

Aufgabe 1:

(23 Punkte)

Ihnen liegt ein Datensatz über Angestellte eines ländlichen mittelständischen Unternehmens vor.

Darin sind u.a. folgende Daten enthalten:

SICK → Krankheitstage; *AGE* → Alter in Jahren; *SEX* → weiblich = 1, männlich = 0; *EDUC* → Hochschulabschluss = 1, kein Hochschulabschluss = 0; *SPORT* → der/die Befragte ist Mitglied in einem Sportverein = 1; ist nicht Mitglied = 0; *DIST* → Entfernung vom Wohnort zum nächsten Sportverein in Kilometern

Sie haben den Auftrag erhalten, die Anzahl der Krankheitstage pro Jahr im Unternehmen zu untersuchen und entscheiden sich für das folgende lineare Modell:

$$SICK = \beta_1 + \beta_2 AGE + \beta_3 SEX + \beta_4 (AGE*SEX) + \beta_5 EDUC + \beta_6 SPORT + e$$

```
Call:
lm(formula = SICK ~ AGE + SEX + AGE*SEX + EDUC + SPORT)

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.58834     0.48408   1.215  0.2281
AGE          ..???.   0.00980   4.466 2.81e-05 ***
SEX          0.72924     0.51790   1.408  0.1633
AGE*SEX      -0.02246     0.01308  -1.717  0.0901 .
EDUC         0.01440     0.15804   0.091  0.9277
SPORT       -1.65159     0.20856  -7.919 1.84e-11 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.6929 on 74 degrees of freedom
Multiple R-Squared:  ???, Adjusted R-squared:  0.715
F-statistic: 40.63 on 5 and 74 DF, p-value: < 2.2e-16
```

- a) Berechnen Sie den fehlenden Wert für b_2 und das multiple Bestimmtheitsmaß (R^2). (5 Punkte)
- b) Interpretieren Sie die geschätzten Parameter (b_1 bis b_6) inhaltlich und statistisch. (7,5 Punkte)
- c) Welche Wahrscheinlichkeit gibt der p-Wert der Koeffizientenschätzer im Regressionsoutput an (rechte Spalte)? (2,5 Punkte)
- d) Eine Kollegin zweifelt Ihr Modell und Ihre Schätzergebnisse an. Sie argumentiert, dass von Hause aus gesunde Menschen häufiger Mitglied in einem Sportverein sind als kränkliche.
 - di) Worauf zielt der Einwand Ihrer Kollegin ab und was wäre die Konsequenz für Ihre KQ-Schätzung, wenn sie Recht hat? Erläutern Sie. (3 Punkte)
 - dii) Ein anderer Kollege schlägt vor, eine Instrumentvariablenschätzung durchzuführen. Gegeben den vorliegenden Datensatz, welche Möglichkeit haben Sie, den Zusammenhang der Variable *SPORT* mit den Krankheitstagen verlässlich zu schätzen? Erläutern Sie kurz. Zählen Sie die Anforderungen auf, die allgemein an ein Instrument gestellt werden und diskutieren Sie, inwieweit diese Anforderungen im vorliegenden Fall erfüllt sind. (5 Punkte)

Hinweis: R liefert Ihnen folgende Korrelationskoeffizienten:

	SPORT	DIST
SPORT	1	-0,499659
DIST	-0,499659	1

Aufgabe 2:

(24 Punkte)

Ihnen liegt das Ergebnis einer Untersuchung mit 50 Beobachtungen vor, in der das Jahreseinkommen (in 1000 €) als lineare Funktion von Geschlecht (SEX), Ausbildung in Jahren (EDUC), Alter (AGE) und Alter zum Quadrat (AGE2) geschätzt wurde. Das Modell lautet:

$$\text{EINK} = \beta_1 + \beta_2 \text{SEX} + \beta_3 \text{EDUC} + \beta_4 \text{AGE} + \beta_5 \text{AGE}^2 + e$$

```
Call:
lm(formula = EINK ~ SEX + EDUC + AGE + AGE2)
```

Coefficients:

	Estimate	Std. Error	t value
Intercept	16.27950	15.13120	1.076
SEX	4.72105	1.22198	3.863
EDUC	0.03725	0.00960	3.880
AGE	0.96152	0.01308	73.511
AGE2	-0.00920	0.21615	-0.043

