

Biewen, Elena

**Working Paper**

## Within-Schätzung bei anonymisierten Paneldaten

IAW Diskussionspapiere, No. 34

**Provided in Cooperation with:**

Institute for Applied Economic Research (IAW)

*Suggested Citation:* Biewen, Elena (2007) : Within-Schätzung bei anonymisierten Paneldaten, IAW Diskussionspapiere, No. 34, Institut für Angewandte Wirtschaftsforschung (IAW), Tübingen

This Version is available at:

<https://hdl.handle.net/10419/36625>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*



# IAW-Diskussionspapiere

Discussion Paper

| 34 |

## Within-Schätzung bei anonymisierten Paneldaten

Elena Biewen

November 2007

ISSN: 1617-5654

INSTITUT FÜR  
ANGEWANDTE  
WIRTSCHAFTSFORSCHUNG

Ob dem Himmelreich 1  
72074 Tübingen  
T: (0 70 71) 98 96-0  
F: (0 70 71) 98 96-99  
E-Mail: [iaw@iaw.edu](mailto:iaw@iaw.edu)  
Internet: [www.iaw.edu](http://www.iaw.edu)



# IAW-Diskussionspapiere

Das Institut für Angewandte Wirtschaftsforschung (IAW) Tübingen ist ein unabhängiges außeruniversitäres Forschungsinstitut, das am 17. Juli 1957 auf Initiative von Professor Dr. Hans Peter gegründet wurde. Es hat die Aufgabe, Forschungsergebnisse aus dem Gebiet der Wirtschafts- und Sozialwissenschaften auf Fragen der Wirtschaft anzuwenden. Die Tätigkeit des Instituts konzentriert sich auf empirische Wirtschaftsforschung und Politikberatung.

Dieses **IAW-Diskussionspapier** können Sie auch von unserer IAW-Homepage als pdf-Datei herunterladen:

<http://www.iaw.edu/Publikationen/IAW-Diskussionspapiere>

## ISSN 1617-5654

Weitere Publikationen des IAW:

- IAW-News (erscheinen 4x jährlich)
- IAW-Report (erscheinen 2x jährlich)
- IAW-Wohnungsmonitor Baden-Württemberg (erscheint 1x jährlich kostenlos)
- IAW-Forschungsberichte

Möchten Sie regelmäßig eine unserer Publikationen erhalten, dann wenden Sie sich bitte an uns:

IAW Tübingen, Ob dem Himmelreich 1, 72074 Tübingen,  
Telefon 07071 / 98 96-0  
Fax 07071 / 98 96-99  
E-Mail: [iaw@iaw.edu](mailto:iaw@iaw.edu)

Aktuelle Informationen finden Sie auch im Internet unter: <http://www.iaw.edu>

---

Der Inhalt der Beiträge in den IAW-Diskussionspapieren liegt in alleiniger Verantwortung der Autorinnen und Autoren und stellt nicht notwendigerweise die Meinung des IAW dar.



# Within-Schätzung bei anonymisierten Paneldaten

Elena Biewen\*

## Zusammenfassung

Dieser Beitrag befasst sich mit der Untersuchung der Auswirkungen der datenverändernden Anonymisierungsverfahren 'variablenspezifische abstandsorientierte Mikroaggregation' (vaabMA, im Englischen auch als Individual Ranking bezeichnet) und 'multiplikative stochastische Überlagerung' auf die 'Within'-Schätzung eines linearen Panelmodells mit Individualeffekten. Es wird gezeigt, dass der 'Within'-Schätzer auf der Grundlage der mittels vaabMA anonymisierten Daten konsistent bleibt. Bei der multiplikativen stochastischen Überlagerung wird neben der allgemeinen Form der Überlagerung eine spezielle Variante analysiert, bei der die Variablen zuerst mit einem konstantem Grundüberlagerungsfaktor und danach zusätzlich additiv überlagert werden. Es wird weiterhin gezeigt, dass beide Varianten der Überlagerung zu Inkonsistenz der 'Within'-Schätzer führen. Anschließend werden korrigierte Schätzer hergeleitet.

## 1 Einleitung

Dieser Beitrag entstand im Rahmen des vom Bundesministerium für Bildung und Forschung (BMBF) geförderten Projektes "Wirtschaftsstatistische Paneldaten und faktische Anonymisierung" und beschäftigt sich mit zwei datenverändernden Verfahren, die für die Anonymisierung von Paneldaten verwendet werden: variablenspezifische abstandsorientierte Mikroaggregation (fortan vaabMA) und multiplikative stochastische Überlagerung. Das Ziel ist die Analyse des Einflusses der beiden Anonymisierungsstrategien auf die 'Within'-Schätzung eines linearen Panelmodells mit Individualeffekten.

