteilung:

$$\tau_{IV} = \frac{n\hat{u}_{IV}^\top P_W \hat{u}_{IV}}{\hat{u}_{IV}^\top \hat{u}_{IV}} \xrightarrow{d} \chi^2_{I-k}. \quad (7.16)$$

Der Test ist auch bekannt als **Sargan-Test**.

dwh Hansen IV.R

Diese Tests auf Überidentifikation werden ebenfalls **Tests für die Validität der Instrumente** genannt. Wie sollten wir eine Ablehnung interpretieren?

- Sind Instrumente mit dem Fehlerterm korreliert und daher unzulässig?
- Wurden einige Variablen fälschlicherweise ausgelassen?

Im Folgenden betrachten wir einen Spezialfall genauer.

Ein Spezialfall: „Test auf Exogenität“

Betrachten wir ein Modell, bei dem wir nicht sicher sind, ob einige erklärende Variablen exogen sind. Sei

$$\mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \mathbf{u}, \quad (7.17)$$

wobei wir annehmen, dass $E(\mathbf{X}_{1t}^\top \mathbf{u}_t) = \mathbf{0}$, aber wir uns nicht sicher sind, ob $E(\mathbf{X}_{2t}^\top \mathbf{u}_t) = \mathbf{0}$. Daher sind die Variablen in \mathbf{X}_2 potentiell endogen.

Sei \mathbf{X}_1 $n \times k_1$ und \mathbf{X}_2 $n \times k_2$. Daher ist $k = k_1 + k_2$. Wir haben l Instrumente \mathbf{W} , bestehend aus den k_1 Variablen in $\mathbf{W}_1 = \mathbf{X}_1$ und k_2 Instrumente \mathbf{W}_2 , die verwendet werden könnten, wenn \mathbf{X}_2 tatsächlich endogen ist. Außerdem haben wir noch $l - k$ überschüssige Instrumente \mathbf{W}_3 .

Sei

$$\mathbf{X}_2 = \mathbf{W}\Pi + \mathbf{E}. \quad (7.18)$$

Wenn \mathbf{X}_2 endogen ist, dann sind die Elemente von \mathbf{E} mit \mathbf{u} korreliert.
Modelliere \mathbf{u} durch \mathbf{E} ,

$$\mathbf{u} = \mathbf{E}\boldsymbol{\lambda} + \varepsilon,$$

wobei wir annehmen können (lineare Projektion), dass \mathbf{E} und ε unkorreliert sind. Schreibe (7.17) dann als

$$\mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \mathbf{E}\boldsymbol{\lambda} + \varepsilon, \quad (7.19)$$

wobei somit nun \mathbf{X}_2 und ε unkorreliert sind.

Um

$$E(u_t|\mathbf{X}_{2t}) = E(u_t|\mathbf{E}_t) = \mathbf{E}_t\boldsymbol{\lambda} = 0,$$

$t = 1, \dots, n$ zu testen, testen wir also $\boldsymbol{\lambda} = \mathbf{0}$.

Allerdings ist \mathbf{E} nicht beobachtbar. Wir können dennoch weiter machen.

Betrachte zunächst den exakt identifizierten Fall $l = k$. Für jede Spalte in \mathbf{X}_2 haben wir ein Instrument in \mathbf{W}_2 und es gibt keine weiteren Instrumente in \mathbf{W}_3 .

Schreibe das Modell (7.19) mit $\mathbf{E} = \mathbf{X}_2 - \mathbf{W}\Pi$ als

$$\mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2(\beta_2 + \lambda) - \mathbf{W}\Pi\lambda + \varepsilon, \quad (7.20)$$

wobei

$$\mathbf{W}\Pi\lambda = (\mathbf{X}_1, \mathbf{W}_2) \begin{pmatrix} \Pi_1 \\ \Pi_2 \end{pmatrix} \lambda = \mathbf{X}_1\Pi_1\lambda + \mathbf{W}_2\Pi_2\lambda. \quad (7.21)$$

Also kann (7.20) geschrieben werden als

$$\begin{aligned} \mathbf{y} &= \mathbf{X}_1(\beta_1 - \Pi_1\lambda) + \mathbf{X}_2(\beta_2 + \lambda) - \mathbf{W}_2\Pi_2\lambda + \varepsilon \\ &= \mathbf{X}_1\beta_1^* + \mathbf{X}_2\beta_2^* + \mathbf{W}_2\gamma + \varepsilon, \end{aligned} \quad (7.22)$$

wobei γ nicht restringiert ist, da Π_2 für $l = k$ quadratisch und nichtsingulär ist.

