

Aufgabe 1

[46 Punkte]

Ein Marktforschungsunternehmen wurde beauftragt, die Einflussfaktoren auf die tägliche Fernsehdauer zu untersuchen. Dazu wurden T=100 Personen zufällig ausgewählt und zu folgenden Aspekten befragt:

- TV:** durchschnittliche Fernsehdauer pro Tag in Stunden
- EK:** monatliches Nettoeinkommen in €
- Alter:** Alter der befragten Person in Jahren
- Abi:** Person hat Abitur (ja = 1, nein = 0)
- G:** Geschlecht (weiblich = 1, männlich = 0)
- Allein:** Person lebt alleine (ja = 1, nein = 0)

Die Marktforscher unterstellen folgendes Modell:

$$TV_t = \beta_1 + \beta_2 \cdot EK_t + \beta_3 \cdot Alter_t + \beta_4 \cdot Abi_t + \beta_5 \cdot G_t + \beta_6 \cdot Allein_t + e_t$$

Die Auswertung der Daten mit R ergab folgenden Output:

```
Call:
lm(formula = TV ~ EK + Alter + Abi + G + Allein)

Coefficients:
              Estimate      Std. Error    t value    Pr(>|t|)
(Intercept)  4.0020803    0.6649798     6.018    3.37e-08 ***
EK           -0.0010280    0.0002175      ?      8.01e-06 ***
Alter        0.0333995      ?             3.817    0.000242 ***
Abi          -0.2571470    0.2308090    -1.114    0.268072
G            -0.1293437    0.2011797    -0.643    0.521837
Allein       -0.5251771    0.2122891    -2.474    0.015160 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9752 on 94 degrees of freedom
Multiple R-Squared:  ?,    Adjusted R-squared: 0.5526
F-statistic: 25.46 on 5 and 94 DF, p-value: 3.768e-16
```

a) Berechnen Sie unter Angabe des Rechenwegs die im Output fehlenden Werte für (5)

a1) den t-Wert für b_2

$$- \quad t_{b_2} = \frac{b_2}{se(b_2)} = \frac{-0,001028}{0,0002175} = -4,7264$$

a2) den Standardfehler von b_3

$$- \quad se(b_3) = \frac{b_3}{t_{b_3}} = \frac{0,0334}{3,817} = 0,00875$$

a3) die Quadratsumme der Fehlerterme (SSE)

$$- \quad \hat{\sigma}^2 = \frac{\sum \hat{e}_t^2}{T - K} = \frac{SSE}{T - K}$$

$$- \quad \hat{\sigma} = 0,9752 \Rightarrow \hat{\sigma}^2 = 0,951015$$

$$\Rightarrow SSE = 0,951015 \cdot 94 = 89,3954$$

a4) das Bestimmtheitsmaß

$$- \bar{R}^2 = 1 - \frac{SSE / (T - K)}{SST / (T - 1)}$$

$$\bar{R}^2 = 0,5526$$

$$T = 100$$

$$K = 6$$

$$\bar{R}^2 = 0,5526 = 1 - \frac{SSE / (T - K)}{SST / (T - 1)} \Rightarrow \frac{SSE \cdot (T - 1)}{SST \cdot (T - K)} = 0,4474$$

$$\Rightarrow \frac{SSE}{SST} = 0,4474 \cdot \frac{T - K}{T - 1} = 0,4474 \cdot \frac{94}{99} = 0,4248$$

$$R^2 = 1 - \frac{SSE}{SST} = 1 - 0,4248 = 0,5752$$

- ODER:

$$\frac{SSE}{SST} = 0,4248 \Rightarrow SST = \frac{SSE}{0,4248} = \frac{89,3954}{0,4248} = 210,4411$$

$$\Rightarrow R^2 = 1 - \frac{SSE}{SST} = \frac{89,3954}{210,4411} = 0,5752$$

b) Betrachten Sie die Koeffizienten b_2 , b_3 und b_4 .

(4 Punkte)

b1) Beurteilen Sie die Signifikanz dieser Koeffizienten.

- Einkommen und Alter haben einen signifikanten Einfluss auf die Fernsehdauer, da die entsprechenden Koeffizienten b_2 und b_3 statistisch signifikant von Null verschieden sind. Bei b_4 ist dies nicht der Fall, Abi hat also keinen signifikanten Einfluss.

b2) Um wie viele Minuten ändert sich die Fernsehdauer, wenn das Einkommen um 100,- € steigt?

- Einkommenseffekt: Mit steigendem Einkommen nimmt die Fernsehdauer tendenziell ab ($b_2 < 0$). Steigt das Einkommen um 100,- €, dann nimmt die Fernsehdauer im Durchschnitt um $0,001 \cdot 100 = 0,1$ Stunden = 6 Minuten ab.

b3) Wie unterscheidet sich die Fernsehdauer in Minuten für Personen mit und ohne Abitur?

