

Simultangleichungsmodelle und Pfadanalyse

Prof. Dr. Martin Elff

E-mail: martin.elff@zu.de

Fall Semester 2015

1 Einführung

Strukturgleichungen in der Soziologie, Ökonometrie und Psychometrie

- Strukturgleichungen wurden unter der Bezeichnung Pfadanalyse bereits in den 1960er Jahren in die Soziologie eingeführt, u.a. von Soziologen, die sich mit Berufsprestige oder sozialer Mobilität befasst haben, wie z.B. Duncan und Blalock.
- In der empirischen Wirtschaftswissenschaften werden auch Strukturgleichungsmodelle verwendet, insbesondere in der Makro-Ökonomie. Hier wird auch noch zwischen spezifischen Modelltypen unterschieden, z.B. *seeming unrelated regression* (SUR)-Modelle. Auch sind hier spezielle Verfahren geläufig: z.B. Two-Step-Least-Squares oder Three-Step-Least-Squares – ähnlich wie wir es bei der Diskussion der Technik der Instrumental-Variablen gesehen haben.
- Weit entwickelt wurde die Methodologie der Strukturgleichungsmodelle in der Psychometrie, der psychologischen Teildisziplin die sich mit der Messung von psychologischen Konzepten befasst. Hier werden häufig auch Modelle mit latenten Variablen verwendet (dazu mehr in späteren Sitzungen).
- In der psychometrisch inspirierten Umfrageforschung finden sich auch häufig Modelle mit Strukturgleichungen.

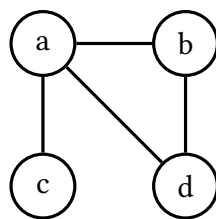
2 Graphen und Kausalbeziehungen

Repräsentation von Kausalbeziehungen durch Graphen

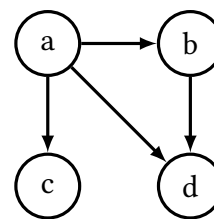
- Graphen allgemein – werden verwendet um kausale oder statistische Beziehungen zwischen Variablen darzustellen. Sie bestehen aus
 - Ecken (*vertices*) oder Knoten (*nodes*): Die Variablen, zwischen denen Beziehungen bestehen können
 - Kanten (*edges*) oder Verbindungen: Verbindung zwischen Ecken, stellen dar, ob eine kausale oder statistische Beziehung besteht
- Gerichtete Graphen (*directed graphs*) – werden verwendet, um kausale oder statistische *Einflussbeziehungen* darzustellen. Verbindungen können
 - einseitig gerichtet sein: Eindeutige Beziehungen von Ursache und Wirkung
 - beidseitig gerichtet sein: kausale Richtung ist wechselseitig, oder beide Variablen werden gemeinsam von einer (unbeobachteten) Drittvariablen beeinflusst (konfundierender Einfluss)

Ein ungerichteter und ein gerichteter Graph

Ein ungerichteter Graph



Ein gerichteter Graph

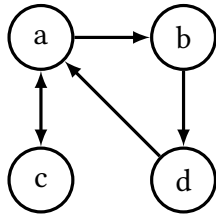


Pfade in gerichteten Graphen

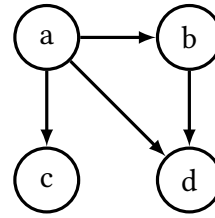
- Ein *Pfad* eine Verkettung von Kanten/Verbindungen mit gemeinsamen Knoten, so dass jeder Endpunkt einer Verbindung der Anfangspunkt der nächsten Verbindung ist.
- Ein Pfad enthält einen Zyklus (oder eine Schleife), wenn man von einem durch ihn verbundenen Knoten entlang den Richtungen der Verbindungen wieder zu diesem zurück gelangen kann.
- Ein *gerichteter azyklischer Graph* (*directed acyclical graph* – DAG) ist ein Graph, in dem kein Pfad einen Zyklus enthält.

