

Kapitel 6

Das OLS Regressionsmodell in Matrixnotation

“What I cannot create, I do not understand.”
(Richard P. Feynman)

Dieses Kapitel bietet im wesentlichen eine Wiederholung der früheren Kapitel. Wir gehen nur in einem einzigen Punkt über die früheren Kapitel hinaus, indem wir die Matrixnotation einführen. Dies vereinfacht die Darstellung der multiplen Regression ganz erheblich.

6.1 OLS-Schätzung in Matrixschreibweise

Das multiple Regressionsmodell kann in Matrixschreibweise deutlich einfacher dargestellt werden als unter Verwendung der bisherigen Summennotation.

Wir gehen wieder von einem einfachen linearen Zusammenhang in der Grundgesamtheit aus, d.h. von einem datengenerierenden Prozess, der durch folgende PRF beschrieben werden kann

$$y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + \cdots + x_{ih}\beta_h + \cdots + x_{ik}\beta_k + \varepsilon_i$$

wobei der erste Index i wie üblich die Beobachtung und der zweite Index $h = 1, \dots, k$ die Variable bezeichnet. Darüber hinaus wollen wir in diesem Kapitel weiterhin annehmen, dass alle Gauss Markov Annahmen erfüllt seien.

Für die n Beobachtungen der Stichprobe können wir die SRF (*‘Sample Regression Function’*) ausführlicher schreiben, nämlich für jede einzelne Beobachtung i

$$\begin{aligned} y_1 &= x_{11}\hat{\beta}_1 + x_{12}\hat{\beta}_2 + \cdots + x_{1k}\hat{\beta}_k + \hat{\varepsilon}_1 \\ y_2 &= x_{21}\hat{\beta}_1 + x_{22}\hat{\beta}_2 + \cdots + x_{2k}\hat{\beta}_k + \hat{\varepsilon}_2 \\ \vdots &= \quad \vdots \quad + \quad \vdots \quad + \cdots + \quad \vdots \quad + \quad \vdots \\ y_i &= x_{i1}\hat{\beta}_1 + x_{i2}\hat{\beta}_2 + \cdots + x_{ik}\hat{\beta}_k + \hat{\varepsilon}_i \\ \vdots &= \quad \vdots \quad + \quad \vdots \quad + \cdots + \quad \vdots \quad + \quad \vdots \\ y_n &= x_{n1}\hat{\beta}_1 + x_{n2}\hat{\beta}_2 + \cdots + x_{nk}\hat{\beta}_k + \hat{\varepsilon}_n \end{aligned}$$

wobei n die Größe der Stichprobe und h die Laufvariable über die Anzahl der erklärenden Variablen (inklusive Regressionskonstante) bezeichnet ($h = 1, \dots, k$). Wie in

der Matrixschreibweise üblich wird mit dem ersten Subindex die Beobachtung – also Zeile – bezeichnet ($i = 1, \dots, n$), und der zweite Subindex ($h = 1, \dots, k$) bezeichnet die erklärende Variable, also Spalte; x_{ih} bezeichnet also die i -te Beobachtung von Variable h , bzw. das Element in Zeile i und Spalte h .

Man beachte, dass wir hier keine speziellen Symbole für Regressionskonstante und Interzept eingeführt haben. Dies ist nicht erforderlich, denn das Interzept ist einfach der Koeffizient der Regressionskonstanten, d.h. eines Einsen-Vektors. Sollte die Regression ein Interzept enthalten ist einfach eine der x -Variablen ein Einsen-Vektor.

Dieses Gleichungssystem kann einfach in Matrixschreibweise geschrieben werden

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} + \begin{pmatrix} \hat{\varepsilon}_1 \\ \hat{\varepsilon}_2 \\ \vdots \\ \hat{\varepsilon}_n \end{pmatrix}$$

oder etwas kompakter

$$\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\varepsilon}}$$

Die PRF schreiben wir analog $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$.

Zur Notation: Üblichen Konventionen folgend bezeichnen fettgedruckte Kleinbuchstaben Vektoren und fettgedruckte Großbuchstaben Matrizen.

6.1.1 Alternative Schreibweisen

Wenn wir aus der i -ten Zeile der \mathbf{X} Matrix einen $k \times 1$ Spaltenvektor formen (für eine ausführlichere Darstellung siehe Appendix 6.A.1, Seite 40).

$$\mathbf{x}_{i\cdot} = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ik} \end{pmatrix}$$

kann das Modell $\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\varepsilon}}$ alternativ für eine einzelne Beobachtung auch elementenweise geschrieben werden als

$$\begin{aligned} y_i &= \mathbf{x}'_{i\cdot}\hat{\boldsymbol{\beta}} + \hat{\varepsilon}_i = \begin{pmatrix} x_{i1} & x_{i2} & \cdots & x_{ik} \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} + \hat{\varepsilon}_i \\ &= x_{i1}\hat{\beta}_1 + x_{i2}\hat{\beta}_2 + \cdots + x_{ik}\hat{\beta}_k + \hat{\varepsilon}_i, \quad i = 1, \dots, n \end{aligned}$$

Wenn keine Gefahr von Missverständnissen besteht schreiben wir im Folgenden statt $\mathbf{x}_{i\cdot}$ kürzer \mathbf{x}_i , also insgesamt n Spaltenvektoren, von denen jeder die Dimension $k \times 1$ hat.

Mit Hilfe des Spaltenvektors \mathbf{x}_i , kann auch die \mathbf{X} alternativ geschrieben werden als

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}'_1 \\ \mathbf{x}'_2 \\ \vdots \\ \mathbf{x}'_n \end{pmatrix}$$

Die $k \times k$ Matrix $\mathbf{X}'\mathbf{X}$ kann auch geschrieben werden als $\sum_{i=1}^n \mathbf{x}_i \mathbf{x}'_i$, da

$$\begin{aligned} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}'_i &= \sum_{i=1}^n \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ik} \end{pmatrix} (x_{i1} \ x_{i2} \ \cdots \ x_{ik}) \\ &= \begin{pmatrix} \sum_{i=1}^n x_{i1}^2 & \sum_{i=1}^n x_{i1}x_{i2} & \cdots & \sum_{i=1}^n x_{i1}x_{ik} \\ \sum_{i=1}^n x_{i2}x_{i1} & \sum_{i=1}^n x_{i2}^2 & \cdots & \sum_{i=1}^n x_{i2}x_{ik} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^n x_{ik}x_{i1} & \sum_{i=1}^n x_{ik}x_{i2} & \cdots & \sum_{i=1}^n x_{ik}^2 \end{pmatrix} \\ &= \begin{pmatrix} x_{11} & x_{21} & \cdots & x_{n1} \\ x_{12} & x_{22} & \cdots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1k} & x_{2k} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix} \\ &= \mathbf{X}'\mathbf{X} \end{aligned}$$

Genauso ist $\mathbf{X}'\mathbf{y} = \sum_i \mathbf{x}_i y_i$, weil

$$\begin{aligned} \sum_{i=1}^n \mathbf{x}_i y_i &= \sum_{i=1}^n \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ik} \end{pmatrix} y_i \\ &= \begin{pmatrix} x_{11}y_1 + x_{21}y_2 + \cdots + x_{n1}y_n \\ x_{12}y_1 + x_{22}y_2 + \cdots + x_{n2}y_n \\ \vdots \\ x_{1k}y_1 + x_{2k}y_2 + \cdots + x_{nk}y_n \end{pmatrix} \\ &= \begin{pmatrix} x_{11} & x_{21} & \cdots & x_{n1} \\ x_{12} & x_{22} & \cdots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1k} & x_{2k} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \\ &= \mathbf{X}'\mathbf{y} \end{aligned}$$

Also kann der OLS Schätzer auch in Vektornotation geschrieben werden

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}'_i \right)^{-1} \sum_{i=1}^n \mathbf{x}_i y_i$$

Diese Vektor-Schreibweise findet sich häufig in Lehrbüchern.

6.2 Die OLS Schätzfunktion

Wie üblich gehen wir wieder davon aus, dass der Parametervektor der Grundgesamtheit β nicht beobachtbar ist, aber wenn wir aus der Grundgesamtheit eine Stichprobe mit dem Umfang n ziehen erwarten wir darin einen ähnlichen Zusammenhang wie in der Grundgesamtheit zu finden. Unser Problem besteht also darin, eine möglichst gute Schätzfunktion $\hat{\beta}$ für den ‘wahren’ Koeffizientenvektor β zu finden.

Die Grundidee der OLS-Schätzung (*Ordinary Least Squares*) besteht darin, den Residuenvektor $\hat{\varepsilon}$ so zu wählen, dass die *Summe der quadrierten Abweichungen* in der Stichprobe (d.h. $\sum_{i=1}^n \hat{\varepsilon}_i^2$) so klein wie möglich wird. Dies ist wieder die übliche Minimierungsaufgabe.

Zur Erinnerung, die SRF $y = X\beta + \varepsilon$ ist ausführlich geschrieben

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} + \begin{pmatrix} \hat{\varepsilon}_1 \\ \hat{\varepsilon}_2 \\ \vdots \\ \hat{\varepsilon}_n \end{pmatrix}$$

Die einzelnen Spalten der X -Matrix sind die erklärenden Variablen, und der $n \times 1$ Vektor $\hat{\varepsilon}$ sind die zu minimierenden Stichprobenresiduen.

Die Summe der quadrierten Residuen kann in Matrixschreibweise einfach als inneres Produkt geschrieben werden

$$\hat{\varepsilon}'\hat{\varepsilon} = \begin{pmatrix} \hat{\varepsilon}_1 & \hat{\varepsilon}_2 & \cdots & \hat{\varepsilon}_n \end{pmatrix} \begin{pmatrix} \hat{\varepsilon}_1 \\ \hat{\varepsilon}_2 \\ \vdots \\ \hat{\varepsilon}_n \end{pmatrix} = \sum_{i=1}^n \hat{\varepsilon}_i^2$$

Da per Definition $\hat{\varepsilon} = y - X\hat{\beta}$ suchen wir den Vektor $\hat{\beta}$, der die Quadratsumme der Residuen minimiert. Dies wird manchmal geschrieben als

$$\hat{\beta} = \arg \min \left[(y - X\hat{\beta})'(y - X\hat{\beta}) \right]$$

wobei die Funktion $\arg \min$ der Bestimmung der Stelle dient, an der die Funktion $(y - X\hat{\beta})'(y - X\hat{\beta})$ ihr Minimum annimmt.

Die Quadratsumme der Residuen ist

$$\begin{aligned} \hat{\varepsilon}'\hat{\varepsilon} &= (y - X\hat{\beta})'(y - X\hat{\beta}) \\ &= (y' - \hat{\beta}'X')(y - X\hat{\beta}) \\ &= y'y - \hat{\beta}'X'y - y'X\hat{\beta} + \hat{\beta}'X'X\hat{\beta} \\ &= y'y - 2\hat{\beta}'X'y + \hat{\beta}'X'X\hat{\beta} \end{aligned}$$

Man beachte, dass $\hat{\beta}'X'y = (\hat{\beta}'X'y)' = y'X\hat{\beta}$, weil beide Terme die Dimension (1×1) haben, also Skalare sind, und die Transponierte eines Skalars der Skalar selbst ist.

Um die Quadratsumme der Residuen zu minimieren müssen wir den obigen Ausdruck nach dem Vektor $\hat{\beta}$ ableiten und diese Ableitungen Null setzen. Dazu benötigen wir zwei Rechenregeln für das Differenzieren von Matrizen.

Diese beiden Rechenregeln sind

1.

$$\frac{\partial \hat{\beta}' \mathbf{X}' \mathbf{y}}{\partial \hat{\beta}} = \mathbf{X}' \mathbf{y}$$

2.

$$\frac{\partial \hat{\beta}' \mathbf{X}' \mathbf{X} \hat{\beta}}{\partial \hat{\beta}} = 2 \mathbf{X}' \mathbf{X} \hat{\beta}$$

und werden im Appendix zu diesem Kapitel ausführlich erläutert (siehe Appendix 6.A.2 Seite 40).

Diese zwei Rechenregeln können wir nun für unsere Minimierungsaufgabe verwenden.

$$\begin{aligned} \min_{\hat{\beta}} (\hat{\varepsilon}' \hat{\varepsilon}) &= \min_{\hat{\beta}} (\mathbf{y} - \mathbf{X} \hat{\beta})' (\mathbf{y} - \mathbf{X} \hat{\beta}) \\ &= \min_{\hat{\beta}} [\mathbf{y}' \mathbf{y} - 2 \hat{\beta}' \mathbf{X}' \mathbf{y} + \hat{\beta}' \mathbf{X}' \mathbf{X} \hat{\beta}] \\ \frac{\partial (\hat{\varepsilon}' \hat{\varepsilon})}{\partial \hat{\beta}} &= -2 \mathbf{X}' \mathbf{y} + 2 \mathbf{X}' \mathbf{X} \hat{\beta} \stackrel{!}{=} 0 \\ \mathbf{X}' \mathbf{X} \hat{\beta} &= \mathbf{X}' \mathbf{y} \end{aligned}$$

Diesen Ausdruck können wir nun einfach nach $\hat{\beta}$ lösen, indem wir mit der Inversen $(\mathbf{X}' \mathbf{X})^{-1}$ vormultiplizieren. Als Ergebnis erhalten wir den **OLS-Punktschätzer**

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

Die Bedingung 2. Ordnung für ein Minimum verlangt, dass die Matrix $\mathbf{X}' \mathbf{X}$ positiv definit ist. Diese Bedingung ist aufgrund der Eigenschaften der Matrix $\mathbf{X}' \mathbf{X}$ unter sehr allgemeinen Bedingungen erfüllt, wenn \mathbf{X} vollen Spaltenrang hat.

Übung: Schreiben Sie die Matrix $\mathbf{X}' \mathbf{X}$ mit Hilfe des Summennotation ausführlich an. Ist die Matrix $\mathbf{X}' \mathbf{X}$ symmetrisch?

Frage: wie sieht $\mathbf{X}' \mathbf{X}$ aus, wenn die erste Spalte von \mathbf{X} die Regressionskonstante ist?

Beispiel: Die Matrix \mathbf{X} sei ein $n \times 1$ Vektor mit lauter Einsen, also einer Regression auf die Regressionskonstante

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \\ &= \left[\begin{pmatrix} 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \right]^{-1} \begin{pmatrix} 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \\ &= (1^2 + 1^2 \dots + 1^2)^{-1} (y_1 + y_2 + \dots + y_n) \\ &= n^{-1} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}\end{aligned}$$

Die Anwendung des OLS-Schätzers liefert also wie erwartet den Mittelwert.

Beispiel: Gegeben sei wieder folgende Stichprobe mit 5 Beobachtungen für x und y :

| | | | | | |
|-------|-----|-----|-----|-----|-----|
| y : | 2.6 | 1.6 | 4.0 | 3.0 | 4.9 |
| x : | 1.2 | 3.0 | 4.5 | 5.8 | 7.2 |

In Matrixschreibweise ist das Modell $\mathbf{y} = \mathbf{X}\hat{\beta} + \varepsilon$ oder ausführlicher

$$\begin{pmatrix} 2.6 \\ 1.6 \\ 4.0 \\ 3.0 \\ 4.9 \end{pmatrix} = \begin{pmatrix} 1 & 1.2 \\ 1 & 3.0 \\ 1 & 4.5 \\ 1 & 5.8 \\ 1 & 7.2 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} + \begin{pmatrix} \hat{\varepsilon}_1 \\ \hat{\varepsilon}_2 \\ \hat{\varepsilon}_3 \\ \hat{\varepsilon}_4 \\ \hat{\varepsilon}_5 \end{pmatrix}$$

wobei der Koeffizient des Einsenvektors (1. Spalte) das zu berechnende Interzept ist. Der OLS-Schätzer ist

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \\ &= \left(\begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ 1.2 & 3.0 & 4.5 & 5.8 & 7.2 \end{pmatrix} \begin{pmatrix} 1 & 1.2 \\ 1 & 3.0 \\ 1 & 4.5 \\ 1 & 5.8 \\ 1 & 7.2 \end{pmatrix} \right)^{-1} \begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ 1.2 & 3.0 & 4.5 & 5.8 & 7.2 \end{pmatrix} \begin{pmatrix} 2.6 \\ 1.6 \\ 4.0 \\ 3.0 \\ 4.9 \end{pmatrix} \\ &= \begin{pmatrix} 5 & 21.7 \\ 21.7 & 116.17 \end{pmatrix}^{-1} \begin{pmatrix} 16.1 \\ 78.6 \end{pmatrix} = \begin{pmatrix} 1.0565 & -0.1973 \\ -0.1973 & 0.0455 \end{pmatrix} \begin{pmatrix} 16.1 \\ 78.6 \end{pmatrix} = \begin{pmatrix} 1.498 \\ 0.397 \end{pmatrix}\end{aligned}$$

Dies sind die Punktschätzer $\hat{\beta}_1$ und $\hat{\beta}_2$. In der üblichen Schreibweise

$$\hat{y}_i = 1.498 + 0.397x_i$$

Hinweis: Die Inverse kann z.B. einfach mit Excel berechnet werden, siehe dazu die Hinweise zu Übungsaufgabe 1 auf Seite 32.

6.2.1 Das Bestimmtheitsmaß R^2 in Matrixschreibweise

Wir erinnern uns

$$R^2 = 1 - \frac{\text{SSR}}{\text{TSS}} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{\sum \hat{\varepsilon}_i^2}{\sum (y_i - \bar{y})^2}$$

wobei SSR für *Sum of Squared Residuals* und TSS für *Total Sum Squared* steht.

Dies lässt sich in Matrixschreibweise auch schreiben als

$$R^2 = 1 - \frac{\hat{\varepsilon}'\hat{\varepsilon}}{\sum (y_i - \bar{y})^2} = 1 - \frac{\hat{\varepsilon}'\hat{\varepsilon}}{\mathbf{y}'\mathbf{y} - n\bar{y}^2}$$

Übungsbeispiel: Zeigen Sie, dass $\sum (y_i - \bar{y})^2 = \mathbf{y}'\mathbf{y} - n\bar{y}^2$.

Das korrigierte Bestimmtheitsmaß \bar{R}^2 (*adjusted* R^2)

Wenn in eine Regressionsgleichung zusätzliche x Variablen aufgenommen werden kann das Bestimmtheitsmaß R^2 nie kleiner, sondern nur größer werden. Deshalb eignet sich das Bestimmtheitsmaß nicht für einen Vergleich von Regressionen mit einer unterschiedlichen Anzahl von erklärenden x Variablen. Eine für diesen Zweck besser geeignete Kennziffer erhält man, wenn man im üblichen R^2 die SSR und TSS um die Anzahl der Freiheitsgrade ‘korrigiert’

$$\bar{R}^2 = 1 - \frac{\frac{\hat{\varepsilon}'\hat{\varepsilon}}{n-k}}{\frac{\mathbf{y}'\mathbf{y} - n\bar{y}^2}{n-1}} = 1 - \left(\frac{\hat{\varepsilon}'\hat{\varepsilon}}{\mathbf{y}'\mathbf{y} - n\bar{y}^2} \right) \left(\frac{n-1}{n-k} \right)$$

Mit einer zunehmenden Zahl erklärender Variablen k wird der Faktor $(n-1)/(n-k)$ größer und kompensiert damit dafür, dass $\hat{\varepsilon}'\hat{\varepsilon}$ tendenziell kleiner wird. Deshalb eignet sich das korrigierte Bestimmtheitsmaß \bar{R}^2 eher für einen Vergleich zweier Regressionen mit einer unterschiedlichen Anzahl erklärender Variablen.

Rechenbeispiel

Gegeben sei wieder folgende Stichprobe mit 5 Beobachtungen für x und y :

| | | | | | |
|-------|-----|-----|-----|-----|-----|
| y : | 2.6 | 1.6 | 4.0 | 3.0 | 4.9 |
| x : | 1.2 | 3.0 | 4.5 | 5.8 | 7.2 |

Die OLS-Punktschätzer $\hat{\beta}_1$ und $\hat{\beta}_2$ haben wir bereits früher berechnet

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \begin{pmatrix} 1.0565 & -0.1973 \\ -0.1973 & 0.0455 \end{pmatrix} \begin{pmatrix} 16.1 \\ 78.6 \end{pmatrix} = \begin{pmatrix} 1.498 \\ 0.397 \end{pmatrix}$$

Um die Standardfehler dieser Punktschätzer $\hat{\sigma}_{\hat{\beta}_1}$ und $\hat{\sigma}_{\hat{\beta}_2}$ berechnen zu können benötigen wir zuerst einen Schätzer $\hat{\sigma}^2$ für das σ^2 der Grundgesamtheit. Wie gezeigt ist

$$\hat{\sigma}^2 = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n-k}$$

ein erwartungstreuer Schätzer für σ^2 .

Die Residuen $\hat{\varepsilon}$ berechnen wir aus

$$\hat{\varepsilon} = \mathbf{y} - \mathbf{X}\hat{\beta} = \begin{pmatrix} 2.6 \\ 1.6 \\ 4.0 \\ 3.0 \\ 4.9 \end{pmatrix} - \begin{pmatrix} 1 & 1.2 \\ 1 & 3.0 \\ 1 & 4.5 \\ 1 & 5.8 \\ 1 & 7.2 \end{pmatrix} \begin{pmatrix} 1.498 \\ 0.397 \end{pmatrix} = \begin{pmatrix} 0.626 \\ -1.088 \\ 0.717 \\ -0.799 \\ 0.545 \end{pmatrix}$$

und

$$\hat{\varepsilon}'\hat{\varepsilon} = \begin{pmatrix} 0.626 & -1.088 & 0.717 & -0.799 & 0.545 \end{pmatrix} \begin{pmatrix} 0.626 \\ -1.088 \\ 0.717 \\ -0.799 \\ 0.545 \end{pmatrix} = 3.0257$$

Der Standardfehler der Regression ist also

$$\hat{\sigma} = \sqrt{\frac{\hat{\varepsilon}'\hat{\varepsilon}}{n-k}} = \sqrt{\frac{3.0257}{3}} = 1.004273$$

Der Standardfehler des Koeffizienten h ist

$$\hat{\sigma}_{\hat{\beta}_h} = \hat{\sigma}\sqrt{v_{hh}}$$

wobei v_{hh} das h -te Diagonalelement der Matrix $(\mathbf{X}'\mathbf{X})^{-1}$ ist. Für dieses Beispiel erhält man (siehe z.B. den mathematischen Appendix http://www.uibk.ac.at/econometrics/einf/app2_matrizen.pdf, Abschnitt B.6)

$$(\mathbf{X}'\mathbf{X})^{-1} = \begin{pmatrix} 1.0565 & -0.1973 \\ -0.1973 & 0.0455 \end{pmatrix}$$

Also sind

$$\hat{\sigma}_{\hat{\beta}_1} = \hat{\sigma}\sqrt{v_{11}} = 1.004273\sqrt{1.0565} = 1.0322$$

$$\hat{\sigma}_{\hat{\beta}_2} = \hat{\sigma}\sqrt{v_{22}} = 1.004273\sqrt{0.0455} = 0.2137$$

Diese Standardfehler der Koeffizienten werden uns die Berechnung der Konfidenzintervalle und Hypothesentests ermöglichen.