- Personen mit Abitur sehen im Durchschnitt ca. 0,26 Stunden \approx 15 Minuten weniger fern.

c) Die Marktforscher vermuten, dass der Einkommenseffekt für Männer und Frauen gleich ist. Wie gehen Sie vor, um diese Hypothese zu testen? Unter welchen Umständen verwerfen Sie sie? Beschreiben Sie detailliert Ihre Vorgehensweise. (4 Punkte)

- Frage: Ist der Einfluss des Einkommens auf die durchschnittliche Fernsehdauer für Männer und Frauen gleich?
- Aufnahme eines Interaktionsterms ins Modell: $\gamma_1 \cdot (EK_t \cdot G_t)$
- Schätzung der Koeffizienten, wenn Modell den Interaktionsterm enthält
- t-Test der $H_0: \gamma_1 = 0$; Wenn γ_1 signifikant ist, wird H_0 verworfen, es besteht ein geschlechtsspezifischer EK-Effekt, d. h. dann ist der Einfluss des Einkommens auf die durchschnittliche Fernsehdauer für Männer und Frauen verschieden. Wenn γ_1 nicht signifikant von Null verschieden ist, wird H_0 nicht verworfen, es besteht kein statistisch relevanter Unterschied zwischen den Einkommenseffekten der Geschlechter.

d) Die Marktforscher wollen wissen, ob ihr Modell korrekt spezifiziert ist.

d1) Nennen Sie drei mögliche Arten von Fehlspezifikation.

(3 Punkte)

- Auslassen relevanter Variablen
- Einbeziehen nicht relevanter Variablen
- falsche funktionale Form

d2) Schreiben Sie das „künstliche“ Modell hin, das einem RESET-Test zu Grunde liegt, den man in R mit folgendem Befehl aufruft: `>reset(lm(TV ~ EK + Alter + Abi + G + Allein),2:3)`. (2 Punkte)

Das künstliche Modell lautet:

$$TV_t = \beta_0 + \beta_1 \cdot EK_t + \beta_2 \cdot Alter_t + \beta_3 \cdot Abi_t + \beta_4 \cdot G_t + \beta_5 \cdot Allein_t + \gamma_1 \cdot \widehat{TV}_t^2 + \gamma_2 \cdot \widehat{TV}_t^3 + e_t$$

d3) Ein RESET-Test in R liefert folgendes Ergebnis.

```
RESET test
data: lm(TV ~ EK + Alter + Abi + G + Allein)
RESET = 0.1156, df1 = 1, df2 = 93, p-value = 0.7346
```

Interpretieren Sie dieses Ergebnis.

(2 Punkte)

- p-Wert = 0.7346 > 5%
 ⇒ Der RESET-Test zeigt, dass der zusätzlich berücksichtigte Koeffizient nicht signifikant ist. Es wird keine (statistisch relevante) Fehlspezifikation angezeigt.

e) Ein häufig auftretendes Problem bei der Schätzung von Regressionsmodellen ist Heteroskedastie.

e1) Erläutern Sie verbal, was man unter Heteroskedastie versteht.

(2 Punkte)

- Heteroskedastie bedeutet, dass die Varianz des Fehlerterms nicht konstant ist, sondern über die Beobachtungseinheiten variiert.

e2) Nennen Sie zwei Konsequenzen von Heteroskedastie für KQ-Schätzer.

(2 Punkte)

- Mögliche Konsequenzen von Heteroskedastie
 - KQ-Schätzer sind zwar weiterhin unverzerrt, aber nicht mehr BLUE
 - Standardfehler werden falsch geschätzt
 - herkömmliche Testverfahren und Intervallschätzung sind ungültig

e3) Sie sollen nun für die Variable „Alter“ testen, ob Heteroskedastie vorliegt. Führen Sie den entsprechenden einseitigen Test ($\alpha=5\%$) mit Hilfe der folgenden ANOVA-Tabellen für die beiden Teilstichproben durch (Hinweis: Der Datensatz wurde nach dem Alter sortiert und anschließend in die 2 Gruppen aufgeteilt.). Geben Sie dabei auch die Null- und Alternativhypothese, die Teststatistik, die Verteilung der Teststatistik und die Ablehnungsregion an. Interpretieren Sie das Ergebnis. (8 Punkte)

Analysis of Variance Table					
Response: TV[1:50]					
	Df	Sum Sq	Mean Sq	F value	Pr(>F)
EK[1:50]	1	197899	19.7899	31.5932	1.214e-06 ***
Alter[1:50]	1	24.3355	24.3355	38.8499	1.533e-07 ***
Abi[1:50]	1	0.5054	0.5054	0.8068	0.373962
G[1:50]	1	0.8491	0.8491	1.3556	0.250576
Allein[1:50]	1	7.6950	7.6950	12.2846	0.001063 **
Residuals	44	27.5615	0.6264		

