



Einsatz von SPSS-Trends bei der Regressionsanalyse mit seriell abhängigen Daten

1 EINLEITUNG	3
2 OLS-REGRESSION TROTZ AUTOKORRELIERTER RESIDUEN	7
2.1 Schätz- und Testergebnisse	7
2.2 Untersuchung der OLS Residuen auf serielle Abhängigkeit	8
2.2.1 Der Durbin-Watson-Test	9
2.2.2 Die (partielle) Autokorrelationsfunktion	11
3 REGRESSION MIT INTEGRIERTEM ARIMA(1,0,0)-MODELL FÜR DIE FEHLER	13
3.1 Das erweiterte Modell	13
3.2 Schätzergebnisse	13
3.3 Prüfung der Modellgültigkeit	14
3.4 Zur Modellgüte	17

4 REGRESSION MIT INTEGRIERTEM ARIMA(P,D,Q)-MODELL FÜR DIE FEHLER	19
5 SONSTIGE HINWEISE ZU REGRESSIONSANALYSEN MIT ZEITREIHENDATEN	20
5.1 Warnung vor Regressionsanalysen mit trendbelasteten Zeitreihen	20
5.2 Zeitreihen mit Saisonfigur	20
6 LITERATUR	21
7 STICHWORTVERZEICHNIS	22

Herausgeber: Universitäts-Rechenzentrum Trier
 Universitätsring 15
 D-54286 Trier
 Tel.: (0651) 201-3417, Fax.: (0651) 3921

Leiter: Prof. Dr.-Ing. Manfred Paul

Autor: Bernhard Baltes-Götz
 Mail: baltes@uni-trier.de

Copyright © 1997; URT

Vorwort

Das Manuskript beschreibt die im SPSS-Zusatzmodul **Trends** verfügbaren Optionen, die serielle Abhängigkeit der Fälle einer Regressionsanalyse durch ein ARIMA-Modell für die Residuen zu berücksichtigen, wobei häufig bereits ein AR(1)-Modell genügt.

Als Software kommt SPSS 6.1 für Windows zum Einsatz, jedoch können praktisch alle vorgestellten Verfahren auch mit jüngeren SPSS-Versionen unter Windows, MacOS oder Linux realisiert werden.

Das Manuskript ist als PDF-Dokument zusammen mit den im Kurs benutzten Dateien auf dem Webserver der Universität Trier von der Startseite (<http://www.uni-trier.de/>) ausgehend folgendermaßen zu finden:

[Rechenzentrum](#) > [Studierende](#) > [EDV-Dokumentationen](#) > [Statistik](#) >
[Einsatz von SPSS-Trends bei der Regression mit seriell abhängigen Daten](#)

Hinweise auf Unzulänglichkeiten im Manuskript werden mit Dank entgegen genommen

1 Einleitung

Die Inferenzstatistik in der üblichen OLS(=ordinary least squares)-Regressionsanalyse setzt u.a. voraus, daß die Modellresiduen unabhängig und identisch normalverteilt sind (mit Mittelwert Null und homogener Varianz). In einer Untersuchung mit K Regressoren und T Fällen (Beobachtungseinheiten), kann man das Modell der multiplen OLS-Regression z.B. folgendermaßen aufschreiben:

$$Y_t = b_0 + \sum_{k=1}^K b_k X_{k_t} + U_t, \quad t = 1, \dots, T \quad (1)$$

Im einzelnen bedeuten:

- X_{k_t} Der Wert des k -ten Regressors bei der Beobachtung t
- U_t Der Anteil von Y_t , der durch die Regressoren nicht erklärt werden kann
 U_t muß ein **unkorrelierter** (und wegen der Normalverteilung damit unabhängiger), zufälliger Fehler mit Erwartungswert Null und homogener Varianz σ^2 sein.

Wenn die Beobachtungen auf einer Zufallsstichprobe aus einer quasi unendlichen Population beruhen, kann die Unabhängigkeit der Residuen als gesichert gelten. Häufig möchte man jedoch regressive Beziehungen anhand von (zeitlich oder räumlich) geordneten Beobachtungen untersuchen. Dabei hat jede Beobachtung einen eindeutigen Vorgänger, von dem sie durchaus über die im Regressionsmodell enthaltenen Variableneffekte hinaus beeinflusst sein könnte.

Als Beispiel wollen wir einen von Durbin & Watson (1951) berichteten Datensatz betrachten, der drei logarithmisch transformierte Zeitreihen mit jährlich in England vorgenommenen Messungen aus dem Beobachtungszeitraum von 1870 bis 1938 enthält:

- Alkoholverbrauch (Variablenname: CONSUMP)
- Pro-Kopf-Einkommen (Variablenname: INCOME)
- Inflationsbereinigter Preisindex (Variablenname: PRICE)

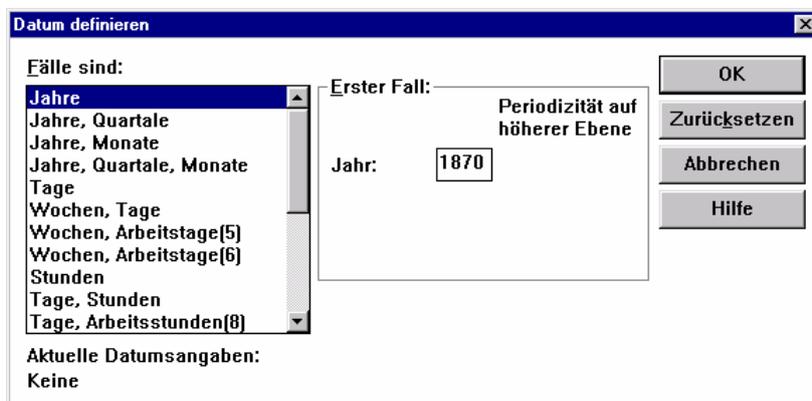
Sie finden den Datensatz in der Datei **CHAP9.SAV** an der im Vorwort vereinbarten Stelle.

Es soll versucht werden, den Alkoholverbrauch durch die beiden anderen Variablen zu erklären. Dabei ist zu befürchten, daß zahlreiche im Modell nicht berücksichtigte Einflüsse auf den Alkoholverbrauch in benachbarten Jahren relativ ähnlich ausgeprägt waren, so daß zwischen zeitlich benachbarten Residuen Korrelationen bestehen, die im OLS-Regressionsmodell (1) verboten sind.

Zunächst wollen wir uns einen optischen Eindruck vom Verlauf der drei Zeitreihen verschaffen. Dazu öffnen wir die oben angegebene SPSS-Datendatei und vereinbaren nach dem Menübefehl

Daten > Datum definieren...