Testen von $\gamma = \mathbf{0}$ in (7.22) ist lediglich ein Test für die Validität der Instrumente, welche nun gegeben sind durch $(\mathbf{W}_1, \mathbf{X}_2, \mathbf{W}_2)$.

Die Teststatistik kann ähnlich zu (7.16) formuliert werden, wobei nun der restringierte IV-Schätzer der OLS-Schätzer einer Regression auf \mathbf{X}_1 und \mathbf{X}_2 ist.

Die Teststatistik und deren asymptotische Verteilung unter $H_0 : \lambda = \mathbf{0}$ oder äquivalent $\gamma = \mathbf{0}$ sind gegeben durch

$$\tau = \frac{\hat{\mathbf{u}}_{\text{OLS}}^{\top} \mathbf{P}_{(\mathbf{X}, \mathbf{W}_2)} \hat{\mathbf{u}}_{\text{OLS}}}{\hat{\sigma}^2} = \frac{\hat{\mathbf{u}}_{\text{OLS}}^{\top} \mathbf{P}_{\{(I - \mathbf{P}_X) \mathbf{W}_2\}} \hat{\mathbf{u}}_{\text{OLS}}}{\hat{\sigma}^2} \xrightarrow{d} \chi_{k_2}^2. \quad (7.23)$$

Die zweite Gleichung folgt mit dem FWL-Theorem 2.4 und (2.9), da

$$\mathbf{P}_{(\mathbf{X}, \mathbf{W}_2)} \hat{\mathbf{u}}_{\text{OLS}} = \mathbf{P}_{\{(I - \mathbf{P}_X) \mathbf{W}_2\}} (I - \mathbf{P}_X) \hat{\mathbf{u}}_{\text{OLS}} = \mathbf{P}_{\{(I - \mathbf{P}_X) \mathbf{W}_2\}} \hat{\mathbf{u}}_{\text{OLS}}$$

Die Anzahl der Freiheitsgrade entspricht der Anzahl der zusätzlichen Instrumente in $(\mathbf{W}_1, \mathbf{X}_2, \mathbf{W}_2)$, also k_2 .

dwh Hansen IV.R

Durbin-Wu-Hausman-Test

Im Fall von Überidentifikation mit $l > k$ enthält \mathbf{W}_3 zusätzliche Instrumente.

Wir können \mathbf{E} in (7.19) durch $\hat{\mathbf{E}}$, die Residuen der Regression der ersten Stufe aus (7.18), ersetzen, da diese in Wahrscheinlichkeit gegen \mathbf{E} konvergieren.

Wir können $\gamma = \mathbf{0}$ in der Regression

$$\mathbf{y} = \mathbf{X}_1\beta_1 + \mathbf{X}_2\beta_2 + \hat{\mathbf{E}}\gamma + \text{residuals}. \quad (7.24)$$

testen. Die Teststatistik und die asymptotische Verteilung sind gegeben durch

$$\tau = \frac{\hat{\mathbf{u}}_{\text{OLS}}^\top \mathbf{P}_{(\mathbf{X}, \hat{\mathbf{E}})} \hat{\mathbf{u}}_{\text{OLS}}}{\hat{\sigma}^2} = \frac{\hat{\mathbf{u}}_{\text{OLS}}^\top \mathbf{P}_{(\mathbf{X}, \mathbf{W}\hat{\Pi})} \hat{\mathbf{u}}_{\text{OLS}}}{\hat{\sigma}^2} = \frac{\hat{\mathbf{u}}_{\text{OLS}}^\top \mathbf{P}_{\{(I - \mathbf{P}_\mathbf{X})\mathbf{W}\hat{\Pi}\}} \hat{\mathbf{u}}_{\text{OLS}}}{\hat{\sigma}^2} \xrightarrow{d} \chi_{k_2}^2. \quad (7.25)$$

Im exakt identifizierten Fall entspricht diese Teststatistik der in (7.23), da $\hat{\Pi}$ dann quadratisch und invertierbar ist, vgl. (2.13). dwh Hansen IV.R

Der Kontrastvektor

Betrachten wir den Fall, ähnlich zum vorherigen Abschnitt, bei dem alle Variablen in \mathbf{X} potenziell endogen sind und wir für jede ein Instrument haben, das in der $n \times k$ Matrix \mathbf{W} enthalten ist. Da $l = k$, ist das Modell exakt identifiziert, wenn nicht davon ausgegangen werden kann, dass die erklärenden Variablen exogen sind.