Analysis of Variance Table					
Response: TV[51:100]					
	Df	Sum Sq	Mean Sq	F value	Pr(>F)
EK[51:100]	1	17.527	17.527	14.3133	0.0004639 ***
Alter[51:100]	1	6.811	6.811	5.5624	0.0228572 *
Abi[51:100]	1	4.649	4.649	3.7965	0.0577545
G[51:100]	1	0.001	0.001	0.0008	0.9771001
Allein[51:100]	1	1.392	1.392	1.1369	0.2921332
Residuals	44	53.880	1.225		

- $H_0: \sigma_t^2 = \sigma^2$ (Homoskedastie)

- $H_1: \sigma_t^2 = \sigma^2 \cdot x_t$ (Heteroskedastie)

- $GQ = \frac{\hat{\sigma}_1^2}{\hat{\sigma}_2^2} \sim F_{m_1=T_1-K, m_2=T_2-K}$ (unter H_0)

$m_1 = 50 - 6 = 44$

$m_2 = 50 - 6 = 44$

Da $m_1 = m_2 = 44$ nicht tabelliert ist, unterstellen wir:

$F_{0,95;44;44} \approx F_{0,95;40;40} = 1,69$

\Rightarrow Ablehnungsbereich = $[1,69; \infty[$

- Stichprobenbefund:

$\hat{\sigma}_1^2 = 1,225$

$\hat{\sigma}_2^2 = 0,6264$

- $\Rightarrow GQ = \frac{1,225}{0,6264} = 1,9556$; liegt im Ablehnungsbereich $\Rightarrow H_0$ ablehnen

- Interpretation: Es liegt (statistisch relevante) Heteroskedastie in Abhängigkeit vom Alter vor.

e4) Aus früheren Studien ist bekannt, dass $\text{var}(e_t) = \sigma^2 \cdot \text{Alter}_t^2$. Formulieren Sie ein GLS-Modell, das das Heteroskedastie-Problem löst, und zeigen Sie, warum im transformierten Modell keine Heteroskedastie mehr vorliegt. (5 Punkte)

- $\text{var}(e_t) = \sigma^2 \cdot \text{Alter}_t^2 = \sigma^2 \cdot h_t$

- GLS-Modell: Transformation des Ursprungsmodells durch Division des Ursprungsmodells durch $\sqrt{h_t} = \text{Alter}_t$

$$\frac{TV_t}{\text{Alter}_t} = \beta_1 \cdot \frac{1}{\text{Alter}_t} + \beta_2 \cdot \frac{EK_t}{\text{Alter}_t} + \beta_3 \cdot \frac{\text{Alter}_t}{\text{Alter}_t} + \beta_4 \cdot \frac{Abi_t}{\text{Alter}_t} + \beta_5 \cdot \frac{G_t}{\text{Alter}_t} + \beta_6 \cdot \frac{\text{Allein}_t}{\text{Alter}_t} + \frac{e}{\text{Alter}_t}$$

- Nachweis, dass das transformierte Modell homoskedastisch ist:

$$\text{var}\left(\frac{e_t}{\text{Alter}_t}\right) = \frac{1}{\text{Alter}_t^2} \cdot \text{var}(e_t) = \frac{1}{\text{Alter}_t^2} \cdot \sigma^2 \cdot \text{Alter}_t^2 = \sigma^2$$

- f) Einer der Marktforscher vermutet, dass als wichtiges Kriterium für die Fernsehdauer auch berücksichtigt werden muss, ob jemand arbeitslos ist (Variable „al“, ja=1, nein=0). Dazu wird das obige Modell separat für Arbeitslose und Nicht-Arbeitslose geschätzt. Die jeweiligen ANOVA-Tabellen sind im Folgenden angegeben (Tabelle 1: nur Arbeitslose, Tabelle 2: nur Nicht-Arbeitslose). Führen Sie einen Chow-Test auf dem 5%-Niveau durch. Geben Sie dabei auch die Null- und Alternativhypothese, die Teststatistik, die Verteilung der Teststatistik und die Ablehnungsregion an. Interpretieren Sie das Ergebnis. (9 Punkte)

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
EK[al==0]	1	15.992	15.992	15.4978	0.0003050 ***
Alter[al==0]	1	30.729	30.729	29.7801	2.387e-06 ***
Abi[al==0]	1	8.088	8.088	7.8381	0.0076964 **
G[al==0]	1	0.002	0.002	0.0024	0.9612339
Allein[al==0]	1	2.432	2.432	2.3565	0.1322614
Residuals	42	43.339	1.032		

