

Methoden der Ökonometrie - Übung 9

Aufgabe 1:

Sei $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ ein lineares Modell, sodass \mathbf{X} vollen Spaltenrang hat (mit Wahrscheinlichkeit 1). Erinnern Sie sich an die Aussage und den Beweis des Satzes von Gauß-Markov. Nehmen Sie jetzt an, dass \mathbf{X} zufällig ist und nicht fest vorgegeben. Die Annahmen aus der Vorlesung lauten jetzt

- $E(\mathbf{u}|\mathbf{X}) = \mathbf{0}$,
- $Var(\mathbf{u}|\mathbf{X}) = \sigma^2 \mathbf{I}_n$.

Zeigen Sie mit Hilfe der Beweise aus der Vorlesung für den OLS Schätzer $\hat{\boldsymbol{\beta}}_{OLS}$, dass

- $E(\hat{\boldsymbol{\beta}}_{OLS}|\mathbf{X}) = \boldsymbol{\beta}$,
- $Var(\hat{\boldsymbol{\beta}}_{OLS}|\mathbf{X}) = \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$,
- für jeden weiteren linearen erwartungstreuen Schätzer $\tilde{\boldsymbol{\beta}}$ gilt $Var(\tilde{\boldsymbol{\beta}}|\mathbf{X}) \geq Var(\hat{\boldsymbol{\beta}}_{OLS}|\mathbf{X})$.
- Wie können Sie daraus folgern, dass die Aussagen aus der Vorlesung gelten, also $E(\hat{\boldsymbol{\beta}}_{OLS}) = \boldsymbol{\beta}$,
- sowie $Var(\tilde{\boldsymbol{\beta}}) \geq Var(\hat{\boldsymbol{\beta}}_{OLS})$? *Hinweis:* Nutzen Sie Aufgabe 3a) aus der zweiten Übung.
- Zeigen Sie, dass die Erwartungstreue von $\tilde{\boldsymbol{\beta}}$ notwendig ist als Voraussetzung im Satz von Gauß-Markov, also dass es ein $\tilde{\boldsymbol{\beta}}$ gibt welches nicht erwartungstreu ist, sodass $Var(\tilde{\boldsymbol{\beta}}|\mathbf{X}) < Var(\hat{\boldsymbol{\beta}}_{OLS}|\mathbf{X})$. *Hinweis:* Was ist mit einem konstanten Schätzer?
- Berechnen Sie $Var(\hat{\mathbf{u}}|\mathbf{X})$.

Aufgabe 2:

Betrachten Sie auf Folie 3-15 den OLS-Schätzer $\hat{\boldsymbol{\beta}} = \{\mathbf{X}^T(\mathbf{I} - \mathbf{P}_Z)\mathbf{X}\}^{-1}\mathbf{X}^T(\mathbf{I} - \mathbf{P}_Z)\mathbf{y}$ unter Einschluss von irrelevanten Variablen \mathbf{Z} . Zeigen Sie, dass $\hat{\boldsymbol{\beta}}$ diese Form hat und berechnen Sie die Kovarianzmatrix von $\hat{\boldsymbol{\beta}}$. Warum kann diese nicht minimal sein (im Matrixsinn), wenn das wahre lineare Modell $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ ist?

Aufgabe 3:

Betrachten Sie den Datensatz „wage2“, dieser enthält Daten zum logarithmierten Lohn, zum gemessenen IQ und zu den gearbeiteten Stunden. Betrachten Sie den folgenden R-Output.

```
> summary(lm(log(wage)~IQ,data=wage2))

Call:
lm(formula = log(wage) ~ IQ, data = wage2)

Residuals:
    Min       1Q   Median       3Q      Max
-2.09324 -0.25547  0.02261  0.27544  1.21486

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.8869943   0.0890206   66.13  <2e-16 ***
IQ           0.0088072   0.0008694   10.13  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3999 on 933 degrees of freedom
Multiple R-squared:  0.09909, Adjusted R-squared:  0.09813
F-statistic: 102.6 on 1 and 933 DF, p-value: < 2.2e-16
```

- Interpretieren Sie die Eingabe. Was wurde hier gemacht?
- Interpretieren Sie die geschätzten Koeffizienten.
- Sind die Koeffizienten signifikant? Was bedeutet Signifikanz?
- Was bedeuten die Angaben „Residual standard error“, „Multiple R-squared“, „Adjusted R-Squared“ und „F-statistic“? Berechnen Sie das *SSR*.
- Lässt sich aus (c) bereits ein kausaler Zusammenhang zwischen dem Log-Lohn und der Intelligenz folgern?
- Betrachten Sie den nächsten Output. Was wurde hier gemacht? Wieso hat das Vorteile, auch wenn Sie sich nur für den Zusammenhang zwischen dem Log-Lohn und der Intelligenz interessieren?
- Warum kann man jetzt annehmen, dass bereits im ersten Modell unverzerrt geschätzt wurde? Aus welchen zwei Gründen könnte dies passiert sein? Welcher trifft eher zu?
- Wie erklären Sie sich inhaltlich den signifikant negativen Regressionskoeffizienten von „hours“?

```
> summary(lm(log(wage)~IQ+hours,data=wage2))

Call:
lm(formula = log(wage) ~ IQ + hours, data = wage2)

Residuals:
    Min       1Q   Median       3Q      Max
-2.06915 -0.24533  0.02365  0.27646  1.28406

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  6.0536074   0.1150416   52.621  <2e-16 ***
IQ           0.0089535   0.0008698   10.293  <2e-16 ***
hours       -0.0041302   0.0018124   -2.279   0.0229 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3991 on 932 degrees of freedom
Multiple R-squared:  0.1041, Adjusted R-squared:  0.1022
F-statistic: 54.14 on 2 and 932 DF, p-value: < 2.2e-16
```

Aufgabe 4:

- a) In einer Studie des National Women's Hospital in Auckland wurden bei Föten Ultraschallmessungen der Leber durchgeführt. Laden Sie den Datensatz "liver" aus der Datei "liverdata.Rda" von der Githubseite. Die Variable "Age" beschreibt die Schwangerschaftswoche und "Length" die gemessene Leberlänge der Föten.
- b) Plotten Sie die Leberlänge gegen das Alter. Bevor Sie weitermachen, überlegen Sie sich wie eine gute Regressionskurve aussehen könnte.
- c) Betrachten Sie die folgenden drei Regressionsmodelle:

i) $\log(Y_i) = \alpha + \beta x_i + \gamma x_i^2 + u_i$

ii) $\frac{Y_i}{x_i^{3/2}} = \alpha + \beta x_i + \gamma x_i^2 + u_i$

iii) $Y_i = \alpha + \beta x_i + u_i$

wobei x dem Alter und Y der Lebergröße entspricht. Ermitteln Sie die jeweiligen Bestimmtheitsmaße und vergleichen Sie. Für welches Modell würden Sie sich demzufolge entscheiden?

- d) Zeichnen Sie die drei Regressionskurven in den Scatterplot ein. Für welches Modell würden Sie sich jetzt entscheiden?