- Berechnen Sie den marginalen Effekt des Alters auf das Jahreseinkommen und interpretieren sie ihn. Wann ist er Null? Wie können Sie die Signifikanz des Alterseffekts testen? Geben Sie die Nullhypothese und die Alternativhypothese an. *(5 Punkte)*
- Was ist der Effekt *(2 Punkte)*
 - der Logarithmierung des Jahreseinkommens auf die Interpretation von β_3 ?
 - auf die Ausprägung des Koeffizienten β_3 wenn das Jahreseinkommen in € statt in 1.000 € gemessen wurde?
- Betrachten Sie nun den Ausbildungseffekt.
 - Unterscheidet sich im vorliegenden Schätzmodell der Ausbildungseffekt für Männer und Frauen? Begründen Sie kurz. *(1 Punkt)*
 - Wie könnte man prüfen, ob es einen statistisch signifikanten Geschlechterunterschied im Ausbildungseffekt gibt? *(2 Punkte)*
 - Welche Hypothese testet in diesem Zusammenhang der Chow Test? Beschreiben Sie am obigen Beispiel die Vorgehensweise des Tests. Geben Sie die Teststatistik, die Freiheitsgrade, die Null- und Alternativhypothese sowie die Interpretation möglicher Testergebnisse an. *(6 Punkte)*
- Wie würden Sie vorgehen, um die Hypothese zu testen, dass sich die Varianz des Störterms für Männer und Frauen unterscheidet? Beschreiben Sie Ihre Vorgehensweise, die Nullhypothese und die genaue Teststatistik. *(4 Punkte)*
- Berechnen Sie das 95%-Konfidenzintervall für den Effekt der Ausbildungsjahre und interpretieren Sie es. *(4 Punkte)*

Aufgabe 3:**(22 Punkte)**

Da Sie sich mit dem Gedanken tragen, in Nürnberg ein Eiscafé zu eröffnen, interessieren Sie sich für Einflussfaktoren, die den örtlichen Eiskonsum bestimmen. Sie haben Daten für 30 Beobachtungen erhoben, wobei jede Beobachtung einem Zeitraum von 4 Wochen entspricht. Der letzte Beobachtungszeitraum ging gestern zu Ende.

Die abgefragten Variablen sind:

- Q: Eiskonsum pro Kopf (in Litern)
- P: Preis pro Liter Eis (in Euro)
- E: Durchschnittliches wöchentliches Haushaltseinkommen (in Euro)
- T: Durchschnittliche Temperatur (in Grad Celsius)

Die Regressionsgleichung lautet: $Q_i = \beta_0 + \beta_1 \cdot P_i + \beta_2 \cdot E_i + \beta_3 \cdot T_i + e_i$

Sie führen in R eine Kleinstquadrateschätzung (KQ) durch und erhalten folgenden Output :

Coefficients:				
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.1973151	0.2702162	0.730	0.47179
P	-1.0444140	0.8343573	-1.252	0.22180
E	0.0033078	0.0011714	2.824	0.00899 **
T	0.0034584	0.0004455	7.762	3.1e-08 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1				
Residual standard error: 0.03683 on 26 degrees of freedom				
Multiple R-Squared: 0.719, Adjusted R-squared: 0.6866				
F-statistic: 22.17 on 3 and 26 DF, p-value: 2.451e-07				

- a) Sie zeigen die Ergebnisse einem Kommilitonen, der sich die KQ-Residuen in einem Plot genauer ansieht und Sie auf das Problem möglicher Autokorrelation hinweist. Er kritisiert Ihre Ergebnisse und sagt, die KQ-Schätzer seien verzerrt und die Standardfehler zu groß. Nehmen Sie zu der Kritik Stellung. (2 Punkte)
- b) Was ist mit Autokorrelation gemeint? Geben Sie im Anschluss für den allgemeinen Fall eine formale Darstellung eines AR(1) Prozesses, in der Sie die einzelnen Formelelemente kurz beschreiben. (4 Punkte)
- c) Um zu überprüfen, ob Ihr Kommilitone mit seiner Autokorrelationsvermutung Recht hat, führen Sie in R einen Durbin-Watson Test auf positive Autokorrelation erster Ordnung auf dem 5%-Signifikanzniveau durch. Sie erhalten folgenden R-Output: (4 Punkte)

Durbin-Watson test	
data:	kq
DW =	1.0212, p-value = 0.0003024

Geben Sie die getestete Null- und Alternativhypothese, die kritischen Werte für den Durbin-Watson Test sowie die darauf basierende Testentscheidung an. Erläutern Sie, ob der p-Wert Ihre Testentscheidung unterstützt. Zu welcher Schlussfolgerung kämen Sie, wenn $DW=1,307$ wäre?

