

Einführung in die Empirische Wirtschaftsforschung

Übungsaufgaben 7

KQ Methode in der linearen Einfachregression - Lineare Mehrfachregression - Multikollinearität - Ausgelassene Prediktoren

1. Schreiben Sie die Zielfunktion der KQ-Methode für die lineare Einfachregression auf. Charakterisieren Sie die Bedingungen erster Ordnung für ein Minimum.

Optional: Leiten Sie die Schätzer für den Achsenabschnitt und die Steigung her.

2. Wenn eine in einem Duopol operierende Firma den Preis um einen Franken senkt, und der Preis der Konkurrenz gleich bleibt, dann steigt der Umsatz um 200. Wenn die Konkurrenz den Preis ebenfalls um einen Franken senkt, dann bleibt der Umsatz unverändert. In 2 von 3 Fällen senkt die Firma den Preis um einen Franken (in den Anderen Fällen lässt sie den Preis unverändert). Im Falle einer Preissenkung, kopiert die Konkurrenz diese Preissenkung in 50% der Fälle. Was sind dann b_1 , b_2 und g_1 in den zwei Gleichungen

$$\begin{aligned}\widehat{Umsatz} &= b_0 + b_1 Preis_{eigen} + b_2 Preis_{Konkurrenz} \\ \widetilde{Umsatz} &= g_0 + g_1 Preis_{eigen} \quad ?\end{aligned}$$

Hinweis: R kann bei der Lösung der Aufgabe helfen.

3. Betrachten Sie das Datenset “miete”. Wir betrachten die Variable “rent” als Zielvariable und die Variablen “size” und “rooms” als Prediktoren.

- Regressieren Sie in einer linearen Einfachregression “rent” auf “rooms”. Interpretieren Sie die Koeffizienten.
- Erweitern Sie die lineare Einfachregression um den Prediktor “size”. Was fällt Ihnen im Vergleich zu Aufgabe a) im Bezug auf b_1 auf? Wie kann es sein, dass b_1 nun nicht mehr signifikant ist?
- In der Regression aus b), um wie viel ändert sich die erwartete Miete ceteris paribus, wenn bei einer Wohnung ein zusätzliches Schlafzimmer hinzukommt?
Macht es Sinn, in diesem Beispiel eine ceteris paribus Aussage zu treffen?
- Gibt es weitere Variablen, welche man dem Regressionsmodel hinzufügen sollte?

4. Wir betrachten eine lineare Mehrfachregression mit p “interessanten” Prediktoren (abgesehen von der Konstanten, die auch im Modell enthalten ist). Die zugehörige R^2 -Statistik sei R_p^2 . Jetzt nehmen wir zusätzlich noch einen weiteren Regressor hinzu, d.h., es gibt jetzt $p + 1$ “interessante” Variablen. Die zugehörige R^2 -Statistik sei R_{p+1}^2 . Zeigen Sie mathematisch, dass notwendigerweise $R_{p+1}^2 \geq R_p^2$.

5. Nehmen Sie an, (Y, X_1, X_2) genüge den Annahmen der linearen Mehrfachregression

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon,$$

wobei $E(\epsilon) = 0$ (beziehungsweise $E(\epsilon X_1) = E(\epsilon X_2) = 0$). β_1, β_2 sind die interessierenden Effekte von X_1, X_2 auf Y . Alternativ wird auch das folgende lineare Einfachregressionsmodell geschätzt

$$Y = \alpha_0 + \alpha_1 X_1 + u.$$

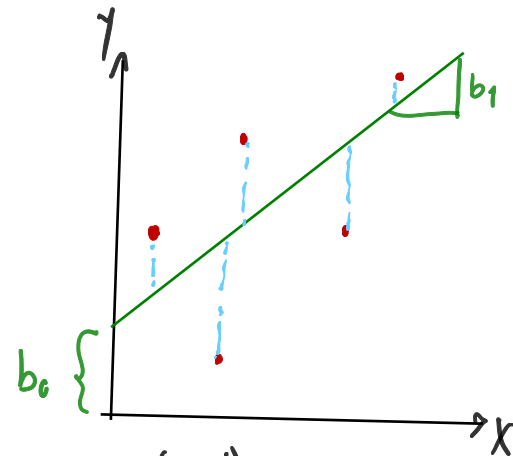
Zeigen Sie, dass der Schätzer a_1 im kleinen Modell i.d.R. nicht unverzerrt (unbiased) ist für den Koeffizienten β_1 im grossen Modell: Leiten Sie eine Formel für den zugehörigen Omitted Variable Bias $E(a_1) - \beta_1$ her.