Die vaabMA ist eine Variante der getrennten Mikroaggregation, bei der die zu Gruppen zusammengefassten Beobachtungen durch die Gruppenmittelwerte ersetzt werden. Die Mikroaggregation bei Paneldaten wird getrennt sowohl für alle Merkmale als auch für alle Zeitpunkte ausgeführt. Andere Varianten der Mikroaggregation, allerdings für Querschnittsdaten, werden beispielsweise bei Schmid et al. (2005), Schmid (2006b) behandelt.

---

\*Elena Biewen, Institut für Angewandte Wirtschaftsforschung, Ob dem Himmelreich 1, 72074 Tübingen, Deutschland. Tel.: (07071) 98 96-36, Fax -99, E-Mail: elena.biewen@iaw.edu. Ich danke Gerd Ronning und Martin Rosemann für viele sehr hilfreiche Hinweise.

Bei der stochastischen Überlagerung wird ein stochastisch erzeugter Fehler zu den Originaldaten hinzugefügt. Grundsätzlich unterscheidet man zwischen der additiven und der multiplikativen Überlagerung. Die multiplikative Überlagerung hat den Vorteil, dass sie alle Originalwerte proportional verändert, so dass für die Untersuchungseinheiten ein besserer Schutz gewährleistet ist. Außerdem bleiben die Nullen im Datensatz erhalten (Rosemann 2006). Aus diesen Gründen betrachtet die vorliegende Arbeit ausschließlich den multiplikativen Fall. Während der additive Fehler in der Literatur eingehend untersucht wurde, findet man zur multiplikativen Variante noch wenig. Stefanski (1985) leitet unter anderem einen korrigierten KQ-Schätzer für den multiplikativen Messfehler her. Iturria et al. (1999) schlagen einen konsistenten M-Schätzer für ein polynomiales Regressionsmodell vor. Lechner (2007), Ronning et al. (2005) und Rosemann (2006) untersuchen die SIMEX-Methode bei Vorliegen der multiplikativ überlagerten Daten.

In diesem Aufsatz wird zuerst die allgemeine Form der multiplikativen Überlagerung und anschließend eine spezielle Variante untersucht, bei der die Variablen zuerst mit einem konstanten Grundüberlagerungsfaktor multipliziert und danach zusätzlich additiv überlagert werden. Das letztere Verfahren schlug Höhne (2004) mit dem Ziel vor, die Verhältnisse zwischen unterschiedlichen Variablen nach der Anonymisierung zu erhalten (Ronning 2007).

Der Beitrag ist wie folgt gegliedert. Im Kapitel 2 wird das zu schätzende lineare Panelmodell definiert. Kapitel 3 untersucht nach einer kurzen Erläuterung des Verfahrens vaabMA die Auswirkungen desselben auf die 'Within'-Schätzung. Dies geschieht sowohl theoretisch als auch mithilfe eines Simulationsexperimentes. Kapitel 4 beschäftigt sich mit der Analyse der beiden Varianten der multiplikativen stochastischen Überlagerung. Eine Zusammenfassung der wichtigsten Ergebnisse schließt den Beitrag.

## 2 Das Modell

Das lineare Panelmodell mit Individualeffekten ist definiert als (Mundlak 1978)

$$\mathbf{Y} = \underset{(NT \times N)(N \times 1)}{\mathbf{D}} \boldsymbol{\alpha} + \underset{(NT \times K)(K \times 1)}{\mathbf{X}} \boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (2-1)$$

$\boldsymbol{\alpha}$  enthält  $N$  Individualeffekte, die für jede Untersuchungseinheit nicht über die Zeit variieren, je nach Individuum jedoch verschieden sind.  $\alpha$ 's werden als stochastische Effekte mit Erwartungswert  $\mu_\alpha$  und Varianz  $\sigma_\alpha^2$  behandelt. Die Korrelation zwischen Regressoren und Individualeffekten ist daher möglich.  $\mathbf{D}$  ist definiert als  $\mathbf{D} = \mathbf{I}_N \otimes \boldsymbol{\nu}_T$ , wobei  $\mathbf{I}_N$  eine Einheitsmatrix der Dimension  $(N \times N)$  und  $\boldsymbol{\nu}_T$  ein  $T$ -dimensionaler Spaltenvektor mit Einsen sind.  $\mathbf{X}$  enthält die  $K$ -Regressoren und  $\boldsymbol{\beta}$  ist der Vektor mit entsprechenden Koeffizienten. Der Störterm des Modells  $\epsilon_{it}$  hat die Eigenschaft:

$$\epsilon_{it} \sim iid(0, \sigma_\epsilon^2), \quad i = 1, \dots, N, \quad t = 1, \dots, T \quad (2-2)$$