Die Teststatistik auf Exogenität kann wie in (7.23) formuliert werden

$$\tau = \frac{\hat{\mathbf{u}}_{\text{OLS}}^{\top} \mathbf{P}_{(\mathbf{X}, \mathbf{W})} \hat{\mathbf{u}}_{\text{OLS}}}{\hat{\sigma}^2} \xrightarrow{d} \chi_k^2.$$

Diese Teststatistik kann ebenfalls auf Grundlage des **Kontrastvektors**, der Differenz zwischen $\hat{\beta}_{\text{OLS}}$ und $\hat{\beta}_{\text{IV}}$, konstruiert werden.

Da

$$\begin{aligned}\hat{\beta}_{IV} &= (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{y}, \\ \hat{\beta}_{OLS} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y},\end{aligned}$$

erhalten wir

$$\begin{aligned}\hat{\beta}_{IV} - \hat{\beta}_{OLS} &= (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W (\mathbf{I} - \mathbf{P}_X) \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \hat{\mathbf{u}}_{OLS}.\end{aligned}\tag{7.26}$$

Wir sehen

$$\mathbf{P}_{(\mathbf{X}, \mathbf{W})} = \mathbf{P}_{(\mathbf{X}, \mathbf{P}_W \mathbf{X})}$$

Dies gilt, da

$$(\mathbf{X}, \mathbf{P}_W \mathbf{X}) = (\mathbf{X}, \mathbf{W}) \mathbf{A},$$
$$\mathbf{A} = \begin{pmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{X} \end{pmatrix},$$

wobei \mathbf{A} eine nichtsinguläre $2k \times 2k$ Matrix ist.

Daher kann $\hat{\sigma}^2 \tau$ umgeschrieben werden in

$$\begin{aligned} \hat{\sigma}^2 \tau &= \hat{\mathbf{u}}_{\text{OLS}}^\top \mathbf{P}_{(\mathbf{X}, \mathbf{P}_W \mathbf{X})} \hat{\mathbf{u}}_{\text{OLS}} \\ &= \hat{\mathbf{u}}_{\text{OLS}}^\top (\mathbf{X}, \mathbf{P}_W \mathbf{X}) \begin{pmatrix} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{P}_W \mathbf{X} \\ \mathbf{X}^\top \mathbf{P}_W \mathbf{X} & \mathbf{X}^\top \mathbf{P}_W \mathbf{X} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}^\top \\ \mathbf{X}^\top \mathbf{P}_W \end{pmatrix} \hat{\mathbf{u}}_{\text{OLS}}, \end{aligned}$$

Und da $\mathbf{X}^\top \hat{\mathbf{u}}_{\text{OLS}} = \mathbf{0}$,

$$\hat{\sigma}^2_\tau = \hat{\mathbf{u}}_{\text{OLS}}^\top \mathbf{P}_W \mathbf{X} (\mathbf{0}, \mathbf{I}_k) \left(\begin{array}{cc} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{P}_W \mathbf{X} \\ \mathbf{X}^\top \mathbf{P}_W \mathbf{X} & \mathbf{X}^\top \mathbf{P}_W \mathbf{X} \end{array} \right)^{-1} \left(\begin{array}{c} \mathbf{0} \\ \mathbf{I}_k \end{array} \right) \mathbf{X}^\top \mathbf{P}_W \hat{\mathbf{u}}_{\text{OLS}} \quad (7.27)$$

Hierbei ist

$$(\mathbf{0}, \mathbf{I}_k) \left(\begin{array}{cc} \mathbf{X}^\top \mathbf{X} & \mathbf{X}^\top \mathbf{P}_W \mathbf{X} \\ \mathbf{X}^\top \mathbf{P}_W \mathbf{X} & \mathbf{X}^\top \mathbf{P}_W \mathbf{X} \end{array} \right)^{-1} \left(\begin{array}{c} \mathbf{0} \\ \mathbf{I}_k \end{array} \right)$$

das südöstliche Element der Inverse, mit Matrixreader, Resultat 23 gleich

$$\begin{aligned} & [\mathbf{X}^\top \mathbf{P}_W \mathbf{X} - \mathbf{X}^\top \mathbf{P}_W \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{X}]^{-1} \\ &= (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} [(\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} - (\mathbf{X}^\top \mathbf{X})^{-1}]^{-1} (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \end{aligned}$$