1. Schreiben Sie die Zielfunktion der KQ-Methode für die lineare Einfachregression auf. Charakterisieren Sie die Bedingungen erster Ordnung für ein Minimum.

Optional: Leiten Sie die Schätzer für den Achsenabschnitt und die Steigung her.

$$(b_0, b_1) = \underset{b_0', b_1'}{\operatorname{argmin}} SQR(b_0', b_1') = \underset{b_0', b_1'}{\operatorname{argmin}} \sum_{i=1}^n e_i^2$$

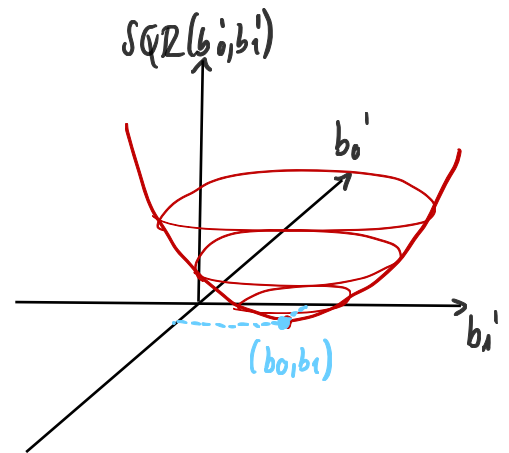
$$= \underset{b_0', b_1'}{\operatorname{argmin}} \sum_{i=1}^n (y_i - b_0' - b_1' x_i)^2$$



• Wir minimieren diese Funktion in Bezug auf b_0, b_1 :

$$\frac{\partial SQR(b_0', b_1')}{\partial b_0'} = 2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) \stackrel{!}{=} 0$$

$$\frac{\partial SQR(b_0', b_1')}{\partial b_1'} = 2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) (-x_i) \stackrel{!}{=} 0$$



- Gleichungssystem mit zwei Gleichungen und zwei Unbekannten (siehe Lösungen für Auflösung)
- Verbindung Mathematik: Wir können das obige schreiben als

$$\nabla SQR(b_0', b_1') = (0, 0)^T$$

2. Wenn eine in einem Duopol operierende Firma den Preis um einen Franken senkt, und der Preis der Konkurrenz gleich bleibt, dann steigt der Umsatz um 200. Wenn die Konkurrenz den Preis ebenfalls um einen Franken senkt, dann bleibt der Umsatz unverändert. In 2 von 3 Fällen senkt die Firma den Preis um einen Franken (in den Anderen Fällen lässt sie den Preis unverändert). Im Falle einer Preissenkung, kopiert die Konkurrenz diese Preissenkung in 50% der Fälle. Was sind dann b_1 , b_2 und g_1 in den zwei Gleichungen

$$\widehat{Umsatz_1} = b_0 + b_1 \text{Preis}_{\text{eigen}} + b_2 \text{Preis}_{\text{Konkurrenz}}$$

$$\widehat{Umsatz_2} = g_0 + g_1 \text{Preis}_{\text{eigen}} \quad ?$$

Hinweis: R kann bei der Lösung der Aufgabe helfen.

$$\left. \begin{array}{l} b_1 = -200 \\ b_2 = 200 \end{array} \right\} \widehat{Umsatz} = b_0 - 200 \text{Preis}_{\text{eigen}} + 200 \text{Preis}_{\text{Konkurrenz}}$$

Was ist der erwartete Umsatz i.e. \widehat{Umsatz} ?