Multiplikation der beiden Seiten der Gleichung (2-1) mit der Matrix  $\mathbf{Q} = \mathbf{I}_{NT} - \mathbf{D}(\mathbf{D}'\mathbf{D})^{-1}\mathbf{D}'$  bewirkt das Abziehen der Mittelwerte über die Zeit von den Beobachtungen und führt damit zur Eliminierung der Individualeffekte ('Within'-Transformation).<sup>1</sup>

Der 'Within'-Schätzer<sup>2</sup> wird geschrieben als (Hsiao 2003)

$$\widehat{\beta}_w = [\mathbf{X}'\mathbf{Q}\mathbf{X}]^{-1} \mathbf{X}'\mathbf{Q}\mathbf{Y} \quad (2-3)$$

Der naive Schätzer, d.h. der Schätzer unter Verwendung der anonymisierten Daten, ist dann<sup>3</sup>

$$\widehat{\beta}_w^a = [\mathbf{X}^{a'}\mathbf{Q}\mathbf{X}^a]^{-1} \mathbf{X}^{a'}\mathbf{Q}\mathbf{Y}^a \quad (2-4)$$

bzw.

$$\widehat{\beta}_w^a = \left[ \sum_i^N \sum_t^T (\mathbf{x}_{it}^a - \bar{\mathbf{x}}_i^a) (\mathbf{x}_{it}^a - \bar{\mathbf{x}}_i^a)' \right]^{-1} \sum_i^N \sum_t^T (\mathbf{x}_{it}^a - \bar{\mathbf{x}}_i^a) (y_{it}^a - \bar{y}_i^a) \quad (2-5)$$

mit

$$\bar{\mathbf{x}}_i^a = \frac{1}{T} \sum_t^T \mathbf{x}_{it}^a, \quad \bar{y}_i^a = \frac{1}{T} \sum_t^T y_{it}^a$$

Aus (2-5) wird deutlich, dass der Schätzer mithilfe der empirischen 'Within'-Varianz und 'Within'-Kovarianz ausgedrückt werden kann (Greene 2000, S.563-564):

$$\begin{aligned} \widehat{\beta}_w^a &= (\mathbf{S}_{x^a x^a}^w)^{-1} \mathbf{S}_{x^a y^a}^w \\ &= \begin{pmatrix} S_{x_1^a}^{w2} & S_{x_1^a x_2^a}^w & \cdots & S_{x_1^a x_K^a}^w \\ \vdots & & & \vdots \\ S_{x_K^a x_1^a}^w & \cdots & & S_{x_K^a}^{w2} \end{pmatrix}^{-1} \begin{pmatrix} S_{x_1^a y^a}^w \\ \vdots \\ S_{x_K^a y^a}^w \end{pmatrix} \end{aligned} \quad (2-6)$$

<sup>1</sup>Nicht nur Individualeffekte, sondern auch andere zeitkonstante Variablen werden durch diese Transformation gelöscht und ihr Effekt ist dann nicht mehr schätzbar.  $\alpha$  könnte man jedoch mittels einer Teilschätzung schätzen (siehe z.B. Greene 2000, Ronning 2007 (Abschnitt D.2)). Ein weiteres Schätzverfahren der Panelökonometrie, die sogenannte Random-Effects-Schätzung, erlaubt die Schätzung der zeitkonstanten Variablen, unterstellt jedoch die Annahme der Unkorreliertheit der Regressoren und der individuellen Effekte.

<sup>2</sup>Die Bezeichnung 'Within' kommt daher, dass dieser Schätzer nur die Variation innerhalb der Gruppen ausnutzt (Hsiao 2003, S.33).

<sup>3</sup>Mit dem Index  $a$  wird im Folgenden die anonymisierte Variable bzw. der naive Schätzer bezeichnet.



$\widehat{\beta}_w^a$  ist dann konsistent, wenn die einzelnen 'Within'-Varianzen und 'Within'-Kovarianzen in (2-6) konsistent sind, genauer gesagt, wenn sie in Wahrscheinlichkeit gegen die entsprechenden wahren theoretischen Momente konvergieren.<sup>4</sup>

Im Folgenden wird untersucht, wie die Anonymisierung die Konsistenz der 'Within'-Schätzer beeinflusst. Dabei wird nur der Einfluss auf  $\beta$  untersucht.