Ein zyklischer und ein nicht-zyklischer gerichteter Graph

Ein nicht azyklischer Graph



Ein azyklischer Graph



Endogenität und Exogenität in Pfadmodellen

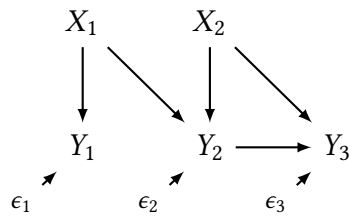
- Eine Variable in einem Pfadmodell heißt *endogen*, wenn sie der Endpunkt mindestens einer Verbindung ist.
- andernfalls heißt sie *exogen*.
- Das ist so zu interpretieren, dass endogene Variablen von mindestens einer anderen Variablen im Modell beeinflusst wird (Einfluss kann auch wechselseitig sein.)
- Exogene Variablen treten im Modell(!) nur als Einflussfaktoren (Ursachen) für andere Variablen auf, ohne selbst von Variablen im Modell beeinflusst zu werden.

Mediation, direkte und indirekte Effekte

- Eine Variable X_1 ist ein Mediator, wenn sie den Einfluss von einer Variablen X_0 zu einer anderen Variable Y vermittelt, in dem Sinne, dass sie
 - der Endpunkt der Verbindung zwischen X_0 und X_1 ist
 - und der Anfangspunkt der Verbindung zwischen X_1 und Y
- Wenn es auch eine Verbindung zwischen X_0 und Y gibt, dann
 - repräsentiert diese den *direkten Effekt* von X_0 auf Y
 - während der Pfad der Verbindungen von X_0 zu X_1 und X_1 zu Y den (durch X_1 vermittelten) *indirekten Effekt* von X_0 auf Y repräsentiert.

Illustration der eingeführten Begriffe

Gegeben sei das Pfadmodell



- Die Variablen X_1 und X_2 sind *exogen*.
- Die Variablen Y_1 , Y_2 , und Y_3 sind *endogen*.
- Die Variable Y_2 ist ein Mediator für Effekte von X_1 und X_2 auf Y_3 .
- Die Variable X_1 hat nur einen indirekten Effekt auf Y_3 , während die Variable X_2 sowohl einen direkten als auch einen indirekten Effekt auf Y_3 hat.

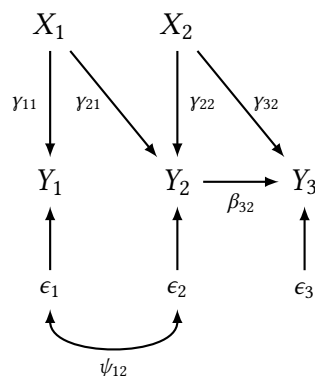
3 Strukturgleichungssysteme

Pfadmodelle und Strukturgleichungen

- Pfadmodelle und lassen sich in Systeme von Strukturgleichungen übersetzten und umgekehrt.
- Dabei entspricht jede gerichtete Verbindung zwischen zwei Variablen einem Strukturkoeffizienten.
- Weiterhin können Strukturgleichungen auch weitere Parameter enthalten: Kovarianzen zwischen den endogenen Variablen bzw. deren Fehlertermen.

Ein Beispiel für ein Pfadmodell und das zugehörige Strukturgleichungssystem

Das Pfadmodell:



Die zugehörigen Strukturgleichungen:

$$\begin{aligned} Y_1 &= \alpha_1 + \gamma_{11}X_1 + \epsilon_1 \\ Y_2 &= \alpha_2 + \gamma_{21}X_1 + \gamma_{22}X_2 + \epsilon_2 \\ Y_3 &= \alpha_3 + \beta_{32}Y_2 + \gamma_{32}X_2 + \epsilon_3 \end{aligned}$$

$$\text{Var}(\epsilon_1) = \psi_{11}$$

$$\text{Cov}(\epsilon_1, \epsilon_2) = \psi_{12} \quad \text{Var}(\epsilon_2) = \psi_{22}$$

$$\text{Var}(\epsilon_3) = \psi_{33}$$

$$\text{Var}(X_1) = \phi_{11}$$

$$\text{Cov}(X_1, X_2) = \phi_{12} \quad \text{Var}(X_2) = \phi_{22}$$

Strukturgleichungen in Matrixform

Strukturgleichungen in Matrixform:

$$\mathbf{y} = \boldsymbol{\alpha} + \mathbf{B}\mathbf{y} + \boldsymbol{\Gamma}\mathbf{x} + \boldsymbol{\epsilon} \quad \text{Cov}(\boldsymbol{\epsilon}) = \boldsymbol{\Psi} \quad \text{Cov}(\mathbf{x}) = \boldsymbol{\Phi}$$