$$E[\widehat{Umsatz}] = E[\widehat{Umsatz} | \text{Preissenkung}] P(\text{Preissenkung}) \quad E[\widehat{Umsatz} | P] \quad E[\widehat{Umsatz} | k.P.]$$

$$+ E[\widehat{Umsatz} | \text{keine Preissenkung}] P(\text{keine Preissenkung})$$

$$= E E[\widehat{Umsatz} | \text{Preissenkung}] P(\text{Preissenkung})$$

$$= 0.5 \cdot 0 + 0.5 \cdot 200 \quad = 2/3$$

, wobei 0.5 :=
Wahrscheinlichkeit,
dass Konk. preis senkt

$$+ E E[\widehat{Umsatz} | \text{keine Preissenkung}] P(\text{keine Preissenkung})$$

$$= 0$$

$$= 1/3$$

→ Der erwartete Umsatz (\widehat{Umsatz}) liegt bei $2/3 \cdot (0.5 \cdot 0 + 0.5 \cdot 200)$,
wobei der erwartete Umsatz nach einer eigenen Preissenkung bei
 $0.5 \cdot 0 + 0.5 \cdot 200 = 100$ liegt.

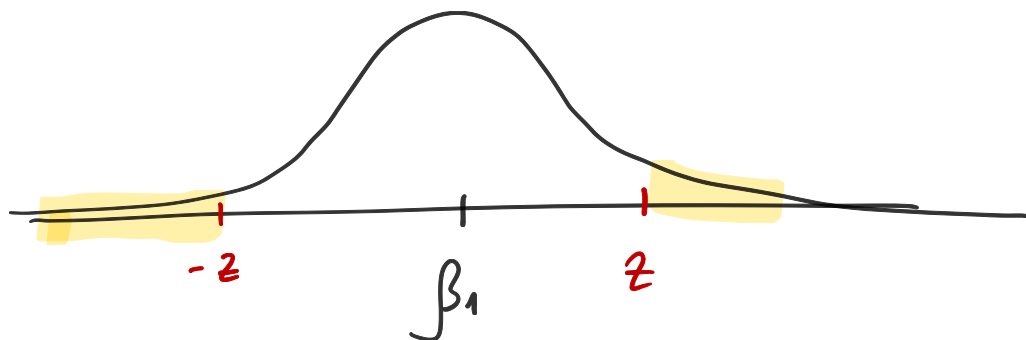
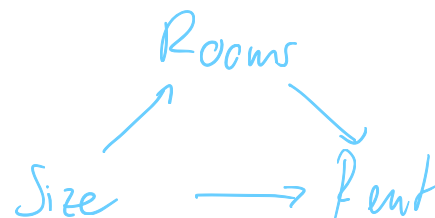
⇒ Somit haben wir:

$$\tilde{\text{Umsatz}} = g_0 - 100 \text{ Preiseigen}$$

3. Betrachten Sie das Datenset "miete". Wir betrachten die Variable "rent" als Zielvariable und die Variablen "size" und "rooms" als Prediktoren.

- Regressieren Sie in einer linearen Einfachregression "rent" auf "rooms". Interpretieren Sie die Koeffizienten.
- Erweitern Sie die lineare Einfachregression um den Prediktor "size". Was fällt Ihnen im Vergleich zu Aufgabe a) im Bezug auf b_1 auf? Wie kann es sein, dass b_1 nun nicht mehr signifikant ist?
- In der Regression aus b), um wie viel ändert sich die erwartete Miete ceteris paribus, wenn bei einer Wohnung ein zusätzliches Schlafzimmer hinzukommt?
Macht es Sinn, in diesem Beispiel eine ceteris paribus Aussage zu treffen?
- Gibt es weitere Variablen, welche man dem Regressionsmodel hinzufügen sollte?