Die t -Statistiken für die Nullhypothese $\beta_1 = 0$, bzw. $\beta_2 = 0$, sind

$$t\text{-Stat}(\hat{\beta}_1) = \frac{1.498}{1.0322} = 1.4512, \quad t\text{-Stat}(\hat{\beta}_2) = \frac{0.397}{0.2137} = 1.8577$$

$$p\text{-Wert}(\hat{\beta}_1) = [1 - \Phi_3^t(1.4512)] \times 2 = 0.2426$$

$$p\text{-Wert}(\hat{\beta}_2) = [1 - \Phi_3^t(1.8577)] \times 2 = 0.1602$$

Für die Berechnung des Bestimmtheitsmaßes R^2 benötigen wir noch $\mathbf{y}'\mathbf{y} = 58.33$ und $\bar{y}^2 = 10.3684$. Damit erhalten wir

$$R^2 = 1 - \frac{\hat{\varepsilon}'\hat{\varepsilon}}{\mathbf{y}'\mathbf{y} - n\bar{y}^2} = 1 - \frac{3.0257}{58.33 - 5 \times 10.3684} = 0.5336$$

und ein korrigiertes Bestimmtheitsmaß

$$\bar{R}^2 = 1 - \left(\frac{\hat{\varepsilon}'\hat{\varepsilon}}{\mathbf{y}'\mathbf{y} - n\bar{y}^2} \right) \left(\frac{n-1}{n-k} \right) = 1 - \frac{3.0257}{58.33 - 5 \times 10.3684} \frac{4}{3} = 0.3782$$

Für einen Vergleich hier der R Output für dieses Beispiel:

$$\begin{aligned} y &= 1.498 + 0.3968 x \\ &\quad (1.0322) \quad (0.2142) \\ R^2 &= 0.5336, \quad \text{adj. } R^2 = 0.3782, \quad n = 5 \end{aligned}$$

6.3 Die Annahmen des klassischen linearen Regressionsmodells

Eine der wichtigsten Schlussfolgerungen der früheren Kapitel war, dass der OLS-Schätzer unter den *Gauss-Markov Annahmen* effizient – d.h. unverzerrt und varianzminimal – ist. Dies gilt selbstverständlich auch für das multiple Regressionsmodell. Wir erinnern uns, dass wir die Gauss-Markov Annahmen einteilen in Annahmen, die die Spezifikation betreffen (z.B. Funktionsform & ‘richtige’ Variablenselektion), in solche die die \mathbf{X} Matrix betreffen (z.B. Exogenität & lineare Unabhängigkeit), und solche über die Störterme der Grundgesamtheit ε_i ($\varepsilon_i \sim \text{i.i.d.}(0, \sigma^2)$).

Ein Großteil der weiteren Veranstaltung wird sich damit beschäftigen, was zu tun ist, wenn eine oder mehrere dieser Annahmen nicht erfüllt sind. Deshalb wollen wir die *Gauss-Markov Annahmen* nun für das allgemeinere multiple Regressionsmodell unter Zuhilfenahme der Matrixschreibweise wiederholen und etwas kompakter darstellen.

Der eigentliche Gauss Markov Beweis wird im Abschnitt 6.8 skizziert.

A1: Linearität der PRF: Die Beobachtungen \mathbf{y} sind eine lineare Funktion der erklärenden Variablen \mathbf{X} und der Störterme $\boldsymbol{\varepsilon}$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

oder ausführlicher

$$y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} + \varepsilon_i \quad \text{mit } i = 1, \dots, n$$

Diese Linearitätsannahme bezieht sich nur auf die Parameter, nicht auf die Variablen. So ist z.B.

$$y_i = \beta_1 + \beta_2 \ln(x_{i2}) + \beta_3 x_{i3}^2 + \varepsilon_i$$

linear in den Parametern, aber nicht linear in den Variablen.

Modelle, die linear in den Parametern sind können mit OLS geschätzt werden, auch wenn sie *nicht* linear in den Variablen sind.

Andererseits ist z.B.

$$y_i = \beta_1 + \ln(\beta_2)x_{i2} + \beta_3^2 x_{i3} + \varepsilon_i$$

zwar linear in den Variablen, aber *nicht* linear in den Parametern.

Solche Modelle können *nicht* mit OLS geschätzt werden, dazu werden andere Schätzverfahren benötigt (z.B. *Maximum Likelihood* Schätzer).

Insbesondere impliziert diese Annahme, dass das Modell *richtig spezifiziert* ist, das heißt, dass wir nicht nur die ‘richtige’ (lineare) Funktionsform unterstellt haben, sondern auch, dass im Modell keine relevanten erklärenden Variablen fehlen (*‘omitted variable bias’*), und dass keine irrelevanten erklärenden Variablen vorkommen (Ineffizienz).

A2: Voller Spaltenrang: Die $n \times k$ Matrix \mathbf{X} hat vollen Spaltenrang,¹ d.h.

$$\text{rk}(\mathbf{X}) = k$$

wobei rk den Rang bezeichnet.²

Das bedeutet *erstens*, dass die Zahl der Beobachtungen nicht kleiner sein darf als die Zahl der zu schätzenden Parameter, $n \geq k$, und *zweitens*, dass zwischen den einzelnen erklärenden x Variablen (den Spalten der Matrix \mathbf{X}) keine exakte lineare Abhängigkeit besteht.³ Besteht zwischen den Spalten der \mathbf{X} Matrix eine exakte lineare Abhängigkeit, spricht man von *perfekter Multikollinearität* und der OLS-Schätzer kann nicht berechnet werden, da die $\mathbf{X}'\mathbf{X}$ Matrix singulär ist. Die meisten statistischen Programme brechen in diesem Fall entweder ab, oder entfernen mit einem entsprechenden Warnhinweis die linear abhängigen Variablen.

Zur Erinnerung: Angenommen, $y_i = \hat{\beta}_1 + \hat{\beta}_2 x_{i2} + \hat{\beta}_3 x_{i3} + \hat{\varepsilon}_i$. Wenn z.B. $x_{i3} = \alpha x_{i2}$, dann ist $y_i = \hat{\beta}_1 + (\hat{\beta}_2 + \hat{\beta}_3 \alpha) x_{i2} + \hat{\varepsilon}_i$, d.h. es kann nur die Summe $(\hat{\beta}_2 + \hat{\beta}_3 \alpha)$ geschätzt werden, $\hat{\beta}_2$ ist nicht identifiziert!

Die Annahme $\text{rk}(\mathbf{X}) = k$ stellt sicher, dass der Koeffizientenvektor β *eindeutig* ist. Diese Bedingung kann alternativ auch folgendermaßen geschrieben werden:

$$\mathbf{X}\beta^{(1)} = \mathbf{X}\beta^{(2)} \quad \text{wenn und nur wenn} \quad \beta^{(1)} = \beta^{(2)}$$

In anderen Worten, falls \mathbf{X} keinen vollen Spaltenrang hat ist der Koeffizientenvektor *nicht identifiziert*, die gemeinsame Verteilung von \mathbf{y} und \mathbf{X} ist diesem Fall mit vielen möglichen Koeffizientenvektoren $\hat{\beta}$ kompatibel, es kann kein eindeutiger Koeffizientenvektor berechnet werden.

¹Unter dem Spaltenrang versteht man die Anzahl linear unabhängiger Spaltenvektoren einer Matrix; in einer grafischen Interpretation ist der Spaltenrang die Dimension des von den Spaltenvektoren aufgespannten Vektorraums.

²Wenn man aus den Zeilen der $n \times k$ Matrix \mathbf{X} eine nichtsinguläre $k \times k$ Matrix bilden kann (d.h. eine Matrix, deren Determinante nicht Null ist), hat \mathbf{X} vollen Spaltenrang.

³Eine lineare Abhängigkeit zwischen den Spaltenvektoren besteht, wenn mindestens ein Spaltenvektor als Linearkombination der restlichen Spaltenvektoren dargestellt werden kann.

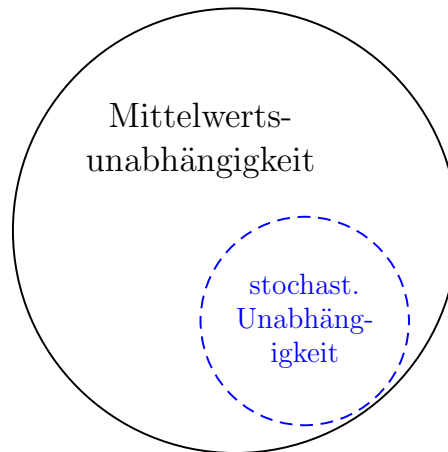


Abbildung 6.1: Mittelwerts- und stochastische Unabhängigkeit. Mittelwertsunabhängigkeit reicht aus (u.a.), damit der OLS-Schätzer erwartungstreu ist. Volle stochastische Unabhängigkeit ist meist nicht erforderlich.

A3: Exogenität der erklärenden Variablen: Die erklärenden Variablen sind *mittelwertsunabhängig* (*‘mean independent’*) von den Störtermen

$$E(\varepsilon_i | \mathbf{X}) = 0$$

Aufgrund des Gesetzes der iterierten Erwartungen impliziert dies auch $E(\varepsilon_i) = 0$.

Mittelwertsunabhängigkeit ist eine schwächere Bedingung als volle stochastische Unabhängigkeit⁴, sie bezieht sich nur auf den Erwartungswert, nicht auf die gesamte Verteilung (z.B. Varianzen und höhere Momente).

Man beachte, dass sich diese Annahme auf die Störterme des datengenerierenden Prozesses bezieht, nicht auf die Residuen der Stichprobe! Die Stichprobenresiduen sind aufgrund der Bedingungen erster Ordnung *immer* (d.h. per Konstruktion) unkorreliert mit den x Variablen.

Diese Annahme steht im Zentrum der Ökonometrie und wird uns noch ausführlicher beschäftigen.

Annahme A3 fordert, dass die Störterme ε_i nicht nur von den individualspezifischen x_{ih} stochastisch unabhängig sein müssen, sondern von allen x_{jh} , mit $j = 1, \dots, n$ und $h = 1, \dots, k$. Für Querschnittsdaten impliziert dies, dass der Störterm eines Individuums i stochastisch unabhängig sein muss von den Regressoren dieses Individuums i sowie von den Regressoren aller anderen Individuen j . Für Zeitreihendaten bedeutet diese Annahme, dass der Störterm der Periode t stochastisch unabhängig sein muss von allen x -Variablen vergangener, gegenwärtiger und zukünftiger Zeitperioden!

⁴Stochastische Unabhängigkeit bedeutet, dass die gemeinsame Dichte gleich dem Produkt der Randdichten ist, d.h. $f(x, y) = f(x)f(y)$. Stochastische Unabhängigkeit impliziert Mittelwertsunabhängigkeit, aber nicht umgekehrt! Mittelwertsunabhängigkeit bezieht sich nur auf lineare Abhängigkeiten, stochastische Unabhängigkeit umfasst auch nichtlineare Abhängigkeiten.

Exkurs: Bedingter Erwartungswert

Seien X und U zwei Zufallsvariablen mit der gemeinsamen Dichte $f(x, u)$. Die Randdichte von X und U ist definiert als

$$f(x) = \int f(x, u) du \quad \text{bzw.} \quad f(u) = \int f(x, u) dx$$

Die auf X bedingte Dichte von U ist definiert

$$f(u|x) = \frac{f(x, u)}{f(x)}$$

Zwei Zufallsvariablen sind stochastisch unabhängig, wenn

$$f(x, u) = f(x)f(u)$$

Der Erwartungswert von U ist definiert als $E(U) = \int uf(u) du$, und der auf X bedingte Erwartungswert von U ist

$$E(U|X = x) = \int uf(u|x) du$$

Wenn U und X *stochastisch unabhängig* sind folgt daraus unter bestimmten Regularitätsbedingungen, dass der auf X bedingte Erwartungswert von U gleich dem unbedingten Erwartungswert von U ist, denn

$$\begin{aligned} E(U|X = x) &= \int uf(u|x) du = \int u \frac{f(x, u)}{f(x)} du \\ &= \int u \frac{f(x)f(u)}{f(x)} du \quad (\text{Unabhängigkeit}) \\ &= \int uf(u) du \equiv E(U) \end{aligned}$$

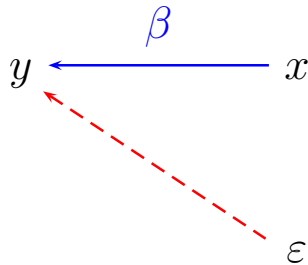
wenn $f(x) > 0$.

Dies zeigt, dass stochastische Unabhängigkeit Mittelwertsunabhängigkeit impliziert, aber wie schon gesagt gilt dies nicht umgekehrt!



Eine intuitive Vorstellung vom Problem vermittelt Abbildung 6.2. Durch eine Korrelation zwischen x und ε wird der mit OLS gemessene Einfluss von x auf y gewissermaßen ‘verschmutzt’, OLS liefert deshalb verzerrte Ergebnisse.

$$E(\varepsilon_i|\mathbf{X}) = 0:$$



$$E(\varepsilon_i|\mathbf{X}) \neq 0:$$

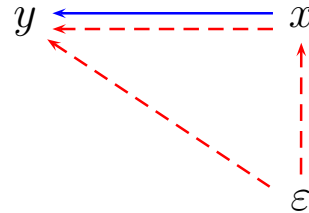


Abbildung 6.2: Wenn $E(\varepsilon_i|\mathbf{X}) = 0$ misst der OLS Schätzer den direkten Einfluss von x auf y (linke Grafik), der OLS Schätzer ist erwartungstreu. Falls ein Regressor x mit dem Störterm korreliert ist ($E(\varepsilon_i|\mathbf{X}) \neq 0$, rechte Grafik), wird dadurch die Messung des Einflusses von x auf y gewissermaßen ‘verschmutzt’, denn in diesem Fall misst der OLS Schätzer den direkten Einfluss von x auf y sowie den indirekten Einfluss von ε über x auf y gemeinsam, der alleinige Einfluss von x auf y ist mit OLS nicht mehr identifiziert, d.h. der OLS Schätzer liefert verzerrte Ergebnisse.

Vereinfacht gesprochen erfordert diese Annahme, dass der datenerzeugende Prozess, der die x erzeugt, unabhängig vom datenerzeugenden Prozess sein muss, der die Störterme ε erzeugt. Regressoren, die diese Annahme verletzen, werden *endogene Regressoren* genannt.

Wenn die Annahme $E(\varepsilon_i|\mathbf{X}) = E(\varepsilon_i) = 0$ sowie die beiden vorhergehenden Annahmen erfüllt sind, ist die OLS Schätzfunktion erwartungstreu!

Leider ist diese Annahme A3 in der Praxis ziemlich häufig verletzt, wichtige Beispiele sind

- *Fehlende relevante Variablen* (‘omitted variables’): wenn eine *nicht* in der Regression vorkommende Variable mit y und mindestens einer in der Regression vorkommenden x Variable korreliert ist;⁵
- *Simultane Kausalität*: wenn z.B. eine einzelne Gleichung eines Mehrgleichungssystems (wie z.B. eine Konsumfunktion) mit OLS geschätzt wird sind die Ergebnisse systematisch verzerrt;
- *Messfehler in den erklärenden Variablen*: selbst ein nicht systematischer Messfehler in einer oder mehreren x Variablen führt zu einer Korrelation des Störterms ε mit den x Variablen, und somit zu verzerrten Schätzungen.

All diese Probleme werden in dem späteren Kapitel über *Instrumentvariablen* ausführlich diskutiert werden.

⁵Ein *Selektionsbias* kann als ein Spezialfall eines ‘omitted variables bias’ angesehen werden.

A4: Störterme: Die Störterme der Grundgesamtheit ε_i sind alle unabhängig und identisch verteilt (*independent and identically distributed*) mit Erwartungswert Null und Varianz σ^2

$$\varepsilon_i \sim \text{i.i.d.}(0, \sigma^2)$$

Diese Annahme umfasst drei einzelne Annahmen:

1. Der Erwartungswert der Störterme in der Grundgesamtheit ist Null

$$E(\varepsilon_i) = 0$$

Wie schon erwähnt ist diese Annahme deutlich weniger streng als die vorhin getroffenen Annahme $E(\varepsilon_i|\mathbf{X}) = 0$, denn $E(\varepsilon_i|\mathbf{X}) = 0$ impliziert $E(\varepsilon_i) = 0$, aber $E(\varepsilon_i) = 0$ impliziert nicht $E(\varepsilon_i|\mathbf{X}) = 0$! Für sich allein genommen ist die Annahme $E(\varepsilon_i) = 0$ nicht besonders schwerwiegend, denn man kann einfach zeigen, dass sich die Verletzung dieser Annahme nur auf die Schätzung des Interzepts auswirkt, was in den meisten Fällen keine große Bedeutung hat.

Dazu nehmen wir z.B. an, dass $\varepsilon_i = \eta + v_i$, wobei η eine Konstante (ungleich Null) und v_i ein üblicher Störterm mit $v_i \sim \text{i.i.d.}(0, \sigma^2)$ ist. Dann kann das Modell $y_i = \beta_1 + \beta_2 x_i + \varepsilon_i$ einfach umgeschrieben werden zu

$$y_i = (\beta_1 + \eta) + \beta_2 x_i + v_i$$

Der Koeffizient β_2 kann in diesem Modell mit OLS effizient geschätzt werden, aber es gibt keine Möglichkeit β_1 oder η einzeln zu schätzen, man kann nur die Summe $(\beta_1 + \eta)$ schätzen. Die Daten enthalten nicht genügend Information um β_1 und η einzeln zu schätzen, man sagt in solchen Fällen, β_1 und η sind nicht *identifiziert*. Das wichtige Konzept der *Identifikation* werden wir später im Zusammenhang mit Systemschätzungen allgemeiner und ausführlicher diskutieren.

2. Homoskedastizität: Jeder Störterm ε_i hat die gleiche, endliche bedingte Varianz σ^2

$$\text{var}(\varepsilon_i|\mathbf{X}) = \sigma^2 \quad \text{für alle } i$$

Ist diese Annahme *nicht* erfüllt, d.h. $\sigma_i^2 \neq \sigma_j^2$ für $i \neq j$, spricht man von *Heteroskedastizität*. Da die PRF in den meisten Fällen nur eine lineare Approximation an die bedingte Erwartungswertfunktion (CEF) ist wird diese Annahme ziemlich häufig verletzt sein. Im Kapitel zur Heteroskedastizität werden wir sehen, dass *robuste Standardfehler* in solchen Fällen manchmal geeigneter sind als OLS Standardfehler.

3. Keine Autokorrelation: Die Störterme ε_i sind untereinander stochastisch unabhängig,

$$E(\varepsilon_i \varepsilon_j | \mathbf{X}) = 0 \quad \text{für } i \neq j$$

Ist diese Annahme verletzt spricht man von *Autokorrelation*.

In Vektornotation können die beiden letzten Annahmen über die Störterme kompakter geschrieben werden als

$$\text{var}(\boldsymbol{\varepsilon}|\mathbf{X}) = \sigma^2 \mathbf{I}$$

Diese vier Annahmen werden für den Gauss-Markov Beweis benötigt, d.h. *wenn* diese vier Annahmen erfüllt sind ist der OLS Schätzer ein **BLUE** (*best linear unbiased estimator*, oder in anderen Worten, effizient).

A5: Normalverteilung: Die Störterme ε sind in der Grundgesamtheit normalverteilt.

Dies ist *keine* Gauss-Markov Annahme, wenn die Gauss-Markov Annahmen A1 – A4 erfüllt sind ist der OLS-Schätzer auch dann effizient, wenn die Normalverteilungsannahme *nicht* erfüllt ist.⁶ Außerdem sind die Stichprobenkennwerte (wie z.B. die geschätzten Koeffizienten) in großen Stichproben aufgrund des zentralen Grenzwertsatzes auch dann asymptotisch normalverteilt, wenn die Störterme ε nicht normalverteilt sind. Die Normalverteilungsannahme wird vor allem für Hypothesentests in kleinen Stichproben benötigt.

6.4 Erwartungstreue

Offensichtlich wird sich der vorhin ermittelte Vektor $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ von Stichprobe zu Stichprobe unterscheiden, wenn wir eine neue Stichprobe ziehen, werden wir auch – hoffentlich nur geringfügig – andere Schätzungen für $\hat{\beta}$ erhalten. Deshalb sind die so berechneten Koeffizienten $\hat{\beta}$ im Gegensatz zu den Parametern der Grundgesamtheit β selbst wieder Zufallsvariablen, genauso wie die damit berechneten Stichprobenresiduen $\hat{\varepsilon}$. Von Zufallsvariablen kann man aus der Stichprobenkennwertverteilung Erwartungswert und Varianz berechnen, und genau dies wollen wir im nächsten Schritt tun.

Insbesondere interessiert uns, ob der OLS-Schätzer $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ erwartungstreue Schätzungen für β liefert, d.h., ob die aus der Stichprobe berechneten Werte $\hat{\beta}$ ‘im Mittel’ den unbeobachtbaren β entsprechen. Dazu setzen wir wieder den ‘wahren’ Zusammenhang der Grundgesamtheit $\mathbf{y} = \mathbf{X}\beta + \varepsilon$ in die Formel unseres Schätzers $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ ein und bilden anschließend den Erwartungswert:

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\underbrace{(\mathbf{X}\beta + \varepsilon)}_{\mathbf{y}} \\ &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon \\ \hat{\beta} &= \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon\end{aligned}$$

Wir bilden den Erwartungswert

$$E(\hat{\beta}) = \beta + E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon]$$

Wenn die erklärenden x Variablen deterministisch (bzw. ‘fixed in repeated sampling’) sind, ist $E((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\varepsilon) = 0$, also ist der OLS-Schätzer erwartungstreu

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

⁶Allerdings gibt es nicht sehr viele Verteilungen außer der Normalverteilung, die die Annahme der Homoskedastizität erfüllen.

Das gilt auch für stochastische x , allerdings ist es der Beweis nicht mehr so trivial. Für ein intuitives Verständnis empfiehlt es sich, einen genaueren Blick auf den $k \times 1$ Vektor $\mathbf{X}'\boldsymbol{\varepsilon}$ zu werfen

$$\mathbf{X}'\boldsymbol{\varepsilon} = \begin{pmatrix} x_{11} & x_{21} & \cdots & x_{n1} \\ x_{12} & x_{22} & \cdots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1k} & x_{2k} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} = \begin{pmatrix} \sum_i x_{i1}\varepsilon_i \\ \sum_i x_{i2}\varepsilon_i \\ \vdots \\ \sum_i x_{ik}\varepsilon_i \end{pmatrix}$$

Man kann zeigen, dass $E(\sum_i x_{ih}\varepsilon_i) = 0$ (für $h = 1, \dots, k$) wenn die Störterme stochastisch unabhängig von den erklärenden x Variablen sind. Wenn die Störterme also unkorreliert mit den x Variablen sind sollte also $E(\mathbf{X}'\boldsymbol{\varepsilon}) = 0$ sein.⁷

Man beachte, dass wir für die Erwartungstreue nur Annahmen A1 bis A3 benötigt haben, d.h. korrekte Spezifikation, voller Spaltenrang sowie Exogenität der Regressoren ($E(\boldsymbol{\varepsilon}|\mathbf{X}) = 0$), nicht aber die Annahmen A4 und A5 (d.h. $\varepsilon_i \sim \text{i.i.d.}(0, \sigma^2)$ und die Normalverteilungsannahme). Deshalb ist der OLS-Schätzer selbst dann erwartungstreu, wenn die Annahmen A4 über die Störterme verletzt sind.

Man beachte auch, dass $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}$ große Ähnlichkeit mit dem Schätzer für $\hat{\boldsymbol{\beta}}$ hat. Der einzige Unterschied zu $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ besteht darin, dass \mathbf{y} durch $\boldsymbol{\varepsilon}$ ersetzt wurde. Intuitiv kann man sich dies folgendermaßen vorstellen: wenn im systematischen Teil des Erklärungsansatzes mit den Variablen \mathbf{X} alles erklärt wurde, sollten die Störterme $\boldsymbol{\varepsilon}$ keine nutzbare Information mehr enthalten. Deshalb sollte die Anwendung des Schätzers auf die Störterme im Erwartungswert Null liefern.