### 3 Variablenspezifische abstandsorientierte Mikroaggregation

#### 3.1 Das Verfahren

Mit der vaabMA werden die Variablen wie folgt anonymisiert (Schmid 2006a). Den Datensatz sortiert man zuerst nach der ersten Variable, es kann sowohl  $x$  als auch  $y$  sein. Nachdem die Anzahl der Elemente je Gruppe ( $A$ ) festgelegt ist, werden alle Beobachtungen dieser ersten Variable in Gruppen zusammengefasst und ihre Werte durch die Gruppenschritte ersetzt. Die anderen Variablen bleiben dabei unverändert. Die gleiche Prozedur wird anschließend mit der zweiten, dritten u.s.w. Variable durchgeführt. Im Fall von Paneldaten wird die Mikroaggregation auf jeden Zeitpunkt getrennt angewandt.

Wir haben beispielsweise zwei Variablen, die zu zwei Zeitpunkten beobachtet werden (Tabelle a). Zuerst wird nach der ersten Variable (z.B.  $x$ ) im Zeitpunkt  $t = 1$  sortiert. Danach wird diese Variable mikroaggregiert (Tabelle b). Im nächsten Schritt erfolgt die Sortierung nach der zweiten Variable ( $y$ ) in  $t = 1$ , und wiederum wird nur diese Variable mikroaggregiert (Tabelle c).

a				b				c			
t=1		t=2		t=1		t=2		t=1		t=2	
$x$	$y$	$x$	$y$	$x$	$y$	$x$	$y$	$x$	$y$	$x$	$y$
0,5	0,3	0,6	5,4	0,5	0,1	0,7	3,1	0,5	0,2	0,7	3,1
0,9	0,2	2,3	1,2	0,5	0,3	0,6	5,4	1,2	0,2	2,3	1,2
0,7	0,7	4,2	3,2	0,5	0,7	4,2	3,2	0,5	0,2	0,6	5,4
1,4	1,1	0,2	1,5	1,2	0,2	2,3	1,2	1,2	0,8	2,2	0,3
1,3	0,6	2,2	0,3	1,2	0,6	2,2	0,3	0,5	0,8	4,2	3,2
0,3	0,1	0,7	3,1	1,2	1,1	0,2	1,5	1,2	0,8	0,2	1,5

Die gleiche Vorgehensweise wird weiter für den zweiten Zeitpunkt vorgenommen. Anschließend erhalten wir den Datensatz, in dem alle Variablen mittels vaabMA anonymisiert wurden (Tabelle d).

<sup>4</sup>Da bei den Datensätzen, die im Projekt untersucht werden, die Anzahl der Unternehmen sehr groß, die Anzahl der Perioden hingegen gering ist, wird in diesem Beitrag der Wahrscheinlichkeitsgrenzwert für  $N \rightarrow \infty$  untersucht.

d			
t=1		t=2	
x	y	x	y
1,2	0,8	2,9	1,0
1,2	0,2	2,9	1,0
1,2	0,8	0,5	1,0
0,5	0,2	0,5	3,9
0,5	0,8	2,9	3,9
0,5	0,2	0,5	3,9

### 3.2 Auswirkungen auf die 'Within'-Schätzung

Dieses Kapitel befasst sich mit den Auswirkungen von vaabMA auf die 'Within'-Schätzung. Die folgenden Darstellungen basieren auf den Ausführungen von Matthias Schmid (Schmid 2006a), der die Konsistenz des Kleinstquadrat-Schätzers im Fall der Querschnittsdaten bewiesen hat.

Mit  $S_{x_k^a x_l^a}^w$  wird weiter die empirische 'Within'-Kovarianz zwischen zwei mikroaggregierten Regressoren  $x_k^a$  und  $x_l^a$  bezeichnet.<sup>5</sup>  $S_{x_k^a}^w$  ist die Varianz des  $k$ -ten anonymisierten Regressors.