Siehe 2



$$2P(Z > |z|)$$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$$

$$\hookrightarrow \hat{X}_1 = a_0 + a_1 X_2 + a_2 X_3 \rightarrow R^2 \quad \hat{X}_3 = a'_0 + a'_1 X_2 + a'_2 X_1 \rightarrow R'^2$$

$$\hat{X}_2 = \hat{a}_0 + \hat{a}_1 X_1 + \hat{a}_2 X_3 \rightarrow \tilde{R}^2$$

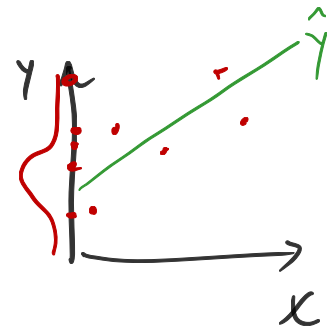
$$VIF = \frac{1}{1 - R^2}$$

4. Wir betrachten eine lineare Mehrfachregression mit p "interessanten" Prediktoren (abgesehen von der Konstanten, die auch im Modell enthalten ist). Die zugehörige R^2 -Statistik sei R_p^2 . Jetzt nehmen wir zusätzlich noch einen weiteren Regressor hinzu, d.h., es gibt jetzt $p+1$ "interessante" Variablen. Die zugehörige R^2 -Statistik sei R_{p+1}^2 . Zeigen Sie mathematisch, dass notwendigerweise $R_{p+1}^2 \geq R_p^2$.

Definition $R^2 := 1 - \frac{SSR}{TSS} = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^N e_i^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$

Wir sehen, dass $TSS = \sum_{i=1}^N (y_i - \bar{y})^2$ von p unabhängig ist. Wir müssen also zeigen, dass:

$$\sum_{i=1}^N (y_i - \hat{y}_{i,p})^2 \geq \sum_{i=1}^N (y_i - \hat{y}_{i,p+1})^2$$



Die zusätzliche Flexibilität durch die $(p+1)$ -te Variable kann die Anpassungsgüte der Regression nur verbessern. Falls die $(p+1)$ -te Variable rein gar nichts von y erklären kann, dann gilt $b_{p+1} = 0$ und somit $\hat{y}_{i,p} = \hat{y}_{i,p+1}$ und $R_p^2 = R_{p+1}^2$.

↳ Das deckt aber natürlich ein Problem von R^2 auf, da es mit zusätzlichen Regressoren nie schlechter wird. Deshalb gibt es Größen wie das adjusted R^2_{adj} :

$$R^2_{adj} := 1 - (1 - R^2) \frac{N-1}{N-p}$$

5. Nehmen Sie an, (Y, X_1, X_2) genüge den Annahmen der linearen Mehrfachregression

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon,$$

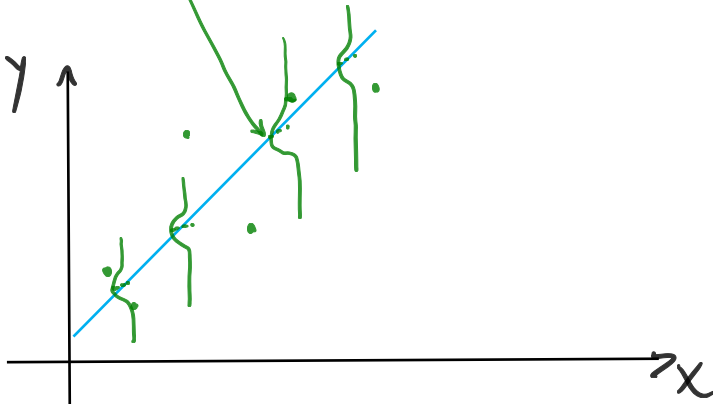
wobei $E(\epsilon) = 0$ (beziehungsweise $E(\epsilon X_1) = E(\epsilon X_2) = 0$). β_1, β_2 sind die interessierenden Effekte von X_1, X_2 auf Y . Alternativ wird auch das folgende lineare Einfachregressionsmodell geschätzt

$$Y = \alpha_0 + \alpha_1 X_1 + u.$$

Zeigen Sie, dass der Schätzer a_1 im kleinen Modell i.d.R. nicht unverzerrt (unbiased) ist für den Koeffizienten β_1 im grossen Modell: Leiten Sie eine Formel für den zugehörigen Omitted Variable Bias $E(a_1) - \beta_1$ her.

$$E(\epsilon) = 0$$

$$E(\epsilon X_1) = E(\epsilon X_2) = 0, \text{ da } \text{Cov}(\epsilon, X_i) = E(\epsilon X_i) - \underbrace{E(\epsilon)E(X_i)}_{=0}, i=1,2$$