Gleichzeitig ist, wiederum mit Resultat 23, $[(\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} - (\mathbf{X}^\top \mathbf{X})^{-1}]^{-1}$ das südöstliche Element von

$$\left(\begin{array}{cc} \mathbf{X}^\top \mathbf{X} & \mathbf{I}_k \\ \mathbf{I}_k & (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \end{array} \right)^{-1}$$

Demnach können wir (7.27) schreiben als

$$\begin{aligned}
 & \hat{\mathbf{u}}_{\text{OLS}}^{\top} \mathbf{P}_W \mathbf{X} (\mathbf{0}, (\mathbf{X}^{\top} \mathbf{P}_W \mathbf{X})^{-1}) \times \\
 & \quad \left(\begin{array}{cc} \mathbf{X}^{\top} \mathbf{X} & \mathbf{I}_k \\ \mathbf{I}_k & (\mathbf{X}^{\top} \mathbf{P}_W \mathbf{X})^{-1} \end{array} \right)^{-1} \left(\begin{array}{c} \mathbf{0} \\ (\mathbf{X}^{\top} \mathbf{P}_W \mathbf{X})^{-1} \end{array} \right) \mathbf{X}^{\top} \mathbf{P}_W \hat{\mathbf{u}}_{\text{OLS}} \\
 (7.26) \quad & \stackrel{=}{=} (\hat{\beta}_{\text{IV}} - \hat{\beta}_{\text{OLS}})^{\top} (\mathbf{0}, \mathbf{I}_k) \left(\begin{array}{cc} \mathbf{X}^{\top} \mathbf{X} & \mathbf{I}_k \\ \mathbf{I}_k & (\mathbf{X}^{\top} \mathbf{P}_W \mathbf{X})^{-1} \end{array} \right)^{-1} \left(\begin{array}{c} \mathbf{0} \\ \mathbf{I}_k \end{array} \right) (\hat{\beta}_{\text{IV}} - \hat{\beta}_{\text{OLS}}) \\
 & = (\hat{\beta}_{\text{IV}} - \hat{\beta}_{\text{OLS}})^{\top} \{ (\mathbf{X}^{\top} \mathbf{P}_W \mathbf{X})^{-1} - (\mathbf{X}^{\top} \mathbf{X})^{-1} \}^{-1} (\hat{\beta}_{\text{IV}} - \hat{\beta}_{\text{OLS}})
 \end{aligned}$$

Wir können τ kann also mithilfe des Kontrastvektors, (4.61) und (7.12) schreiben als

$$\tau = n(\hat{\beta}_{IV} - \hat{\beta}_{OLS})^\top \left\{ \hat{\mathbf{V}}_{IV} - \hat{\mathbf{V}}_{OLS} \right\}^{-1} (\hat{\beta}_{IV} - \hat{\beta}_{OLS}),$$

(plus einem Term, der aufgrund der Wahl von $\hat{\sigma}^2$ in Wahrscheinlichkeit gegen null konvergiert, welche sich von der gewöhnlichen, für $\hat{\mathbf{V}}_{OLS}$ verwendeten, unterscheidet—wegen (4.38) und unter der Nullhypothese würde dieser Unterschied die asymptotische Verteilung nicht beeinflussen.)

Wenn die Variablen in \mathbf{X} tatsächlich exogen sind, dann ist OLS asymptotisch effizient und IV konsistent. Wenn die erklärenden Variablen endogen sind, dann ist OLS inkonsistent, jedoch ist IV immer noch konsistent.

Dies ist ein Beispiel für einen sogenannten **Hausman-Test**.

Hausman-Tests

Ein Hausman-Test kann in einem generellen Kontext mit zwei Schätzern $\hat{\theta}$ und $\tilde{\theta}$ formuliert werden, welche die folgenden Bedingungen erfüllen:

- (i) $\hat{\theta}$ ist asymptotisch effizient unter H_0 und inkonsistent unter H_1 ,
- (ii) $\tilde{\theta}$ ist sowohl unter H_0 als auch unter H_1 konsistent.

Da $\hat{\theta}$ unter H_0 asymptotisch effizient ist, ist es asymptotisch unkorreliert mit allen konsistenten Schätzer von $\mathbf{0}$.