Achtung: Wir können die Stichprobenresiduen $\hat{\boldsymbol{\varepsilon}}$ nicht dazu benützen um zu testen, ob die Störterme der Grundgesamtheit $\boldsymbol{\varepsilon}$ tatsächlich mit den erklärenden Variablen \mathbf{X} unkorreliert sind. Wir haben bereits gezeigt, dass die Bedingungen erster Ordnung für die Herleitung des OLS-Schätzers $\mathbf{X}'\hat{\boldsymbol{\varepsilon}} = 0$ implizieren, d.h. die Stichprobenkorrelation zwischen \mathbf{X} und $\hat{\boldsymbol{\varepsilon}}$ ist aufgrund der Konstruktion unseres Schätzers immer Null, auch wenn in der Grundgesamtheit eine Korrelation zwischen $\boldsymbol{\varepsilon}$ und \mathbf{X} bestehen sollte!

Falls die Regression ein Interzept enthält ist eine Spalte von \mathbf{X} ein Einsen-Vektor. Die dazugehörige Bedingung erster Ordnung der Minimierung der Quadratsumme der Residuen stellt in diesem Fall sicher, dass die Summe von $\hat{\boldsymbol{\varepsilon}}$ immer Null ist, d.h. $\sum_{i=1}^n \hat{\varepsilon}_i = 0$. Dies muss bei einer Schätzung ohne Interzept nicht gelten!

Damit haben wir die Erwartungstreue der OLS Schätzer unter den Annahmen A1 – A3 bewiesen. Als nächstes werden wir die Varianzen der OLS Schätzer berechnen, um wieder Hypothesentests durchführen zu können.

6.4.1 Nichtberücksichtigung relevanter Variablen (*‘omitted variables’*) in Matrixschreibweise

Das ‘wahre’ Modell sei

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$$

⁷Dies beweist bei stochastischen x natürlich nicht, dass $E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}] = 0$, dazu müssten wir weiter ausholen, aber es soll den Blick auf die richtige Stelle lenken.

aber wir schätzen das ‘falsche’ Modell

$$\mathbf{y} = \mathbf{X}_1 \boldsymbol{\beta}_1^* + \boldsymbol{\varepsilon}^*$$

wobei \mathbf{X}_1 eine $n \times k_1$ und \mathbf{X}_2 eine $n \times k_2$ Matrix ist.

Die Erwartungstreue prüfen wir wie üblich, indem wir das ‘wahre’ Modell in den Schätzer für $\hat{\boldsymbol{\beta}}$ einsetzen und den Erwartungswert bilden.

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{y} \\ &= (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' (\mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{X}_2 \boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}) \\ &= \boldsymbol{\beta}_1 + (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{X}_2 \boldsymbol{\beta}_2 + (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \boldsymbol{\varepsilon} \end{aligned}$$

Selbst wenn die erklärenden Variablen wieder exogen sind ist

$$\begin{aligned} E(\hat{\boldsymbol{\beta}}) &= \boldsymbol{\beta}_1 + (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{X}_2 \boldsymbol{\beta}_2 + (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' E(\boldsymbol{\varepsilon}) \\ &= \boldsymbol{\beta}_1 + (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{X}_2 \boldsymbol{\beta}_2 \\ &\neq \boldsymbol{\beta}_1 \end{aligned}$$

Im Unterschied dazu erzeugt die irrtümliche Berücksichtigung irrelevanter Variablen keinen Bias (siehe 6.B.6, Seite 51), allerdings ist in diesem Fall OLS nicht mehr effizient, da Freiheitsgrade ‘verschwendet’ werden. Darüber hinaus erzeugt die Berücksichtigung irrelevanter Variablen häufig ein ‘*overfitting*’ in der Stichprobe, mit einhergehender schlechte Prognosegüte bei neuen Daten (Out-of-sample). Eventuell vorliegende Multikollinearität vergrößert diese Probleme in der Regel noch. Außerdem kann die Interpretierbarkeit der Ergebnisse unter zu vielen irrelevanten Variablen Schaden nehmen.

6.5 Schätzung der Varianz von $\hat{\boldsymbol{\beta}}$

Die Varianz einer skalaren Zufallsvariable $\hat{\beta}_h$ ist definiert als $\text{var}(\hat{\beta}_h) = E[\hat{\beta}_h - E(\hat{\beta}_h)]^2$. Die geschätzten Koeffizienten $\hat{\boldsymbol{\beta}}$ bilden einen $k \times 1$ Spaltenvektor von Zufallsvariablen. Für Hypothesentests sowie die Berechnung von Konfidenzintervallen benötigen wir unter anderem die Standardfehler der geschätzten Koeffizienten, die einfach die Wurzel der Hauptdiagonalelemente der Varianz-Kovarianzmatrix von $\hat{\boldsymbol{\beta}}$ sind. Diese Varianz-Kovarianzmatrix des Vektors von Zufallsvariablen $\hat{\boldsymbol{\beta}}$ ist definiert als

$$\text{var}(\hat{\boldsymbol{\beta}}) = E \left[\hat{\boldsymbol{\beta}} - E(\hat{\boldsymbol{\beta}}) \right] \left[\hat{\boldsymbol{\beta}} - E(\hat{\boldsymbol{\beta}}) \right]'$$

Beispiel Wenn $\hat{\beta}$ der 2×1 Koeffizientenvektor des bivariaten Modells ist erhalten wir die 2×2 Varianz-Kovarianzmatrix

$$\begin{aligned} \text{var}(\hat{\beta}) &= E \left[\hat{\beta} - E(\hat{\beta}) \right] \left[\hat{\beta} - E(\hat{\beta}) \right]' \\ &= E \left[\begin{pmatrix} \hat{\beta}_1 - E(\hat{\beta}_1) \\ \hat{\beta}_2 - E(\hat{\beta}_2) \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 - E(\hat{\beta}_1) & \hat{\beta}_2 - E(\hat{\beta}_2) \end{pmatrix} \right] \\ &= \begin{pmatrix} E[\hat{\beta}_1 - E(\hat{\beta}_1)]^2 & E[\hat{\beta}_1 - E(\hat{\beta}_1)][\hat{\beta}_2 - E(\hat{\beta}_2)] \\ E[\hat{\beta}_1 - E(\hat{\beta}_1)][\hat{\beta}_2 - E(\hat{\beta}_2)] & E[\hat{\beta}_2 - E(\hat{\beta}_2)]^2 \end{pmatrix} \\ &= \begin{pmatrix} \text{var}(\hat{\beta}_1) & \text{cov}(\hat{\beta}_1, \hat{\beta}_2) \\ \text{cov}(\hat{\beta}_1, \hat{\beta}_2) & \text{var}(\hat{\beta}_2) \end{pmatrix} \end{aligned}$$

Die Elemente auf der Hauptdiagonale sind die Varianzen und die Elemente auf der Nebendiagonale die Kovarianzen von $\hat{\beta}$.

Wir haben im vorhergehenden Abschnitt bereits gezeigt, dass unter den Gauss-Markov Annahmen $\hat{\beta} = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon$ und $E(\hat{\beta}) = \beta$. Deshalb ist

$$\begin{aligned} E[\hat{\beta} - E(\hat{\beta})] &= E[\beta + ((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon - \beta)] \\ &= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon] \end{aligned}$$

Unter Verwendung dieses Ergebnisses erhalten wir die Varianz-Kovarianzmatrix von $\hat{\beta}$ als

$$\begin{aligned} \text{var}(\hat{\beta}) &= E \left[\hat{\beta} - E(\hat{\beta}) \right] \left[\hat{\beta} - E(\hat{\beta}) \right]' \\ &= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon][(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon]' \\ &= E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon\varepsilon'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}] \end{aligned}$$

da $(\mathbf{AB})' = \mathbf{B}'\mathbf{A}'$ sowie $(\mathbf{A}')' = \mathbf{A}$ (siehe mathematischen Appendix).

Wenn die \mathbf{X} deterministisch sind, sind die einzigen Zufallsvariablen in diesem Ausdruck die Elemente des Vektors der Störterme ε . Deshalb können wir den Erwartungswertoperator direkt vor die Matrix der Zufallsvariablen $\varepsilon\varepsilon'$ schreiben

$$\text{var}(\hat{\beta}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\varepsilon\varepsilon')\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \quad (6.1)$$

Zur Berechnung der eigentlich interessierenden Varianz-Kovarianzmatrix von $\hat{\beta}$ benötigen wir also die Matrix $E(\varepsilon\varepsilon')$, die wir uns nun etwas näher ansehen wollen.

Die Elemente des $n \times 1$ Spaltenvektor der Störterme ε haben einen Erwartungswert von Null, d.h. $E(\varepsilon) = \mathbf{O}$, bzw. ausführlicher

$$E(\varepsilon) = E \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} = \begin{pmatrix} E(\varepsilon_1) \\ E(\varepsilon_2) \\ \vdots \\ E(\varepsilon_n) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \mathbf{O}$$

deshalb ist $E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}')$ einfach die Varianz-Kovarianzmatrix der Störterme.

$$\begin{aligned}
 \text{var}(\boldsymbol{\varepsilon}) &= E \begin{pmatrix} \varepsilon_1 - E(\varepsilon_1) \\ \varepsilon_2 - E(\varepsilon_2) \\ \vdots \\ \varepsilon_n - E(\varepsilon_n) \end{pmatrix} \begin{pmatrix} \varepsilon_1 - E(\varepsilon_1) & \varepsilon_2 - E(\varepsilon_2) & \cdots & \varepsilon_n - E(\varepsilon_n) \end{pmatrix} \\
 &= E \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} \begin{pmatrix} \varepsilon_1 & \varepsilon_2 & \cdots & \varepsilon_n \end{pmatrix} = E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}') \\
 &= \begin{pmatrix} E(\varepsilon_1)^2 & E(\varepsilon_1\varepsilon_2) & \cdots & E(\varepsilon_1\varepsilon_n) \\ E(\varepsilon_2\varepsilon_1) & E(\varepsilon_2)^2 & \cdots & E(\varepsilon_2\varepsilon_n) \\ \vdots & \vdots & \ddots & \vdots \\ E(\varepsilon_n\varepsilon_1) & E(\varepsilon_n\varepsilon_2) & \cdots & E(\varepsilon_n)^2 \end{pmatrix} \\
 &= \begin{pmatrix} \text{var}(\varepsilon_1) & \text{cov}(\varepsilon_1\varepsilon_2) & \cdots & \text{cov}(\varepsilon_1\varepsilon_n) \\ \text{cov}(\varepsilon_2\varepsilon_1) & \text{var}(\varepsilon_2) & \cdots & \text{cov}(\varepsilon_2\varepsilon_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(\varepsilon_n\varepsilon_1) & \text{cov}(\varepsilon_n\varepsilon_2) & \cdots & \text{var}(\varepsilon_n) \end{pmatrix}
 \end{aligned}$$

Diese Varianz-Kovarianzmatrix der Grundgesamtheit erhält eine sehr einfache Form, wenn man *Homoskedastizität* ($\text{var}(\varepsilon_i) = \sigma^2$) und *Abwesenheit von Autokorrelation* ($\text{cov}(\varepsilon_i, \varepsilon_j) = 0$ für $i \neq j$) unterstellt, also

$$\text{var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$$

denn in diesem Fall ist die Varianz-Kovarianzmatrix eine einfache Diagonalmatrix.

$$E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}') = \begin{pmatrix} \sigma^2 & 0 & \cdots & 0 \\ 0 & \sigma^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma^2 \end{pmatrix} = \sigma^2 \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} = \sigma^2 \mathbf{I}_n$$

wobei \mathbf{I}_n die $n \times n$ Einheitsmatrix ist.

Wie schon gesagt, diese Varianz-Kovarianzmatrix der Störterme ε_i benötigen wir, um die eigentlich interessierende *Varianz-Kovarianzmatrix der geschätzten OLS Koeffizienten* zu berechnen, d.h. $\text{var}(\hat{\boldsymbol{\beta}})$.

Dies ist unter den Annahmen $E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}') = \sigma^2 \mathbf{I}$ sehr einfach, wir brauchen nur in Gleichung (6.1) einzusetzen

$$\begin{aligned}
 \text{var}(\hat{\boldsymbol{\beta}}) &= E \left[\hat{\boldsymbol{\beta}} - E(\hat{\boldsymbol{\beta}}) \right] \left[\hat{\boldsymbol{\beta}} - E(\hat{\boldsymbol{\beta}}) \right]' = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \underbrace{E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}')}_{=\sigma^2 \mathbf{I}} \mathbf{X}' (\mathbf{X}'\mathbf{X})^{-1} \\
 &= \sigma^2 [(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}] \\
 &= \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}
 \end{aligned}$$

da σ^2 ein Parameter der Grundgesamtheit ist (d.h. eine fixe, aber unbekannte Zahl) und die \mathbf{X} exogen sind. Damit haben wir einen Ausdruck für die Varianz-Kovarianzmatrix des *Koeffizientenvektors*, der allerdings noch von der unbeobachtbaren Varianz der Störterme σ^2 abhängt. Wie früher benötigen wir wieder einen Schätzer $\hat{\sigma}^2$ für σ^2 .

6.6 Eine Schätzfunktion für die Varianz von ε (σ^2)

Für Hypothesentests benötigen wir die Standardfehler der Koeffizienten. Da der Ausdruck für die Varianz von $\hat{\beta}$, $\text{var}(\hat{\beta}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$, noch die unbeobachtbare Varianz σ^2 der Grundgesamtheit enthält, benötigen wir einen Schätzer für σ^2 . Den Ansatz, den wir zur Herleitung des Schätzers für β wählten, nämlich die Minimierung von $\hat{\varepsilon}'\hat{\varepsilon}$, können wir hier nicht anwenden, da die Zielfunktion die Varianz σ^2 nicht enthält.

Ein möglicher Kandidat für einen solchen Schätzer wäre die Varianz der Stichprobenresiduen, aber wie schon im bivariaten Fall kann man zeigen, dass diese einen verzerrten Schätzer darstellt. Wir werden aber im folgenden Beweis zeigen, dass

$$\hat{\sigma}^2 = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n-k} = \frac{\sum_i \hat{\varepsilon}_i^2}{n-k}$$

ein unverzerrter Schätzer für die Varianz σ^2 der Störterme ist, wobei k die Anzahl der erklärenden Variablen (inkl. Interzept) ist; $(n-k)$ bezeichnet wieder die Anzahl der Freiheitsgrade. Die Wurzel daraus,

$$\hat{\sigma} = \sqrt{\frac{\hat{\varepsilon}'\hat{\varepsilon}}{(n-k)}}$$

wird *Standardfehler der Regression* genannt.

Dies wollen wir nun ausführlich zeigen und beweisen.

Die Störterme ε_i sind Zufallsvariablen, und unter den früher gemachten Annahmen ist die Varianz dieser Zufallsvariablen ist die Zahl σ^2 . Aber natürlich ist $\varepsilon = \mathbf{y} - \mathbf{X}\beta$ nicht beobachtbar. Da $\mathbf{y} = \mathbf{X}\hat{\beta} + \hat{\varepsilon}$ und $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ kann man für die Stichproben-Residuen schreiben

$$\hat{\varepsilon} = \mathbf{y} - \mathbf{X}\hat{\beta} = \mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}$$

Zur Vereinfachung der Schreibweise führen wir für die $n \times n$ Matrix $[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']$ das Symbol \mathbf{M} ein, d.h. wir definieren

$$\mathbf{M} := \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

Die Matrix \mathbf{M} spielt in der Ökonometrie eine wichtige Rolle und wird **residuenerzeugende Matrix** genannt, da – wie soeben gezeigt – $\mathbf{M}\mathbf{y}$ den Residuenvektor $\hat{\varepsilon}$ ergibt

$$\hat{\varepsilon} = \mathbf{M}\mathbf{y} \quad (6.2)$$

d.h. Vormultiplikation des $n \times 1$ Vektors \mathbf{y} mit der $n \times n$ Matrix \mathbf{M} erzeugt den Residuenvektor.

Die residuenerzeugende Matrix \mathbf{M} hat mehrere wichtige Eigenschaften

- \mathbf{M} ist symmetrisch, d.h. $\mathbf{M} = \mathbf{M}'$

Warum? Falls \mathbf{M} symmetrisch ist muss gelten $\mathbf{M} = \mathbf{M}'$.

$$\text{Da } \mathbf{M}' = [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']' = \mathbf{I}' - (\mathbf{X}')'((\mathbf{X}'\mathbf{X})^{-1})'\mathbf{X}' = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \mathbf{M}$$

- \mathbf{M} ist idempotent, d.h. $\mathbf{M}\mathbf{M} = \mathbf{M}$. Dies kann einfach durch Ausmultiplizieren überprüft werden.⁸
- Die Matrizen \mathbf{M} und \mathbf{X} sind orthogonal, d.h. $\mathbf{M}\mathbf{X} = \mathbf{O}$, weil

$$\mathbf{M}\mathbf{X} = [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{O}$$

Daraus folgt gemeinsam mit Gleichung (6.2) auch

$$\hat{\varepsilon} = \mathbf{M}\varepsilon$$

weil

$$\hat{\varepsilon} = \mathbf{M}\mathbf{y} = \mathbf{M}(\underbrace{\mathbf{X}\beta + \varepsilon}_{\mathbf{y}}) = \mathbf{O}\beta + \mathbf{M}\varepsilon = \mathbf{M}\varepsilon$$

- \mathbf{M} ist singulär mit Rang $n-k$ (idempotente Matrizen haben mit Ausnahme der Einheitsmatrix nie vollen Rang; siehe Appendix zur Matrixalgebra). Deshalb besitzt \mathbf{M} keine Inverse und der obige Zusammenhang kann z.B. nicht benutzt werden, um aus den Residuen $\hat{\varepsilon}$ die Störterme ε zu berechnen.

Exkurs: Die Verteilung der Residuen bei normalverteilten Störtermen

Wenn die Störterme ε normalverteilt sind mit $\varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ und die \mathbf{X} Matrix deterministisch ist erlaubt der Zusammenhang $\hat{\varepsilon} = \mathbf{M}\varepsilon$ die Verteilung der Residuen zu bestimmen.

Da jede Linearkombination normalverteilter Zufallsvariablen wieder normalverteilt ist, sind unter obigen Annahmen auch die Residuen $\hat{\varepsilon}$ normalverteilt, weil $\mathbf{M}\varepsilon$ ein linearer Zusammenhang ist.

⁸ \mathbf{M} ist idempotent weil

$$\begin{aligned} \mathbf{M}\mathbf{M} &= (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')(\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \\ &= \mathbf{I}\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ &= \mathbf{I} - 2\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ &= \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \mathbf{M} \end{aligned}$$

Der Erwartungswert und die Varianz der Residuen können ebenfalls einfach bestimmt werden:

$$\begin{aligned} E(\hat{\varepsilon}) &= E(\mathbf{M}\varepsilon) \\ &= \mathbf{M} E(\varepsilon) \\ &= \mathbf{0} \end{aligned}$$

und

$$\begin{aligned} \text{var}(\hat{\varepsilon}) &= E(\hat{\varepsilon}\hat{\varepsilon}') \\ &= E(\mathbf{M}\varepsilon\varepsilon'\mathbf{M}') \\ &= \mathbf{M} E(\varepsilon\varepsilon')\mathbf{M}' \\ &= \mathbf{M}(\sigma^2\mathbf{I}_n)\mathbf{M}' \quad (\text{wenn } \varepsilon_i \sim \text{i.i.d.}(0, \sigma^2)) \\ &= \sigma^2\mathbf{M}\mathbf{M}' \\ &= \sigma^2\mathbf{M} \end{aligned}$$

Deshalb ist

$$\hat{\varepsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{M})$$

Man beachte, dass für dieses Resultat alle Annahmen des klassischen linearen Regressionsmodells (A1 – A4) erforderlich waren, und dass die Varianz-Kovarianzmatrix der Residuen singulär ist!

Selbst wenn die Störterme homoskedastisch und nicht autokorreliert sind, sind die Residuen im allgemeinen heteroskedastisch und autokorreliert, da die Nebendiagonalelemente von $\text{var}(\hat{\varepsilon}) = \sigma^2\mathbf{M}$ ungleich Null sind (d.h. $E(\hat{\varepsilon}_i\hat{\varepsilon}_j) \neq 0$ für $i \neq j$ und die Hauptdiagonalelemente von \mathbf{M} ungleich groß sind ($\text{var}(\hat{\varepsilon}_i) \neq \sigma^2$ für alle i)).

¶

Die residuenerzeugende Matrix \mathbf{M} ist eng verwandt mit der **Projektionsmatrix** \mathbf{P}

$$\mathbf{P} := \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

da $\mathbf{M} = [\mathbf{I} - \mathbf{P}]$.

Wenn man den \mathbf{y} Vektor mit der Projektionsmatrix \mathbf{P} vormultipliziert erhält man die gefitteten Werte, denn $\hat{\mathbf{y}}$

$$\hat{\mathbf{y}} = \mathbf{y} - \hat{\varepsilon} = \mathbf{y} - \mathbf{M}\mathbf{y} = (\mathbf{I} - \mathbf{M})\mathbf{y} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{P}\mathbf{y}$$

Die Projektionsmatrix \mathbf{P} wird manchmal auch ‘Hutmatrix’ (*‘hat matrix’*) genannt, da sie dem \mathbf{y} Vektor ‘einen Hut aufsetzt’.

Ebenso wie die residuenerzeugende Matrix \mathbf{M} ist auch \mathbf{P} eine symmetrische, idempotente und singuläre $n \times n$ Matrix mit Rang k .

Die Bedeutung der Projektionsmatrix \mathbf{P} und der residuenerzeugenden Matrix \mathbf{M} wird im Abschnitt über die geometrische Interpretation des OLS-Schätzers (Abschnitt 6.11) ausführlich diskutiert.