Im vorangegangenen Beispiel:

$$\mathbf{y} = \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \end{bmatrix}$$

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & \beta_{32} & 0 \end{bmatrix} \quad \boldsymbol{\Gamma} = \begin{bmatrix} \gamma_{11} & 0 \\ \gamma_{21} & \gamma_{22} \\ 0 & \gamma_{32} \end{bmatrix} \quad \boldsymbol{\Psi} = \begin{bmatrix} \psi_{11} & \psi_{12} & 0 \\ \psi_{12} & \psi_{22} & 0 \\ 0 & 0 & \psi_{33} \end{bmatrix}$$

- Wie zu sehen ist, entspricht fast jeder möglichen gerichteten Verbindung im Pfadmodell ein Element in den Koeffizienten-Matrizen \mathbf{B} und $\boldsymbol{\Gamma}$.
- Im Pfadmodell nicht „realisierte“ Verbindungen entsprechen Elementen in den Matrizen, die gleich Null gesetzt sind, d.h. auf Null festgelegt sind.
- Solche Festlegungen nennt man die *Parameter-Restriktionen* des Modells.
- Alle Elemente von $\boldsymbol{\alpha}$ und der Matrizen \mathbf{B} , $\boldsymbol{\Gamma}$, und $\boldsymbol{\Psi}$, die nicht auf bestimmte Werte festgelegt sind, heißen die *freien Parameter* des Modells.
- Die Kovarianzmatrix $\boldsymbol{\Phi}$ der exogenen Variablen wird *nie* restringiert.
- Das Strukturgleichungsmodell lässt sich auch in eine sogenannte *reduzierte* Form bringen, in der die endogenen Variablen in \mathbf{y} nur auf der linken Seite der Gleichung stehen:

$$\mathbf{y} = \boldsymbol{\alpha} + \mathbf{B}\mathbf{y} + \boldsymbol{\Gamma}\mathbf{x} + \boldsymbol{\epsilon} \Leftrightarrow$$

$$(\mathbf{I} - \mathbf{B})\mathbf{y} = \boldsymbol{\alpha} + \boldsymbol{\Gamma}\mathbf{x} + \boldsymbol{\epsilon} \Leftrightarrow$$

$$\mathbf{y} = (\mathbf{I} - \mathbf{B})^{-1}\boldsymbol{\alpha} + (\mathbf{I} - \mathbf{B})^{-1}\boldsymbol{\Gamma}\mathbf{x} + (\mathbf{I} - \mathbf{B})^{-1}\boldsymbol{\epsilon}$$

- Aus der reduzierten Form lassen sich Kriterien für die Schätzung der Parameter eines Strukturgleichungs-Modells herleiten.

Schätzung von Strukturgleichungsmodellen

- Die Parameter können unter gewissen Voraussetzungen aus empirischen Kovarianz-Matrizen geschätzt werden, die sich aus einem Datensatz gewinnen lassen.
- Mit Hilfe der Matrix-Algebra lässt sich eine Verbindung herstellen zwischen den Parametern des Modells und der von dem Modell implizierten Varianzen und Kovarianzen der beobachtbaren Variablen \mathbf{x} und \mathbf{y} . (Die Fehlerterme in ϵ sind nicht direkt beobachtbar.)
- Die freien Parameter des Modells können dann dadurch geschätzt werden, dass der (in bestimmter Weise gewichtete) Unterschied zwischen den vom Modell implizierten Varianzen und Kovarianzen und den empirischen Varianzen und Kovarianzen so klein wie möglich gemacht wird. (Zu den Details mehr in einer späteren Sitzung.)
- Sei Σ_x die Varianz-Kovarianz-Matrix von \mathbf{x} , Σ_{xy} die vom Modell implizierte Kovarianz-Matrix von \mathbf{x} und \mathbf{y} , sowie Σ_y die vom Modell implizierte Varianz-Kovarianz-Matrix von \mathbf{y} . Dann sind die Modellimplikationen:

$$\begin{aligned}\Sigma_x &= \Phi \\ \Sigma_{xy} &= \Phi \Gamma' (I - B)^{-1'} \\ \Sigma_y &= (I - B)^{-1} (\Gamma \Phi \Gamma' + \Psi) (I - B)^{-1'}\end{aligned}$$