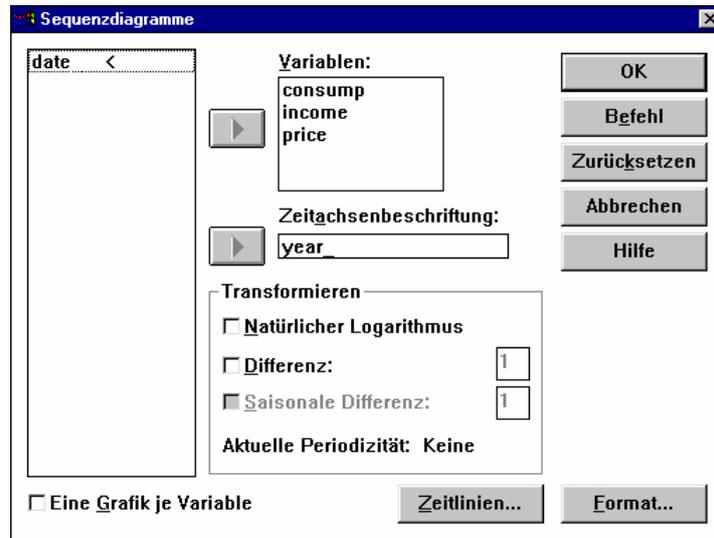
mit folgender Dialogbox eine neue Variable mit den Jahreszahlen ab 1870:



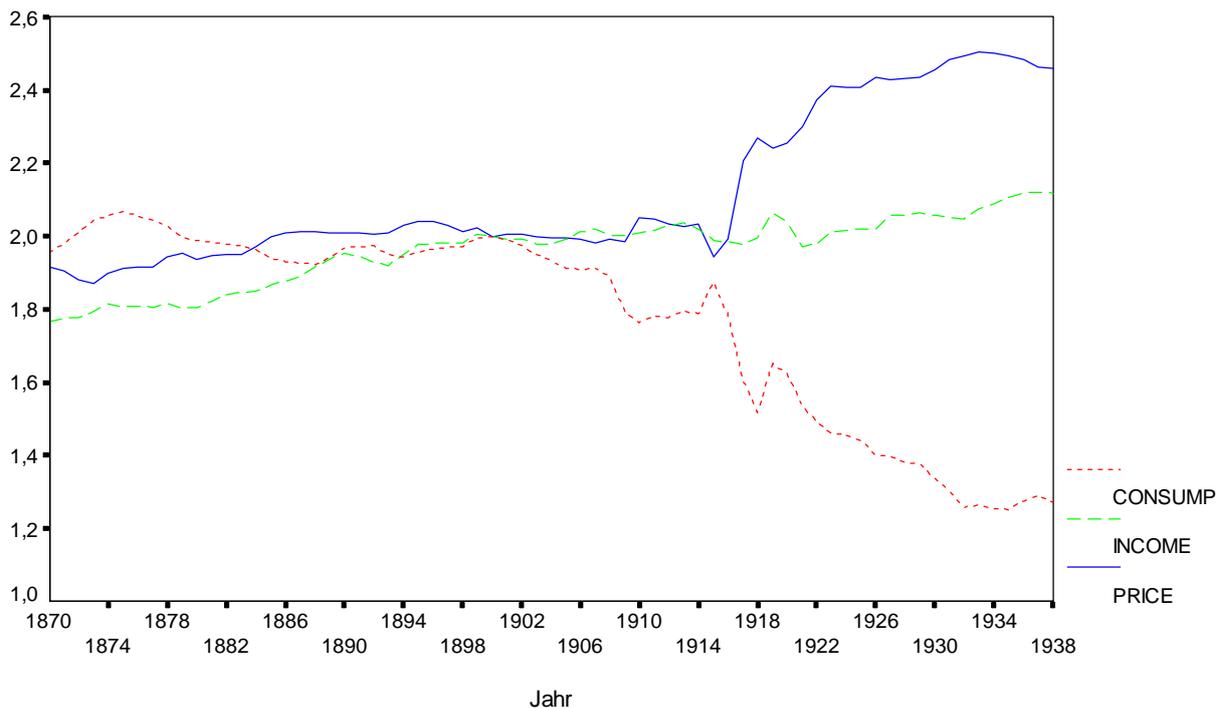
Danach fordern über den Menübefehl

Grafik > Sequenz

mit der Dialogbox



die folgende Abbildung an:



Abgesehen von den Jahren während des ersten Weltkriegs zeigt die Zeitreihe CONSUMP einen relativ glatten Verlauf, d.h. jede Beobachtung liegt ziemlich nahe bei ihrem Vorgänger, woraus wir direkt auf eine hohe positive Autokorrelation (erster Ordnung, s.u.) schließen können. Diese ist an sich aber noch nicht unverträglich mit dem klassischen Regressionsmodell, das die Unkorreliertheit der *Residuen* aus der geplanten Regressionsgleichung mit den Prädiktoren INCOME und PRICE voraussetzt.

Sie werden in diesem Kurs lernen, bei einer Regressionsanalyse mit solchen Zeitreihendaten die kritische Voraussetzung unkorrelierter Residuen zu prüfen und gegebenenfalls die OLS-Regression durch ein adäquateres Verfahren zu ersetzen. Dabei wird die Regressionsmethodologie durch Komponenten aus der Zeitreihenanalyse erweitert. Kenntnisse dieser nicht ganz trivialen Methodenfamilie werden zwar in diesem Kurs nicht vorausgesetzt, sind aber recht nützlich, weil unweigerlich einige Begriffe der Zeitreihenanalyse auftreten (z.B. Autokorrelationsfunktion). Eine kurze Einführung in elementare Begriffe der Zeit-

reihenanalyse und in die allgemeinen Bedienungsmerkmale des zuständigen SPSS-Moduls Trends, zu dem auch die in diesem Kurs vorgestellten Regressionsverfahren für seriell abhängige Daten gehören, finden Sie im URT-Umdruck „Zeitreihenanalyse mit SPSS-Trends“.

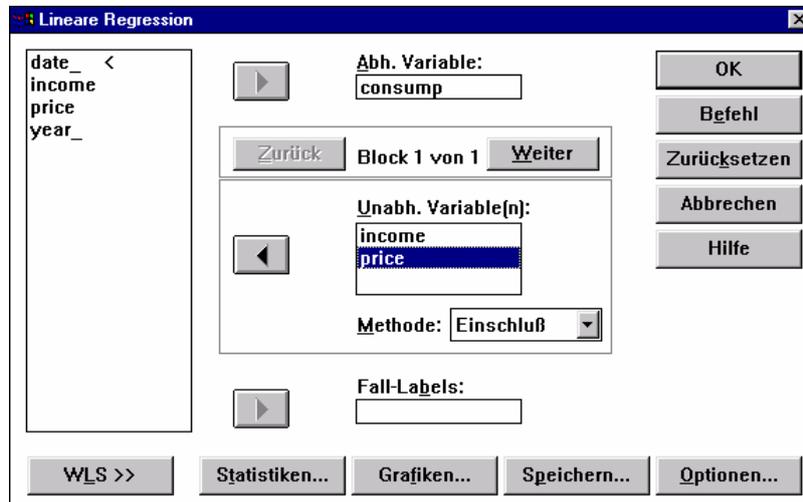
2 OLS-Regression trotz autokorrelierter Residuen

2.1 Schätz- und Testergebnisse

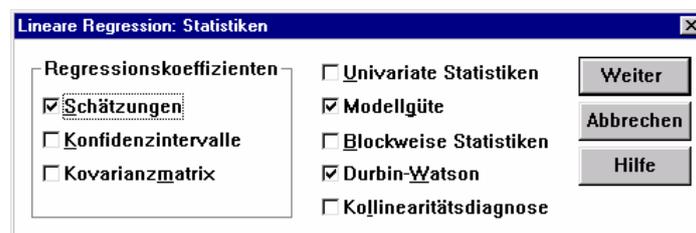
Wir fordern zunächst mit

Statistik > Regression > Linear...

die Dialogbox zur multiplen OLS-Regression an, die wir folgendermaßen ausfüllen:



In der "**Statistiken...**"-Subdialogbox fordern wir über die Standardausgabe hinaus den Durbin-Watson-Koeffizienten an:



und in der "**Speichern...**"-Subdialogbox veranlassen wir, daß SPSS die unstandardisierten Residuen als neue Variable sichert. Wir erhalten u.a. die folgenden Ergebnisse:

Multiple R	,97766				
R Square	,95581				
Adjusted R Square	,95447				
Standard Error	,05786				
Analysis of Variance					
	DF	Sum of Squares	Mean Square		
Regression	2	4,77917	2,38959		
Residual	66	,22095	,00335		
F =	713,78788	Signif F =	,0000		
----- Variables in the Equation -----					
Variable	B	SE B	Beta	T	Sig T
INCOME	-,120141	,108436	-,042713	-1,108	,2719
PRICE	-1,227648	,050052	-,945573	-24,527	,0000
(Constant)	4,606734	,152035		30,301	,0000
Durbin-Watson Test =	,24878				

Das adjustierte R^2 ist mit 0,95 sehr hoch, und die Global-Nullhypothese:

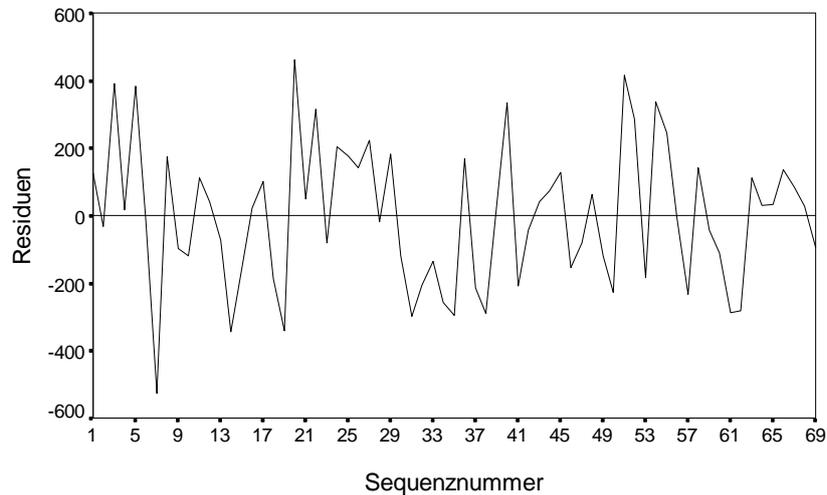
$$H_0: b_{INCOME} = 0 \text{ und } b_{PRICE} = 0$$

wird im F-Test deutlich verworfen. Weniger erfreulich ist die Tatsache, daß für die Einkommensvariable kein signifikanter Einfluß ermittelt wurde. Allerdings sind vor einer Ergebnisinterpretation zunächst die

Residuen auf Indizien für verletzte Modellannahmen zu untersuchen. Da wir Zeitreihendaten analysiert haben, müssen wir vor allem die Frage einer möglichen seriellen Abhängigkeit der Residuen klären.

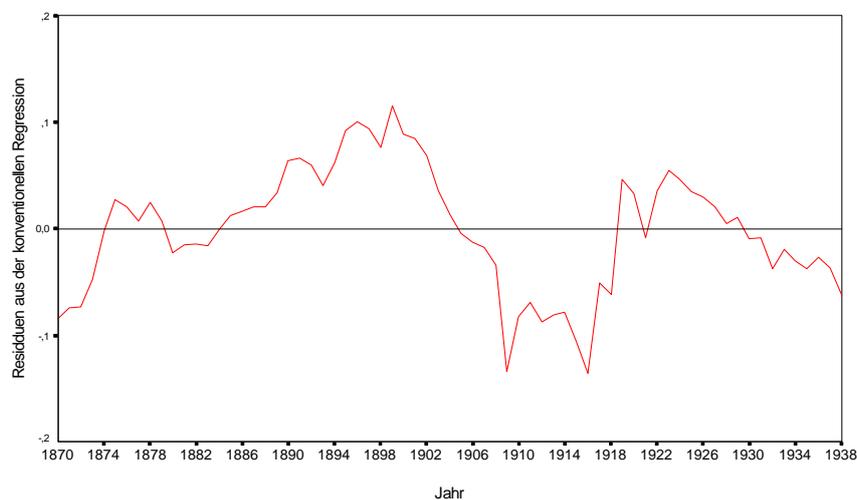
2.2 Untersuchung der OLS Residuen auf serielle Abhängigkeit

Die vom OLS-Regressionsmodell geforderte Unabhängigkeit der Residuen impliziert, daß die Zeitreihe der Residuen ein weißes Rauschen ist. Ein weißes Rauschen "zuckt wild" um seinen Mittelwert, wie es z.B. in der folgenden Abbildung mit den Residuen aus einer Regression mit Querschnittsdaten zu sehen ist:



Bei Querschnittsdaten ist in der Regel aus versuchsplanerischen Gründen eine serielle Abhängigkeit auszuschließen. Einige überbesorgte Zeitgenossen, denen der Durbin-Watson-Koeffizient bei *jedem* Datensatz den Schlaf raubt, können also beruhigt werden.

Bei den Residuen zu unserer OLS-Regressionsgleichung für die Durbin-Watson-Daten zeigt sich (bei gleicher Stichprobengröße) ein deutlich anderes Bild:



Dieses Sequenzdiagramm wurde über den Befehl **Grafik > Sequenz** unter Verwendung der Variablen RES_1 erstellt, die wir im Subdialog **Speichern** zur linearen Regressionsanalyse angefordert haben. Unsere Residualzeitreihe schlängelt sich gemächlich um die Nullage, d.h. jeder Wert liegt relativ nahe bei seinem Vorgänger. Er kann folglich durch seinen Vorgänger gut vorhergesagt werden, und wir erwarten daher einen hohen Wert für die sogenannte **Autokorrelation erster Ordnung** zwischen der Residualzeitreihe und ihrer um einen Zeittakt verschobenen Variante.

Diese Größe wird aus den Stichprobendaten folgendermaßen geschätzt:

$$\hat{\rho} := \frac{\sum_{t=2}^T (u_t - \bar{u})(u_{t-1} - \bar{u})}{\sum_{t=1}^T (u_t - \bar{u})^2}, \quad \text{mit} \quad \bar{u} := \frac{1}{T} \sum_{t=1}^T u_t$$

Da wir in diesem Kurs im wesentlichen nur die Autokorrelation *erster* Ordnung betrachten, bezeichnen wir diese mit einem einfachen ρ ohne Index für die Ordnung der Autokorrelation. Die Schätzformel für die erste Autokorrelation ähnelt offenbar sehr stark der Stichprobenformel für gewöhnliche Korrelationen.

Weil die OLS-Residuen den Mittelwert Null haben, vereinfacht sich der obige Ausdruck zu:

$$\hat{\rho} = \frac{\sum_{t=2}^T u_t u_{t-1}}{\sum_{t=1}^T u_t^2} \quad (2)$$

Bei der Analyse zeitlich geordneter Daten sind häufig positive Residual-Autokorrelation erster Ordnung zu beobachten. Negative Residual-Autokorrelationen erster Ordnung sind hier eher selten.

2.2.1 Der Durbin-Watson-Test

Bei einer Anwendung der OLS-Regressionsanalyse auf zeitlich geordneten Daten muß wenigstens geprüft werden, ob die Autokorrelation 1. Ordnung gleich Null ist. Dies ist eine notwendige, jedoch keine hinreichende Anwendungsvoraussetzung.

Weil in der Regel eine positive Autokorrelation zu befürchten ist, sollte der Signifikanztest für die Autokorrelation 1. Ordnung einseitig angelegt werden:

H_0 : Die Autokorrelation 1. Ordnung der Modellresiduen ist *nicht* positiv.
versus

H_1 : Die Autokorrelation 1. Ordnung der Modellresiduen ist positiv.