Diese Ergebnisse wollen wir nun verwenden um einen erwartungstreuen Schätzer $\hat{\sigma}^2$ für das unbekannte σ^2 zu ermitteln. Wir schreiben zuerst die Quadratsumme der Residuen $\hat{\varepsilon}'\hat{\varepsilon}$ mit Hilfe der residuenerzeugenden Matrix \mathbf{M} an

$$\begin{aligned}\hat{\varepsilon} &= \mathbf{M}\mathbf{y} = \mathbf{M}\mathbf{X}\boldsymbol{\beta} + \mathbf{M}\boldsymbol{\varepsilon} = \mathbf{M}\boldsymbol{\varepsilon} \quad (\text{da } \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \text{ und } \mathbf{M}\mathbf{X} = \mathbf{O}) \\ \hat{\varepsilon}'\hat{\varepsilon} &= \boldsymbol{\varepsilon}'\mathbf{M}'\mathbf{M}\boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}'\mathbf{M}\mathbf{M}\boldsymbol{\varepsilon} = \boldsymbol{\varepsilon}'\mathbf{M}\boldsymbol{\varepsilon} \quad (\text{da } \mathbf{M} \text{ sym. \& idempotent ist})\end{aligned}$$

Für die weiteren Ausführungen benötigen wir den Spur-Operator (*trace*, tr): Die Spur einer symmetrischen Matrix ist definiert als die Summe der Hauptdiagonalelemente $\text{tr}(\mathbf{A}) = a_{11} + a_{22} + \dots + a_{nn}$. Man kann zeigen, dass für geeignet dimensionierte Matrizen gilt: $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$. Da $\boldsymbol{\varepsilon}'\mathbf{M}\boldsymbol{\varepsilon}$ ein Skalar ist kann man davon die Spur nehmen ohne das Resultat zu verändern. Unter Verwendung der obigen Regel können wir also schreiben: $\hat{\varepsilon}'\hat{\varepsilon} = \boldsymbol{\varepsilon}'\mathbf{M}\boldsymbol{\varepsilon} = \text{tr}(\boldsymbol{\varepsilon}'\mathbf{M}\boldsymbol{\varepsilon}) = \text{tr}(\mathbf{M}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}')$. Zur Überprüfung der Erwartungstreue müssen wir davon den Erwartungswert bilden:

$$\begin{aligned}\text{E}(\hat{\varepsilon}'\hat{\varepsilon}) &= \text{E}[\text{tr}(\mathbf{M}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}')] \\ &= \text{tr}[\text{E}(\mathbf{M}\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}')] \quad (\text{da die Spur ein linearer Operator ist}) \\ &= \text{tr}[\mathbf{M}\sigma^2\mathbf{I}] \quad (\text{da } \mathbf{M} \text{ deterministisch ist und wenn } \text{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}') = \sigma^2\mathbf{I}) \\ &= \sigma^2 \text{tr}(\mathbf{M}) \\ &= \sigma^2 \text{tr}[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] \quad (\text{da } \text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})) \\ &= \sigma^2 [\text{tr}(\mathbf{I}) - \text{tr}(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')] \quad (\text{da } \text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})) \\ &= \sigma^2 [\text{tr}(\mathbf{I}) - \text{tr}(\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1})] \\ &= \sigma^2 [\text{tr}(\mathbf{I}_{(n \times n)}) - \text{tr}(\mathbf{I}_{(k \times k)})] \\ &= \sigma^2(n - k)\end{aligned}$$

Die Varianz der Stichprobenresiduen $\tilde{\sigma}_{\hat{\varepsilon}_i}^2$ ist⁹

$$\tilde{\sigma}_{\hat{\varepsilon}_i}^2 := \frac{1}{n} \sum_i \hat{\varepsilon}_i^2 = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n}$$

Bildet man den Erwartungswert, so folgt

$$\text{E}(\tilde{\sigma}_{\hat{\varepsilon}_i}^2) = \frac{\text{E}(\hat{\varepsilon}'\hat{\varepsilon})}{n} = \frac{(n - k)\sigma^2}{n}$$

Somit unterschätzt die Stichprobenvarianz $\tilde{\sigma}_{\hat{\varepsilon}_i}^2$ in ihrem Erwartungswert die Populationsvarianz σ^2 um den Faktor $(n - k)/n$. Multiplizieren wir den Erwartungswert der Stichprobenvarianz mit dem Kehrwert $n/(n - k)$, wird der ‘bias’ korrigiert und wir erhalten eine erwartungstreue Schätzung der Populationsvarianz σ^2 . Wir definieren also ein $\hat{\sigma}^2$

$$\hat{\sigma}^2 = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n - k}$$

für das gilt

$$\text{E}(\hat{\sigma}^2) = \frac{\text{E}(\hat{\varepsilon}'\hat{\varepsilon})}{n - k} = \sigma^2$$

d.h. dieses $\hat{\sigma}^2$ ist ein erwartungstreuer Schätzer für die Varianz der Störterme σ^2 . ■

⁹Diese so definierte Stichprobenvarianz $\tilde{\sigma}_{\hat{\varepsilon}_i}^2$ ist *kein* Schätzer, sondern die (Populations-) Varianz der Zufallsvariablen $\hat{\varepsilon}_i$. Deshalb wird durch n und nicht durch $n - k$ dividiert.

Die Wurzel dieses Schätzers für die Varianz der Störterme $\hat{\sigma}^2$, d.h.

$$\hat{\sigma} = \sqrt{\frac{\hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}}}{n - k}}$$

wird Standardfehler der Schätzung (*Standard Error of the Regression*) genannt und wird von fast allen ökonometrischen Software-Paketen ausgegeben.¹⁰

Damit ist unser Problem gelöst. Zur Erinnerung, wir hatten im letzten Abschnitt die Varianz-Kovarianz Matrix der Regressionskoeffizienten berechnet

$$\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}' \mathbf{X})^{-1}$$

konnten damit aber nicht viel anfangen, da σ^2 ein Parameter der Grundgesamtheit – und deshalb nicht beobachtbar – ist.

Nun haben wir einen erwartungstreuen Schätzer für σ^2 hergeleitet, und dies erlaubt uns endlich eine erwartungstreue Schätzfunktion für die Varianz-Kovarianz Matrix der OLS Schätzer $\hat{\boldsymbol{\beta}}$ anzugeben

$$\widehat{\text{var}}(\hat{\boldsymbol{\beta}}) := \hat{\boldsymbol{\sigma}}_{\hat{\boldsymbol{\beta}}}^2 = \hat{\sigma}^2 (\mathbf{X}' \mathbf{X})^{-1} = \frac{\hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}}}{n - k} (\mathbf{X}' \mathbf{X})^{-1}$$

Man beachte, dass $\hat{\sigma}^2$ ein Schätzer für die Varianz der Grundgesamtheit σ^2 ist (also ein Skalar), während $\hat{\boldsymbol{\sigma}}_{\hat{\boldsymbol{\beta}}}^2$ eine $k \times k$ Matrix mit den Schätzern für die Varianzen und Kovarianzen der Koeffizienten $\hat{\boldsymbol{\beta}}$ ist. Die Wurzel der Hauptdiagonalelemente der Matrix $\hat{\boldsymbol{\sigma}}_{\hat{\boldsymbol{\beta}}}^2$ sind die Standardfehler der Koeffizienten.

Der Standardfehler des h -ten Koeffizienten $\hat{\beta}_h$ ist $\hat{\sigma}_{\hat{\beta}_h} = \hat{\sigma} \sqrt{v_{hh}}$, wobei v_{hh} das h -te Diagonalelement der Matrix $(\mathbf{X}' \mathbf{X})^{-1}$ ist. Da wir die Varianz aus der Stichprobe schätzen ist die entsprechende Teststatistik wieder t-verteilt mit $n - k$ Freiheitsgraden.

Daraus folgt

$$t\text{-Stat}(\hat{\beta}_h) = \frac{\hat{\beta}_h - \beta_h}{\hat{\sigma}_{\hat{\beta}_h}} = \frac{\hat{\beta}_h - \beta_h}{\hat{\sigma} \sqrt{v_{hh}}} \sim t_{n-k}$$

bzw. das Konfidenzintervall

$$\hat{\beta}_h \pm t_{\alpha/2}^c (\hat{\sigma} \sqrt{v_{hh}})$$

Wir können die bisherigen Ergebnisse also zusammenfassen:

Unter den Gauss-Markov Annahmen gilt

¹⁰In Stata kann mit dem postestimation Befehl `e(rmse)` (für Root MSE) auf den Standardfehler und mit `e(rss)` auf die Quadratsumme der Residuen zugreifen; in R erhält man die Quadratsumme der Residuen mit `deviance(eqname)` und den Standardfehler der Regression mit `summary(eqname)$sigma`.

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \\ \widehat{\text{var}}(\hat{\beta}) &:= \hat{\sigma}_{\hat{\beta}}^2 = \hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1} \\ \hat{\sigma}^2 &= \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n-k}\end{aligned}$$

Für $\varepsilon_i \sim N(0, \sigma^2)$ und $\hat{\beta}$ statistisch unabhängig von $\hat{\varepsilon}$ folgt:

$$\begin{aligned}\hat{\beta} &\sim N(\beta, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}) \\ t &= \frac{\hat{\beta}_h - \beta_h}{\sqrt{\hat{\sigma}^2 (\mathbf{X}'\mathbf{X})_{hh}^{-1}}} \sim t_{n-k}\end{aligned}$$

wobei $(\mathbf{X}'\mathbf{X})_{hh}^{-1}$ das h -te Diagonalelement der Matrix $(\mathbf{X}'\mathbf{X})^{-1}$ bezeichnet.

6.7 Die Varianz einer Linearkombination von Koeffizienten*

Um die Varianz einer Linearkombination von OLS Koeffizienten zu berechnen definieren wir einen $k \times 1$ Vektor \mathbf{r} . Durch entsprechende Wahl der Werte von \mathbf{r} können wir beliebige Linearkombinationen der Parameter berechnen; wenn z.B. $\mathbf{r} = \mathbf{1}$ erhalten wir mit $\mathbf{r}'\beta$ die Summe der β ; wenn \mathbf{r} in der Zeile h eine Eins enthält und sonst nur Nullen pickt $\mathbf{r}'\beta$ Element β_h heraus

$$\mathbf{r}'\beta = \begin{pmatrix} 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \sum_{i=1}^3 \beta_i \quad \text{oder} \quad \begin{pmatrix} 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \beta_2$$

Definieren wir einen Skalar $\gamma = \mathbf{r}'\beta$. Wir interessieren uns für die Varianz von $\hat{\gamma} = \mathbf{r}'\hat{\beta}$, diese ist

$$\text{var}(\hat{\gamma}) = \mathbf{r}' \text{var}(\hat{\beta}) \mathbf{r} = \sigma^2 \mathbf{r}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{r}$$

weil

$$\begin{aligned}\text{var}(\mathbf{r}'\hat{\beta}) &= E \left(\mathbf{r}'(\hat{\beta} - \beta)(\hat{\beta} - \beta)' \mathbf{r} \right) \\ &= \mathbf{r}' E \left((\hat{\beta} - \beta)(\hat{\beta} - \beta)' \right) \mathbf{r} \\ &= \mathbf{r}' (\sigma^2 (\mathbf{X}'\mathbf{X})^{-1}) \mathbf{r}\end{aligned}$$

Beispiel: Für $y = \hat{\beta}_1 + \hat{\beta}_2 x_2 + \hat{\beta}_3 x_3 + \hat{\varepsilon}$ ist die Varianz von $\hat{\gamma} = 2\hat{\beta}_2 - \hat{\beta}_3$ gleich $\text{var}(\hat{\gamma}) = 4 \text{var}(\hat{\beta}_2) + \text{var}(\hat{\beta}_3) - 4 \text{cov}(\hat{\beta}_2, \hat{\beta}_3)$ (binomische Formel). Um dies in Ma-

trixnotation zu zeigen definieren wir $\mathbf{r} = (0 \ 2 \ -1)'$

$$\begin{aligned} \text{var}(\mathbf{r}'\hat{\boldsymbol{\beta}}) &= \begin{pmatrix} 0 & 2 & -1 \end{pmatrix} \begin{pmatrix} \text{var}(\hat{\beta}_1) & \text{cov}(\hat{\beta}_1, \hat{\beta}_2) & \text{cov}(\hat{\beta}_1, \hat{\beta}_3) \\ \text{cov}(\hat{\beta}_1, \hat{\beta}_2) & \text{var}(\hat{\beta}_2) & \text{cov}(\hat{\beta}_2, \hat{\beta}_3) \\ \text{cov}(\hat{\beta}_1, \hat{\beta}_3) & \text{cov}(\hat{\beta}_3, \hat{\beta}_2) & \text{var}(\hat{\beta}_3) \end{pmatrix} \begin{pmatrix} 0 \\ 2 \\ -1 \end{pmatrix} \\ &= \begin{pmatrix} 2 \text{cov}(\hat{\beta}_1, \hat{\beta}_2) - \text{cov}(\hat{\beta}_1, \hat{\beta}_3) & 2 \text{var}(\hat{\beta}_2) - \text{cov}(\hat{\beta}_3, \hat{\beta}_2) & 2 \text{cov}(\hat{\beta}_2, \hat{\beta}_3) - \text{var}(\hat{\beta}_3) \end{pmatrix} \begin{pmatrix} 0 \\ 2 \\ -1 \end{pmatrix} \\ &= 4 \text{var}(\hat{\beta}_2) + \text{var}(\hat{\beta}_3) - 4 \text{cov}(\hat{\beta}_2, \hat{\beta}_3) \end{aligned}$$

6.8 Effizienz des OLS Schätzers*

Das Gauss Markov Theorem für die Effizienz des OLS Schätzers kann natürlich auch unter Zuhilfenahme der Matrixnotation bewiesen werden.

Das Gauss Markov Theorem besagt, dass bei Gültigkeit der Gauss Markov Annahmen der OLS Schätzer $\hat{\boldsymbol{\beta}}$ eine kleinere Varianz hat als alle alternativen linearen unverzerrten Schätzfunktionen $\tilde{\boldsymbol{\beta}}$.

Für den Beweis definieren wir eine beliebige lineare Schätzfunktion $\tilde{\boldsymbol{\beta}}$ und untersuchen, unter welchen Bedingungen diese Schätzfunktion unverzerrt ist. Sodann berechnen wir für diese Schätzfunktion $\tilde{\boldsymbol{\beta}}$ unter der Bedingung der Unverzerrtheit die Varianz-Kovarianzmatrix. Man kann dann relativ einfach zeigen, dass die Differenz zwischen dieser Varianz-Kovarianzmatrix und der Varianz-Kovarianzmatrix des OLS Schätzers immer eine positiv-semidefinite Matrix ist.

Dies impliziert, dass für jedes einzelne Element der Schätzfunktionen gilt

$$\text{var}(\tilde{\beta}_h) - \text{var}(\hat{\beta}_h) \geq d \quad \text{mit } d \geq 0$$

für $h = 1, \dots, k$. Analoges gilt auch für jede beliebige Linearkombination der Schätzfunktionen.

Beweis:* Wir definieren eine beliebige lineare Schätzfunktion¹¹

$$\tilde{\boldsymbol{\beta}} = \mathbf{C}\mathbf{y}$$

wobei \mathbf{C} eine $k \times n$ Matrix mit Konstanten ist, die eine Funktion der \mathbf{X} sein können, ähnlich wie die $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ Matrix für den OLS Schätzer.

Der Schätzer $\tilde{\boldsymbol{\beta}}$ ist unverzerrt, wenn

$$\text{E}(\tilde{\boldsymbol{\beta}}) = \boldsymbol{\beta}$$

Dies impliziert Restriktionen auf die Matrix \mathbf{C}

$$\begin{aligned} \text{E}(\tilde{\boldsymbol{\beta}}) &= \text{E}(\mathbf{C}\mathbf{y}) \\ &= \text{E}(\mathbf{C}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon})) \\ &= \mathbf{C}\mathbf{X}\boldsymbol{\beta} + \text{E}(\mathbf{C}\boldsymbol{\varepsilon}) \\ &= \mathbf{C}\mathbf{X}\boldsymbol{\beta} \quad (\text{wenn } \text{E}(\mathbf{C}\boldsymbol{\varepsilon}) = \mathbf{C}\text{E}(\boldsymbol{\varepsilon}) = \mathbf{0}) \end{aligned}$$

¹¹Dieser Beweis folgt eng Vogelvang (2005, 72f).

Deshalb ist $\tilde{\beta}$ nur eine unverzerrte Schätzfunktion wenn

$$\mathbf{C}\mathbf{X} = \mathbf{I}_k$$

Die Varianz-Kovarianzmatrix von $\tilde{\beta}$ ist

$$\begin{aligned} \text{var}(\tilde{\beta}) &= \text{var}(\mathbf{C}\mathbf{y}) = \text{var}(\mathbf{C}(\mathbf{X}\beta + \varepsilon)) \\ &= \text{var}(\mathbf{C}\mathbf{X}\beta + \mathbf{C}\varepsilon) \\ &= \text{var}(\mathbf{C}\varepsilon) \quad (\text{wenn } \tilde{\beta} \text{ unverzerrt ist}) \\ &= \mathbf{C} \text{var}(\varepsilon) \mathbf{C}' \\ &= \mathbf{C}\sigma^2 \mathbf{I}_n \mathbf{C}' \quad (\text{wenn } \varepsilon_i \sim \text{i.i.d.}(0, \sigma^2)) \\ &= \sigma^2 \mathbf{C}\mathbf{C}' \end{aligned}$$

Zum Vergleich, die Varianz-Kovarianzmatrix des OLS Schätzers ist

$$\text{var}(\hat{\beta}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$$

Für den Gauss Markov Beweis benötigen wir die Differenz

$$\left[\text{var}(\tilde{\beta}) - \text{var}(\hat{\beta}) \right]$$

und erinnern uns, dass $\tilde{\beta} = \mathbf{C}\mathbf{y}$ und $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$.

Deshalb definieren wir eine $k \times n$ Matrix \mathbf{D} mit den Differenzen

$$\mathbf{D} = \mathbf{C} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$$

die wir umschreiben zu

$$\mathbf{C} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' + \mathbf{D}$$

Wir multiplizieren beide Seiten mit \mathbf{X} und verwenden die Bedingung für Unverzerrtheit $\mathbf{C}\mathbf{X} = \mathbf{I}_k$

$$\mathbf{C}\mathbf{X} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X} + \mathbf{D}\mathbf{X}$$

Daraus folgt $\mathbf{I}_k = \mathbf{I}_k + \mathbf{D}\mathbf{X}$ oder

$$\mathbf{D}\mathbf{X} = \mathbf{O}$$

Damit ist's fast geschafft, denn unter Verwendung dieser Eigenschaft erhalten wir

$$\begin{aligned} \mathbf{C}\mathbf{C}' &= [(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' + \mathbf{D}] [(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' + \mathbf{D}]' \\ &= [(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' + \mathbf{D}] [\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} + \mathbf{D}'] \\ &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} + \mathbf{D}\mathbf{D}' \quad (\text{weil } \mathbf{D}\mathbf{X} = \mathbf{O}) \\ &= (\mathbf{X}'\mathbf{X})^{-1} + \mathbf{D}\mathbf{D}' \end{aligned}$$

Der Unterschied zwischen den beiden positiv definiten $k \times k$ Matrizen $\mathbf{C}\mathbf{C}'$ und $(\mathbf{X}'\mathbf{X})^{-1}$ ist die positiv semidefinite Matrix $\mathbf{D}\mathbf{D}'$.

Wir haben bereits gezeigt, dass $\text{var}(\tilde{\beta}) = \sigma^2 \mathbf{C}\mathbf{C}'$ und, unter den Gauss Markov Annahmen, $\text{var}(\hat{\beta}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$.

Also folgt

$$\begin{aligned}\text{var}(\tilde{\beta}) &= \sigma^2 \mathbf{C}\mathbf{C}' \\ &= \sigma^2 [(\mathbf{X}'\mathbf{X})^{-1} + \mathbf{D}\mathbf{D}'] \\ &= \text{var}(\hat{\beta}) + \sigma^2 \mathbf{D}\mathbf{D}'\end{aligned}$$

oder

$$\text{var}(\tilde{\beta}) - \text{var}(\hat{\beta}) = \sigma^2 \mathbf{D}\mathbf{D}'$$

Dies impliziert, dass sowohl jedes einzelne Element von $\tilde{\beta}$ als auch jede beliebige Linearkombination dieser Elemente kleinere (oder höchstens gleich große) Varianz haben als die der entsprechenden Elemente von $\hat{\beta}$.

Dies kann einfach gezeigt werden: \mathbf{D} ist eine $k \times n$ Matrix, deshalb hat die symmetrische Matrix $\mathbf{D}\mathbf{D}'$ die Dimension $k \times k$. Aus der linearen Algebra ist bekannt, dass Matrizen der Form $\mathbf{D}\mathbf{D}'$ immer positiv semidefinit sind, d.h., dass für jeden beliebigen $k \times 1$ Vektor \mathbf{r} (für alle $\mathbf{r} \neq \mathbf{0}$) gilt $\mathbf{r}'\mathbf{D}\mathbf{D}'\mathbf{r} \geq 0$ (man beachte, dass dies ein Skalar ist).¹²

Deshalb muss gelten

$$\mathbf{r}' (\text{var}(\tilde{\beta}) - \text{var}(\hat{\beta})) \mathbf{r} = \sigma^2 \mathbf{r}'\mathbf{D}\mathbf{D}'\mathbf{r} \geq 0$$

bzw.

$$\mathbf{r}' \text{var}(\tilde{\beta}) \mathbf{r} = \mathbf{r}' \text{var}(\hat{\beta}) \mathbf{r} + \sigma^2 \mathbf{r}'\mathbf{D}\mathbf{D}'\mathbf{r}$$

mit $\sigma^2 \mathbf{r}'\mathbf{D}\mathbf{D}'\mathbf{r} \geq 0$.

Wir können den Vektor \mathbf{r} beliebig wählen, also z.B. auch, dass er in der Zeile h eine Eins hat und sonst überall Null, dann wird dadurch genau die Varianz des Koeffizienten $\tilde{\beta}_h$ (bzw. $\hat{\beta}_h$) 'ausgeschnitten'. Aber durch geeignete Wahl von \mathbf{r} kann auch jede beliebige Linearkombination von $\tilde{\beta}$ erzeugt werden.

Damit ist gezeigt, dass sowohl jedes einzelne Element von $\tilde{\beta}$ als auch die jede beliebige Linearkombination der einzelnen Elemente von $\tilde{\beta}$ immer eine kleinere (oder höchstens gleich große) Varianz hat, als die entsprechende Varianz jeder beliebigen anderen linearen unverzerrten Schätzfunktion, sofern die Gauss Markov Annahmen erfüllt sind! Der OLS Schätzer ist also ein BLUE ('best linear unbiased estimator').

Dieser Beweis erfolgte für deterministische \mathbf{X} , man kann aber zeigen, dass das Gauss Markov Theorem auch für stochastische Regressoren gilt, wenn alle Gauss Markov Annahmen erfüllt sind (für einen ebenso eleganten wie kurzen Beweis siehe z.B. White and Cho (2012)).

6.9 Konsistenz des OLS Schätzers*

Wir erinnern uns, dass eine Schätzfunktion $\hat{\theta}$ konsistent ist, wenn sowohl der Bias als auch die Varianz von $\hat{\theta}_n$ gegen Null konvergieren, wenn $n \rightarrow \infty$. Dies wird geschrieben als

$$\hat{\theta} \xrightarrow{p} \theta \quad \text{oder} \quad \text{plim}(\hat{\theta}) = \theta$$

¹²Dies ist einfach zu sehen, sei \mathbf{X} eine $n \times k$ Matrix und \mathbf{r} ein $k \times 1$ Spaltenvektor. Dann gilt für den Skalar $\mathbf{r}'\mathbf{X}'\mathbf{X}\mathbf{r} = (\mathbf{X}\mathbf{r})'\mathbf{X}\mathbf{r} = \|\mathbf{X}\mathbf{r}\|^2 \geq 0$.

Die Konsistenz des OLS Schätzers kann auch für *stochastische Regressoren* relativ einfach bewiesen werden.

Dazu benötigen wir zusätzliche Annahmen, die aber nicht allzu streng sind. Im Wesentlichen verlangen diese Annahmen, der datengenerierende Prozess durch $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ beschrieben werden kann, wobei die einzelnen Paare $(y_i, x_{2i}, \dots, x_{ki})$ (mit $i = 1, \dots, n$) Zufallsziehungen aus einer gemeinsamen Verteilung sind (i.i.d. Stichprobe), wobei große Ausreißer unwahrscheinlich sein sollen (die ersten vier Momente existieren und sind nicht unendlich groß). Außerdem müssen die Regressoren exogen sein ($E(\varepsilon_i | \mathbf{X}) = 0$).