Wie berechnen wir a_1 (kleines Modell):

$$\begin{aligned} a_1 &= \frac{\hat{\text{Cov}}(X_1, Y)}{\hat{\text{Var}}(X_1)} = \frac{\frac{1}{N-1} \sum_{i=1}^N (X_{i1} - \bar{X}_1)(Y_i - \bar{Y})}{\frac{1}{N-1} \sum_{i=1}^N (X_{i1} - \bar{X}_1)^2} \stackrel{(*)}{=} \frac{\sum_{i=1}^N (X_{i1} - \bar{X}_1) Y_i}{\underbrace{\sum_{i=1}^N (X_{i1} - \bar{X}_1)^2}_{:= S_X^2}} \\ &= \frac{\sum_{i=1}^N (X_{i1} - \bar{X}_1) (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i)}{S_X^2} \end{aligned}$$

$$= \beta_0 \sum_{i=1}^N (x_{i1} - \bar{x}_1) + \beta_1 \sum_{i=1}^N (\cancel{x_{i1} - \bar{x}_1}) x_{i1} + \beta_2 \sum_{i=1}^N (x_{i1} - \bar{x}_1) x_{i2} + \sum_{i=1}^N (x_{i1} - \bar{x}_1) \varepsilon_i$$

$$S_x^2 = \sum_{i=1}^N (\cancel{x_{i1} - \bar{x}_1}) x_{i1}$$

$$= \beta_1 + \beta_2 \frac{\sum (x_{i1} - \bar{x}_1)^2 x_{i2}}{\sum (x_{i1} - \bar{x}_1) x_{i1}} + \frac{\sum (x_{i1} - \bar{x}_1) \varepsilon_i}{\sum (x_{i1} - \bar{x}_1) x_{i1}}$$

↳ Wie schätzen wir a_1 in Erwartung, also $E(a_1)$:

$$E(a_1) = \beta_1 + \beta_2 \frac{\text{Cor}(x_1, x_2)}{\text{Var}(x_1)} + \frac{\sum (x_{i1} - \bar{x}_1) E(\varepsilon_i)}{\underbrace{\sum (x_{i1} - \bar{x}_1) x_{i1}}_{=0}} \neq \beta_1$$

(*) Alternative form um folgendes zu schreiben: $\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$:

$$= \sum_{i=1}^N x_i y_i - \bar{y} \sum_{i=1}^N x_i - \bar{x} \sum_{i=1}^N y_i + N \bar{x} \bar{y}$$

$$= \sum_{i=1}^N x_i y_i - \bar{x} \sum_{i=1}^N y_i - N \bar{x} \bar{y} + N \bar{x} \bar{y}$$

$$= \sum_{i=1}^N (x_i - \bar{x}) y_i$$

↳ Using the same trick, we can write S_x^2 as $\sum_{i=1}^N (x_{i1} - \bar{x}_1) x_{i1}$

Einfacher:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

$$\begin{aligned} a_1 &= \frac{\hat{\text{Cov}}(X_1, Y)}{\hat{\text{Var}}(X_1)} = \frac{\hat{\text{Cov}}(X_1, \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon)}{\hat{\text{Var}}(X_1)} \\ &= \frac{\hat{\text{Cov}}(X_1, \beta_0) + \beta_1 \hat{\text{Cov}}(X_1, X_1) + \beta_2 \hat{\text{Cov}}(X_1, X_2) + \hat{\text{Cov}}(X_1, \varepsilon)}{\hat{\text{Var}}(X_1)} \\ &= 0 + \beta_1 \frac{\hat{\text{Var}}(X_1)}{\hat{\text{Var}}(X_1)} + \beta_2 \frac{\hat{\text{Cov}}(X_1, X_2)}{\hat{\text{Var}}(X_1)} + 0 \\ &= \beta_1 + \beta_2 \frac{\hat{\text{Cov}}(X_1, X_2)}{\hat{\text{Var}}(X_1)} \end{aligned}$$

$\beta_2 :=$ Populationsparameter von
Size