Die 'Within'-Momente kann man umschreiben als<sup>6 7</sup>

$$S_{x_k^a}^w = \frac{1}{T} \sum_t \left( S_{x_{tk}^a}^2 + (\bar{x}_{tk}^a)^2 \right) - \frac{1}{T^2} \sum_t \sum_s (S_{x_{tk}^a x_{sk}^a} + \bar{x}_{tk}^a \bar{x}_{sk}^a) \quad (3-7)$$

$$S_{x_k^a x_l^a}^w = \frac{1}{T} \sum_t (S_{x_{tk}^a x_{tl}^a} + \bar{x}_{tk}^a \bar{x}_{tl}^a) - \frac{1}{T^2} \sum_t \sum_s (S_{x_{tk}^a x_{sl}^a} + \bar{x}_{tk}^a \bar{x}_{sl}^a) \quad (3-8)$$

$$S_{x_k^a y^a}^w = \frac{1}{T} \sum_t (S_{x_{tk}^a y_t^a} + \bar{x}_{tk}^a \bar{y}_t^a) - \frac{1}{T^2} \sum_t \sum_s (S_{x_{tk}^a y_s^a} + \bar{x}_{tk}^a \bar{y}_s^a) \quad (3-9)$$

$$t, s = 1, \dots, T, \quad k = 1, \dots, K, \quad l = 1, \dots, L$$

Die 'Within'-Kovarianzen/-Varianzen lassen sich durch die Gesamtvarianzen/-kovarianzen und Mittelwerte in den einzelnen Zeitpunkten  $t$  und  $s$  darstellen.

Die Mittelwerte bleiben nach der Durchführung der vaabMA erhalten. Für die Kovarianzen/Varianzen lässt sich der Beweis von Schmid (2006a, Lemma 1 und 2) auf den Fall von Paneldaten übertragen. Bei Gültigkeit der Annahmen von Schmid gilt für jedes feste Paar  $(t, s)$  asymptotisch:

<sup>5</sup>Analog für  $S_{x_k^a y^a}^w$

<sup>6</sup>Herleitung im Appendix A

<sup>7</sup>Ich danke Hans Schneeweiß für wichtige Hinweise bei der Herleitung der Formel.

$$\underset{N \rightarrow \infty}{plim} S_{x_{tk}^a x_{sk}^a} = \sigma_{x_{tk} x_{sk}} \quad (3-10)$$

$$\underset{N \rightarrow \infty}{plim} S_{x_{tk}^a x_{sl}^a} = \sigma_{x_{tk} x_{sl}} \quad (3-11)$$

$$\underset{N \rightarrow \infty}{plim} S_{x_{tk}^a y_s^a} = \sigma_{x_{tk} y_s} \quad (3-12)$$

$$t, s = 1, \dots, T, \quad k = 1, \dots, K, \quad l = 1, \dots, L$$

Daraus folgt, dass (3-7)-(3-9) ebenfalls konsistent bleiben. Folglich konvergieren  $\mathbf{S}_{x^a x^a}^w$ ,  $\mathbf{S}_{x^a y^a}^w$  in Wahrscheinlichkeit gegen die wahren  $\Sigma_{xx}^w$ ,  $\Sigma_{xy}^w$ .

Als ein wichtiges Ergebnis lässt sich somit festhalten, dass der 'Within'-Schätzer auf der Grundlage der mit vaabMA anonymisierten Daten konsistent bleibt. Dabei hat der geführte Beweis keine Unabhängigkeitsannahme von  $\mathbf{X}$  unterstellt. Die Konsistenz gilt auch dann, wenn die Regressoren über die Zeit korreliert sind.

### Simulationsergebnisse

Die Auswirkungen von vaabMA auf die 'Within'-Schätzung wurden anhand einer Monte-Carlo-Studie überprüft. Es wurde das lineare Panelmodell (2-1) mit zwei Regressoren geschätzt. Die Regressoren erfüllen die iid-Annahme und folgen der Lognormalverteilung mit

$$X_1 \sim L(4, 35; 1, 75^2), \quad X_2 \sim L(3, 45; 1, 4^2)$$

Der Störterm des Modells  $\epsilon_{it}$  ist standardnormalverteilt. Der Individual-effekt  $\alpha$  ist mit dem zweiten Regressor korreliert und wird nach der Formel von Biørn (1996) generiert:

$$\alpha_i = (\bar{x}_{i,2} - E[x_2]) \lambda + \epsilon_{\alpha,i} \quad i = 1, \dots, N$$

mit  $\bar{x}_{i,2} = \frac{1}{T} \sum_t x_{it2}$  (Mittelwert des zweiten Regressors über die Zeit für das  $i$ -te Unternehmen),  $\lambda = 1$ ,  $\epsilon_{\alpha,i} \sim N(0, 1)$ .