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
EK[al==1]	1	27.894	27.894	41.6959	5.957e-08 ***
Alter[al==1]	1	36.032	36.032	53.8592	2.834e-09 ***
Abi[al==1]	1	1.120	1.120	1.6737	0.202218
G[al==1]	1	0.241	0.241	0.3608	0.551005
Allein[al==1]	1	7.043	7.043	10.5282	0.002194 **
Residuals	46	30.774	0.669		

- $TV_t = \text{Ursprungsmodell} + \delta_1 \cdot al_t + \delta_2 \cdot EK_t \cdot al_t + \delta_3 \cdot Alter_t \cdot al_t + \dots + \delta_6 \cdot Allein_t \cdot al_t + e_t$

$H_0: \delta_1 = \delta_2 = \dots = \delta_6 = 0$

$H_1: \text{mindestens ein } \delta_i \neq 0$

$$F = \frac{\frac{SSE_R - SSE_U}{J}}{\frac{SSE_U}{T - K}} \sim F_{m_1=J; m_2=T-K} \quad (\text{unter } H_0)$$

$J = 6, \quad T = 100, \quad K = 2 \cdot J = \text{Anzahl der Parameter im unrestringierten Modell} = 12$

$\Rightarrow m_1 = 6, \quad m_2 = 100 - 12 = 88$

$F_{0,95; 6; 88} \approx F_{0,95; 6; 60} = 2,25$

$\Rightarrow \text{Ablehnungsbereich} = [2,25; \infty[$

- $SSE_U = 43,339 + 30,774 = 74,113$

- $SSE_R = 89,3954$ (aus Teilaufgabe a3)

$$\Rightarrow F = \frac{(89,3954 - 74,113)}{\frac{6}{74,113}} = 3,0243 ; \text{ liegt im Ablehnungsbereich } \Rightarrow H_0 \text{ ablehnen!}$$

$$88$$

- Interpretation: Es bestehen signifikante Unterschiede zwischen den Schätzmodellen für Arbeitslose und Nicht-Arbeitslose, die Daten sollten also nicht gepoolt werden.

Aufgabe 2

[10 Punkte]

Welche Antwort ist richtig? Bitte kreuzen Sie die zutreffende Antwort an. Zu jeder Frage gibt es nur eine richtige Antwort. Für jede korrekt angekreuzte Antwort gibt es 1 Punkt, für jede falsch angekreuzte Antwort wird 1 Punkt abgezogen. Die Gesamtpunktzahl kann nicht negativ werden.

1.	Was macht R, wenn Sie den Befehl „>abline(5,10)“ eingeben?	
	<input type="checkbox"/>	Es wird eine Linie erzeugt, die vom x-Wert 5 zum y-Wert 10 verläuft.
	<input type="checkbox"/>	Es wird eine Linie mit der Steigung 5 und dem Achsenabschnitt 10 erzeugt.
	<input checked="" type="checkbox"/>	Es wird eine Linie mit dem Achsenabschnitt 5 und der Steigung 10 erzeugt.

2.	Sie wollen sich die Häufigkeitsverteilung des Merkmals X anzeigen lassen. Welcher R-Befehl ist dazu nicht geeignet?	
	<input type="checkbox"/>	> table(X)
	<input type="checkbox"/>	> hist(X)
	<input checked="" type="checkbox"/>	> quant(X)

3.	Mit welchem der folgenden Parameter der Funktion „>plot()“ ändern Sie die Skalierung der y-Achse?	
	<input checked="" type="checkbox"/>	ylim
	<input type="checkbox"/>	ylab
	<input type="checkbox"/>	yscale

4.	Mit welchem der folgenden Befehle schätzen Sie ein Regressionsmodell ohne Konstante?	
	<input type="checkbox"/>	> lm(y ~ -const+x1+x2+x3)
	<input type="checkbox"/>	> lm(y ~ x1+x2+x3,const=F)
	<input checked="" type="checkbox"/>	> lm(y ~ 0+x1+x2+x3)

5.	Sie wollen die Werte der folgenden Variablen in einem R-Objekt zusammenfassen: Vorname, Alter, Geschlecht und Einkommen. Welchen R-Objektyp müssen Sie dafür wählen?	
	<input type="checkbox"/>	Vektor
	<input checked="" type="checkbox"/>	Dataframe
	<input type="checkbox"/>	Matrix

6.	Welchen Wert berechnen Sie mit folgender Formel (y ist die abhängige Variable einer linearen Regression): >sum((y - predict(y))^2) ?	
	<input checked="" type="checkbox"/>	SSE
	<input type="checkbox"/>	SSR