Die Problematik der Identifikation

- Damit ein Strukturgleichungsmodell überhaupt schätzbar sein kann, muss es *identifiziert* sein:
- Für jede mögliche Kovarianzmatrix muss es genau einen Satz von möglichen Werten der freien Parameter des Modells, die diese Kovarianzmatrix implizieren.
- Wenn eine Kovarianzmatrix mit mehr einem Satz von Parameterwerten vereinbar ist, dann ist das entsprechende Modell *unteridentifiziert*.
- Probleme für die Identifikation eines Strukturgleichungsmodells stellen insbesondere von Null verschiedene Parameter in der Matrix \mathbf{B} und außerhalb der Diagonalen der Matrix Ψ dar.
- In der Regel müssen die meisten Parameter in \mathbf{B} und außerhalb der Diagonalen von Ψ gleich Null gesetzt werden, um die Identifikation des Modells sicher zu stellen.
- Wenn allerdings alle diese Parameter gleich Null gesetzt werden, dann hat man kein „interessantes“ Strukturgleichungsmodell mehr, sondern einfach ein System von linearen Regressionsgleichungen.

- Der Nachweis, dass ein konkretes Strukturgleichungsmodell identifiziert ist, kann im Einzelfall schwierig sein. Allerdings gibt es verschiedene Regeln für notwendige oder hinreichende Bedingungen.
- Die Darstellung des Modells in Form eines Graphs kann dabei hilfreich sein.
- Insbesondere dann, wenn der Graph des Strukturgleichungsmodells ein gerichteter azyklischer Graph ist, dann ist das Modell identifiziert.

Strukturkoeffizienten und Pfadkoeffizienten

- In der soziologischen Pfadanalyse wurden früher häufig nur standardisierte Variablen in Strukturgleichungsmodellen verwendet.
- Koeffizienten von standardisierten Variablen in Strukturgleichungsmodellen werden auch *Pfadkoeffizienten* genannt.
- Der Name rührt daher, dass die „Stärke des kausalen Einflusses“ entlang unterschiedlicher Pfade einfach quantitativ ausgedrückt werden kann.
- Beispiel:
 - Drei standardisierte Variablen Z_1 , Z_2 , und Z_3 .
 - Zwei Strukturgleichungen:

$$Z_2 = \rho_{21}Z_1 + \epsilon_1$$

$$Z_3 = \rho_{31}Z_1 + \rho_{32}Z_2 + \epsilon_1$$

- Direkter Effekt von Z_1 auf Z_3 : ρ_{31}
- Indirekter (von Z_2 vermittelter) Effekt von Z_1 auf Z_3 : $\rho_{32}\rho_{21}$
- Totaler Effekt (direkter und indirekter Effekt) von Z_1 auf Z_3 : $\rho_{31} + \rho_{32}\rho_{21}$

4 Anwendungsbeispiel: Sympathie für die Linkspartei

Datengrundlage: Daten aus der *German Longitudinal Election Study* (GLES) von 2009

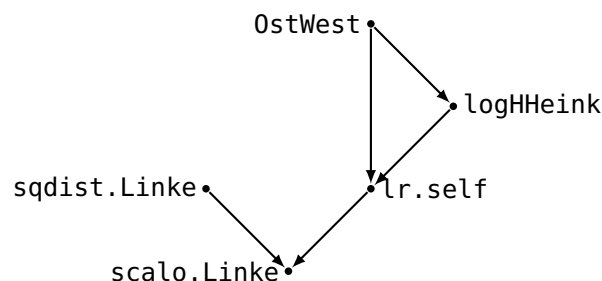
- Datensatz **"work-sem1.RData"**
- Relevante Variablen
 - Links-Rechts-Selbsteinordnung: **lr.self**
 - Wahrgenommene Distanz zwischen der LINKEN und der eigenen Position: **sqdist.Linke**
 - Sympathie für die LINKE: **scal0.Linke**
 - Ostdeutsche und Westdeutsche: Dummy-Variable **OstWest**
 - Haushaltseinkommen, logarithmiert: **logHHeink**