In der vorliegenden Studie wird mit  $N = 1035$  Beobachtungen und  $t = 4$  Zeitpunkten gearbeitet. Sowohl die abhängige Variable als auch alle Regressoren werden mittels vaabMA anonymisiert. Die Gruppen werden mit  $A = 3$  und 5 Elementen gebildet. Die Anzahl der Monte-Carlo-Replikationen beträgt 500. Die wahren Parameter sind auf  $\beta_1 = 1,0$  und  $\beta_2 = -2,5$  gesetzt.

Der 'Within'-Schätzer berechnet mit Originaldaten (Original) wird mit dem naiven 'Within'-Schätzer (Naiv) verglichen. Neben den durchschnittlichen Schätzern<sup>8</sup> und empirischen Standardabweichungen<sup>9</sup> werden noch weitere Größen berechnet. Die durchschnittliche Verzerrung zeigt an, wie stark

<sup>8</sup>Durchschnitt der Schätzer über alle Monte-Carlo-Wiederholungen

<sup>9</sup>Standardabweichung der Schätzer über alle Monte-Carlo-Wiederholungen

der mittlere Schätzer von dem wahren Parameter abweicht. Eine geringe Verzerrung kann jedoch mit einer hohen Varianz des Schätzers einhergehen. Root Mean Squared Error (RMSE: Wurzel des mittleren quadratischen Fehlers) berücksichtigt sowohl die Verzerrung als auch die Varianz (Greene 2000, S.104). Ein weiteres Maß ist die relative Standardabweichung (RELSE). Sie wird berechnet als "Durchschnitt der Standardfehler aller Monte-Carlo-Wiederholungen geteilt durch die empirische Standardabweichung der 'Within'-Schätzer über alle Monte-Carlo-Replikationen". Die Idee ist, dass mit steigender Anzahl der Replikationen die empirische Standardabweichung gegen den wahren Standardfehler konvergiert (Lechner et al. 2003). Wenn RELSE kleiner/größer 1 ist, ist der geschätzte Standardfehler nach unten/oben verzerrt. Der relative Standardfehler macht damit die Aussagen über die Genauigkeit der geschätzten asymptotischen Standardfehler der Schätzer.<sup>10</sup>

Tabelle 3/1 zeigt die Ergebnisse für die vaabMA. Es ist zu sehen, dass die Anonymisierung der Daten mit vaabMA die 'Within'-Schätzung nicht verändert. Bei beiden Gruppengrößen sind die durchschnittliche Verzerrung und RMSE gering. Weiter fällt auf, dass RELSE sehr nahe bei 1 liegt. Das würde implizieren, dass die vaabMA zu keinem Effizienzverlust führt.

Tabelle 3/1: Variablenspezifische abstandsorientierte Mikroaggregation

<b>A = 3</b>	Ø Schätzer	Std.Abweichung	Ø Verzerrung	RMSE	RELSE
Original	1,000	0,011	0,000	0,011	0,987
	-2,501	0,012	-0,001	0,013	1,045
Naiv	0,999	0,011	-0,001	0,011	0,982
	-2,499	0,013	0,001	0,015	1,031

<b>A = 5</b>	Ø Schätzer	Std.Abweichung	Ø Verzerrung	RMSE	RELSE
Original	1,000	0,010	0,000	0,010	1,051
	-2,500	0,012	0,000	0,012	1,056
Naiv	0,998	0,010	-0,002	0,010	1,038
	-2,496	0,013	0,004	0,013	1,028

<sup>10</sup>zum RELSE siehe auch Lechner et al. 2003

## 4 Multiplikative stochastische Überlagerung

### 4.1 Allgemeine Überlagerung

Im Folgenden erfolgt die Analyse des Effekts der multiplikativ überlagerten Variablen auf die 'Within'-Schätzung.<sup>11</sup> Das lineare Panelmodell ist (2-1). Alle Regressoren und die abhängige Variable werden multiplikativ stochastisch überlagert:

$$\mathbf{X}^a = \mathbf{X} \odot \mathbf{U} \quad (4-13)$$

und

$$\mathbf{Y}^a = \mathbf{Y} \odot \mathbf{V} \quad (4-14)$$

$\odot$  bezeichnet das Hadamard-Produkt (elementenweise Multiplikation).  $\mathbf{U}$  und  $\mathbf{V}$  sind  $(NT \times K)$ -dimensionale Zufallsmatrizen.

Der Regressor ist generiert als<sup>12</sup>

$$x_{itk} \sim iid(\mu_k, \sigma_k^2), \quad i = 1, \dots, N, \quad t = 1, \dots, T, \quad k = 1, \dots, K \quad (4-15)$$

Die einzelnen  $x_{itk}$  und  $y_{it}$  sind von den Fehlervariablen für alle  $i, t, k$  stochastisch unabhängig.