<input type="checkbox"/>	SST
--------------------------	-----

7.	Mit welchem R-Befehl bestimmen Sie den kritischen Wert einer t-Verteilung mit 15 Freiheitsgraden bei einem Signifikanzniveau von 10%?	
	<input type="checkbox"/>	> pt(0.1, 15)
	<input checked="" type="checkbox"/>	> qt(0.9, 15)
	<input type="checkbox"/>	> pt(0.9, 15)

8.	Bei welchem der folgenden Befehle bekommen Sie eine Fehlermeldung?	
	<input type="checkbox"/>	> pf(0.214, 12, 2)
	<input type="checkbox"/>	> df(50, 1, 1)
	<input checked="" type="checkbox"/>	> qf(49.99, 1, 1)

9.	Sie wollen die Wahrscheinlichkeit dafür berechnen, dass eine standardnormalverteilte Zufallsvariable zwischen 0.8 und 1.5 liegt. Mit welchem R-Befehl erhalten Sie das richtige Ergebnis?	
	<input type="checkbox"/>	> dnorm(1.5, 0, 1) – dnorm(0.8, 0, 1)
	<input checked="" type="checkbox"/>	> pnorm(1.5, 0, 1) – pnorm(0.8, 0, 1)
	<input type="checkbox"/>	> pnorm(0.7, 0, 1)

10.	Der t-Wert für den Koeffizienten b_2 beträgt 5.25 (Zahl der Freiheitsgrade = 15, Signifikanzniveau = 5%). Wie berechnen Sie mit R den entsprechenden p-Wert, der auch im R-Output erscheinen würde?	
	<input checked="" type="checkbox"/>	> 2*(1-pt(5.25, 15))
	<input type="checkbox"/>	> 1-(2*pt(5.25, 15))
	<input type="checkbox"/>	> 1-(2*qt(0.05, 15))

Aufgabe 3

[6 Punkte]

Wo sind die Fehler? Die folgende Funktion enthält 6 Fehler. Markieren Sie diese deutlich auf dem Aufgabenblatt und schreiben Sie die korrekte Formulierung darüber! Das Erkennen eines Fehlers gibt 0,5 Punkte, die richtige Verbesserung gibt ebenfalls 0,5 Punkte. Es werden 0,5 Punkte abgezogen, wenn

- entweder eine richtige Stelle als Fehler gekennzeichnet wurde oder
- der Fehler zwar richtig erkannt wurde, die Verbesserung aber falsch ist.

Die Gesamtpunktzahl kann nicht negativ werden.

Die folgende Funktion soll einen Vorhersagewert für y berechnen, wenn x_{1t} den Wert x_0 annimmt. Es gilt das folgende Regressionsmodell: $y_t = \beta_0 + \beta_1 * x_{1t} + \beta_2 * x_{1t}^2 + e_t$

Es werden zunächst die Koeffizienten geschätzt und anschließend der y -Wert berechnet, der sich ergibt, wenn man für x_{1t} einen beliebigen Wert x_0 einsetzt:

```
> Vorhersage == function(y,x1,x0)
{
  reg = lm(y~x1+x1^2)
  neu = dataframe(x1=c(X0))
  Vorhersage = Predict(neu,reg)
```

```
return(Vorhersage)
}
```

- Korrigierte Funktion:

```
> Vorhersage = function(y,x1,x0)
{
  reg = lm(y~x1+I(x1^2))
  neu = data_frame(x1=c(x0))
  Vorhersage = predict(reg,neu)
  return(Vorhersage)
}
```

Aufgabe 4

[6 Punkte]

Beantworten Sie die folgenden Fragen zu R. Für jede korrekte Antwort bekommen Sie 2 Punkte. Für falsche Antworten wird nichts abgezogen.

a) Mit welchem Befehl weisen Sie dem Vektor x die natürlichen Logarithmen der ganzen Zahlen 20 bis 120 zu? Welches Ergebnis würde im Anschluss daran der Befehl $>length(x)$ liefern?

- Zuweisung der natürlichen Logarithmen der ganzen Zahlen zwischen 20 und 120 auf den Vektor x : $> x = c(log(20:120))$
- $> length(x)$ gibt die „Länge“ des Vektors x an, also die Zahl der in ihm enthaltenen Elemente. Das sind hier 100.

b) Der Vektor x enthält die ganzen Zahlen von 5 bis 8. Welches Ergebnis liefert der Befehl $> sum(x[1:3]*x[2:4])$?

- Ergebnis des Befehls $> sum(x[1:3]*x[2:4])$, wobei $x = c(5:8)$: $5 \cdot 6 + 6 \cdot 7 + 7 \cdot 8 = 128$

c) Mit welchem Befehl erhalten Sie die Fläche unter der Dichtefunktion der Standardnormalverteilung im Bereich von $-\infty$ bis 0.8?