Insbesondere ist $\hat{\theta}$ unter der Nullhypothese asymptotisch unkorreliert mit dem Kontrastvektor $\hat{\theta} - \tilde{\theta}$, das bedeutet

$$n^{1/2} \begin{pmatrix} \hat{\theta} - \theta \\ \hat{\theta} - \tilde{\theta} \end{pmatrix} \xrightarrow{d} \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} \text{AVar}(\hat{\theta}) & \mathbf{0} \\ \mathbf{0} & \text{AVar}(\hat{\theta} - \tilde{\theta}) \end{pmatrix} \right).$$

Daher gilt mit (3.6)

$$\begin{aligned} n^{1/2} \begin{pmatrix} \hat{\boldsymbol{\theta}} - \boldsymbol{\theta} \\ \tilde{\boldsymbol{\theta}} - \boldsymbol{\theta} \end{pmatrix} &= \begin{pmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{I} & -\mathbf{I} \end{pmatrix} n^{1/2} \begin{pmatrix} \hat{\boldsymbol{\theta}} - \boldsymbol{\theta} \\ \hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}} \end{pmatrix} \\ &\xrightarrow{d} \mathcal{N} \left(\mathbf{0}, \begin{pmatrix} \text{AVar}(\hat{\boldsymbol{\theta}}) & \text{AVar}(\hat{\boldsymbol{\theta}}) \\ \text{AVar}(\hat{\boldsymbol{\theta}}) & \text{AVar}(\hat{\boldsymbol{\theta}}) + \text{AVar}(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}) \end{pmatrix} \right). \end{aligned}$$

Daraus folgt

$$\text{AVar}(\tilde{\boldsymbol{\theta}}) = \text{AVar}(\hat{\boldsymbol{\theta}}) + \text{AVar}(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}),$$

so dass $\text{AVar}(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}) = \text{AVar}(\tilde{\boldsymbol{\theta}}) - \text{AVar}(\hat{\boldsymbol{\theta}})$ bzw.

$$n^{1/2}(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}) \xrightarrow{d} \mathcal{N} \left(\mathbf{0}, \text{AVar}(\tilde{\boldsymbol{\theta}}) - \text{AVar}(\hat{\boldsymbol{\theta}}) \right).$$

Wir erhalten also unter der Nullhypothese

$$n(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}})^\top \left\{ \text{AVar}(\tilde{\boldsymbol{\theta}}) - \text{AVar}(\hat{\boldsymbol{\theta}}) \right\}^{-1} (\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}}) \xrightarrow{d} \chi_k^2,$$

wenn $\boldsymbol{\theta}$ ein k Vektor und $\text{AVar}(\hat{\boldsymbol{\theta}} - \tilde{\boldsymbol{\theta}})$ nichtsingulär ist.

Heteroskedastie-robuste Tests

Heteroskedastie-konsistente Kovarianzmatrix-Schätzung in Gegenwart von Instrumenten erfolgt ähnlich wie für KQ-Schätzung.

D.h., wenn $\mathbf{u}|\mathbf{W} \sim (\mathbf{0}, \mathbf{D})$ und \mathbf{D} diagonal, jedoch nicht notwendigerweise gleich $\sigma^2 \mathbf{I}$ ist, dann gilt unter den üblichen anderen Annahmen

$$n^{1/2}(\hat{\beta}_{IV} - \beta) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{V}_{IV,rob}),$$

wobei, analog zu (4.52),

$$\mathbf{V}_{IV,rob} = \text{plim}_{n \rightarrow \infty} n(\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_W \mathbf{D} \mathbf{P}_W \mathbf{X} (\mathbf{X}^\top \mathbf{P}_W \mathbf{X})^{-1}.$$

Um $\mathbf{W}^\top \mathbf{D} \mathbf{W}$ in $\mathbf{P}_W \mathbf{D} \mathbf{P}_W$ zu schätzen, können wir, analog zu (4.65),

$$\mathbf{W}^\top \hat{\mathbf{D}} \mathbf{W} = \sum_{t=1}^n \mathbf{w}_t^\top \mathbf{w}_t \hat{u}_{IV,t}^2$$

verwenden, wobei \mathbf{W}_t eine Reihe von \mathbf{W} ist. Dann können Teststatistiken mit der robusten Schätzung $\hat{\mathbf{V}}_{IV,rob}$ formuliert werden. `IV robust.R`