Um die Konsistenz des OLS Schätzers zu beweisen müssen wir zeigen, dass

$$\text{plim}_{n \rightarrow \infty}(\hat{\boldsymbol{\beta}}) = \text{plim}_{n \rightarrow \infty}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}) = \boldsymbol{\beta}$$

Für den plim von stochastischen Matrizen und Vektoren gilt, wenn die Multiplikation möglich ist und \mathbf{X} nicht singulär ist

$$\begin{aligned} \text{plim}_{n \rightarrow \infty}(\mathbf{X}^{-1}) &= \text{plim}_{n \rightarrow \infty}(\mathbf{X})^{-1} \\ \text{plim}_{n \rightarrow \infty}(\mathbf{X}\mathbf{y}) &= \text{plim}_{n \rightarrow \infty} \mathbf{X} \text{plim}_{n \rightarrow \infty} \mathbf{y} \end{aligned}$$

Wenden wir uns zuerst der Matrix $\mathbf{X}'\mathbf{X}$ zu

$$\begin{aligned} \mathbf{X}'\mathbf{X} &= \begin{pmatrix} x_{11} & x_{21} & \cdots & x_{n1} \\ x_{12} & x_{22} & \cdots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1k} & x_{2k} & \cdots & x_{nk} \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix} \\ &= \begin{pmatrix} \sum_{i=1}^n x_{i1}^2 & \sum_{i=1}^n x_{i1}x_{i2} & \cdots & \sum_{i=1}^n x_{i1}x_{ik} \\ \sum_{i=1}^n x_{i2}x_{i1} & \sum_{i=1}^n x_{i2}^2 & \cdots & \sum_{i=1}^n x_{i2}x_{ik} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^n x_{ik}x_{i1} & \sum_{i=1}^n x_{ik}x_{i2} & \cdots & \sum_{i=1}^n x_{ik}^2 \end{pmatrix} \end{aligned}$$

Ein typisches Element ist $\sum_{i=1}^n x_{ih}x_{ig}$ (mit $h, g = 1, \dots, k$), eine Summe von n Produkten. Wenn $n \rightarrow \infty$ können wir deshalb nicht erwarten, dass diese Summe von Produkten gegen eine endliche Zahl konvergiert, deshalb existiert im allgemeinen für die Matrix $\mathbf{X}'\mathbf{X}$ kein plim.

Hingegen sollte unter obigen Annahmen der Durchschnitt von n Zufallsvariablen, d.h. $\frac{1}{n} \sum_{i=1}^n x_{ih}x_{ig} = q_{hg}$, gegen eine feste Zahl konvergieren, bzw. in Matrixschreibweise

$$\text{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}'\mathbf{X} = \mathbf{Q}_{\mathbf{X}'\mathbf{X}}$$

wobei $\mathbf{Q}_{\mathbf{X}'\mathbf{X}}$ eine endliche nichtstochastische Matrix mit vollem Rang k ist.

Beim Beweis der Erwartungstreue haben wir bereits gezeigt, dass

$$\hat{\boldsymbol{\beta}} = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\varepsilon}$$

Also ist

$$\begin{aligned}\text{plim}_{n \rightarrow \infty}(\hat{\beta}) &= \text{plim}_{n \rightarrow \infty}(\beta) + \text{plim}_{n \rightarrow \infty}((\mathbf{X}'\mathbf{X})^{-1}) \text{plim}_{n \rightarrow \infty}(\mathbf{X}'\varepsilon) \\ &= \beta + \left(\text{plim}_{n \rightarrow \infty} \frac{1}{n} \mathbf{X}'\mathbf{X}\right)^{-1} \text{plim}_{n \rightarrow \infty} \left(\frac{1}{n} \mathbf{X}'\varepsilon\right) \\ &= \beta + (\mathbf{Q}_{\mathbf{X}'\mathbf{X}})^{-1} \text{plim}_{n \rightarrow \infty} \left(\frac{1}{n} \mathbf{X}'\varepsilon\right)\end{aligned}$$

Der Schätzer $\hat{\beta}$ ist also konsistent, wenn $\text{plim}_{n \rightarrow \infty}(\frac{1}{n} \mathbf{X}'\varepsilon) = \mathbf{0}$. Wir schreiben diesen $k \times 1$ Vektor noch einmal ausführlich an

$$\frac{1}{n} \mathbf{X}'\varepsilon = \begin{pmatrix} \frac{1}{n} \sum_i \varepsilon_i \\ \frac{1}{n} \sum_i x'_{i2} \varepsilon_i \\ \vdots \\ \frac{1}{n} \sum_i x'_{ik} \varepsilon_i \end{pmatrix}$$

wobei hier die erste Spalte der \mathbf{X} Matrix ein Einsenvektor ist ($\hat{\beta}_1$ ist also das Interzept). Man beachte, dass die k Elemente dieses Vektors den Kovarianzen zwischen den Regressoren und den Störtermen ähneln. Falls die Regressoren exogen sind sollten die Elemente dieses Vektors gegen Null konvergieren.

Daraus folgt

$$\begin{aligned}\text{plim}_{n \rightarrow \infty}(\hat{\beta}) &= \beta + (\mathbf{Q}_{\mathbf{X}'\mathbf{X}})^{-1} \text{plim}_{n \rightarrow \infty} \left(\frac{1}{n} \mathbf{X}'\varepsilon\right) \\ &= \beta + (\mathbf{Q}_{\mathbf{X}'\mathbf{X}})^{-1} \mathbf{0} \\ &= \beta\end{aligned}$$

Das bedeutet, wann immer $\text{plim}_{n \rightarrow \infty}(\frac{1}{n} \mathbf{X}'\varepsilon) = \mathbf{0}$ und einige weniger strenge Annahmen erfüllt sind ist der OLS Schätzer konsistent.

Zum Beweis der Konsistenz ist (im Unterschied zur Effizienz) keine Homoskedastizität erforderlich, aber die Annahmen A1 – A3 und sowie endliche zweite Momente (Varianzen dürfen nicht unendlich groß werden).

6.10 Asymptotische Normalverteilung des OLS Schätzers*

Die asymptotische Normalverteilung des OLS Schätzers wird benötigt, wenn die Verteilung der Grundgesamtheit unbekannt ist und zumindest für große Stichproben Hypothesentests benötigt werden oder asymptotische Konfidenzintervalle werden.¹³

Die vollständige Herleitung wäre etwas aufwändiger, deshalb hier nur eine kurze Skizze.

Wir wollen zeigen, dass

$$\sqrt{n}(\hat{\beta}_n - \beta) \xrightarrow{d} N(\mathbf{0}, \Sigma)$$

¹³Dieser Abschnitt wurde noch nicht überprüft und enthält möglicherweise Fehler :(

Dabei ist $\hat{\beta}_n$ die auf n Beobachtungen beruhende Schätzfunktion für den Parametervektor β . Da die Varianz von $\hat{\beta}_n$ mit zunehmenden n abnimmt muss mit der Konvergenzgeschwindigkeit \sqrt{n} multipliziert (skaliert) werden, damit man eine stabile Grenzverteilung erhält.

Beispiel: Varianz des arithmetischen Mittels $\hat{\mu}$

Um das Kollabieren oder Divergieren der Varianz von $(\hat{\mu}_n - \mu)$ mit zunehmenden Stichprobenumfang n zu verhindern kann man mit der Konvergenzgeschwindigkeit \sqrt{n} multiplizieren. Denn aus

$$\text{var}(\hat{\mu}) := n^{-1} \sigma^2$$

folgt

$$\text{var}(\sqrt{n}(\hat{\mu}_n - \mu)) = n \text{var}(\hat{\mu}_n - \mu) = n \frac{1}{n} \sigma^2 = \sigma^2 = \text{asy var}(\hat{\mu}_n)$$

Das heißt, ist σ^2 die asymptotische Varianz des arithmetischen Mittels und die Konvergenzrate beträgt \sqrt{n} . Analoges gilt für das lineare Regressionsmodell.

Das Symbol \xrightarrow{d} bedeutet Konvergenz der Dichte nach (*density*), und $N(\mathbf{0}, \Sigma)$ bedeutet normalverteilt mit Erwartungswert $\mathbf{0}$ und Kovarianzmatrix Σ .

Das lineare Modell ist

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon$$

mit

$$\hat{\beta}_n = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} = \beta + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\varepsilon$$

Wir dividieren durch n

$$\hat{\beta}_n - \beta = \left(\frac{1}{n} \mathbf{X}'\mathbf{X} \right)^{-1} \left(\frac{1}{n} \mathbf{X}'\varepsilon \right)$$

Unter der Annahme von i.i.d. Daten und endlichen zweiten Momenten können Gesetze der Großen Zahl angewandt werden, d.h.

•

$$\frac{1}{n} (\mathbf{X}'\mathbf{X}) \xrightarrow{p} E(\mathbf{X}'\mathbf{X}) := \mathbf{Q}$$

konvergiert der Wahrscheinlichkeit nach zu einer positiv definiten Matrix.

•

$$\frac{1}{\sqrt{n}} \mathbf{X}'\varepsilon \xrightarrow{p} \mathbf{0}$$

Wir wollen die Verteilung von

$$\frac{1}{\sqrt{n}} \mathbf{X}'\varepsilon = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}'_i \varepsilon_i$$

bestimmen.

Darauf können wir einen Zentralen Grenzwertsatz anwenden, wenn die folgenden Annahmen erfüllt sind:

- Die Datenpaare (x_i, ε_i) sind i.i.d.
- $E(\varepsilon_i | \mathbf{x}_i) = 0 \Rightarrow E(\mathbf{x}_i \varepsilon_i) = 0$
- $E(\varepsilon_i^2 | \mathbf{x}_i) = \sigma^2(\mathbf{x}_i)$ (Störterme können heteroskedastisch sein)
- $E(\mathbf{x}_i' \mathbf{x}_i \varepsilon_i^2) < \infty$

Unter diesen Annahmen folgt aus einem multivariaten Zentralen Grenzwertsatz

$$\frac{1}{\sqrt{n}} \mathbf{X}' \boldsymbol{\varepsilon} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i' \varepsilon_i \xrightarrow{d} N(\mathbf{0}, \boldsymbol{\Omega}) \quad \text{mit } \boldsymbol{\Omega} = E(\mathbf{x}_i' \mathbf{x}_i \varepsilon_i)$$

Daraus folgt unter Anwendung von Slutsky's Theorem

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \xrightarrow{d} N(\mathbf{0}, \mathbf{Q}^{-1} \boldsymbol{\Omega} \mathbf{Q}^{-1})$$

mit $\mathbf{Q} = E(\mathbf{x}_i' \mathbf{x}_i)$ und $\boldsymbol{\Omega} = E(\mathbf{x}_i' \mathbf{x}_i \varepsilon_i^2)$

Daraus folgt

- Wenn $n \rightarrow \infty$ konvergiert die Stichprobenkennwertverteilung gegen die Normalverteilung, zentriert um den wahren Wert $\boldsymbol{\beta}$.
- Die Varianz der Stichprobenkennwertverteilung hängt sowohl von der Störterme als auch von der Verteilung der Regressoren ab.
- Dieses Resultat erlaubt die Durchführung von Hypothesentests und die Bestimmung von asymptotischen Konfidenzintervallen in großen Stichproben, unabhängig von der Verteilung des datengenerierenden Prozesses.

Übungsaufgaben

1. Berechnen Sie *mit Excel* die Koeffizienten, deren Standardabweichungen sowie das R^2 des Modells

$$y_i = \hat{\beta}_1 + \hat{\beta}_2 x_{i2} + \hat{\beta}_3 x_{i3} + \hat{\varepsilon}_i$$

für folgende Daten

| y | x_2 | x_3 |
|-----|-------|-------|
| 2 | 9 | 1 |
| 5 | 4 | 2 |
| 4 | 7 | 3 |
| 8 | 2 | 4 |
| 9 | 3 | 5 |
| 9 | 1 | 6 |

Beachten Sie, dass die Gleichung ein Interzept enthält!

Hinweis: Excel erlaubt u.a. das Transponieren [=MTRANS(*Bereich*)], Multiplizieren [=MMULT(*Bereich1*; *Bereich2*)] und Invertieren [=MINV(*Bereich*)] von Matrizen. Markieren Sie dazu zuerst den Ausgabebereich (in der richtigen Dimension!), geben die Formel ein und schließen Sie mit <Umschalt>+<Strg>+<Eingabe> die Eingabe ab (d.h. bei gedrückter Shift- und Strg-Taste die Eingabetaste drücken).

2. Das bivariate Modell in Matrixschreibweise ist

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} + \begin{pmatrix} \hat{\varepsilon}_1 \\ \hat{\varepsilon}_2 \\ \vdots \\ \hat{\varepsilon}_n \end{pmatrix}$$

(a) Zeigen Sie, dass

$$(\mathbf{X}'\mathbf{X}) = \begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix}$$

(b) Daraus lässt sich leicht die Inverse berechnen

$$(\mathbf{X}'\mathbf{X})^{-1} = \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix}$$

Verwenden Sie dieses Resultat, um $\hat{\beta}_1$, $\hat{\beta}_2$ sowie deren Varianzen $\hat{\sigma}_{\hat{\beta}_1}^2$ und $\hat{\sigma}_{\hat{\beta}_2}^2$ zu berechnen.

(c) Zeigen Sie anhand dieses Beispiels, dass $(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{X}) = \mathbf{I}$.

3. (a) Zeigen Sie, dass die Projektionsmatrix $\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ und die residuenerzeugende Matrix $\mathbf{M} = (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')$ symmetrisch und idempotent sind.

(b) Zeigen Sie, dass

$$\mathbf{MP} = \mathbf{PM} = \mathbf{O}$$

4. Angenommen \mathbf{z} sei ein $n \times 1$ Vektor mit lauter Einsen, d.h. $\mathbf{z} = (1, \dots, 1)'$. Zeigen Sie, dass für $\mathbf{X} = \mathbf{z}$ die residuenerzeugende Matrix $\mathbf{M} = (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')$ gleich

$$\mathbf{M} = \mathbf{I} - \frac{\mathbf{z}\mathbf{z}'}{n}$$

ist, und dass $\mathbf{M}\mathbf{y}$ den Vektor \mathbf{y} in die Abweichungen vom Mittelwert $(y_i - \bar{y})$ überführt (dies ist u.a. in der Panelökonometrie von Bedeutung).

6.11 Eine geometrische Interpretation des OLS-Schätzers

In diesem Abschnitt werden wir zeigen, dass der OLS-Schätzer eine sehr einfache geometrische Interpretation hat, nämlich als *orthogonale Projektion in den Spaltenraum*.¹⁴ Wir werden das Grundprinzip an einem sehr einfachen Beispiel erläutern.

Wir gehen von nur zwei Beobachtungen sowie von einer abhängigen und einer einzigen erklärenden Variable aus. Zur Vereinfachung nehmen wir eine Regression durch den Ursprung an, d.h. ohne Interzept. Die Daten seien

$$\mathbf{y} = \begin{pmatrix} 1 \\ 4 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} 3 \\ 5 \end{pmatrix}$$

¹⁴Zwei Vektoren stehen *orthogonal* aufeinander, wenn sie einen rechten Winkel einschließen.

und das Regressionsmodell ist

$$\mathbf{y} = \mathbf{x}\hat{\beta} + \hat{\varepsilon}$$

Das linke Panel in Abbildung 6.3 zeigt die herkömmliche grafische Darstellung. Jede Beobachtung ist ein Punkt im x, y Raum, und dieser Raum hat zwei Dimensionen, eine für y und eine für x . In dieser Darstellung wird die Regressionsgerade bestimmt, indem die Quadratsumme der Residuen $\hat{\varepsilon}_1^2 + \hat{\varepsilon}_2^2$ minimiert wird. Der Regressionskoeffizient $\hat{\beta}$ ist die Steigung dieser Regressionsgerade. Gäbe es mehr als zwei Beobachtungen wäre jede weitere Beobachtung ein zusätzlicher Punkt in diesem Raum. Gäbe es eine weitere x -Variable würde dafür eine weitere Dimension benötigt.

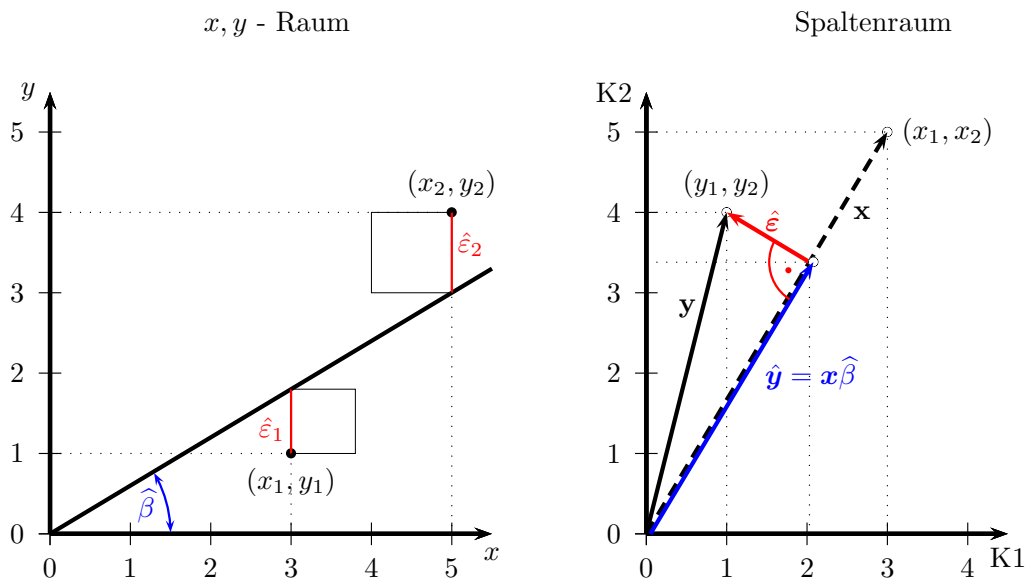


Abbildung 6.3: OLS-Schätzer in herkömmlicher Darstellung im x, y Raum (links) und als orthogonale Projektion im Euklidischen Raum (rechts).
Mit $\mathbf{x}' = (3 \ 5)$, $\mathbf{y}' = (1 \ 4)$

Es gibt eine alternative grafische Darstellungsmöglichkeit dazu. Wir wissen, dass die n Elemente des Vektors \mathbf{y} auch Vektoren im euklidischen Raum \mathbb{E}^n dargestellt werden können. Ebenso können die einzelnen Spalten der \mathbf{X} -Matrix als Vektoren in \mathbb{E}^n dargestellt werden.

Das rechte Panel in Abbildung 6.3 zeigt die Vektor-Darstellung der gleichen Daten im Euklidischen Raum \mathbb{E}^2 , da wir zwei Beobachtungen haben. Man beachte, dass hier die Koordinaten nicht mit x und y beschriftet sind, sondern mit K1 und K2 für Koordinate 1 und 2.

In der Abbildung im rechten Panel sind nicht die einzelnen Beobachtungen als Punkte eingezeichnet, sondern die n -dimensionalen Vektoren. Dabei sind \mathbf{y} und \mathbf{x} unabhängig von der Zahl der Beobachtungen jeweils Vektoren. Da wir in diesem Beispiel nur zwei Beobachtungen haben ($n = 2$) können wir die Vektoren im zweidimensionalen Raum darstellen. Hätten wir hundertausend Beobachtungen wäre \mathbf{y} trotzdem nur ein Vektor (ebenso wie \mathbf{x}), aber im 100 000-dimensionalen Raum, den wir grafisch natürlich nicht darstellen können (im Gegensatz dazu hätten wir im x, y Raum nach wie vor zwei Dimensionen, aber 100,000 Beobachtungspunkte).

Hätten wir mehrere x -Variablen wäre jede weitere x -Variable ein weiterer $n \times 1$ Vektor im n -dimensionalen Euklidischen Raum \mathbb{E}^n .

Ein *Unterraum* (*subspace*) von \mathbb{E}^n hat eine niedrigere Dimension als n und wird von sogenannten *Basisvektoren* aufgespannt.

Ein spezieller Unterraum, der für das Folgende wichtig ist, wird von den k Spalten der \mathbf{X} -Matrix aufgespannt, und wird deshalb *Spaltenraum* genannt. Diesen Unterraum, der von den k Spaltenvektoren aufgespannt wird, bezeichnen wir mit $\mathcal{S}(\mathbf{X})$, bzw. $\mathcal{S}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k)$.

Ein Unterraum $\mathcal{S}(\mathbf{X})$ besteht aus allen Vektoren, die als Linearkombination der Basisvektoren \mathbf{x}_h (mit $h = 1, 2, \dots, k$) gebildet werden können.

Für Mathematik-Freaks noch die formale Definition:

$$\mathcal{S}(\mathbf{X}) \equiv \left\{ \mathbf{z} \in \mathbb{E}^n \left| \mathbf{z} = \sum_{h=1}^k \hat{\beta}_h \mathbf{x}_h, \hat{\beta}_h \in \mathbb{R} \right. \right\}$$

Da wir in diesem einfachen Beispiel nur einen \mathbf{x} Vektor haben hat dieser Spaltenraum Dimension Eins. Er besteht aus allen $\hat{\beta}\mathbf{x}$ mit $\hat{\beta} \in \mathbb{R}$.

Im Falle von k erklärenden Variablen ist jede der erklärenden x -Variablen ein Vektor in \mathbb{E}^n , und diese insgesamt k Vektoren spannen einen k -dimensionalen Unterraum auf.

Durch die OLS-Methode werden die \mathbf{y} so in diesen Spaltenraum $\mathcal{S}(\mathbf{X})$ projiziert, dass die *Länge* des Residuenvektors minimal wird.

Dies ist im rechten Panel in Abbildung 6.3 dargestellt. Die kürzeste Entfernung zwischen \mathbf{y} und dem Spaltenraum besteht dann, wenn der Residuenvektor $\hat{\mathbf{e}}$ rechtwinklig auf den Spaltenraum steht, oder in anderen Worten, wenn $\hat{\mathbf{e}}$ und \mathbf{x} *orthogonal* sind.

Wir erinnern uns, dass zwei Vektoren orthogonal sind, wenn das innere Produkt Null ist, d.h., $\mathbf{x} \perp \hat{\mathbf{e}}$ wenn $\mathbf{x}'\hat{\mathbf{e}} = 0$ (vgl. Abbildung 6.4).

Man beachte, dass $\mathbf{x}'\hat{\mathbf{e}} = 0$ eine Bedingung erster Ordnung bei der Herleitung des OLS-Schätzers war.

$$\mathbf{x}_1 = \begin{pmatrix} 3 \\ 2 \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} -1 \\ 1.5 \end{pmatrix}$$

$$\begin{aligned} \mathbf{x}_1' \mathbf{x}_2 &= \begin{pmatrix} 3 & 2 \end{pmatrix} \begin{pmatrix} -1 \\ 1.5 \end{pmatrix} \\ &= 3 \times (-1) + 2 \times 1.5 \\ &= 0 \end{aligned}$$

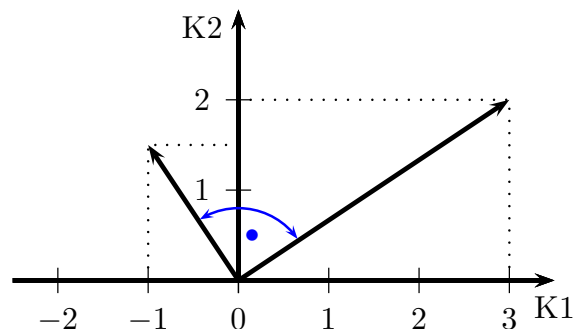


Abbildung 6.4: Orthogonalität: Zwei Vektoren \mathbf{x}_1 und \mathbf{x}_2 stehen rechtwinklig aufeinander, d.h. sie sind *orthogonal* $\mathbf{x}_1 \perp \mathbf{x}_2$, wenn und nur wenn $\mathbf{x}_1' \mathbf{x}_2 = 0$

Da die Länge (bzw. der Betrag oder die Norm) eines Vektors $\hat{\mathbf{e}}$ definiert ist als¹⁵

$$||\hat{\mathbf{e}}|| = \sqrt{\hat{\mathbf{e}}' \hat{\mathbf{e}}}$$

¹⁵Dies folgt aus dem Pythagoräischen Lehrsatz.

minimiert der OLS-Schätzer zugleich die Länge des Residuenvektors $\hat{\varepsilon}$, da $\hat{\varepsilon}'\hat{\varepsilon}$ die Quadratsumme der Residuen ist, die durch den OLS-Schätzer minimiert wird.