Ein für dieses Testproblem geeignetes Verfahren wurde von Durbin & Watson (1951) entwickelt. Die DW-Prüfstatistik für eine Stichprobe der Größe T ist folgendermaßen definiert (vgl. Hartung 1989):

$$DW := \frac{\sum_{t=2}^T (u_t - u_{t-1})^2}{\sum_{t=1}^T u_t^2} = \frac{\frac{1}{T} \sum_{t=2}^T (u_t - u_{t-1})^2}{\frac{1}{T} \sum_{t=1}^T u_t^2} \quad (3)$$

Dabei sind mit u_t , $t=1, \dots, T$, wie im gesamten Manuskript die üblichen Kleinst-Quadrat-Residuen gemeint. Offenbar wird die Statistik um so größer, je mehr sich u_t von seinem Vorgänger u_{t-1} unterscheidet. Im Zähler des rechten Quotienten steht annähernd die mittlere quadrierte Differenz. Diese wird durch den Nenner normiert, der die Stichprobenvarianz enthält, weil die Residuen den Mittelwert Null haben.

Durch Umformungen des DW-Koeffizienten wird seine enge Beziehung zur Stichproben-Autokorrelation erster Ordnung deutlich (vgl. Formel (2)):

$$\begin{aligned}
 \frac{\sum_{t=2}^T (u_t - u_{t-1})^2}{\sum_{t=1}^T u_t^2} &= \frac{\sum_{t=2}^T u_t^2 - 2 \sum_{t=2}^T u_t u_{t-1} + \sum_{t=2}^T u_{t-1}^2}{\sum_{t=1}^T u_t^2} \\
 &= \frac{2 \sum_{t=2}^{T-1} u_t^2 + u_1^2 + u_N^2}{\sum_{t=1}^T u_t^2} - 2 \frac{\sum_{t=2}^T u_t u_{t-1}}{\sum_{t=1}^T u_t^2} \\
 &\approx 2 - 2\hat{\rho} = 2(1 - \hat{\rho})
 \end{aligned} \tag{4}$$

Damit ist klar, daß der Durbin-Watson-Koeffizient zwischen Null und Vier variiert und in Abhängigkeit von der ersten Autokorrelation folgendes Verhalten zeigt:

$$\begin{aligned}
 \hat{\rho} \rightarrow 1 &\Rightarrow DW \rightarrow 0 \\
 \hat{\rho} \rightarrow 0 &\Rightarrow DW \rightarrow 2 \\
 \hat{\rho} \rightarrow -1 &\Rightarrow DW \rightarrow 4
 \end{aligned}$$

In unserer Stichprobe haben wir einen DW-Wert von 0,24878 erhalten, der also für eine deutlich positive Autokorrelation 1. Ordnung spricht.

Zur Durchführung des Durbin-Watson-Tests müssen leider Tabellen herangezogen werden, die z.B. im SPSS-Handbuch zur Zeitreihenanalyse (SPSS 1994, Anhang A) abgedruckt sind. In der oben beschriebenen Testsituation haben wir normalerweise die empirisch ermittelte Teststatistik mit einem kritischen Wert zu vergleichen. Im Fall des DW-Tests kann der kritische Wert leider nicht exakt bestimmt werden. Es ist lediglich möglich, eine untere (dL) und eine obere Schranke (dU) anzugeben. Diese werden in Abhängigkeit von der Stichprobengröße T , der Prädiktorenzahl k (gerechnet *ohne* den konstanten Term) und dem gewünschten Signifikanzniveau α aus der zuständigen Tabelle ermittelt. Dann sind bei dem oben angegebenen einseitigen Test (H_0 versus H_1) die folgenden Entscheidungsregeln anzuwenden:

DW-Wert	Testentscheidung
$DW \leq dL$	H_0 ablehnen, d.h. positive Autokorrelation
$dL < DW < dU$	keine Entscheidung möglich
$dU \leq DW$	H_0 beibehalten

Da unsere Stichprobengröße ($T = 69$) in der Tabelle auf Seite 331 des SPSS-Trends-Handbuchs fehlt, verwenden wir den nächst kleineren Stichprobenumfang und lesen bei $T = 65$, $k = 2$ und $\alpha = 0,05$ ab:

$$dL = 1,536 \quad \text{und} \quad dU = 1,662$$

Aufgrund unseres empirischen DW-Wertes von 0,24878 wird die Nullhypothese deutlich verworfen. Wir gehen also davon aus, daß in den Modellresiduen der OLS-Regression eine starke Autokorrelation erster Ordnung vorliegt, so daß die Voraussetzungen dieses Verfahrens verletzt sind. In dieser Situation sind u.a. die Signifikanzbeurteilungen zu den Regressionskoeffizienten verfälscht.

Bei dem alternativen einseitigen Testproblem:

$$\begin{aligned}
 H'_0 &: \text{Die Autokorrelation 1. Ordnung der Modellresiduen ist } \textit{nicht} \text{ negativ.} \\
 &\text{versus} \\
 H'_1 &: \text{Die Autokorrelation 1. Ordnung der Modellresiduen ist negativ.}
 \end{aligned}$$

sowie bei dem zweiseitigen Testproblem:

H_0 : Die Autokorrelation 1. Ordnung der Modellresiduen ist gleich Null.

versus

H_1 : Die Autokorrelation 1. Ordnung der Modellresiduen ist ungleich Null.

gelten andere Entscheidungsregeln (siehe Hartung 1989, S. 740f).

Weitere Hinweise zum DW-Test (vgl. SPSS 1994, S. 327f):

- Bei Modellen ohne konstantem Term sind andere DW-Tabellen heranzuziehen, die im SPSS-Handbuch zur Zeitreihenanalyse (SPSS 1994, Anhang A) ebenfalls abgedruckt sind.
- Der Test ist nicht anwendbar, wenn zeitversetzte Varianten der abhängigen Variablen als Regressoren verwendet werden.

2.2.2 Die (partielle) Autokorrelationsfunktion

Nachdem wir nun wissen, daß die Residuen aus der OLS-Regressionsanalyse der Durbin-Watson-Daten autokorreliert sind, haben wir u.a. die folgenden Möglichkeiten:

- Wir erweitern unser Modell um einen autoregressiven Fehlerprozeß erster Ordnung, weil sich die Autokorrelation erster Ordnung im Durbin-Watson-Test als signifikant erwiesen hat. In der Tat ist die Annahme recht plausibel, daß ein Residuum in erster Linie von seinem unmittelbaren Vorgänger beeinflusst ist. SPSS bietet daher für diesen wichtigen Spezialfall eine eigene Prozedur, die recht zu bedienen ist.
- Wir sind durch den Durbin-Watson-Befund hellhörig geworden und wollen nun genaue Informationen über den Typ der Abhängigkeit in den Residuen, damit wir diese in einem erweiterten Modell geeignet berücksichtigen können. Bisher wissen wir nur von der Autokorrelation erster Ordnung, weil der Durbin-Watson-Koeffizient für andere Abhängigkeiten blind ist. Üblicherweise beschränkt man sich bei der Suche nach einer Erklärung für die Abhängigkeit in den Residuen auf die ARIMA-Modellfamilie, die unter anderem den autoregressiven Prozeß erster Ordnung als Spezialfall enthält, außerdem aber sehr viele andere Zeitreihenmodelle, die für den signifikanten Durbin-Watson-Test gesorgt haben könnten.