- d) Sie möchten Ihr Ergebnis aus c) überprüfen und führen in R einen Lagrange-Multiplier Test auf dem 5%-Signifikanzniveau durch. Sie erhalten folgenden R-Output: (4,5 Punkte)

```

Coefficients:
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.2588680  0.2571651   1.007  0.32376
P            -1.1920552  0.7918621  -1.505  0.14476
E             0.0031920  0.0011085   2.879  0.00805
T             0.0032551  0.0004328   7.520  7.12e-08
e             0.4282815  0.2112149   2.028  0.05338
---
Residual standard error: 0.03481 on 25 degrees of freedom
Multiple R-Squared:  0.7587,    Adjusted R-squared:  0.7201
F-statistic: 19.65 on 4 and 25 DF,  p-value: 2.009e-07

```

Geben Sie die geschätzte Gleichung für das Beispiel der Aufgabe an und erläutern Sie allgemein die Testidee knapp und präzise. Wird das Ergebnis des ersten Durbin-Watson Tests (aus c)) bestätigt? Warum könnten sich die Testergebnisse unterscheiden?

- e) Sie führen eine Generalized Least Squares (GLS-) Schätzung in R durch und erhalten folgende Schätzergebnisse: (7,5 Punkte)

```

Generalized least squares fit by REML
Model: Q ~ P + E + T

Coefficients:
      Value  Std.Error  t-value  p-value
(Intercept)  0.1570    0.2896    0.5421   0.4725
P            -0.8922    0.8108   -1.1004   0.2215
E             0.0032    0.0015    2.0730   0.0490
T             0.0036    0.0006    6.4175   0.0000

```

Machen Sie auf Basis der Ergebnisse der GLS-Schätzung eine möglichst präzise Vorhersage für die Höhe des Nürnberger pro Kopf Eiskonsums in den kommenden vier Wochen unter Ausnutzung aller Informationen. Gehen Sie davon aus, dass

$$Q_{30} = 0,763$$

$$P_{30} = 4,1 \quad P_{31} = 3,8$$

$$E_{30} = 300 \quad E_{31} = 310$$

$$T_{30} = 18,6 \quad T_{31} = 19,2$$

$$\sum_{t=2}^T \hat{e}_{t-1} \cdot \hat{e}_t = 0,072$$

$$\sum_{t=2}^T \hat{e}_{t-1}^2 = 0,081$$

und geben Sie Ihren Rechenweg an.

Aufgabe 4:**(10 Punkte)**

In R wurde folgende Funktion programmiert:

```

> my.prog <- function(x,y)
{
  s1 <- sum(x)/length(x)
  s2 <- sum(y[-5])/length(y[-5])
  plot(x,y)
  points(s1,s2)
  c <- c(s1,s2)
  return(c)
}

```

Die Funktion soll auf zwei Vektoren v und z angewendet werden, welche die Ziffern von 1 bis 5 bzw. von 2 bis 6 enthalten.

- a) Geben Sie einen möglichen R-Befehl an, mit dem man den Vektor z generieren kann. (1 Punkt)
- b) Mit welchem R-Befehl rufen Sie die Funktion für die beiden Vektoren auf? (1 Punkt)
- c) Stellen Sie alle Ausgaben so dar, wie sie mit dieser Funktion für diese beiden Vektoren erzeugt werden. (8 Punkte)