- $> pnorm(0.8,0,1)$

Aufgabe 5

[21 Punkte]

Wahr oder falsch? Tragen Sie für jede der folgenden Aussagen ein „w“ für „wahr“ oder ein „f“ für „falsch“ ein. Für jede richtige Antwort gibt es 1 Punkt, für jede falsche Antwort wird 1 Punkt abgezogen. Die Gesamtpunktzahl kann nicht negativ werden.

F	Relative Konzentrationsmaße beschreiben, auf wie viele Merkmalsträger ein bestimmter Anteil der Merkmalssumme entfällt.
W	Mit dem Lagrange-Multiplier Test kann getestet werden, ob der Störterm eines Modells durch einen autoregressiven Prozess zweiter Ordnung determiniert wird.
F	Die Summe quadrierter χ^2 -verteilter Zufallsvariablen ist F-verteilt, wenn diese Zufallsvariablen statistisch unabhängig voneinander sind.
F	Zur Durchführung eines Instrumentvariablenschätzverfahrens benötigt man eine erklärende Variable, die mit keiner anderen erklärenden Variablen korreliert ist.

F	Der Whiteschätzer wird auf Modelle mit transformierten Variablen angewendet.
F	Der Koeffizient einer Dummyvariable gibt an, ob sich die Steigungsparameter eines Modells für Teilgruppen unterscheiden.
W	Der Herfindahlindex wird berechnet als die Summe der quadrierten Anteile der Merkmalsträger an der Merkmalssumme.
W	Multikollinearitätsprobleme lassen sich über eine Erhöhung der Beobachtungszahl reduzieren.
W	Bei heteroskedastischen Störtermen bleiben die Kleinstquadrateschätzer der Steigungsparameter unverzerrt.
W	Bei am Niveau α statistisch insignifikanten Koeffizienten liegt der Wert „Null“ innerhalb ihres $(1-\alpha) \cdot 100\%$ -Konfidenzintervalls.
W	Der Goldfeldt-Quandt Test kann sowohl als einseitiger als auch als zweiseitiger Test durchgeführt werden.
W	Je höher der Ginikoeffizient ist, umso ungleicher ist die unterliegende Verteilung.
W	Nach dem Method of Moments Schätzverfahren, lassen sich Bevölkerungsparameter durch analoge Parameter der Stichprobe schätzen.
W	Es ist für die Qualität der Schätzergebnisse günstiger, zu viele als zu wenige erklärende Variablen im Modell zu berücksichtigen.
W	Der Paascheindex ist kommensurabel.
W	Ein Schätzer ist konsistent, wenn bei steigender Stichprobengröße die Wahrscheinlichkeit gegen 1 konvergiert, dass der Schätzer in der Nähe des wahren Wertes liegt.
W	Saisonkoeffizienten beschreiben die saisonale Abweichung eines Zeitreihenmittelwerts von \bar{s} .
W	Nicht alle Parameter des logistischen Trendmodells können im linearen Regressionsmodell geschätzt werden.
F	Im Rahmen des Kleinstquadrateschätzers können ausschließlich lineare Beziehungen zwischen erklärenden und abhängigen Variablen geschätzt werden.
W	Der Jarque-Bera Test nutzt eine χ^2 -verteilte Teststatistik.
F	Im einfachen linearen Modell führt eine Multiplikation der erklärenden Variable mit 100 zu einem um den Faktor 100 größeren Steigungsparameter.

Aufgabe 6

[10 Punkte]

Welche Antwort ist richtig? Bitte kreuzen Sie die zutreffende Antwort an. Zu jeder Frage gibt es nur eine richtige Antwort. Für jede korrekt angekreuzte Antwort gibt es 1 Punkt, für jede falsch angekreuzte Antwort wird 1 Punkt abgezogen. Die Gesamtpunktzahl kann nicht negativ werden.

1.	Mit welcher Wahrscheinlichkeit fällt eine mit 3 Freiheitsgraden t-verteilte Zufallsvariable in das Intervall [1.638 ; 3.182]?	
	<input checked="" type="checkbox"/>	0.075
	<input type="checkbox"/>	0.15
	<input type="checkbox"/>	0.925