Als hochmotivierte Wissenschaftler(innen) beginnen wir natürlich mit der Suche nach einem geeigneten ARIMA-Modell für die Residuen, wobei die **Autokorrelationsfunktion** (ACF) sowie die **partielle Autokorrelationsfunktion** (PACF) unsere entscheidenden Hilfsmittel sind (siehe z.B. Pindyck & Rubinfeld 1990; Schlittgen & Streitberg 1989; SPSS 1994, Abschnitt 6). Die empirische **Autokorrelationsfunktion** gibt für jede natürliche Zahl τ die Korrelation einer Zeitreihe mit der um τ Takte verschobenen Version an. Für $\tau = 1$ erhält man also gerade die Autokorrelation erster Ordnung, die wir oben im Zusammenhang mit dem Durbin-Watson-Test schon kennengelernt haben. Analog dazu gibt die geschätzte **partielle Autokorrelationsfunktion** für jede natürliche Zahl τ die partielle Korrelation einer Zeitreihe mit der um τ Takte verschobenen Version bei Kontrolle der dazwischenliegenden Verschiebungsvarianten an.

Um für die Residualvariable RES_1 aus der OLS-Regression die ACF und die PACF anzufordern zu können, wählen wir den Menübefehl:

Grafik > Zeitreihen > Autokorrelationen...

und erhalten folgendes Bild:

```

Autocorrelations:  RES_1  Residual

   Auto- Stand.
Lag  Corr.  Err. -1  -.75  -.5  -.25  0  .25  .5  .75  1  Box-Ljung  Prob.
-----+-----+-----+-----+-----+-----+-----+-----+-----+
 1  ,851  ,118          .          .          .          .          .          .          .          .          .          .
 2  ,738  ,117          .          .          .          .          .          .          .          .          .          .
 3  ,629  ,116          .          .          .          .          .          .          .          .          .          .
 4  ,500  ,115          .          .          .          .          .          .          .          .          .          .
 5  ,392  ,114          .          .          .          .          .          .          .          .          .          .
 6  ,284  ,113          .          .          .          .          .          .          .          .          .          .
 7  ,149  ,112          .          .          .          .          .          .          .          .          .          .
 8  ,005  ,112          .          .          .          .          .          .          .          .          .          .
 9  -,081  ,111          .          .          .          .          .          .          .          .          .          .
10  -,199  ,110          .          .          .          .          .          .          .          .          .          .
11  -,255  ,109          .          .          .          .          .          .          .          .          .          .
12  -,308  ,108          .          .          .          .          .          .          .          .          .          .
13  -,395  ,107          .          .          .          .          .          .          .          .          .          .
14  -,432  ,106          .          .          .          .          .          .          .          .          .          .
15  -,431  ,105          .          .          .          .          .          .          .          .          .          .
16  -,389  ,104          .          .          .          .          .          .          .          .          .          .

Plot Symbols:  Autocorrelations *  Two Standard Error Limits .

Total cases:  69  Computable first lags:  68
    
```

Die Autokorrelationsfunktion zeigt ein exponentiell ausschwingendes Verhalten, während die partielle Autokorrelationsfunktion nach dem ersten Lag abbricht:

```

Partial Autocorrelations:  RES_1  Residual

   Pr-Aut- Stand.
Lag  Corr.  Err. -1  -.75  -.5  -.25  0  .25  .5  .75  1
-----+-----+-----+-----+-----+-----+-----+-----+
 1  ,851  ,120          .          .          .          .          .          .          .          .          .
 2  ,049  ,120          .          .          .          .          .          .          .          .          .
 3  -,035  ,120          .          .          .          .          .          .          .          .          .
 4  -,133  ,120          .          .          .          .          .          .          .          .          .
 5  -,027  ,120          .          .          .          .          .          .          .          .          .
 6  -,064  ,120          .          .          .          .          .          .          .          .          .
 7  -,177  ,120          .          .          .          .          .          .          .          .          .
 8  -,180  ,120          .          .          .          .          .          .          .          .          .
 9  ,074  ,120          .          .          .          .          .          .          .          .          .
10  -,174  ,120          .          .          .          .          .          .          .          .          .
11  ,089  ,120          .          .          .          .          .          .          .          .          .
12  -,076  ,120          .          .          .          .          .          .          .          .          .
13  -,189  ,120          .          .          .          .          .          .          .          .          .
14  ,020  ,120          .          .          .          .          .          .          .          .          .
15  ,059  ,120          .          .          .          .          .          .          .          .          .
16  ,131  ,120          .          .          .          .          .          .          .          .          .

Plot Symbols:  Autocorrelations *  Two Standard Error Limits .

Total cases:  69  Computable first lags:  68
    
```

Dieses Bild paßt relativ gut zu einem autoregressiven Prozeß erster Ordnung, also zu einem ARIMA (1,0,0)-Prozeß. Allerdings bin ich mir nicht ganz sicher, ob man aus den Residuen aus einem offenbar falschen Modell hinreichend sicher auf das Modell für die „echten“ Residuen schließen kann.

3 Regression mit integriertem ARIMA(1,0,0)-Modell für die Fehler

3.1 Das erweiterte Modell

Wir erweitern unser Modell nun so, daß die Autokorrelation erster Ordnung in den Residuen explizit berücksichtigt und somit eine verzerrte Schätzung der Regressionskoeffizienten verhindert wird. Dabei wird die abhängige Variable Y_t zum Zeitpunkt t folgendermaßen erklärt:

$$Y_t = b_0 + \sum_{i=1}^k b_i X_{i,t} + U_t \quad (5)$$

$$t = 1, \dots, T$$

$$U_t = \rho U_{t-1} + \varepsilon_t \quad (6)$$

Hier handelt es sich offenbar um ein normales multiples Regressionsmodell (5), das um einen ARIMA(1,0,0)-Fehlerprozeß (6) erweitert wurde. Im einzelnen bedeuten:

- $X_{k,t}$ Der Wert des k -ten Regressors bei der Beobachtung t
- U_t Der Anteil von Y_t , der durch die "Fremdregressoren" $X_{k,t}$ nicht erklärt werden kann
- ρ Man kann zeigen, daß der autoregressive Parameter erster Ordnung zum Fehlerprozeß gerade mit der ersten Autokorrelation der U_t -Zeitreihe identisch ist.
- ε_t Ein unkorrelierter, zufälliger Fehler, der normalverteilt ist mit Erwartungswert Null und homogener Varianz σ^2

Für ε_t (nicht aber für U_t) werden also die klassischen Annahmen der OLS-Regressionsanalyse gemacht.

3.2 Schätzergebnisse

SPSS beherrscht zum obigen Modell die folgenden Berechnungsmethoden:

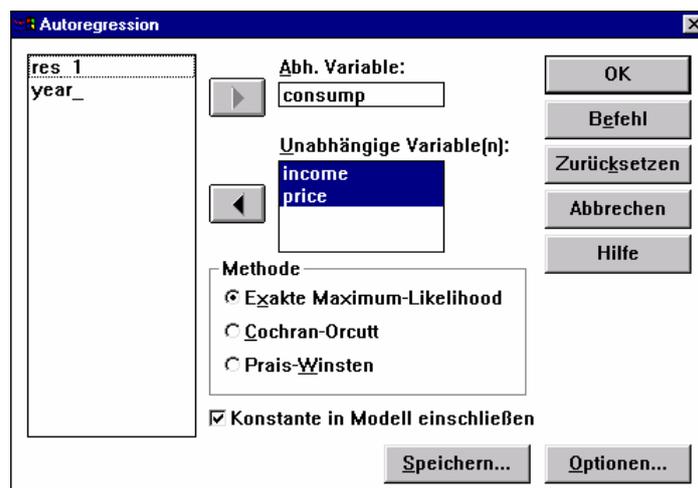
- Prais-Winsten (GLS)
- Cochrane-Orcutt
- Maximum Likelihood

Das Maximum Likelihood - Verfahren benötigt zwar die meiste Rechenzeit, liefert aber auch die genauesten Ergebnisse und toleriert außerdem fehlende Daten. Die beiden anderen Verfahren schätzen außerdem nicht das obige Modell (5,6), sondern transformieren die Regressionsgleichung, um die Autokorrelation der Residuen zu beseitigen.