Aufgabe 5:**(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 welchem R-Befehl erzeugen Sie einen Vektor x , der ungerade Zahlen zwischen 1 und 99 enthält?
	<input type="checkbox"/> <code>> x <- seq(1,99, by=2)</code>
	<input type="checkbox"/> <code>> x <- seq(1,99, x[2]-1)</code>
2.	Welchen Grafik-Bestandteil erzeugen Sie mit dem R-Befehl <code>> abline(h=0, lty=2)</code> ?
	<input type="checkbox"/> Abszisse in Fettdruck
	<input type="checkbox"/> Horizontale gestrichelte Linie mit Ordinatenabschnitt Null
3.	Mit welchem R-Befehl können Sie die Standardfehler eines als Objekt <code>mod.kq</code> vorliegenden Modelloutputs auslesen?
	<input type="checkbox"/> <code>> mod.kq\$coef[,1]</code>
	<input type="checkbox"/> <code>> mod.kq\$coef[,2]</code>
4.	Bei welchem der folgenden R-Befehle erhalten Sie keine Fehlermeldung?
	<input type="checkbox"/> <code>> pf(.95;12,2)</code>
	<input type="checkbox"/> <code>> pf(.95,12;2)</code>
	<input type="checkbox"/> <code>> pf(.95,12,2)</code>

5.	Welchen R-Befehl müssen Sie verwenden, um ein KQ-Modell zu schätzen, welches eine Variable enthält, die die Interaktion zwischen x und z abbildet?
	<input type="checkbox"/> <code>> lm(y ~ x + z + Int{x*z})</code>
	<input type="checkbox"/> <code>> lm(y ~ x + z + IA[x^z])</code>
6.	Mit welcher Option des R-Befehls <code>read.table</code> kann man einzulesenden Variablen Namen zuweisen?
	<input type="checkbox"/> <code>> col.names=c("var1","var2","var3")</code>
	<input type="checkbox"/> <code>> var.names=c("var1","var2","var3")</code>
7.	Welchen R-Befehl müssen Sie verwenden, um ein KQ-Modell mit der ersten Hälfte eines Datensatzes mit 50 Beobachtungen durchzuführen?
	<input type="checkbox"/> <code>> lm(y ~ x, data.frame[1:25])</code>
	<input type="checkbox"/> <code>> lm(y ~ x, lower.half=T)</code>
8.	Welchen R-Befehl kann man verwenden, um einen Chow-Test durchzuführen?
	<input type="checkbox"/> <code>> chow.test(mod1.kq, mod2.kq)</code>
	<input type="checkbox"/> <code>> anova(mod1.kq, mod2.kq)</code>
9.	Welche Kennzahl berechnet man mit dem Befehl <code>sum.mod\$sigma^2</code> (<code>sum.mod</code> sei der Modelloutput)?
	<input type="checkbox"/> Autokorrelationskoeffizient eines AR(1)-Fehlers
	<input type="checkbox"/> quadriertes Vorhersagewert bei Heteroskedastie
10.	Der t -Wert eines geschätzten Koeffizienten sei 0.67 (bei 48 Freiheitsgraden und einem Signifikanzniveau von 5%). Mit welchem R-Befehl kann man nicht den im Output ausgewiesenen p -Wert manuell berechnen?
	<input type="checkbox"/> <code>> 2*(1-pt(0.67,df=48))</code>
	<input type="checkbox"/> <code>> 2*pt(0.67,df=48,lower.tail=F)</code>
	<input type="checkbox"/> <code>> 1-(2*pt(0.67,df=48,lower.tail=T))</code>

Aufgabe 6:

(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.

	Eine mögliche Interpretation des F-Tests ist, dass er den Erklärungsgehalt unterschiedlicher Modelle vergleicht.
	Um einen Chow Test durchzuführen, sind zwei Schätzungen erforderlich.
	Beim einseitigen t-Test liegt die Ablehnungsregion im Bereich positiver t-Werte.
	Zur Bestimmung des Kleinstquadrateschätzers ist die Annahme normalverteilter Fehlerterme nicht erforderlich.