Modellspezifikation

- Die Sympathie für die Linke (gemessen mit der Variable **scal0.Linke**) wird beeinflusst
 - von der Position der Wähler auf der Links-Rechts-Skala (enthalten in der Variablen **lr.self**)
 - von der wahrgenommenen Distanz zwischen der „Linken“ und der eigenen Position der Wähler auf der Links-Rechts-Skala (enthalten in der Variablen **sqdist.Linke**)
 - Die Position der Wähler auf der Links-Rechts-Skala wird wiederum beeinflusst
 - davon, ob Sie in Ostdeutschland oder in Westdeutschland leben (repräsentiert durch die Dummy-Variable **OstWest**) und
 - vom logarithmierten Gesamteinkommen des Haushalts des Wählers/der Wählerin (in der Variablen **logHHeink**).
 - Das Haushaltseinkommen wird wiederum davon beeinflusst, ob man im Osten oder im Westen Deutschlands lebt.
- ⇒ Welche Variablen sind in dem Modell *endogen* und welche *exogen*?

Das verbal spezifizierte Modell als Graph und als Gleichungssystem

- Die beschriebenen Zusammenhänge lassen sich graphisch folgendermaßen darstellen:



- Das graphisch dargestellte Kausalmodell entspricht folgendem Strukturgleichungsmodell:

$$\begin{aligned}\log HHeink &= \alpha_1 + \gamma_{11} OstWest + \epsilon_1 \\ lr.self &= \alpha_2 + \beta_{21} \log HHeink + \gamma_{21} OstWest + \epsilon_2 \\ scal0.Linke &= \alpha_3 + \beta_{32} lr.self + \gamma_{31} sqdist.Linke + \epsilon_3\end{aligned}$$

Schätzung des Modells mit systemfit

- **systemfit** ist ein Zusatzpaket für R, das vor allem für Simultangleichungsmodelle geeignet ist, wie sie in der Ökonomie und Ökonometrie vor kommen.
- Angeboten werden vor allem Varianten des Kleinstquadrateschätzers: Ordinary Least Squares, Two-Step Least Squares etc. für nicht-rekursive Gleichungssysteme
- Code für die Schätzung des Modells:

```
load ("work-sem1.RData") # Der Datensatz
library(systemfit) # Das R-Zusatzpaket
sfit.Linke <- systemfit(
  list( # Mehrere Strukturgleichungen kommen in eine Liste
    logHHeink ~ OstWest, # Kommas nicht vergessen!
    lr.self ~ logHHeink + OstWest,
    scalo.Linke ~ lr.self + sqdist.Linke + OstWest # Hier kein Komma!
  ), # Komma nicht vergessen!
  data=work.sem1)
summary(sfit.Linke)
```

Schätzergebnisse mit systemfit

- Zusammenfassende Statistiken für die Gleichungen:

Output

```
systemfit results
method: OLS
```

	N	DF	SSR	detRCov	OLS-R2	McElroy-R2
system	8449	8440	29106	6.51	0.298	0.186

	N	DF	SSR	MSE	RMSE	R2	Adj R2
eq1	2537	2535	936	0.369	0.608	0.022	0.021
eq2	2308	2305	8093	3.511	1.874	0.039	0.038
eq3	3604	3600	20077	5.577	2.362	0.374	0.374

- Kovarianzen und Korrelationen der Residuen (unerklärten Anteilen der endogenen Variablen)

Output

```
The covariance matrix of the residuals
```

	eq1	eq2	eq3
eq1	0.34487	-0.00186	-0.0321
eq2	-0.00186	3.52070	-0.1474

eq3 -0.03213 -0.14738 5.3691

The correlations of the residuals

	eq1	eq2	eq3
eq1	1.00000	-0.00143	-0.0275
eq2	-0.00143	1.00000	-0.0338
eq3	-0.02749	-0.03377	1.00000

- Erste Gleichung

_____ Output _____

OLS estimates for 'eq1' (equation 1)

Model Formula: logHHeink ~ OstWest

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.0983	0.0421	168.68	<2e-16 ***
OstWest	0.1905	0.0245	7.78	1e-14 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.608 on 2535 degrees of freedom