In der Literatur zu Messfehlern wird häufig unterstellt, dass die Fehlervariablen miteinander unkorreliert sind. In dieser Arbeit wird lediglich die Annahme getroffen, dass die Überlagerungen zu unterschiedlichen Zeitpunkten unkorreliert sind:

$$\begin{aligned} cov(u_{itk}, v_{is}) &= 0, \quad t \neq s \\ cov(u_{itk}, u_{isl}) &= 0, \quad t \neq s, k \neq l \end{aligned} \quad \left\{ \begin{array}{l} i = 1, \dots, N, \\ t = 1, \dots, T, \\ k = 1, \dots, K. \end{array} \right. \quad (4-16)$$

Die Annahme

$$E[u_{itk}] = E[v_{it}] = 1 \quad (4-17)$$

stellt sicher, dass die anonymisierte und die Originalvariable den gleichen Erwartungswert haben.

Der Schätzer auf Basis der anonymisierten Daten wurde in (2-4) bzw. (2-6) definiert. Im Folgenden ist zu untersuchen, ob die einzelnen Elemente

<sup>11</sup>Stefanski (1985, S.586f.) hat bereits gezeigt, dass der KQ-Schätzer bei Vorliegen des multiplikativen Messfehlers inkonsistent ist, und hat einen korrigierten Schätzer hergeleitet.

<sup>12</sup>Diese Annahme kann bei Paneldaten verletzt sein, erleichtert jedoch für den Anfang die Analyse.

in (2-6) im Fall der multiplikativen stochastischen Überlagerung konsistente Schätzer der wahren Momente sind.

Zuerst haben wir die Varianz des  $k$ -ten Regressors  $S_{x_k^a}^{w,2}$  angeschaut. Sie lässt sich schreiben als<sup>13</sup>

$$\underset{N \rightarrow \infty}{plim} S_{x_k^a}^{w,2} = \left(1 - \frac{1}{T}\right) [(\sigma_k^2 + \mu_k^2) \sigma_{u_k}^2 + \sigma_k^2] \quad (4-18)$$

mit  $\sigma_k^2, \mu_k^2$  als Varianz bzw. Erwartungswert des Originalregressors  $k$  und  $\sigma_{u_k}^2$  als Varianz der Fehlervariable von  $x_k$ .

Für die 'Within'-Kovarianz  $S_{x_k^a x_l^a}^w$  zwischen zwei Regressoren  $k$  und  $l$  erhält man<sup>14</sup>

$$\underset{N \rightarrow \infty}{plim} S_{x_k^a x_l^a}^w = \left(1 - \frac{1}{T}\right) [(\sigma_{kl} + \mu_k \mu_l) \sigma_{u_k u_l} + \sigma_{kl}] \quad (4-19)$$

$\sigma_{kl}$  bezeichnet die Kovarianz der beiden wahren Regressoren  $k$  und  $l$ .  $\sigma_{u_k u_l}$  steht für die Kovarianz der Fehlervariablen der jeweiligen Regressoren. Anschließend kann man zeigen, dass<sup>15</sup>

$$\begin{aligned} & \underset{N \rightarrow \infty}{plim} S_{x_k^a y^a}^w \\ &= \left(1 - \frac{1}{T}\right) \left( \sum_l \beta_l \sigma_{kl} + \left( \sigma_{\alpha k} + \mu_\alpha \mu_k + \sum_l \beta_l (\sigma_{kl} + \mu_k \mu_l) \right) \sigma_{u_k v} \right) \end{aligned} \quad (4-20)$$

gilt.

Dabei ist  $\sigma_{u_k v}$  die Kovarianz zwischen der Fehlervariable des Originalregressors  $x_k$  und der Fehlervariable von  $y$ ,  $\sigma_{\alpha k}$  die Kovarianz des  $k$ -ten Regressores und des Individualeffekts.  $\mu_\alpha$  bezeichnet den Erwartungswert des Individualeffekts.

Während bei der 'Within'-Schätzung ohne Anonymisierung sich Individualeffekte rauskürzen lassen, wird aus (4-20) deutlich, dass im Fall der Anonymisierung durch allgemeine Überlagerung Individualeffekte bleiben.