Abbildung 6.5 zeigt die in Lehrbüchern übliche orthogonale Projektion in den Spaltenraum. In diesem Fall ist der Spaltenraum nur eindimensional (man beachte, dass diese Abbildung nur aus dem rechten Panel von Abbildung 6.3 ‘herausgezeichnet’ wurde).

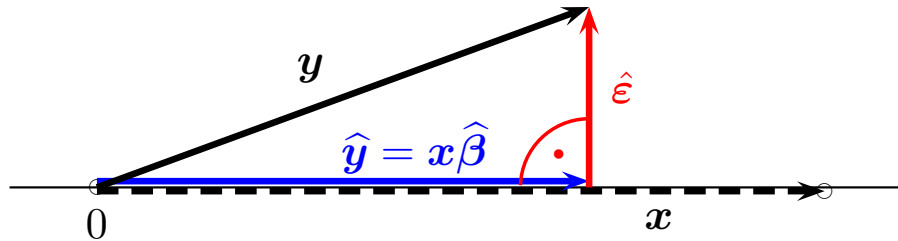


Abbildung 6.5: OLS-Schätzfunktion $\hat{\beta}$ als Ergebnis einer orthogonalen Projektion in den Spaltenraum.

Wie man sieht ist $y = \hat{y} + \hat{\varepsilon}$, wobei $\hat{y} = x\hat{\beta}$ ein Vielfaches des x -Vektors ist.

Wenn es zwei erklärende x Variablen gibt spannen diese zwei Vektoren eine Ebene als Spaltenraum $\mathcal{S}(\mathbf{X})$ auf. Durch die OLS-Methode wird y so in diesen Spaltenraum $\mathcal{S}(\mathbf{X})$ projiziert, dass die Länge des Residuenvektors minimal wird.

Wenn eine weitere Beobachtung dazu kommt wird die Abbildung im Euklidischen Raum dreidimensional (im x, y Raum ändert sich die Dimension nicht, dort kommt dadurch nur ein weiterer Punkt dazu). Abbildung 6.6 zeigt den Fall für drei Beobachtungen und eine erklärende Variable. Die drei Beobachtungen spannen einen 3-dimensionalen euklidischen Raum auf, und y wird orthogonal auf den eindimensionalen Spaltenraum, der durch den Vektor x aufgespannt wird, projiziert.

Abbildung 6.7 zeigt eine Abbildung der gleichen 2 Vektoren wie in Abbildung 6.6, aber mit einem zusätzlichen Vektor x_2 (also $n = 3$ und $k = 2$). Die beiden Vektoren x_1 und x_2 spannen einen zweidimensionalen Spaltenraum auf (d.h. eine Ebene), und y wird wieder orthogonal auf diese Ebene projiziert.

Wenn man den Spaltenraum $\mathcal{S}(\mathbf{X})$ wieder ‘herauszeichnet’ erhält man eine Grafik wie in Abbildung 6.8. Dies ist eine übliche Darstellung einer orthogonalen Projektion in den Spaltenraum, wie sie in Lehrbüchern häufig zu finden ist. In dieser Abbildung wird gut sichtbar, dass \hat{y} eine Linearkombination der Vektoren x_1 und x_2 ist, wobei die ‘Gewichte’ $\hat{\beta}_1$ und $\hat{\beta}_2$ so bestimmt werden, dass der Vektor $\hat{\varepsilon}$ möglichst kurz wird.

Höhere Dimensionen können grafisch nicht mehr dargestellt werden, aber das Prinzip bleibt gleich. Das Interzept ist einfach ein Einsen-Vektor und spannt eine eigene Dimension auf, d.h. es gibt keine Asymmetrie zwischen dem Interzept und den anderen erklärenden Variablen.

Man beachte, dass die Basisvektoren x_1, x_2, \dots, x_k nicht rechtwinklig aufeinander stehen müssen, die x -Variablen können also durchaus untereinander korreliert sein (und sind es normalerweise auch), aber sie müssen *linear unabhängig* sein, d.h. sie

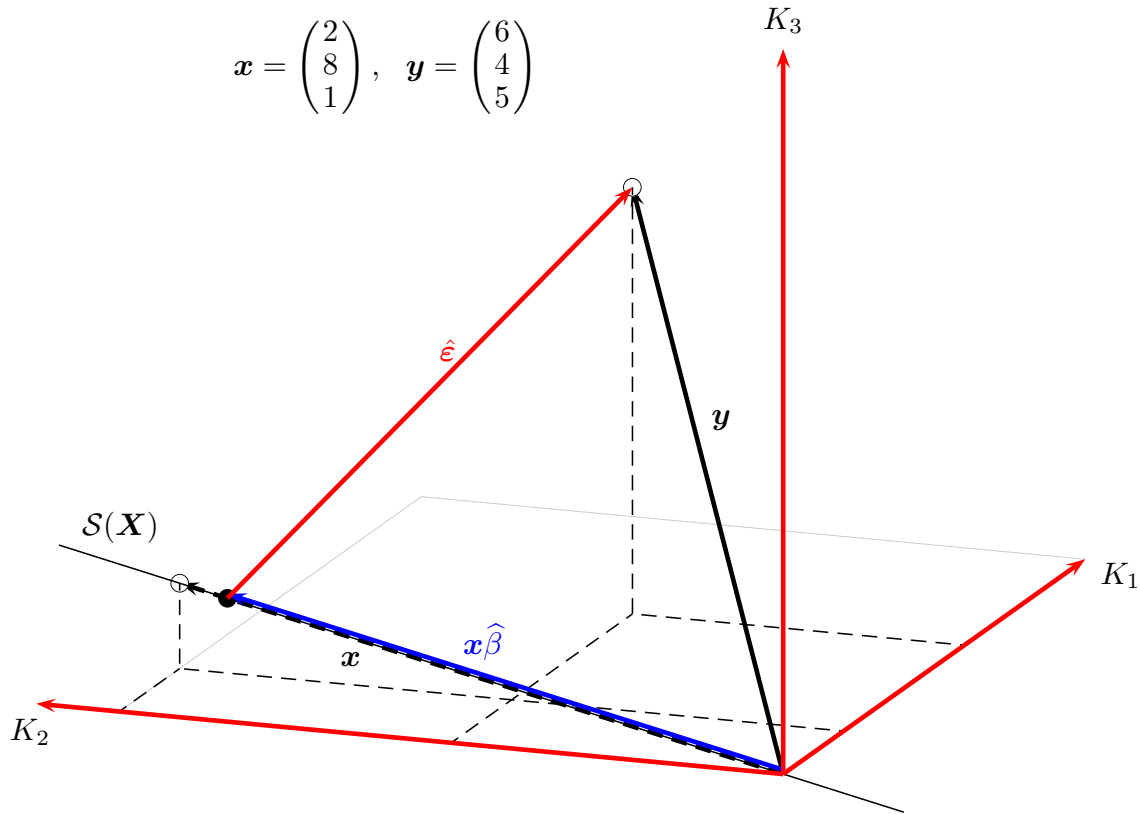


Abbildung 6.6: OLS-Schätzer $\hat{\beta}$ als orthogonale Projektion für $n = 3$ und $k = 1$. Der Spaltenraum $\mathcal{S}(X)$ ist eindimensional (alle Linearkombinationen des Vektors x), und $\hat{\varepsilon}$ steht orthogonal auf x .

dürfen nicht perfekt korreliert sein. Dies ist die Bedingung, dass die X Matrix vollen Spaltenrang haben muss, oder in anderen Worten, dass keine *perfekte Multikollinearität* vorliegen darf.

In Abschnitt 6.6 haben wir die Projektionsmatrix $P = X(X'X)^{-1}X'$ und die residuenerzeugende Matrix $M = I - X(X'X)^{-1}X'$ kennengelernt.

Wir erinnern uns:

$$\hat{\varepsilon} = y - X\hat{\beta} = [I - X(X'X)^{-1}X'] y = My$$

und

$$\hat{y} = y - \hat{\varepsilon} = (I - M)y = X(X'X)^{-1}X'y = Py$$

Durch Vormultiplikation von y mit P wird y in den Spaltenraum $\mathcal{S}(X)$ projiziert, daher der Name Projektionsmatrix.

Wie man einfach überprüfen kann sind P und M orthogonal, d.h.

$$PM = MP = O$$

Ein Punkt, der schon im Spaltenraum liegt, wird nicht transformiert, denn

$$PX = X(X'X)^{-1}X'X = X$$

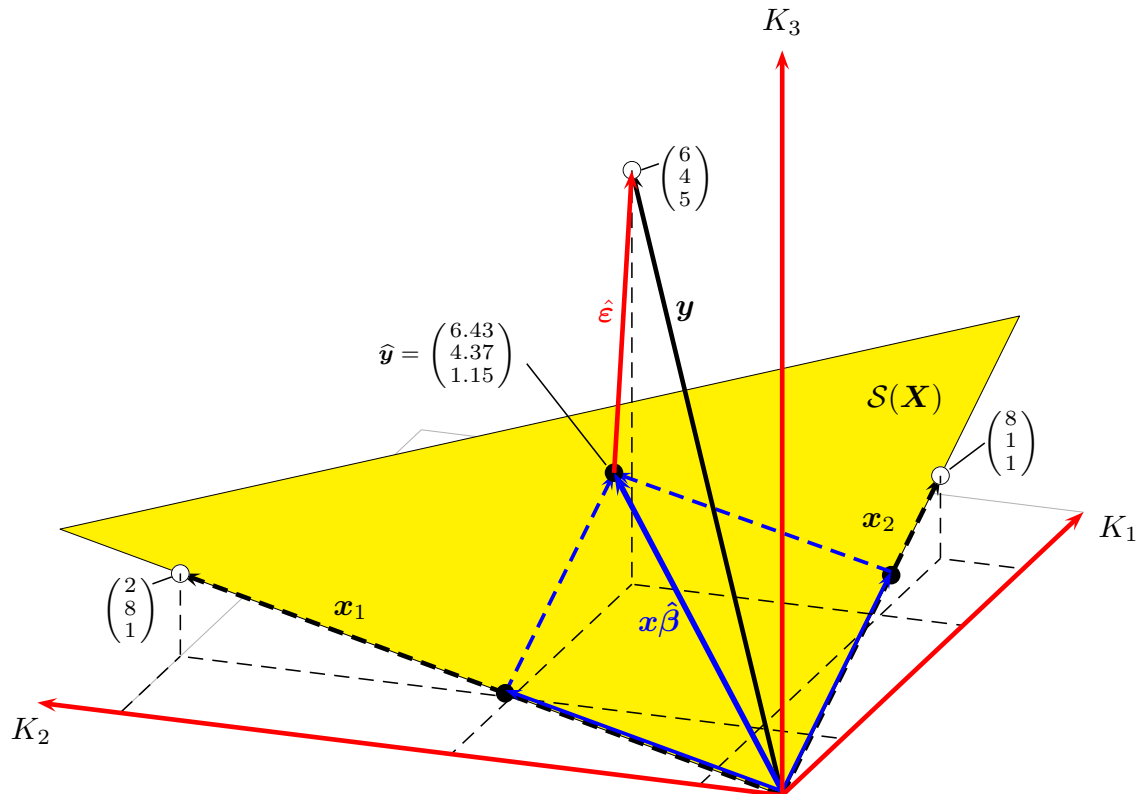


Abbildung 6.7: OLS-Schätzer $\hat{\beta}$ als orthogonale Projektion für $n = 3$ und $k = 2$. Der Spaltenraum $\mathcal{S}(X)$ ist zweidimensional (alle Linearkombinationen der Vektoren x_1 und x_2), und $\hat{\varepsilon}$ steht orthogonal auf diese Ebene (bzw. auf $x\hat{\beta}$).

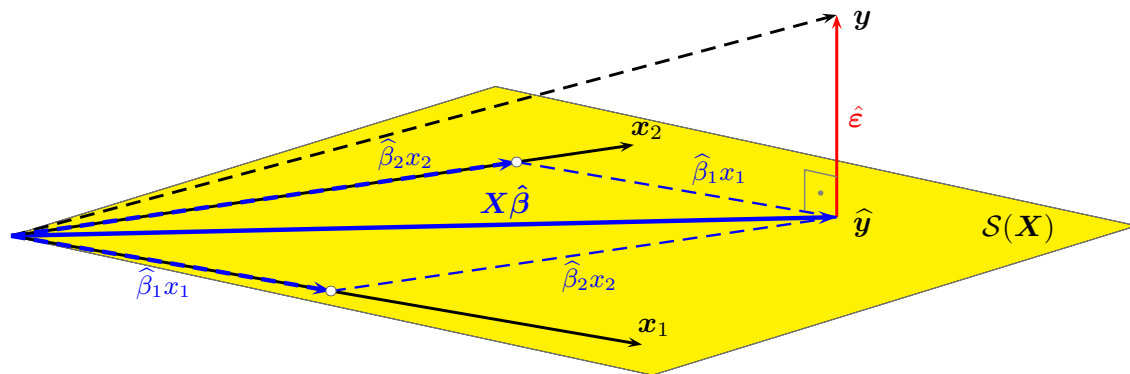


Abbildung 6.8: OLS-Schätzer $\hat{\beta}$ als Ergebnis einer orthogonalen Projektion in den Spaltenraum.

Mit Hilfe der Projektions- und residuenerzeugenden Matrix kann y in zwei orthogonale Teile zerlegt werden, denn

$$\begin{aligned} y &= \hat{y} + \hat{\varepsilon} \\ &= Py + My \\ &= \text{Projektion} + \text{Residuum} \end{aligned}$$

Beispiel: Wir berechnen die Abbildung 6.3 zugrunde liegenden \mathbf{P} und \mathbf{M} Matrizen. Die Vektoren sind

$$\mathbf{y} = \begin{pmatrix} 1 \\ 4 \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} 3 \\ 5 \end{pmatrix}$$

und das Modell ist

$$\mathbf{y} = \mathbf{x}\hat{\beta} + \hat{\varepsilon}$$

$$\mathbf{P} = \mathbf{x}(\mathbf{x}'\mathbf{x})^{-1}\mathbf{x}' = \begin{pmatrix} 3 \\ 5 \end{pmatrix} (34)^{-1} \begin{pmatrix} 3 & 5 \end{pmatrix} = (34)^{-1} \begin{pmatrix} 9 & 15 \\ 15 & 25 \end{pmatrix}$$

$$\hat{\mathbf{y}} = \mathbf{P}\mathbf{y} = (34)^{-1} \begin{pmatrix} 9 & 15 \\ 15 & 25 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} = \frac{1}{34} \begin{pmatrix} 9 + 60 \\ 15 + 100 \end{pmatrix} \approx \begin{pmatrix} 2.03 \\ 3.38 \end{pmatrix}$$

$$\mathbf{M} = \mathbf{I} - \mathbf{P} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} - (34)^{-1} \begin{pmatrix} 9 & 15 \\ 15 & 25 \end{pmatrix} = \frac{1}{34} \begin{pmatrix} 25 & -15 \\ -14 & 9 \end{pmatrix}$$

$$\hat{\varepsilon} = \mathbf{M}\mathbf{y} = \frac{1}{34} \begin{pmatrix} 25 & -15 \\ -15 & 9 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} = \frac{1}{34} \begin{pmatrix} 25 - 60 \\ -15 + 36 \end{pmatrix} \approx \begin{pmatrix} -1.03 \\ 0.62 \end{pmatrix}$$

Kontrolle:

$$\mathbf{y} = \hat{\mathbf{y}} + \hat{\varepsilon} = \begin{pmatrix} 2.03 \\ 3.38 \end{pmatrix} + \begin{pmatrix} -1.03 \\ 0.62 \end{pmatrix} = \begin{pmatrix} 1 \\ 4 \end{pmatrix}$$

Herkömmliche Berechnung als OLS-Schätzer:

$$\hat{\beta} = (\mathbf{x}'\mathbf{x})^{-1}\mathbf{x}'\mathbf{y} = (34)^{-1} \begin{pmatrix} 3 & 5 \end{pmatrix} \begin{pmatrix} 1 \\ 4 \end{pmatrix} = \frac{23}{34} = 0.676$$

$$\hat{\mathbf{y}} = \mathbf{x}\hat{\beta} = \begin{pmatrix} 3 \\ 5 \end{pmatrix} 0.676 \approx \begin{pmatrix} 2.03 \\ 3.38 \end{pmatrix}$$

Diese geometrische Interpretation ist für viele Fälle nützlich, zum Beispiel um zu verstehen, was mit den Koeffizienten eines Regressionsmodells passiert, wenn eine Variable mit einer Konstanten multipliziert wird.

Literaturverzeichnis

Johnston, J. and Dinardo, J. (1996), *Econometric Methods*, 4 edn, McGraw-Hill/Irwin.

Vogelvang, B. (2005), *Econometrics: Theory & Applications With EvIEWS*, Financial Times Management.

White, Jr., H. L. and Cho, J. S. (2012), 'A Three Line Proof that OLS is BLUE', Available at SSRN: <http://ssrn.com/abstract=2050611>.

URL: <http://dx.doi.org/10.2139/ssrn.2050611>

6.A Appendix

6.A.1 Alternative Vektorschreibweise

In der Literatur wird häufig für Beobachtung i der $(k \times 1)$ Spaltenvektor \mathbf{x}_i zu definiert

$$\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ik} \end{pmatrix}$$

sodass die Transponierte \mathbf{x}_i' einfach Zeile i aus Matrix \mathbf{X} darstellt. Mit dieser Definition kann $y_i = \sum_h x_{ih}b_h + \hat{\varepsilon}_i$, d.h. die Gleichung für eine spezifische Beobachtung, auch in Vektornotation geschrieben werden als

$$y_i = \mathbf{x}_i' \hat{\boldsymbol{\beta}} + \hat{\varepsilon}_i$$

was nur eine andere Kurzschreibweise für das komplette Gleichungssystem ist

$$y_i = x_{i1}\hat{\beta}_1 + x_{i2}\hat{\beta}_2 + \cdots + x_{ik}\hat{\beta}_k + \hat{\varepsilon}_i$$

Mit dieser Vektor-Schreibweise kann der OLS Schätzer alternativ in folgender Vektornotation geschrieben werden

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \left(\sum_i \mathbf{x}_i\mathbf{x}_i' \right)^{-1} \sum_i \mathbf{x}_i y_i$$

Diese Vektor-Schreibweise erweist sich für manche spätere Anwendungen als nützlich.

6.A.2 Ableitungen von Matrizen

Beweis für

$$\frac{\partial \hat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{y}}{\partial \hat{\boldsymbol{\beta}}} = \mathbf{X}' \mathbf{y}$$

Da \mathbf{X} eine $n \times k$ Matrix und \mathbf{y} ein $n \times 1$ Vektor ist, ist $\mathbf{X}'\mathbf{y}$ ein $k \times 1$ Spaltenvektor. Zur Vereinfachung der Schreibweise bezeichnen wir diesen Spaltenvektor mit \mathbf{c} , d.h. $\mathbf{c} := \mathbf{X}'\mathbf{y}$.

Wir schreiben $\hat{\boldsymbol{\beta}}'\mathbf{c}$ ausführlich an

$$\hat{\boldsymbol{\beta}}'\mathbf{c} = \begin{pmatrix} \hat{\beta}_1 & \hat{\beta}_2 & \cdots & \hat{\beta}_k \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_k \end{pmatrix} = \hat{\beta}_1 c_1 + \hat{\beta}_2 c_2 + \cdots + \hat{\beta}_k c_k$$

Also ist

$$\frac{\partial \hat{\boldsymbol{\beta}}'\mathbf{c}}{\partial \hat{\beta}_1} = c_1, \quad \frac{\partial \hat{\boldsymbol{\beta}}'\mathbf{c}}{\partial \hat{\beta}_2} = c_2, \quad \dots, \quad \frac{\partial \hat{\boldsymbol{\beta}}'\mathbf{c}}{\partial \hat{\beta}_k} = c_k,$$

oder einfacher

$$\frac{\partial \hat{\beta}' \mathbf{c}}{\partial \hat{\beta}} = \mathbf{c}$$

bzw. in unserem Fall mit $\mathbf{c} := \mathbf{X}'\mathbf{y}$

$$\frac{\partial \hat{\beta}' \mathbf{X}'\mathbf{y}}{\partial \hat{\beta}} = \mathbf{X}'\mathbf{y}$$

■

Beweis für

$$\frac{\partial \hat{\beta}' \mathbf{X}' \mathbf{X} \hat{\beta}}{\partial \hat{\beta}} = 2 \mathbf{X}' \mathbf{X} \hat{\beta}'$$

Wir beachten, dass $\mathbf{X}'\mathbf{X}$ symmetrisch ist und schreiben zur einfacheren Darstellung \mathbf{A} für $\mathbf{X}'\mathbf{X}$ (d.h. $\mathbf{X}'\mathbf{X} := \mathbf{A}$).

$$\begin{aligned} \hat{\beta}' \mathbf{A} \hat{\beta} &= \begin{pmatrix} \hat{\beta}_1 & \hat{\beta}_2 & \cdots & \hat{\beta}_k \end{pmatrix} \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1k} \\ a_{21} & a_{22} & \cdots & a_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ a_{k1} & a_{k2} & \cdots & a_{kk} \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} \\ &= \begin{pmatrix} \sum_{h=1}^k \hat{\beta}_h a_{h1} & \sum_{h=1}^k \hat{\beta}_h a_{h2} & \cdots & \sum_{h=1}^k \hat{\beta}_h a_{hk} \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} \\ &= \hat{\beta}_1^2 a_{11} + \hat{\beta}_1 \hat{\beta}_2 a_{21} + \cdots + \hat{\beta}_1 \hat{\beta}_k a_{k1} + \\ &\quad + \hat{\beta}_2 \hat{\beta}_1 a_{12} + \hat{\beta}_2^2 a_{22} + \cdots + \hat{\beta}_2 \hat{\beta}_k a_{k2} + \\ &\quad \vdots \\ &\quad + \hat{\beta}_k \hat{\beta}_1 a_{1k} + \hat{\beta}_k \hat{\beta}_2 a_{2k} + \cdots + \hat{\beta}_k^2 a_{kk} \end{aligned}$$

wenn \mathbf{A} symmetrisch ist (d.h. $a_{ij} = a_{ji}$) sind die jeweiligen Ableitungen

$$\begin{aligned} \frac{\partial \hat{\beta}' \mathbf{A} \hat{\beta}}{\partial \hat{\beta}_1} &= 2\hat{\beta}_1 a_{11} + \hat{\beta}_2 a_{21} + \cdots + \hat{\beta}_k a_{k1} + \hat{\beta}_2 a_{12} + \cdots + \hat{\beta}_k a_{1k} \\ &= 2(\hat{\beta}_1 a_{11} + \hat{\beta}_2 a_{12} + \cdots + \hat{\beta}_k a_{1k}) \\ \frac{\partial \hat{\beta}' \mathbf{A} \hat{\beta}}{\partial \hat{\beta}_2} &= 2(\hat{\beta}_1 a_{21} + \hat{\beta}_2 a_{22} + \cdots + \hat{\beta}_k a_{2k}) \\ &\vdots \\ \frac{\partial \hat{\beta}' \mathbf{A} \hat{\beta}}{\partial \hat{\beta}_k} &= 2(\hat{\beta}_1 a_{k1} + \hat{\beta}_2 a_{k2} + \cdots + \hat{\beta}_k a_{kk}) \end{aligned}$$

oder kürzer

$$\frac{\partial \hat{\beta}' A \hat{\beta}}{\partial \hat{\beta}} = 2 A \hat{\beta}$$

bzw. in unserem Fall mit $A := X'X$

$$\frac{\partial \hat{\beta}' X' X \hat{\beta}}{\partial \hat{\beta}} = 2 X' X \hat{\beta}$$

Hinweis: $\hat{\beta}' X' X \hat{\beta}$ ist eine *quadratische Form*.