Mit dem Menübefehl:

Statistik > Zeitreihen > Autoregression...

wird die folgende Dialogbox zur Spezifikation der Regressionsanalyse mit ARIMA(1,0,0)-Residualmodell angefordert:



Wir erhalten u.a. folgende Ergebnisse:

FINAL PARAMETERS:

Number of residuals 69
 Standard error ,02266171
 Log likelihood 163,29768
 AIC -318,59536
 SBC -309,65893

Analysis of Variance:

	DF	Adj. Sum of Squares	Residual Variance
Residuals	65	,03554197	,00051355

Variables in the Model:

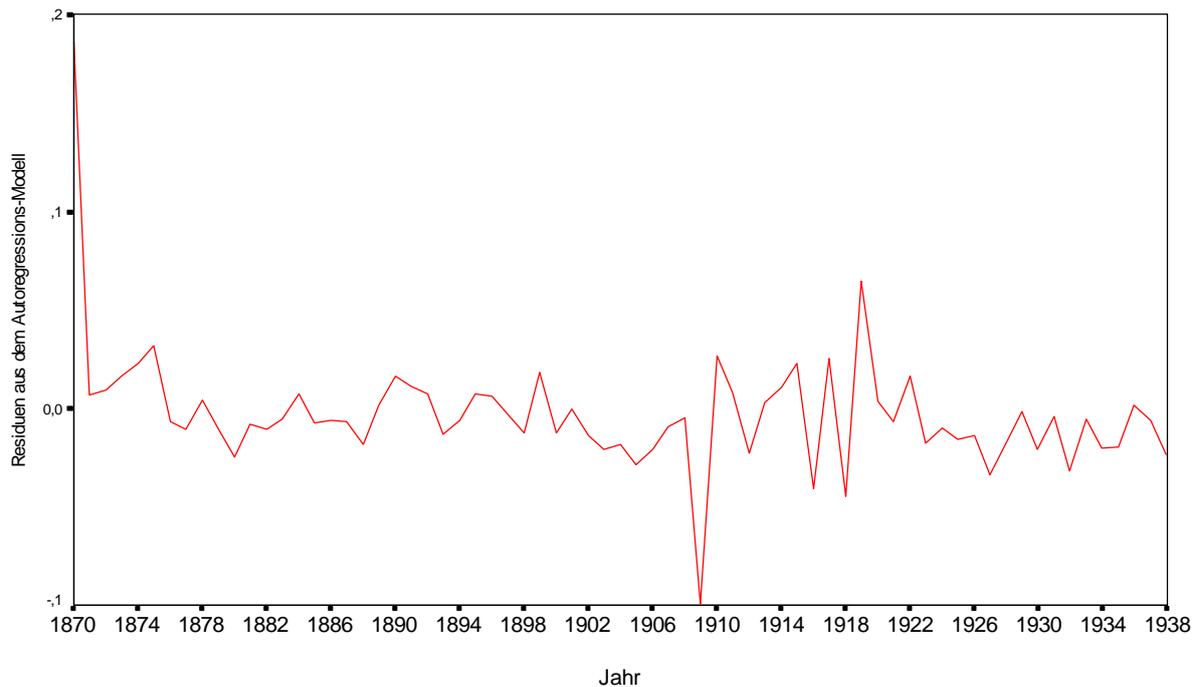
	B	SEB	T-RATIO	APPROX. PROB.
AR1	,9933830	,01167538	85,083568	,00000000
INCOME	,6234051	,14689261	4,243951	,00007119
PRICE	-,9280366	,07816780	-11,872365	,00000000
CONSTANT	2,4485054	,37391717	6,548256	,00000000

Unterschiede zu den Ergebnissen der OLS-Regression ergeben sich vor allem beim Regressor INCOME: Während die OLS-Regression hier ein (nicht-signifikantes) negatives Regressionsgewicht liefert, erhalten wir bei Berücksichtigung der Autokorrelation in den Residuen ein signifikant positives Gewicht (APPROX. PROB < 0,001).

Der Signifikanztest zum Parameter des autoregressiven Modells (6) entscheidet sich ebenfalls deutlich gegen seine Nullhypothese ($\rho = 0$). Leider ist die Höhe der Schätzung (= 0,99) leicht besorgniserregend, weil ein ARIMA(1,0,0)- bzw. AR(1)-Prozeß nur stationär ist bei $|\rho| < 1$ (siehe Schlittgen & Streitberg 1989, S. 99). Dieses Kriterium ist also nur knapp erfüllt.

3.3 Prüfung der Modellgültigkeit

Natürlich müssen wir auch die Residuen des verbesserten Modells (laut SPSS-Ausgabe gespeichert in der Variablen ERR_1) einer kritischen Würdigung unterziehen. Der Verlauf zeigt nun deutliche Ähnlichkeit mit einem weißen Rauschen:



Das Residuum zum ersten Jahr fällt nur deshalb so extrem aus dem Rahmen, weil bei seiner Berechnung kein u_{t-1} gemäß Modellgleichung (6) zur Verfügung stand. Außerdem sind noch die Residuen zu den Jahren 1909 und 1919 auffällig. Möglicherweise hat der dramatische Effekt des ersten Weltkriegs auf die Zeitreihen dafür gesorgt, daß in seinem zeitlichen Umfeld eine relativ schlechte Modellpassung auftritt.

Bei einer Analyse der (mit **Grafik > Zeitreihen > Autokorrelationen...** angeforderten) Autokorrelationsfunktion der Residuen aus dem erweiterten Modell ergeben sich keinerlei Hinweise auf verbliebene serielle Abhängigkeiten:

```
Autocorrelations:   ERR_1   Error for CONSUMP from AREG, MOD_2
```

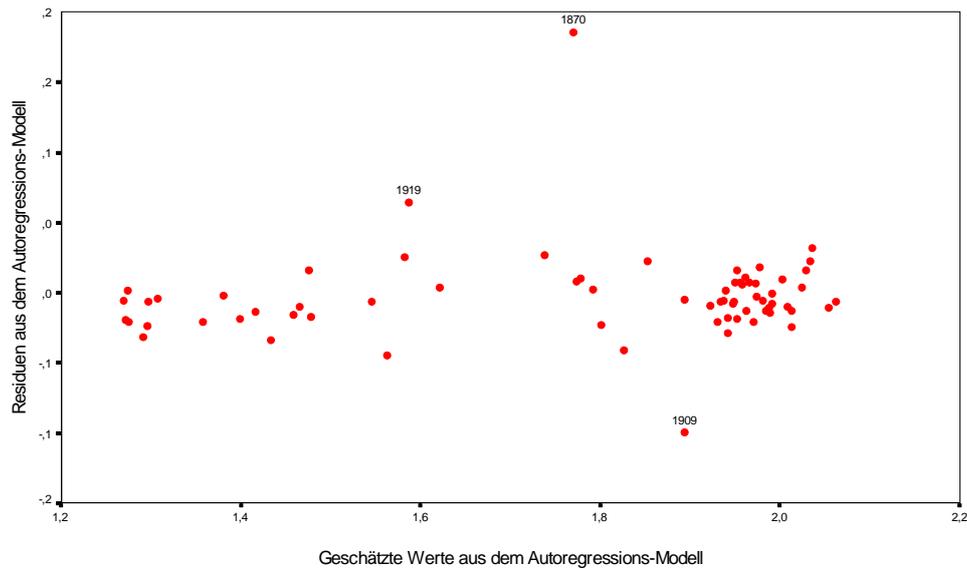
Lag	Auto-Corr.	Stand. Err.	-1	-.75	-.5	-.25	0	.25	.5	.75	1	Box-Ljung	Prob.
1	-,022	,118					*					,034	,854
2	,122	,117					**					1,119	,572
3	,119	,116					**					2,178	,536
4	,126	,115					**					3,377	,497
5	,123	,114					**					4,529	,476
6	-,042	,113					*					4,667	,587
7	,044	,112					*					4,821	,682
8	-,073	,112					*					5,246	,731
9	,086	,111					**					5,856	,754
10	-,216	,110					****					9,719	,466
11	-,012	,109					*					9,730	,555
12	-,017	,108					*					9,755	,637
13	-,134	,107					***					11,333	,583
14	,025	,106					*					11,388	,655
15	-,042	,105					*					11,545	,713
16	,019	,104					*					11,577	,773

Plot Symbols: Autocorrelations * Two Standard Error Limits .