Number of observations: 2537 Degrees of Freedom: 2535

SSR: 935.857 MSE: 0.369 Root MSE: 0.608

Multiple R-Squared: 0.022 Adjusted R-Squared: 0.021

- Zweite Gleichung

_____ Output _____

OLS estimates for 'eq2' (equation 2)

Model Formula: lr.self ~ logHHeink + OstWest

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-2.4286	0.4988	-4.87	0.0000012 ***
logHHeink	0.0731	0.0667	1.09	0.27
OstWest	0.7774	0.0831	9.35	< 2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.874 on 2305 degrees of freedom

Number of observations: 2308 Degrees of Freedom: 2305

SSR: 8093.095 MSE: 3.511 Root MSE: 1.874

Multiple R-Squared: 0.039 Adjusted R-Squared: 0.038

- Dritte Gleichung

_____ Output _____

OLS estimates for 'eq3' (equation 3)

Model Formula: `scalo.Linke ~ lr.self + sqdist.Linke + OstWest`

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.74150	0.16900	4.39	0.00001179 ***
lr.self	-0.68904	0.03365	-20.48	< 2e-16 ***
sqdist.Linke	-0.01882	0.00358	-5.26	0.00000015 ***
OstWest	-1.22888	0.08337	-14.74	< 2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.362 on 3600 degrees of freedom

Number of observations: 3604 Degrees of Freedom: 3600

SSR: 20076.693 MSE: 5.577 Root MSE: 2.362

Multiple R-Squared: 0.374 Adjusted R-Squared: 0.374

Schätzung des Modells mit lavaan

- **lavaan** ist ein Zusatzpaket für R, das für Faktor- und Strukturgleichungsmodelle geeignet sind, die (auch) latente Variablen enthalten. (**lavaan**=latent variable analysis)
- Angeboten werden Maximum-Likelihood Schätzer und Varianten, wie Generalised Least Squares, Distribution-Free Weighted Least Squares, etc.
- Code für die Schätzung des Modells:

```
load ("work-sem1.RData") # Der Datensatz
library(lavaan)
lavaan.Linke <- sem(
  ' # Mehrere Strukturgleichungen kommen in eine Zeichenkette!
    logHHeink ~ OstWest # keine Kommas!
    lr.self ~ logHHeink + OstWest
    scalo.Linke ~ lr.self + sqdist.Linke + OstWest
  ', # Komma nicht vergessen!
  data=work.sem1)
summary(lavaan.Linke)
```

Schätzergebnisse mit lavaan

- Zusammenfassende Statistiken über das Modell als Ganzes:

Output

lavaan (0.5-20) converged normally after 31 iterations

	Used	Total
Number of observations	2246	4288
Estimator	ML	
Minimum Function Test Statistic	2257.425	
Degrees of freedom	3	
P-value (Chi-square)	0.000	

Parameter Estimates:

Information	Expected
Standard Errors	Standard

- Koeffizienten der Strukturgleichungen

Output

Regressions:

	Estimate	Std.Err	Z-value	P(> z)
logHHeink ~				
OstWest	0.183	0.026	7.035	0.000
lr.self ~				
logHHeink	0.069	0.068	1.015	0.310
OstWest	0.809	0.084	9.620	0.000
scalo.Linke ~				
lr.self	-0.723	0.026	-27.804	0.000
sqdist.Linke	-0.020	0.003	-7.042	0.000
OstWest	-1.139	0.105	-10.811	0.000

- Varianzen der Fehlerterme

Output

Variances:

	Estimate	Std.Err	Z-value	P(> z)
logHHeink	0.343	0.010	33.511	0.000
lr.self	3.516	0.105	33.511	0.000
scalo.Linke	5.339	0.159	33.511	0.000

- Die Strukturgleichungen formatiert:

	logHHeink	lr.self	scalo.Linke
OstWest	0.183*** (0.026)	0.809*** (0.084)	-1.139*** (0.105)
logHHeink		0.069 (0.068)	
lr.self			-0.723*** (0.026)
sqdist.Linke			-0.020*** (0.003)
gfi	0.798		
rmsea	0.578		
N	2246		