■

6.A.3 Der OLS-Schätzer mit den Matrixbefehlen von R

Das folgende kleine Program zeigt, wie der OLS-Schätzer etwas komplizierter mit den Matrixbefehlen von R berechnet werden kann.

```
X <- matrix(c(1,1,1,1,1,0,1,2,3,4),nr=5)
y <- c(3, 5, 8, 9, 9)

b <- solve(crossprod(X))%*%t(X)%*%y
e <- y - X%*%b
R2 <- 1 - (crossprod(e)/(crossprod(y)-5*mean(y)^2))
se <- sqrt(crossprod(e)/(5-2))

cov <- as.numeric(crossprod(e)/(5-2)) * solve(crossprod(X))
StdErr <- sqrt(diag(cov))
t.stat <- b/StdErr
p.val <- 2*(1 - pt(abs(b/StdErr), df=5-2))

# Kontrolle
x2 <- 0:4
(eq <- summary(lm(y ~ x2)))
```

6.B Wald-Test auf allgemeine lineare Restriktionen in Matrixschreibweise

Wie schon erwähnt beruhen sehr viele Tests in der Ökonometrie auf einem Vergleich zwischen einem *restringiertem* Modell und einem *nicht-restringierten* Modell, wobei das restringierte Modell als ‘Spezialfall’ des nicht-restringierten Modells dargestellt werden kann (sogenannte ‘*nested*’ Tests).

Im Folgenden wollen wir eine allgemeinere Art lineare Restriktionen abzubilden darstellen. Wir werden gleich sehen, dass sich lineare Hypothesen immer in der Form

$$R\beta = r$$

anschreiben lassen, wobei die Matrix \mathbf{R} aus Konstanten besteht, die aus der Nullhypothese folgen, und die Dimension $q \times k$ hat (q ist die Anzahl der Restriktionen und k die Anzahl der geschätzten Koeffizienten). Der Vektor \mathbf{r} hat die Dimension $q \times 1$ und besteht ebenfalls aus Konstanten, deren Wert(e) aus der Nullhypothese folgen.

Beispiele: Das Modell sei

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i2} + \varepsilon_i$$

- Die Nullhypothese $H_0: \beta_2 = c$ (wobei c eine beliebige Konstante ist) kann mit $\mathbf{R} = [0 \ 1 \ 0]$ und $\mathbf{r} = c$ (also ein Skalar) angeschrieben werden als

$$\mathbf{R}\boldsymbol{\beta} = [0 \ 1 \ 0] \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = c = \mathbf{r}$$

bzw. $\beta_2 = c$ mit $q = 1$ (eine Restriktion).

- $H_0: \beta_2 + \beta_3 = 1$:

$$\mathbf{R} = [0 \ 1 \ 1], \quad \mathbf{r} = 1, \quad q = 1$$

$$\mathbf{R}\boldsymbol{\beta} = [0 \ 1 \ 1] \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = 1$$

bzw.

$$\beta_2 + \beta_3 = 1$$

mit $q = 1$ (eine Restriktion).

- $H_0: \beta_2 = \beta_3$ (bzw. $H_0: \beta_2 - \beta_3 = 0$):

$$\mathbf{R} = [0 \ 1 \ -1], \quad \mathbf{r} = 0, \quad q = 1$$

- $H_0: \beta_2 = \beta_3$ und $\beta_1 = 1$:

$$\mathbf{R} = \begin{bmatrix} 0 & 1 & -1 \\ 1 & 0 & 0 \end{bmatrix}, \quad \mathbf{r} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad q = 2$$

weil

$$\begin{bmatrix} 0 & 1 & -1 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

gleich ist

$$\begin{aligned} \beta_2 - \beta_3 &= 0 \\ \beta_1 &= 1 \end{aligned}$$

Ganz allgemein können lineare Hypothesen als $\mathbf{R}\boldsymbol{\beta} - \mathbf{r} = \mathbf{0}$ geschrieben werden, wobei die Matrix \mathbf{R} die Dimension $q \times k$ hat, wobei q die Anzahl der Restriktionen und k die Anzahl der Koeffizienten ist. Die Anzahl der Restriktionen q darf die Anzahl der Koeffizienten k nicht übersteigen, d.h. es muss $q < k$ gelten.

Wenn wir als (OLS-) Schätzung für $\boldsymbol{\beta}$ wieder $\hat{\boldsymbol{\beta}}$ verwenden ist der Vektor $(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})$ ein Maß für die ‘Abweichung’ der erwarteten Werte von den beobachteten Werten.

Wie üblich können wir nun die Stichprobenkennwertverteilung von $\mathbf{R}\hat{\boldsymbol{\beta}}$ unter der Nullhypothese $\mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ bestimmen.

Man kann zeigen, dass die folgende Teststatistik

$$F\text{-Stat} = (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})' \left[\widehat{\text{Rvar}}(\hat{\boldsymbol{\beta}}) \mathbf{R}' \right]^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}) \frac{1}{q} \sim F_{q, n-k}$$

F -verteilt ist und für einen Test beliebiger linearer Nullhypothesen $H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ herangezogen werden kann.

Bevor wir diese Teststatistik allgemein herleiten wollen wir eine kurze und eher intuitive Erklärung geben.

6.B.1 Eine intuitive Darstellung

Zur Verdeutlichung der Idee gehen wir zunächst induktiv vor, d.h. wir bedienen uns für den einfachen Hypothesentest in einem bivariaten Modell der Matrixschreibweise.

Angenommen, wir möchten die Nullhypothese $H_0: \beta_2 = 0$ testen. Um dieses Problem allgemeiner in Matrixschreibweise anzuschreiben definieren wir

$$\mathbf{R} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \quad \text{und} \quad \mathbf{r} = 0$$

Wenn wir den Koeffizientenvektor mit diesem \mathbf{R} vormultiplizieren erhalten wir den Koeffizienten $\hat{\beta}_2$:

$$\mathbf{R}\boldsymbol{\beta} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = 0 = \mathbf{r}$$

oder $\beta_2 = 0$.

Also können Null- und Alternativhypothese auch geschrieben werden als

$$H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{r} \quad \text{gegen} \quad H_1: \mathbf{R}\boldsymbol{\beta} \neq \mathbf{r}$$

Dementsprechendes gilt natürlich auch für die Stichprobe:

$$\mathbf{R}\hat{\boldsymbol{\beta}} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{bmatrix} = \hat{\beta}_2$$

Als nächstes benötigen wir die Varianz von $\hat{\beta}_2$.

Mit Hilfe der selben Matrix \mathbf{R} kann auch das gewünschte Element der Varianz-Kovarianzmatrix von $\hat{\beta}$ ‘freigestellt’ werden, indem die Varianz- Kovarianzmatrix der Koeffizienten $\widehat{\text{var}}(\hat{\beta})$ mit \mathbf{R} vor- und mit \mathbf{R}' nachmultipliziert wird, also

$$\begin{aligned}\widehat{\text{var}}(\hat{\beta}_2) &= \widehat{\text{var}}(\mathbf{R}\hat{\beta}) = \mathbf{R}\widehat{\text{var}}(\hat{\beta})\mathbf{R}' \\ &= \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \widehat{\text{var}}(\hat{\beta}_1) & \widehat{\text{cov}}(\hat{\beta}_1, \hat{\beta}_2) & \widehat{\text{cov}}(\hat{\beta}_1, \hat{\beta}_3) \\ \widehat{\text{cov}}(\hat{\beta}_1, \hat{\beta}_2) & \widehat{\text{var}}(\hat{\beta}_2) & \widehat{\text{cov}}(\hat{\beta}_2, \hat{\beta}_3) \\ \widehat{\text{cov}}(\hat{\beta}_1, \hat{\beta}_3) & \widehat{\text{cov}}(\hat{\beta}_2, \hat{\beta}_3) & \widehat{\text{var}}(\hat{\beta}_3) \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ &= \widehat{\text{var}}(\hat{\beta}_2)\end{aligned}$$

Wir gehen nun von der üblichen t-Statistik für den Test eines einzelnen Koeffizienten aus und bedienen uns der Matrixschreibweise:

$$t = \frac{\hat{\beta}_2 - \beta_2}{\sqrt{\widehat{\text{var}}(\hat{\beta}_2)}} = \frac{\mathbf{R}\hat{\beta} - \mathbf{R}\beta}{\sqrt{\mathbf{R} \underbrace{\widehat{\text{var}}(\hat{\beta})}_{\hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}} \mathbf{R}'}} = \frac{\mathbf{R}\hat{\beta} - \mathbf{R}\beta}{\sqrt{\mathbf{R}\hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'}}$$

bzw.

$$t = \frac{\mathbf{R}\hat{\beta} - \mathbf{R}\beta}{\hat{\sigma} \sqrt{\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'}}$$

Man kann zeigen, dass das Quadrat einer t-verteilten Zufallsvariable mit n Freiheitsgraden F-verteilt mit einem Zähler und n Nennerfreiheitsgraden ist, d.h. $t^2(n) = F(1, n)$ (siehe z.B. Johnston & DiNardo 1997, p. 489f).

Das Quadrat der obigen t-Statistik wäre demnach F -verteilt und könnte folgendermaßen angeschrieben werden:

$$\begin{aligned}F &= \frac{(\mathbf{R}\hat{\beta} - \mathbf{r})' [\mathbf{R}\widehat{\text{var}}(\hat{\beta})\mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})}{q} \\ &= \frac{(\mathbf{R}\hat{\beta} - \mathbf{r})' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})}{q\hat{\sigma}^2} \\ &= \frac{(\mathbf{R}\hat{\beta} - \mathbf{r})' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})/q}{\hat{\epsilon}'\hat{\epsilon}/(n-k)}\end{aligned}$$

(wir erinnern uns, dass $\widehat{\text{var}}(\hat{\beta}) = \hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}$ und $\hat{\sigma}^2 = \hat{\epsilon}'\hat{\epsilon}/(n-k)$)

Dies gilt tatsächlich, und sogar wenn wir allgemeiner q lineare Hypothesen testen und die Matrix \mathbf{R} deshalb die Dimension $(q \times k)$ hat und \mathbf{r} ein $(q \times 1)$ Vektor ist, allerdings ist der Beweis etwas komplizierter.

6.B.2 Eine etwas allgemeinere Darstellung

Wir möchten die Stichprobenkennwertverteilung von $\mathbf{R}\hat{\beta}$ unter der Nullhypothese $\mathbf{R}\beta = \mathbf{r}$ bestimmen.

Da $E(\hat{\beta}) = \beta$ und \mathbf{R} deterministisch ist folgt

$$E(\mathbf{R}\hat{\beta}) = \mathbf{R}\beta$$

sowie

$$\begin{aligned}
 \text{var}(\mathbf{R}\hat{\boldsymbol{\beta}}) &= \text{E} \left[\mathbf{R}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{R}' \right] \\
 &= \mathbf{R} \text{E} \left[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \right] \mathbf{R}' \\
 &= \mathbf{R} \text{var}(\hat{\boldsymbol{\beta}}) \mathbf{R}' \\
 &= \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'
 \end{aligned}$$

Wenn

$$\boldsymbol{\varepsilon} \sim \text{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

folgt

$$\mathbf{R}\hat{\boldsymbol{\beta}} \sim \text{N} \left[\mathbf{R}\boldsymbol{\beta}, \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}' \right]$$

bzw.

$$\mathbf{R}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sim \text{N} \left[\mathbf{0}, \sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}' \right]$$

Man kann zeigen, dass unter dieser Bedingung die folgende Variable χ^2 -verteilt mit q Freiheitsgraden ist

$$(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})' [\sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}) \sim \chi^2(q)$$

Dies folgt aus der Verteilung quadratischer Formen (siehe z.B. Johnston & DiNardo 1997, S. 493ff). Um eine Intuition dafür zu bekommen starten wir mit einem $(k \times 1)$ Vektor \mathbf{x} , dessen Elemente unabhängig standardnormalverteilt seien (d.h. Mittelwert = 0 und Varianz = 1, bzw. $\mathbf{x} \sim \text{N}(\mathbf{0}, \mathbf{I})$). Wir erinnern uns, dass die Quadratsumme von k unabhängig standardnormalverteilten Variablen χ^2 -verteilt ist, also

$$\mathbf{x}'\mathbf{x} \sim \chi^2(k)$$

Angenommen, \mathbf{x} seien nicht standardnormalverteilt, sondern

$$\mathbf{x} \sim \text{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

Die einzelnen x sind immer noch unabhängig und haben einen Mittelwert Null, aber die Elemente müssen durch σ^2 dividiert werden damit die Varianz Eins wird.

Also

$$\frac{x_1^2}{\sigma^2} + \frac{x_2^2}{\sigma^2} + \dots + \frac{x_k^2}{\sigma^2} = \frac{1}{\sigma^2} \mathbf{x}'\mathbf{x} \sim \chi^2(k)$$

Dies kann auch geschrieben werden als

$$\mathbf{x}'(\sigma^2 \mathbf{I})^{-1} \mathbf{x} \sim \chi^2(k)$$

Wenn die \mathbf{x} zudem nicht unabhängig verteilt sind, also

$$\mathbf{x} \sim \text{N}(\mathbf{0}, \boldsymbol{\Sigma})$$

wäre der äquivalente Ausdruck

$$\mathbf{x}'(\boldsymbol{\Sigma})^{-1} \mathbf{x} \sim \chi^2(k)$$

was in der Struktur dem obigen Ausdruck

$$(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})' [\sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}) \sim \chi^2(q)$$

entspricht, wenn $\mathbf{x} = (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})$ und $\boldsymbol{\Sigma} = [\sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']$.

Dieses Resultat ist tatsächlich richtig, wenngleich der Beweis dafür etwas komplizierter ist¹⁶ (siehe z.B. Johnston & DiNardo 1997, S. 494).

Allgemein gilt für $\mathbf{x} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ und wenn \mathbf{A} symmetrisch und idempotent mit Rang q ist, dass

$$\frac{1}{\sigma^2} \mathbf{x}' \mathbf{A} \mathbf{x} \sim \chi^2(q)$$

Mit Hilfe dieses Resultats kann u.a. gezeigt werden, dass die Quadratsumme der Residuen des Regressionsmodells χ^2 -verteilt ist.

Wir erinnern uns, dass der Residuenvektor $\hat{\boldsymbol{\varepsilon}}$ geschrieben werden kann als

$$\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{M}\mathbf{y} \quad \text{mit} \quad \mathbf{M} = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$$

wobei \mathbf{M} eine symmetrische und idempotente Matrix mit Rang $n-k$ ist und $\mathbf{M}\mathbf{X} = \mathbf{0}$. Daraus folgt

$$\hat{\boldsymbol{\varepsilon}} = \mathbf{M}\mathbf{y} = \mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{M}\boldsymbol{\varepsilon}$$

und

$$\hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}} = \boldsymbol{\varepsilon}' \mathbf{M} \boldsymbol{\varepsilon}$$

Unter der Annahme $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ folgt also

$$\frac{\hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}}}{\sigma^2} \sim \chi^2(n-k)$$

Kehren wir zurück. Wir wissen nun also, dass

$$(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})' [\sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r}) \sim \chi^2(q)$$

Wie im bivariaten Fall haben wir aber wieder das Problem, dass σ^2 unbekannt ist.

Man kann zeigen, dass

$$\frac{\hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}}}{\sigma^2} \sim \chi^2(n-k)$$

unabhängig von $\hat{\boldsymbol{\beta}}$ verteilt ist. Wir wissen außerdem bereits, dass das Verhältnis zweier unabhängig χ^2 -verteilter Zufallsvariablen F -verteilt ist.¹⁷

Also ist unter der Nullhypothese $\mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ die folgende Statistik F -verteilt

$$\frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})' [\sigma^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})/q}{\hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}} / [\sigma^2(n-k)]} \sim \frac{\chi^2(q)}{\chi^2(n-k)} \sim F(q, (n-k))$$

¹⁶Der Beweis beruht u.a. darauf, dass jede positiv definite Matrix $\boldsymbol{\Sigma}$ zerlegt werden kann in $\boldsymbol{\Sigma} = \mathbf{P}\mathbf{P}'$, wobei \mathbf{P} eine nichtsinguläre $(k \times k)$ Matrix ist.

¹⁷Wenn X und Y zwei stochastisch unabhängige χ^2 -verteilte Zufallsvariablen sind mit n_1 bzw. n_2 Freiheitsgraden, dann ist $(X/n_1)/(Y/n_2)$ F -verteilt mit n_1 Zähler- und n_2 Nennerfreiheitsgraden.

Analog zum bivariaten Fall, wo wir eine standardnormalverteilte Variable durch eine χ^2 -verteilte Variable dividierten, wodurch sich die unbekannte Varianz der Grundgesamtheit σ^2 herauskürzte und wir eine t-Statistik erhielten, erhalten wir hier eine F -verteilte Teststatistik, wobei sich die unbekannte Varianz der Grundgesamtheit σ^2 wieder herauskürzt.

Als Resultat erhalten wir deshalb

$$\frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})/q}{\hat{\boldsymbol{\varepsilon}}'\hat{\boldsymbol{\varepsilon}}/(n-k)} \sim F(q, n-k) \quad (6.3)$$

Mit der üblichen Notation

$$\hat{\sigma}^2 \equiv \frac{\hat{\boldsymbol{\varepsilon}}'\hat{\boldsymbol{\varepsilon}}}{(n-k)}$$

kann dies auch folgendermaßen geschrieben werden

$$(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})' [\hat{\sigma}^2 \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})/q \sim F(q, n-k)$$

bzw. weil $\hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1} = \widehat{\text{var}}(\hat{\boldsymbol{\beta}})$

$$(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})' [\mathbf{R}\widehat{\text{var}}(\hat{\boldsymbol{\beta}})\mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})/q \sim F(q, n-k)$$

Wenn also die auf diese Weise berechnete F -Statistik größer ist als der vorher festgelegte kritische Wert kann die Nullhypothese $\mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ verworfen werden!

6.B.3 Test auf gemeinsame Signifikanz aller Steigungskoeffizienten

Als Spezialfall dieses allgemeinen F -Tests für lineare Restriktionen können wir zum Beispiel testen, ob alle Steigungskoeffizienten (d.h. alle Koeffizienten außer dem Interzept) simultan gleich Null sind.

Wir haben diesen Test bereits früher mit Hilfe einer ANOVA-Tafel dargestellt. Aber ebenso gut können diese Restriktionen mit Hilfe des hier vorgestellten F -Tests getestet werden.

Einfachheitshalber beginnen wir wieder mit dem einfachen Modell

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i$$

und testen, ob β_2 und β_3 gemeinsam signifikant von Null verschieden sind (die Verallgemeinerung auf den k Variablen Fall ist einfach).

$$H_0: \boldsymbol{\beta}^s = \begin{pmatrix} \beta_2 \\ \beta_3 \end{pmatrix} = \mathbf{0} \quad \text{gegen} \quad H_1: \boldsymbol{\beta}^s = \begin{pmatrix} \beta_2 \\ \beta_3 \end{pmatrix} \neq \mathbf{0}$$

Die geeignete Matrix \mathbf{R} und der Vektor \mathbf{r} sind

$$\mathbf{R} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \text{und} \quad \mathbf{r} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Also

$$\mathbf{R}\boldsymbol{\beta} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \begin{bmatrix} \beta_2 \\ \beta_3 \end{bmatrix} = \boldsymbol{\beta}^s$$

Ähnlich

$$\mathbf{R}\widehat{\text{var}}(\hat{\boldsymbol{\beta}})\mathbf{R}' = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \widehat{\text{var}}(\hat{\beta}_1) & \widehat{\text{cov}}(\hat{\beta}_1, \hat{\beta}_2) & \widehat{\text{cov}}(\hat{\beta}_1, \hat{\beta}_3) \\ \widehat{\text{cov}}(\hat{\beta}_1, \hat{\beta}_2) & \widehat{\text{var}}(\hat{\beta}_2) & \widehat{\text{cov}}(\hat{\beta}_2, \hat{\beta}_3) \\ \widehat{\text{cov}}(\hat{\beta}_1, \hat{\beta}_3) & \widehat{\text{cov}}(\hat{\beta}_2, \hat{\beta}_3) & \widehat{\text{var}}(\hat{\beta}_3) \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mathbf{R}\widehat{\text{var}}(\hat{\boldsymbol{\beta}})\mathbf{R}' = \begin{bmatrix} \widehat{\text{var}}(\hat{\beta}_2) & \widehat{\text{cov}}(\hat{\beta}_2, \hat{\beta}_3) \\ \widehat{\text{cov}}(\hat{\beta}_2, \hat{\beta}_3) & \widehat{\text{var}}(\hat{\beta}_3) \end{bmatrix} = \widehat{\text{var}}(\mathbf{b}^s)$$

Da in diesem Fall $q = k - 1$ folgt

$$F = \frac{(\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})' [\mathbf{R}\widehat{\text{var}}(\hat{\boldsymbol{\beta}})\mathbf{R}']^{-1} (\mathbf{R}\hat{\boldsymbol{\beta}} - \mathbf{r})}{q} = \frac{(\hat{\boldsymbol{\beta}}^s)' [\widehat{\text{var}}(\hat{\boldsymbol{\beta}}^s)]^{-1} \hat{\boldsymbol{\beta}}^s}{k - 1}$$

wobei $\hat{\boldsymbol{\beta}}^s$ den Vektor der Steigungskoeffizienten (d.h. aller Koeffizienten mit Ausnahme des Interzepts) bezeichnet.

Dies ist der übliche F-Test, inwieweit die erklärenden Variablen gemeinsam einen Erklärungsbeitrag leisten (d.h. ob mindestens eine erklärende Variable außer dem Interzept signifikant von Null verschieden ist).

6.B.4 Ein Schätzer unter linearen Restriktionen

Bisher haben wir Hypothesen der Form $\mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ getestet. Wenn wir diese Nullhypothese nicht verwerfen können, bzw. Informationen außerhalb der Stichprobe haben, dass diese Restriktion erfüllt ist, dann können wir diese Information für eine restriktierte Schätzung nutzen. Dazu minimieren wir die Quadratsumme der Residuen unter dieser Restriktion mit Hilfe des Lagrange Verfahrens.

$$\min(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

unter der Nebenbedingung

$$\mathbf{R}\boldsymbol{\beta} = \mathbf{r}$$

Die Lagrangefunktion für die Stichprobenbeobachtungen ist

$$\mathcal{L} = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_r)'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_r) - \boldsymbol{\lambda}'(\mathbf{r} - \mathbf{R}\hat{\boldsymbol{\beta}}_r)$$

wobei $\boldsymbol{\lambda}$ ein q -Vektor mit den Lagrangemultiplikatoren ist.