Total cases: 69 Computable first lags: 68

Die ersten 16 empirischen Autokorrelationen liegen innerhalb der Konfidenzschranken, und Der Box-Ljung-Test akzeptiert für jedes Lag τ die Nullhypothese, daß die ersten τ Autokorrelationen gleich Null sind.

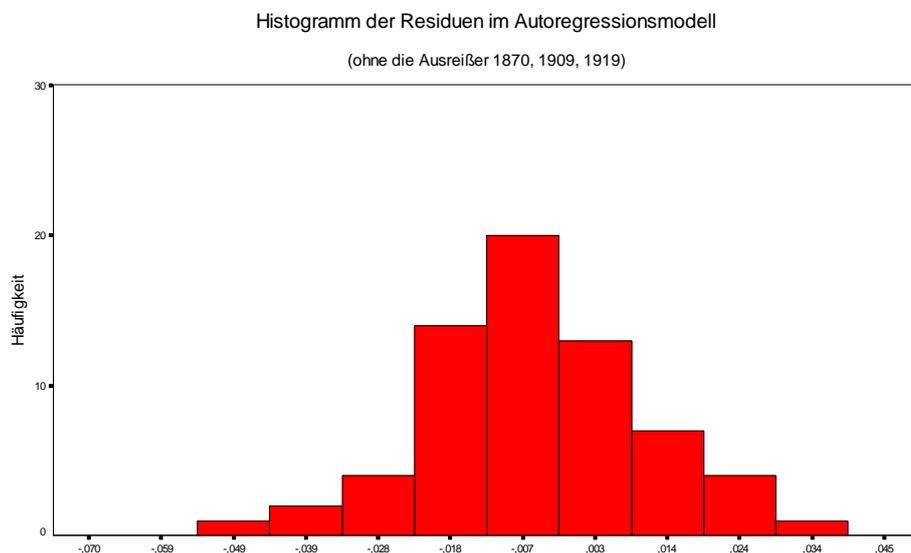
Wir wollen auch den restlichen Annahmen des erweiterten Regressionsmodells noch etwas Aufmerksamkeit schenken. Zur Prüfung der Linearität und der Homoskedastizität sollte man u.a. den Plot der Residuen gegen die geschätzten Werte inspizieren:



Die drei bereits bekannten Ausreißer sind hier mit dem SPSS-Punktauswahlmodus () beschriftet worden. Während der Fall 1870 ein Opfer seiner Randlage und damit kein wirklicher Ausreißer ist, sollte man vielleicht über die beiden anderen nachdenken.

Der Punkteschwarm zeigt einen leichten Aufwärtstrend, was als Argument gegen die Linearitätsannahme gewertet werden könnte. Auch scheint die Varianz mit zunehmendem Schätzwert leicht anzusteigen, was gegen die Homoskedastizität spräche.

Bei einer Untersuchung der Residuen mit der Prozedur zur exploratorischen Datenanalyse (erreichbar über **Statistik > Deskriptive Statistik > Explorative Datenanalyse...**) sorgen die drei Ausreißer dafür, dass der Normalverteilungsanpassungstest Alarm schlägt. Entfernt man die drei, resultiert aber ein ganz passables Bild:



Insgesamt erreicht unser Modell bei den Gültigkeitstests die Note *ausreichend*.

3.4 Zur Modellgüte

Im Unterschied zu den beiden alternativen Schätzverfahren liefert die ML-Methode kein Bestimmtheitsmaß in Analogie zum R^2 -Wert der OLS-Regression. Um die Anpassungsgüte der Regression mit ARIMA(1,0,0)-Fehlermodell mit derjenigen der reinen OLS-Regression zu vergleichen, analysieren wir daher die jeweiligen Modellresiduen (RES_1 bzw. ERR_1) mit der SPSS-Prozedur FIT, die leider nur via Syntax verfügbar ist:

```
fit error = res_1 err_1
/obs = consump, consump
/dfe = 66 65.
```

Die abhängige Variable war in beiden Fällen CONSUMP, die Freiheitsgradzahlen sind den bisherigen SPSS-Ausgaben zu entnehmen. Wegen des zusätzlichen autoregressiven Parameters haben die Residuen des erweiterten Modells nur noch 65 Freiheitsgrade.

FIT Error Statistics			
Error Variable		RES_1	ERR_1
Observed Variable		CONSUMP	CONSUMP
N of Cases	Use	69	69
Deg Freedom	Use	66	65
Mean Error	Use	,0000	-,0027
Mean Abs Error	Use	,0457	,0184
Mean Pct Error	Use	-,1116	-,2226
Mean Abs Pct Err	Use	2,5766	1,0627
SSE	Use	,2210	,0676
MSE	Use	,0033	,0010
RMS	Use	,0579	,0322
Durbin-Watson	Use	,2488	1,4945

In der "RMS"-Zeile steht für jedes Modell die Wurzel aus der geschätzten Fehlervarianz, also der geschätzte Standardfehler. Hier erreicht das erweiterte Modell mit 0,0322 einen erheblich besseren Wert als das OLS-Regressionsmodell (0,0579), d.h. die Anpassungsgüte hat sich deutlich verbessert. In der Ausgabe der Autoregressionsprozedur wird für das erweiterte Modell ein noch günstigerer Standardfehler (0,023) angegeben, weil hier offenbar der erste Fall, dessen Residuum nicht korrekt geschätzt werden kann (s.o.), weggelassen wurde. Um diese Vermutung zu überprüfen, wollen wir den ersten Fall ausschließen und dann die Fit-Analyse wiederholen. Das temporäre Deaktivieren des Jahres 1870 kann z.B. nach dem Befehl **Daten > Fälle auswählen...** mit der folgenden Dialogbox geschehen:

Fälle auswählen: Bereich

Erster Fall: 1871 Letzter Fall: 1938

Jahr: 1871 1938

Weiter
Abbrechen
Hilfe

Nun liefert der Fit-Befehl mit angepaßten Freiheitsgraden:

```
fit error = res_1 err_1
/obs = consump, consump
/dfe = 65 64.
```

Das folgende Ergebnis:

FIT Error Statistics			
Error Variable		RES_1	ERR_1
Observed Variable		CONSUMP	CONSUMP
N of Cases	Use	68	68
Deg Freedom	Use	65	64
Mean Error	Use	,0012	-,0055
Mean Abs Error	Use	,0451	,0159
Mean Pct Error	Use	-,0502	-,3658
Mean Abs Pct Err	Use	2,5515	,9384
SSE	Use	,2139	,0329
MSE	Use	,0033	,0005
RMS	Use	,0574	,0227
Durbin-Watson	Use	,2565	2,0949

Tatsächlich stimmt der RMS-Wert der Fit-Prozedur nun mit dem Standardfehler der Autoregressions-Prozedur überein. Ferner zeigt die letzte Zeile, daß sich auch der Durbin-Watson-Koeffizient perfekt normalisiert hat, wobei aber eine Signifikanzbeurteilung im Modell mit Autoregressor leider nicht möglich ist (siehe oben).