Musterlösung zur Diplomvorprüfung Statistik II – Einf. Ökonometrie im SS 05

2.	Um das lineare Modell $y_t = \beta_0 + \beta_1 \cdot x_{1t} + \beta_2 \cdot x_{2t} + \beta_3 \cdot x_{3t} + e_t$ unter der Restriktion zu schätzen, dass die Summe der Steigungsparameter 5 beträgt, nutzt man folgendes Modell:	
	<input type="checkbox"/>	$y_t - 5 \cdot x_{1t} = \beta_0 \cdot (1 - x_0) + \beta_2 \cdot (x_{2t} - x_{1t}) + \beta_3 \cdot (x_{3t} - x_{1t}) + e_t$
	<input checked="" type="checkbox"/>	$y_t - 5 \cdot x_{1t} = \beta_0 + \beta_2 \cdot (x_{2t} - x_{1t}) + \beta_3 \cdot (x_{3t} - x_{1t}) + e_t$
	<input type="checkbox"/>	$y_t - 5 \cdot x_{1t} - \beta_0 = \beta_2 \cdot (x_{2t} - x_{1t}) + \beta_3 \cdot (x_{3t} - x_{1t}) + e_t$
3.	Welche Interpretation gilt für β_1 im Modell für die Einfachregression $\ln y = \beta_0 + \beta_1 \cdot \ln x + e$?	
	<input type="checkbox"/>	Marginaler Effekt: bei einer Änderung von x um eine Einheit ändert sich y um eine Einheit.
	<input type="checkbox"/>	Semielastizität: bei einer Änderung von x um eine Einheit ändert sich y um ein Prozent.
	<input checked="" type="checkbox"/>	Elastizität: bei einer Änderung von x um ein Prozent ändert sich y um ein (falsch! β_1) Prozent.
4.	Im linearen Modell ist die Varianz des Prognosefehlers umso kleiner,	
	<input type="checkbox"/>	je größer der Wert der erklärenden Variablen für den Prognosezeitpunkt.
	<input checked="" type="checkbox"/>	je kleiner die Streuung der geschätzten Parameter.
	<input type="checkbox"/>	je stärker der Zusammenhang zwischen den geschätzten Parametern.
5.	Das R^2 beschreibt	
	<input type="checkbox"/>	den Anteil der auf Basis des Modells korrekt vorhersagbaren Werte der abhängigen Variable.
	<input type="checkbox"/>	den Korrelationskoeffizienten zwischen der abhängigen und den erklärenden Variablen.
	<input checked="" type="checkbox"/>	den mit dem Modell erklärten Anteil der Variation der abhängigen Variablen.
6.	Interaktionseffekte zwischen erklärenden Variablen	
	<input type="checkbox"/>	sind nötig, wenn die Effekte qualitativer erklärender Variablen geschätzt werden sollen.
	<input type="checkbox"/>	können die Schätzgüte eines Modells reduzieren.
	<input checked="" type="checkbox"/>	bieten die Möglichkeit, für Teilstichproben unterschiedliche Steigungsparameter zu schätzen.
7.	Der Kleinstquadratschätzer	
	<input type="checkbox"/>	sollte nur bei normalverteilten abhängigen Variablen verwendet werden.
	<input type="checkbox"/>	sollte nur bei normalverteilten Residuen verwendet werden.
	<input checked="" type="checkbox"/>	kann auch bei nicht normalverteilten Koeffizienten das Gauss-Markov-Theorem erfüllen.
8.	Die Normalgleichungen des KQ-Schätzers	
	<input checked="" type="checkbox"/>	ergeben sich bei Minimierung der Zielfunktion.
	<input type="checkbox"/>	sind über das Method of Moments Verfahren nicht herleitbar.
	<input type="checkbox"/>	können nur im einfachen KQ-Modell bestimmt werden.
9.	Die Varianz von in erster Ordnung autokorrelierten Störtermen (AR(1))	
	<input type="checkbox"/>	ist immer heteroskedastisch.
	<input checked="" type="checkbox"/>	ist nur definiert für $\rho \neq 1$.
	<input type="checkbox"/>	ist umso größer, je länger die von der Stichprobe beschriebene Zeitspanne ist.

10.	Der Durbin-Watson Test	
<input checked="" type="checkbox"/>	unterstellt eine AR(1)-Verteilung der Fehlerterme.	
<input type="checkbox"/>	kann genutzt werden, wenn sich im Modell die verzögerte abhängige Variablen unter den erklärenden Variablen befindet.	
<input type="checkbox"/>	ist nur für positive Autokorrelation anwendbar.	