Als Lösung für dieses Minimierungsproblem erhält man (siehe z.B. Johnston & Dinardo 1997, S. 96f und 103f)

$$\hat{\beta}_r = \hat{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} (\mathbf{r} - \mathbf{R}\hat{\beta})$$

wobei $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ der übliche OLS-Schätzer ohne Restriktionen ist.

Die Varianz-Kovarianzmatrix der restringierten Schätzer $\hat{\beta}_r$ erhält man aus

$$\widehat{\text{var}}(\hat{\beta}_r) = \hat{\sigma}_r^2 (\mathbf{X}'\mathbf{X})^{-1} \left\{ \mathbf{I}_k - \mathbf{R}' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \right\}$$

mit

$$\hat{\sigma}_r^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\beta}_r)'(\mathbf{y} - \mathbf{X}\hat{\beta}_r)}{n - k + q} = \frac{\hat{\epsilon}'_r \hat{\epsilon}_r}{n - k + q}$$

Der restringierte Schätzer ist unverzerrt und effizient, wenn die Restriktion in der Grundgesamtheit gilt. Wenn die Restriktion in der Grundgesamtheit nicht erfüllt ist, liefert der restringierte Schätzer verzerrte Ergebnisse!

6.B.5 Restringierte Schätzung und F -Test

Wir haben früher den F -Test

$$F = \frac{(\hat{\epsilon}'_r \hat{\epsilon}_r - \hat{\epsilon}' \hat{\epsilon})/q}{\hat{\epsilon}' \hat{\epsilon} / (n - k)}$$

verwendet um lineare Restriktionen zu testen.

Wir werden nun zeigen, dass dieser zum gleichen Ergebnis führt wie der F -Test, den wir vorhin hergeleitet haben, nämlich

$$F = \frac{(\mathbf{R}\hat{\beta} - \mathbf{r})' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})/q}{\hat{\epsilon}' \hat{\epsilon} / (n - k)}$$

Wir beginnen mit den Residuen der restringierten Regression

$$\begin{aligned} \hat{\epsilon}_r &= \mathbf{y} - \mathbf{X}\hat{\beta}_r \\ &= \mathbf{y} - \mathbf{X}\hat{\beta} - \mathbf{X}(\hat{\beta}_r - \hat{\beta}) \\ &= \hat{\epsilon} - \mathbf{X}(\hat{\beta}_r - \hat{\beta}) \end{aligned}$$

Transponieren und ausmultiplizieren gibt

$$\hat{\epsilon}'_r \hat{\epsilon}_r = \hat{\epsilon}' \hat{\epsilon} + (\hat{\beta}_r - \hat{\beta})' \mathbf{X}' \mathbf{X} (\hat{\beta}_r - \hat{\beta})$$

Einsetzen der Formel für den restringierten Schätzer

$$\hat{\beta}_r - \hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} (\mathbf{r} - \mathbf{R}\hat{\beta})$$

gibt nach Vereinfachungen

$$\hat{\epsilon}'_r \hat{\epsilon}_r - \hat{\epsilon}' \hat{\epsilon} = (\mathbf{R}\hat{\beta} - \mathbf{r})' [\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})$$

Dies ist gleich dem Zähler der obigen F -Statistik, weshalb

$$F = \frac{(\hat{\epsilon}'_r \hat{\epsilon}_r - \hat{\epsilon}' \hat{\epsilon})/q}{\hat{\epsilon}' \hat{\epsilon} / (n - k)} \sim F_{(q, n-k)}$$

nur ein anderer Ausdruck für die Teststatistik der Nullhypothese $\mathbf{R}\beta = \mathbf{r}$ ist (vgl. Johnston and Dinardo, 1996, 97).

6.B.6 Irrtümliche Berücksichtigung nicht relevanter Variablen (*‘redundant variables’*) in Matrixschreibweise

Das ‘wahre’ Modell sei

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\varepsilon}$$

aber wir schätzen das ‘falsche’ Modell

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$$

wobei \mathbf{X}_1 wieder eine $n \times k_1$ und \mathbf{X}_2 eine $n \times k_2$ Matrix ist ($k_1 + k_2 = k$).

Das ‘lange’ Modell kann mit Hilfe partitionierter Matrizen geschrieben werden als

$$\mathbf{y} = (\mathbf{X}_1 \ \mathbf{X}_2) \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix} + \boldsymbol{\varepsilon}$$

und der OLS-Schätzer in partitionierter Schreibweise ist

$$\begin{pmatrix} \hat{\boldsymbol{\beta}}_1 \\ \hat{\boldsymbol{\beta}}_2 \end{pmatrix} = ((\mathbf{X}_1 \ \mathbf{X}_2)'(\mathbf{X}_1 \ \mathbf{X}_2))^{-1} (\mathbf{X}_1 \ \mathbf{X}_2)' \mathbf{y}$$

Wir werden gleich zeigen, dass die Subvektoren $\hat{\boldsymbol{\beta}}_1$ und $\hat{\boldsymbol{\beta}}_2$ folgendermaßen dargestellt werden können:

$$\begin{aligned} \hat{\boldsymbol{\beta}}_1 &= (\mathbf{X}_1' \mathbf{M}_{\mathbf{X}_2} \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{M}_{\mathbf{X}_2} \mathbf{y} \\ \hat{\boldsymbol{\beta}}_2 &= (\mathbf{X}_2' \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2)^{-1} \mathbf{X}_2' \mathbf{M}_{\mathbf{X}_1} \mathbf{y} \end{aligned}$$

mit

$$\mathbf{M}_{\mathbf{X}_i} = \mathbf{I}_n - \mathbf{X}_i(\mathbf{X}_i' \mathbf{X}_i)^{-1} \mathbf{X}_i' \quad (i = 1, 2)$$

Beweis: Das mit OLS geschätzte (falsche) Modell ist

$$\mathbf{y} = \mathbf{X}_1\hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2\hat{\boldsymbol{\beta}}_2 + \hat{\boldsymbol{\varepsilon}}$$

Wir multiplizieren diese Gleichung mit $\mathbf{M}_{\mathbf{X}_1}$ von links und erhalten

$$\mathbf{M}_{\mathbf{X}_1} \mathbf{y} = \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1 + \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \mathbf{M}_{\mathbf{X}_1} \hat{\boldsymbol{\varepsilon}}$$

und berücksichtigen, dass $\mathbf{M}_{\mathbf{X}_1} \mathbf{X}_1 = \mathbf{0}$ und $\mathbf{M}_{\mathbf{X}_1} \hat{\boldsymbol{\varepsilon}} = \hat{\boldsymbol{\varepsilon}}$ (der zweite Ausdruck folgt aus $\mathbf{X}_1' \hat{\boldsymbol{\varepsilon}} = \mathbf{0}$). Also

$$\begin{aligned} \mathbf{M}_{\mathbf{X}_1} \mathbf{y} &= \underbrace{\mathbf{M}_{\mathbf{X}_1} \mathbf{X}_1}_{=\mathbf{0}} \hat{\boldsymbol{\beta}}_1 + \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \underbrace{\mathbf{M}_{\mathbf{X}_1} \hat{\boldsymbol{\varepsilon}}}_{=\hat{\boldsymbol{\varepsilon}}} \\ &= \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \hat{\boldsymbol{\varepsilon}} \end{aligned}$$

Als nächstes multiplizieren wir diesen Ausdruck von links mit \mathbf{X}_2'

$$\begin{aligned} \mathbf{X}_2' \mathbf{M}_{\mathbf{X}_1} \mathbf{y} &= \mathbf{X}_2' \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \mathbf{X}_2' \hat{\boldsymbol{\varepsilon}} \\ \mathbf{X}_2' \mathbf{M}_{\mathbf{X}_1} \mathbf{y} &= \mathbf{X}_2' \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 \end{aligned}$$

weil $\mathbf{X}_2' \hat{\boldsymbol{\varepsilon}} = \mathbf{0}$ (Orthogonalität). Wenn die Matrix $\mathbf{X}_2' \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2$ nicht singulär ist folgt

$$\hat{\boldsymbol{\beta}}_2 = (\mathbf{X}_2' \mathbf{M}_{\mathbf{X}_1} \mathbf{X}_2)^{-1} \mathbf{X}_2' \mathbf{M}_{\mathbf{X}_1} \mathbf{y}$$

Ganz ähnlich kann gezeigt werden, dass

$$\hat{\boldsymbol{\beta}}_1 = (\mathbf{X}_1' \mathbf{M}_{\mathbf{X}_2} \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{M}_{\mathbf{X}_2} \mathbf{y}$$



Wir nutzen nun diese Ergebnisse um zu zeigen, dass

$$E \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \mathbf{0} \end{pmatrix}$$

Wir starten mit dem Subvektor $\hat{\beta}_1$ und setzen wieder den ‘wahren’ Zusammenhang ein:

$$\begin{aligned} \hat{\beta}_1 &= (\mathbf{X}'_1 \mathbf{M}_{X_2} \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{M}_{X_2} \mathbf{y} \\ &= (\mathbf{X}'_1 \mathbf{M}_{X_2} \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{M}_{X_2} (\mathbf{X}_1 \beta_1 + \boldsymbol{\varepsilon}) \\ &= \beta_1 + (\mathbf{X}'_1 \mathbf{M}_{X_2} \mathbf{X}_1)^{-1} \mathbf{X}'_1 \mathbf{M}_{X_2} \boldsymbol{\varepsilon} \end{aligned}$$

Wenn die \mathbf{X} wieder deterministisch sind und $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ folgt

$$E(\hat{\beta}_1) = \beta_1$$

Als nächstes untersuchen wir den Subvektor $\hat{\beta}_2$:

$$\begin{aligned} \hat{\beta}_2 &= (\mathbf{X}'_2 \mathbf{M}_{X_1} \mathbf{X}_2)^{-1} \mathbf{X}'_2 \mathbf{M}_{X_1} \mathbf{y} \\ &= (\mathbf{X}'_2 \mathbf{M}_{X_1} \mathbf{X}_2)^{-1} \mathbf{X}'_2 \mathbf{M}_{X_1} (\mathbf{X}_1 \beta_1 + \boldsymbol{\varepsilon}) \\ &= (\mathbf{X}'_2 \mathbf{M}_{X_1} \mathbf{X}_2)^{-1} \mathbf{X}'_2 \underbrace{\mathbf{M}_{X_1} \mathbf{X}_1}_{=0} \beta_1 + (\mathbf{X}'_2 \mathbf{M}_{X_1} \mathbf{X}_2)^{-1} \mathbf{X}'_2 \mathbf{M}_{X_1} \boldsymbol{\varepsilon} \\ &= \mathbf{0} + (\mathbf{X}'_2 \mathbf{M}_{X_1} \mathbf{X}_2)^{-1} \mathbf{X}'_2 \mathbf{M}_{X_1} \boldsymbol{\varepsilon} \end{aligned}$$

Also

$$E(\hat{\beta}_2) = \mathbf{0}$$

Im Falle der Berücksichtigung nicht-relevanter Variablen bleibt der OLS-Schätzer also erwartungstreu und konsistent, aber man kann zeigen, dass die Varianz der Koeffizienten

$$\text{var}(\hat{\beta}_1) = \sigma_{\varepsilon}^2 (\mathbf{X}'_1 \mathbf{M}_{X_2} \mathbf{X}_1)^{-1}$$

größer ist als die Varianz-Kovarianzmatrix der Koeffizienten des ‘wahren’ Modells $\text{var}(\hat{\beta}_1) = \sigma_{\varepsilon}^2 (\mathbf{X}'_1 \mathbf{X}_1)^{-1}$, der Schätzer ist also *nicht effizient*!

6.C Das Frisch-Waugh-Lovell Theorem in Matrixschreibweise(‘Residual Regression’)

Bereits im Kapitel zur deskriptiven Regressionsanalyse haben wir das Frisch-Waugh-Lovell Theorem erläutert, welches in der allerersten Ausgabe der *Econometrica* von Ragnar Frisch und F.V. Waugh¹⁸ publiziert wurde.

¹⁸Frisch, Ragnar and F.V. Waugh (1933), “Partial time regression as compared with individual trends”, *Econometrica*, 1, pp. 387-401.

Lovell¹⁹ hat dieses Ergebnis 1963 verallgemeinert und gezeigt, dass dies für beliebige Partitionierungen gilt.

Für diese Verallgemeinerung partitionieren wir die Datenmatrix \mathbf{X} und den Koeffizientenvektor $\boldsymbol{\beta}$

$$\mathbf{X} = (\mathbf{X}_1 \quad \mathbf{X}_2), \quad \boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix}$$

wobei \mathbf{X}_1 die Dimension $n \times k_1$ und \mathbf{X}_2 die Dimension $n \times k_2$ hat, mit $k_1 + k_2 = k$, und $\boldsymbol{\beta}_1$ (bzw. $\boldsymbol{\beta}_2$) entsprechend die Dimension $k_1 \times 1$ (bzw. $k_2 \times 1$) hat. Das Regressionsmodell kann nun geschrieben werden als

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$$

Dieses Modell kann mit OLS geschätzt werden und gibt

$$\mathbf{y} = \mathbf{X}_1\hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2\hat{\boldsymbol{\beta}}_2 + \hat{\boldsymbol{\varepsilon}}$$

Wir zeigen nun, dass wir $\hat{\boldsymbol{\beta}}_2$ schätzen können, indem wir die Residuen einer Schätzung von \mathbf{y} auf \mathbf{X}_1 auf die Residuen einer Schätzung von \mathbf{X}_2 auf \mathbf{X}_1 regressieren.

Dies folgt im wesentlichen aus zwei Eigenschaften der OLS Methode

1. Die Residuen einer OLS Schätzung ($\hat{\boldsymbol{\varepsilon}}$) sind mit den erklärenden Variablen (\mathbf{X}) unkorreliert (stehen also orthogonal aufeinander).
2. Die Koeffizienten einer Teilmenge der erklärenden x Variablen sind Null, wenn diese Variablen sowohl mit den restlichen x Variablen als auch mit der abhängigen y Variable unkorreliert sind.

Der folgende Beweis folgt eng dem Lehrbuch von Bruce E. Hansen (University of Wisconsin, <http://www.ssc.wisc.edu/~bhansen/econometrics/>).

Wir beginnen damit die residuenerzeugende Matrix für die erste Submatrix \mathbf{X}_1 zu definieren

$$\mathbf{M}_1 = \mathbf{I}_n - \mathbf{X}_1(\mathbf{X}_1'\mathbf{X}_1)^{-1}\mathbf{X}_1'$$

Aus $\mathbf{M} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ und $\mathbf{X}_1'\mathbf{M}_1 = \mathbf{O}$ folgt

$$\mathbf{M}_1\mathbf{M} = \mathbf{M} - \mathbf{X}_1(\mathbf{X}_1'\mathbf{X}_1)^{-1}\mathbf{X}_1'\mathbf{M} = \mathbf{M}$$

Wir haben bereits früher gezeigt, dass $\hat{\boldsymbol{\varepsilon}} = \mathbf{M}\boldsymbol{\varepsilon}$.²⁰ Prämultiplizieren mit \mathbf{M}_1 gibt

$$\mathbf{M}_1\hat{\boldsymbol{\varepsilon}} = \mathbf{M}_1\mathbf{M}\boldsymbol{\varepsilon} = \mathbf{M}\boldsymbol{\varepsilon} = \hat{\boldsymbol{\varepsilon}}$$

Als nächstes prämultiplizieren wir $\mathbf{y} = \mathbf{X}_1\hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2\hat{\boldsymbol{\beta}}_2 + \hat{\boldsymbol{\varepsilon}}$ mit \mathbf{M}_1 und erhalten

$$\begin{aligned} \mathbf{M}_1\mathbf{y} &= \mathbf{M}_1\mathbf{X}_1\hat{\boldsymbol{\beta}}_1 + \mathbf{M}_1\mathbf{X}_2\hat{\boldsymbol{\beta}}_2 + \mathbf{M}_1\hat{\boldsymbol{\varepsilon}} \\ &= \mathbf{M}_1\mathbf{X}_2\hat{\boldsymbol{\beta}}_2 + \hat{\boldsymbol{\varepsilon}} \end{aligned}$$

¹⁹Lovell, Michael C. (1963) "Seasonal Adjustment of Economic Time Series and Multiple Regression Analysis", Journal of the American Statistical Association, December 1963, pp. 993-1010.

²⁰ $\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = (\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{y} = \mathbf{M}\mathbf{y} = \mathbf{M}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \mathbf{M}\boldsymbol{\varepsilon}$ da $\mathbf{M}\mathbf{X} = \mathbf{O}$.

da $\mathbf{M}_1 \mathbf{X}_1 = 0$ und $\mathbf{M}_1 \hat{\boldsymbol{\varepsilon}} = \hat{\boldsymbol{\varepsilon}}$.

Nun prämultiplizieren wir obigen Ausdruck mit \mathbf{X}_2' und erinnern uns, dass aufgrund der Bedingungen erster Ordnung $\mathbf{X}_2' \hat{\boldsymbol{\varepsilon}} = 0$ (d.h. die Spalten der \mathbf{X} Matrix stehen paarweise orthogonal auf den OLS-Residuenvektor $\hat{\boldsymbol{\varepsilon}}$)

$$\mathbf{X}_2' \mathbf{M}_1 \mathbf{y} = \mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \mathbf{X}_2' \hat{\boldsymbol{\varepsilon}} = \mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2$$

und lösen nach $\hat{\boldsymbol{\varepsilon}}$

$$\hat{\boldsymbol{\varepsilon}}_2 = (\mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2' \mathbf{M}_1 \mathbf{y}$$

Wir erinnern uns, dass \mathbf{M}_1 eine residuenerzeugende Matrix ist und definieren die Residuen einer OLS Regression auf \mathbf{X}_1

$$\begin{aligned}\check{\mathbf{X}}_2 &= \mathbf{M}_1 \mathbf{X}_2 \\ \check{\mathbf{y}} &= \mathbf{M}_1 \mathbf{y}\end{aligned}$$

Da \mathbf{M}_1 außerdem symmetrisch und idempotent ist, also $\mathbf{M}_1 = \mathbf{M}_1 \mathbf{M}_1$, folgt

$$\begin{aligned}\hat{\boldsymbol{\beta}}_2 &= (\mathbf{X}_2' \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2' \mathbf{M}_1 \mathbf{y} \\ &= (\mathbf{X}_2' \mathbf{M}_1 \mathbf{M}_1 \mathbf{X}_2)^{-1} \mathbf{X}_2' \mathbf{M}_1 \mathbf{M}_1 \mathbf{y} \\ &= (\check{\mathbf{X}}_2' \check{\mathbf{X}}_2)^{-1} (\check{\mathbf{X}}_2' \check{\mathbf{y}})\end{aligned}$$

Zudem folgt aus $\mathbf{M}_1 \mathbf{y} = \mathbf{M}_1 \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \hat{\boldsymbol{\varepsilon}}$ und den Definitionen $\check{\mathbf{X}}_2 = \mathbf{M}_1 \mathbf{X}_2$ sowie $\check{\mathbf{y}} = \mathbf{M}_1 \mathbf{y}$ dass die ‘kurze’ Regression

$$\check{\mathbf{y}} = \check{\mathbf{X}}_2 \hat{\boldsymbol{\beta}}_2 + \hat{\boldsymbol{\varepsilon}}$$

den gleichen Residuenvektor $\hat{\boldsymbol{\varepsilon}}$ liefert wie die ‘lange’ Regression.²¹

Deshalb erhält man aus einer OLS Schätzung von

$$\mathbf{y} = \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \hat{\boldsymbol{\varepsilon}}$$

den gleichen Koeffizientenvektor $\hat{\boldsymbol{\beta}}_2$ und den gleichen Residuenvektor $\hat{\boldsymbol{\varepsilon}}$, den man durch folgende Vorgangsweise erhält

1. Regressiere \mathbf{y} auf \mathbf{X}_1 und berechne daraus die Residuen $\check{\mathbf{y}}$
2. Regressiere die \mathbf{X}_2 auf \mathbf{X}_1 und berechne daraus die Residuen $\check{\mathbf{X}}_2$
3. Regressiere $\check{\mathbf{y}}$ auf die $\check{\mathbf{X}}_2$ und berechne daraus den OLS Koeffizientenvektor $\hat{\boldsymbol{\beta}}_2$ und die Residuen $\hat{\boldsymbol{\varepsilon}}$.

²¹Der Korrelationskoeffizient zwischen diesen um den Einfluss von \mathbf{X}_1 ‘bereinigten’ Residuen wird übrigens *partieller Korrelationskoeffizient* genannt.

Beispiel

Wir greifen auf die Daten des früheren Beispiels mit der orthogonalen Projektion zurück

| y | x_1 | x_2 |
|-----|-------|-------|
| 6 | 2 | 8 |
| 4 | 8 | 1 |
| 5 | 1 | 1 |

Wenn wir den Einfluss der Variablen x_1 aus y und x_2 eliminieren wollen ist die Matrix \mathbf{X}_1 einfach der \mathbf{x}_1 Vektor

$$\mathbf{x}'_1 = (2 \ 8 \ 1), \quad \mathbf{x}'_1 \mathbf{x}_1 = 69, \quad (\mathbf{x}'_1 \mathbf{x}_1)^{-1} = 1/69 = 0.01449$$

Die \mathbf{M}_1 Matrix ist

$$\begin{aligned} \mathbf{M}_1 &= \mathbf{I}_3 - \mathbf{x}_1(\mathbf{x}'_1 \mathbf{x}_1)^{-1} \mathbf{x}'_1 \\ &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} - \begin{pmatrix} 2 \\ 8 \\ 1 \end{pmatrix} (1/69) (2 \ 8 \ 1) \\ &= \begin{pmatrix} 0.942 & -0.232 & -0.029 \\ -0.232 & 0.072 & -0.116 \\ -0.029 & -0.116 & 0.986 \end{pmatrix} \end{aligned}$$

Die OLS-Residuen der Regressionen von \mathbf{x}_1 auf \mathbf{x}_2 (d.h. $\hat{\varepsilon}^{x_2} = x_2 - b^{x_2} x_1$) erhält man aus $\mathbf{M}_1 \mathbf{x}_2$, und die Residuen der Regressionen von \mathbf{x}_1 auf \mathbf{y} (d.h. $\hat{\varepsilon}^y = y - b^y x_1$) erhält man aus $\mathbf{M}_1 \mathbf{y}$

$$\hat{\varepsilon}^{x_2} = \mathbf{M}_1 \mathbf{x}_2 = \begin{pmatrix} 7.275 \\ -1.899 \\ 0.638 \end{pmatrix} \quad \text{und} \quad \hat{\varepsilon}^y = \mathbf{M}_1 \mathbf{y} = \begin{pmatrix} 4.580 \\ -1.681 \\ 4.290 \end{pmatrix}$$

Den Koeffizienten $\hat{\beta}_2$ der OLS-Regression $y = \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \hat{\varepsilon}$ erhält man alternativ, indem man die Residuen der ersten Stufe aufeinander regressiert, also

$$\hat{\beta}_2 = [(\mathbf{M}_1 \mathbf{x}_2)'(\mathbf{M}_1 \mathbf{x}_2)]^{-1} (\mathbf{M}_1 \mathbf{x}_2)'(\mathbf{M}_1 \mathbf{y}) = 0.689$$

Dies ist der gleiche Koeffizient $\hat{\beta}_2$, den man auch aus der multiplen Regression erhält

$$y = \underset{(0.505)}{0.46} x_1 + \underset{(0.516)}{0.689} x_2$$

$$R^2 = -6.576, \quad s = 3.893, \quad n = 3$$