4 Regression mit integriertem ARIMA(p,d,q)-Modell für die Fehler

Zwar ist der autoregressive Fehlerprozeß erster Ordnung ein in der Praxis außerordentlich wichtiger Spezialfall, doch kommen natürlich prinzipiell für den Fehlerprozeß beliebige ARIMA-Modelle in Frage. SPSS bietet für den allgemeinen Fall die Prozedur ARIMA, die mit

Statistik > Zeitreihen > ARIMA...

aufgerufen und mit folgender Dialogbox gesteuert wird:

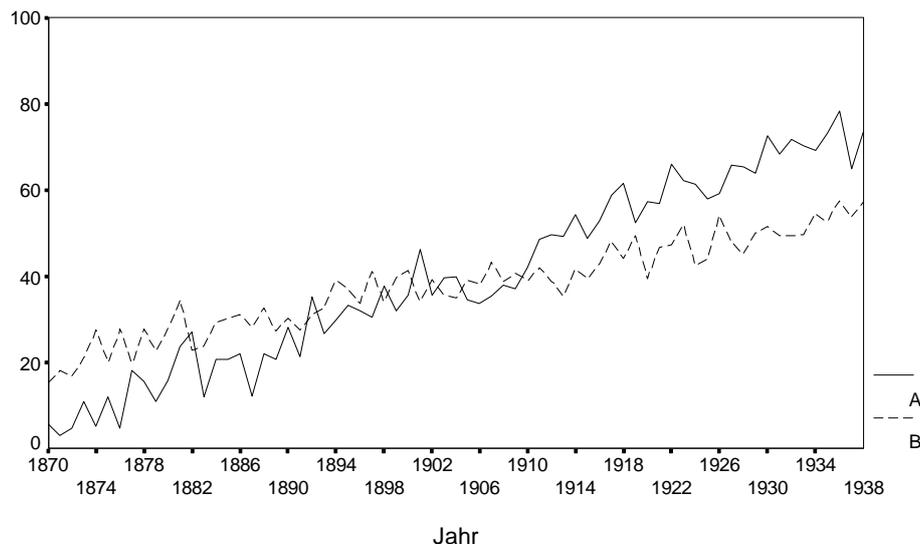
The screenshot shows the SPSS ARIMA dialog box. On the left, a list of variables includes 'err_1', 'fit_1', 'lcl_1', 'res_1', 'sep_1', 'ucl_1', and 'year_'. The 'Abhängige Variable:' field is set to 'consump'. The 'Transformation:' dropdown menu is set to 'Keine'. The 'Unabhängige Variable(n):' list contains 'income' and 'price'. The 'Modell' section has 'Autoregressiv' with p=1, 'Differenz' with d=0, and 'Gleitender Durchschnitt' with q=0. The 'Saisonal' section has 'sp: 0', 'sd: 0', and 'sq: 0'. The checkbox 'Konstante in Modell einschließen' is checked. At the bottom, 'Aktuelle Periodizität:' is 'Keine'. Buttons for 'OK', 'Befehl', 'Zurücksetzen', 'Abbrechen', 'Hilfe', 'Speichern...', and 'Optionen...' are visible.

In diesem Beispiel wird der für unsere Daten passende ARIMA(1,0,0)-Fehlerprozeß spezifiziert, so daß wir (bis auf Rundungsfehler) dieselben Ergebnisse erhalten wie bei der oben verwendeten Prozedur Autoregression. ARIMA benutzt im Unterschied zur Prozedur Autoregression stets die Maximum Likelihood - Schätzmethode.

5 Sonstige Hinweise zu Regressionsanalysen mit Zeitreihendaten

5.1 Warnung vor Regressionsanalysen mit trendbelasteten Zeitreihen

Unabhängig von unserer bisherigen Argumentationslinie soll an dieser Stelle noch vor der Regression einer trendbelasteten Zeitreihe *A* auf eine andere, ebenfalls trendbelastete Zeitreihe *B* gewarnt werden. Hier hat offenbar der Faktor ZEIT einen Effekt auf die beiden Zeitreihen bzw. Variablen *A* und *B*. ZEIT muß als Regressor in das Modell aufgenommen werden, weil sonst ein erheblicher "omitted-variable-error" auftreten kann. Der Regressor-Zeitreihe *A* wird in diesem Fall zu Unrecht ein starker Effekt auf die abhängige Zeitreihe *B* zugeschrieben. Die beiden Zeitreihen in der folgenden Abbildung haben bis auf den Trend keinerlei Gemeinsamkeit, sie korrelieren jedoch zu 0,92 miteinander:



Der "omitted-variable-error" kann mit den oben beschriebenen Verfahren zur Residuenanalyse natürlich *nicht* aufgedeckt werden.

Um trendbereinigte Ergebnisse zu erzielen, können Sie ...

- den Trend beseitigen, z.B. durch Differenzieren der Zeitreihen,
- die Zeit als Prädiktor in das Modell einbringen.

5.2 Zeitreihen mit Saisonfigur

Wenn Sie eine abhängige Zeitreihe mit jahreszeitlichen Schwankungen analysieren wollen, haben Sie u.a. die beiden folgenden Möglichkeiten (siehe Fieger & Toutenburg 1995, S. 71ff; SPSS 1994, S. 153ff):

- Beziehen Sie die Monate mit Kodiervariablen in das Modell mit ein.
- Entfernen Sie Saisonfigur mit dem SPSS-Befehl

Statistik > Zeitreihen > Saisonale Zerlegung...

6 Literatur

- Durbin, J. & Watson, G.S. (1951). Testing for serial correlation in least-squares regression II. *Biometrika*, 38, 159-178.
- Fieger, A. & Toutenburg, H. (1995). *SPSS Trends für Windows*. München: Prentice Hall.
- Hartung, J. (1989). *Statistik*. München: Oldenbourg.
- Pindyck, R. S. & Rubinfeld, D.L. (1990). *Econometric models and econometric forecasts*. Auckland: McGraw-Hill.
- Schlittgen, R. & Streitberg, B.H.J. (1989). *Zeitreihenanalyse*. München: Oldenbourg.
- SPSS, Inc. (1994). *SPSS Trends 6.1*. Chicago, IL.

7 Stichwortverzeichnis

	omitted-variable-error	20
	P	
	PACF	11
	Partielle Autokorrelationsfunktion.....	11
	R	
	Residuen	4
	RMS	17
	S	
	Standardfehler.....	17
	Stationär.....	14
	T	
	Trendbelastete Zeitreihen	20
	Trends	6
	W	
	Weißes Rauschen.....	8
	Z	
	Zeitreihenanalyse.....	5
A		
ACF		11
AR(1)-Prozeß		14
ARIMA(1,0,0)-Modell		12, 13
Autokorrelation erster Ordnung.....		9
Autokorrelationsfunktion.....		11
B		
Box-Ljung-Test		15
D		
Durbin-Watson-Daten.....		4
Durbin-Watson-Koeffizien		7
Durbin-Watson-Test		9
F		
FIT		17
M		
Maximum Likelihood - Schätzungen.....		13
Modellresiduen		4
O		
OLS		4