Aufgabe 7

[21 Punkte]

Sie untersuchen im Rahmen eines einfachen linearen Regressionsmodells den Zusammenhang zwischen der Anzahl verkaufter Personenwagen eines Monats (X_t) und dem monatlichen Börsenschlusswert (Y_t) eines Automobilherstellers. Ihnen liegen Daten zu Verkaufszahlen und Börsenkursen für 60 Monate vor. Ihr Modell lautet

$$Y_t = \beta_0 + \beta_1 \cdot X_t + e_t$$

a) Sie sind besorgt, ob man den Verkaufszahlen trauen kann und vermuten, dass die X-Variable fehlerhaft gemessen ist. Erläutern Sie die Konsequenz für Ihren Schätzer. (3 Punkte)

- Messfehler in der erklärenden Variablen führen im linearen Modell zu inkonsistenten und verzerrten Schätzern für den Steigungsparameter. Die geschätzten Parameter sind darüber hinaus nicht mehr approximativ normalverteilt. "Attenuation bias" führt dazu, dass der Koeffizient zu klein ausgewiesen wird bzw. gegen die 0 hin verzerrt ist.

b) Außerdem sind Sie nicht sicher, ob die Annahme $cov(e_t, e_s) = 0$ in Ihrem Fall zutrifft. Beschreiben Sie die einzelnen Schritte des LM-Tests zur Überprüfung dieser Hypothese. (4 Punkte)

- Wenn $cov(e_t, e_s) \neq 0$ müssen wir eine Annahme zum Störtermprozess treffen, z. B. AR(1):

$$e_t = \rho \cdot e_{t-1} + u_t$$
- Der LM-Test schätzt dann:

$$y_t = \beta_0 + \beta_1 \cdot x_t + \rho \cdot e_{t-1} + e_t$$
 für $t = 2, \dots, T$
- Dazu wird zunächst das Basismodell $y_t = \beta_0 + \beta_1 \cdot x_t + e_t$ geschätzt, dann der verzögerte Wert e_{t-1} gewonnen und im LM-Modell genutzt.
- Bei insignifikantem ρ wird die Hypothese eines AR(1)-Störterms verworfen; ist ρ statistisch signifikant von Null verschieden, kann die Hypothese eines AR(1)-Störterms nicht verworfen werden.

c) Was bedeutet es für die Eigenschaften Ihres Schätzers, wenn die Annahme $cov(e_t, e_s) = 0$ nicht zutrifft? (2 Punkte)

- Wenn $cov(e_t, e_s) \neq 0$ ist der KQ-Schätzer nicht effizient (nicht BLUE), die Standardfehler werden falsch ausgewiesen und somit sind die Hypothesentests und Konfidenzintervalle nicht mehr verlässlich.

d) Ihr Kollege bezweifelt, dass die Annahme $cov(x_t, e_s) = 0$ erfüllt ist und schlägt vor, dass Sie dies mit einem Hausman-Test überprüfen. Beschreiben Sie die Grundidee des Tests und erläutern Sie die notwendigen Schritte, um den Test durchzuführen. (7 Punkte)

- *Grundidee:*
Vergleiche Parameterschätzer aus KQ und IV-Verfahren, da unter $H_0: cov(x_t, e_t) = 0$ beide Schätzer konsistent sind (Differenz ≈ 0) und unter $H_1: cov(x_t, e_t) \neq 0$ nur der IV-Schätzer konsistent ist (Differenz $\neq 0$).
- *Schritte:*
 - Regressiere x_t auf Instrument
 - Bestimme Residuum dieser Schätzgleichung
 - Füge dieses Residuum in die Schätzgleichung für y ein und schätze dessen Koeffizienten.
 - Teste diesen Koeffizienten auf statistische Signifikanz. Falls er signifikant ist, verwirfe H_0 .

- e) **Angenommen, die Schätzung ergibt: $\hat{\beta}_0 = 3$ und $\hat{\beta}_1 = 2$ und es gilt ein AR(1) Modell mit $\hat{\rho} = 0.8$. Sie erwarten, dass die Verkaufszahlen im Monat $t = 61$ mit $X_{61} = 100$ PKW genauso hoch sind, wie im letzten Beobachtungsmonat ($X_{60} = 100$), zu dessen Ende der Börsenkurs $Y_{60} = 205$ betragen hat. Wie lautet Ihre Prognose des Börsenschlusskurses für die Periode $t = 61$? Stellen Sie Ihren Rechenweg dar und erläutern Sie die notwendigen Schritte knapp in verbaler Form. (5 P)**

- $\hat{y}_{61} = \hat{\beta}_0 + \hat{\beta}_1 \cdot x_{61} + e_{61}$

- $e_{61} = e_{60} \cdot \hat{\rho} + v_t$

- $\hat{e}_{61} = \hat{e}_{60} \cdot \hat{\rho}$

- $\hat{e}_{60} = \hat{y}_{60} - \hat{\beta}_0 - \hat{\beta}_1 \cdot x_{60}$
 $= 205 - 3 - 2 \cdot 100 = 2$

- $\hat{y}_{61} = 3 + 2 \cdot 100 + 2 \cdot 0,8 = 204,6$

- Die Prognose enthält eine Fortschreibung des Störprozesses für die Periode $t = 61$ neben dem herkömmlichen vorhergesagten Wert $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 \cdot x_0$.