	Die Annahme, dass der Störterm im multivariaten Regressionsmodell normalverteilt ist, ist für Erwartungstreue und Varianz der Kleinstquadrate-Schätzer unerheblich.
	Im multiplen Regressionsmodell steigt die Varianz eines geschätzten Steigungsparameters β_k , wenn die entsprechende erklärende Variable x_k stark mit anderen erklärenden Variablen im Modell korreliert ist.
	Es ist möglich, den Parameter ρ bei autokorrelierten Störtermen erster Ordnung als Steigungsparameter in einer KQ-Schätzung zu schätzen.
	Der Goldfeld-Quandt Test ist nur für Situationen mit proportionaler Heteroskedastie geeignet.
	Der Herfindahl-Index ist ein absolutes Konzentrationsmaß.
	Die Varianz des Vorhersagefehlers im einfachen Regressionsmodell ist am Mittelwert der erklärenden Variablen am geringsten.
	Ein Typ I Fehler wird wahrscheinlicher, wenn α steigt.
	Das $(1-\alpha)\%$ Konfidenzintervall für den Steigungsparameter β_2 besagt, dass der wahre Wert von β_2 mit einer Wahrscheinlichkeit von $(1-\alpha)$ im beschriebenen Intervall liegt.
	Konsistente Schätzer können verzerrt sein.
	Die gemeinsame Dichtefunktion $f(X,Y)$ zweier unabhängiger Zufallsvariablen X und Y unterscheidet sich von der gemeinsamen Dichtefunktion zweier korrelierter Zufallsvariablen.
	Der White Schätzer korrigiert das Problem stochastischer Fehlertermvarianzen.
	Der LM Test auf Autokorrelation erster Ordnung in den Störtermen besteht aus einem Signifikanztest für den geschätzten Koeffizienten des um eine Periode verzögerten Fehlerterms, der zusätzlich ins ursprüngliche Modell aufgenommen wird.
	Unter Heteroskedastie können bessere Vorhersagen gemacht werden als ohne Heteroskedastie.
	Kategoriale erklärende Variablen werden typischerweise mit Bezug auf eine Referenzgruppe interpretiert.
	Multikollinearitätsprobleme lassen sich über eine Erhöhung der Beobachtungszahl reduzieren.
	Um eine saisonbereinigte Zeitreihe zu erstellen, werden lineare, exponentielle oder logistische Saisonmodelle genutzt.
	Ein hoher Gini-Koeffizient lässt auf eine gleichmäßige Verteilung schließen.

Aufgabe 7:

(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.	Die Präzision der Schätzung eines Steigungsparameters ist umso höher,
	<input type="checkbox"/> je weniger Beobachtungen vorliegen.
	<input type="checkbox"/> je mehr Parameter geschätzt werden.
	<input type="checkbox"/> je größer die Streuung der erklärenden Variable.
2.	Der Typ II Fehler
	<input type="checkbox"/> tritt auf, wenn die Nullhypothese verworfen wird, obwohl sie zutrifft.
	<input type="checkbox"/> ist umso wahrscheinlicher, je größer die Stichprobe ist.
	<input type="checkbox"/> wird unwahrscheinlicher, wenn der Typ I Fehler wahrscheinlicher wird.

3.	Der Two Stage Least Squares Schätzer
	<input type="checkbox"/> schätzt das gleiche lineare Regressionsmodell zweimal.
	<input type="checkbox"/> nutzt vorhergesagte Werte auf der zweiten Stufe.
4.	Eine Division der erklärenden Variable x_k durch den Faktor a führt zu
	<input type="checkbox"/> einem um den Faktor a reduzierten Parameterschätzwert für β_k .
	<input type="checkbox"/> einem um den Faktor a erhöhten Parameterschätzwert für β_k .
5.	Die Normalgleichungen des KQ-Schätzers
	<input type="checkbox"/> ergeben sich bei Minimierung der Zielfunktion.
	<input type="checkbox"/> sind über das Method of Moments Verfahren nicht herleitbar.
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.
7.	Ein RESET Test mit quadrierten und kubischen vorhergesagten Werten (\hat{y}^2 und \hat{y}^3) der abhängigen Variable ergibt eine Teststatistik von 4,8 mit einem p-Wert von 0,067. Dies bedeutet:
	<input type="checkbox"/> Das Modell sollte in logarithmierter Form geschätzt werden.
	<input type="checkbox"/> Am Signifikanzniveau von 10% wird H_0 nicht verworfen.
8.	Bei gegen unendlich konvergierender Stichprobengröße
	<input type="checkbox"/> konvergiert der Intervallschätzer der Steigungsparameter gegen das Signifikanzniveau.
	<input type="checkbox"/> konvergiert die Varianz des KQ-Schätzers gegen Null.
9.	Die Varianz von in erster Ordnung autokorrelierten Störtermen (AR(1))
	<input type="checkbox"/> ist immer heteroskedastisch.
	<input type="checkbox"/> ist nur definiert für $\rho \neq 1$.
10.	Punktschätzer sind
	<input type="checkbox"/> informativer als Intervallschätzer.
	<input type="checkbox"/> umso verlässlicher, je kleiner die geschätzte Fehlervarianz $\hat{\sigma}^2$.