Nachdem die einzelnen Elemente hergeleitet wurden, kann man schreiben:

$$\underset{N \rightarrow \infty}{plim} \mathbf{S}_{x^a x^a}^w \frac{1}{(1 - 1/T)}$$

<sup>13</sup>Herleitung im Appendix B.1

<sup>14</sup>Appendix B.2

<sup>15</sup>Appendix B.3

$$= \underbrace{\begin{pmatrix} \sigma_{u_1}^2(\sigma_1^2 + \mu_1^2) + \sigma_1^2 & \sigma_{u_1 u_2}(\sigma_{12} + \mu_1 \mu_2) + \sigma_{12} & \cdots & \sigma_{u_1 u_K}(\sigma_{1K} + \mu_1 \mu_K) + \sigma_{1K} \\ \sigma_{u_2 u_1}(\sigma_{21} + \mu_2 \mu_1) + \sigma_{21} & \sigma_{u_2}^2(\sigma_2^2 + \mu_2^2) + \sigma_2^2 & & \sigma_{u_2 u_K}(\sigma_{2K} + \mu_2 \mu_K) + \sigma_{2K} \\ \vdots & & \ddots & \\ \sigma_{u_K u_1}(\sigma_{K1} + \mu_K \mu_1) + \sigma_{K1} & & & \sigma_{u_K}^2(\sigma_K^2 + \mu_K^2) + \sigma_K^2 \end{pmatrix}}_{=A}, \quad (4-21)$$

$$\begin{aligned} & \underset{N \rightarrow \infty}{plim} \mathbf{S}_{x^a y^a}^w \frac{1}{(1 - 1/T)} \\ &= \begin{pmatrix} \sum_l \beta_l \sigma_{1l} + (\sigma_{\alpha 1} + \mu_{\alpha} \mu_1 + \sum_l \beta_l (\sigma_{1l} + \mu_1 \mu_l)) \sigma_{u_1 v} \\ \vdots \\ \sum_l \beta_l \sigma_{Kl} + (\sigma_{\alpha K} + \mu_{\alpha} \mu_K + \sum_l \beta_l (\sigma_{Kl} + \mu_K \mu_l)) \sigma_{u_K v} \end{pmatrix} \\ &= \underbrace{\begin{pmatrix} \sigma_{u_1 v} (\sigma_{1\alpha} + \mu_1 \mu_{\alpha}) \\ \vdots \\ \sigma_{u_K v} (\sigma_{K\alpha} + \mu_K \mu_{\alpha}) \end{pmatrix}}_{=B} + \underbrace{\begin{pmatrix} \sigma_1^2 + \sigma_{u_1 v} (\sigma_1^2 + \mu_1^2) & \cdots & \sigma_{1K} + \sigma_{u_K v} (\sigma_{1K} + \mu_1 \mu_K) \\ \vdots & & \vdots \\ \sigma_{K1} + \sigma_{u_K v} (\sigma_{K1} + \mu_K \mu_1) & \cdots & \sigma_K^2 + \sigma_{u_K v} (\sigma_K^2 + \mu_K^2) \end{pmatrix}}_{=C} \underbrace{\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_K \end{pmatrix}}_{=\beta} \end{aligned} \quad (4-22)$$

Der Wahrscheinlichkeitsgrenzwert des naiven 'Within'-Schätzers sieht dann wie folgt aus:

$$\underset{N \rightarrow \infty}{plim} \widehat{\beta}_w^a = \mathbf{A}^{-1} (\mathbf{B} + \mathbf{C}) \boldsymbol{\beta} \quad (4-23)$$

Aus (4-23) lässt sich ein korrigierter 'Within'-Schätzer herleiten:

$$\widehat{\beta}_w^{korr} = \widehat{\mathbf{C}}^{-1} (\widehat{\mathbf{A}} \widehat{\beta}_w^a - \mathbf{B}) \quad (4-24)$$

Mit Dach werden die geschätzten Matrizen/Vektoren bezeichnet.

Jetzt spielt es eine Rolle, wie die Fehlervariablen generiert werden. Bei korrelierten Fehlern, d.h.  $\sigma_{u_k v} \neq 0$ , erscheint in der Korrekturformel die Kovarianz zwischen Originalregressoren und Individualeffekten. Die Korrektur ist nur dann möglich, wenn diese Kovarianz bekannt bzw. schätzbar ist. Bei Unkorreliertheit der Fehlervariablen verschwindet  $\mathbf{B}$  und  $\mathbf{C}$  kürzt sich zur Varianz-Kovarianz-Matrix der Regressoren. Die Korrektur der Verzerrung ist damit ohne Kenntnis von  $\sigma_{\alpha k}$  möglich.

### Schätzung der Momente

In der Praxis stehen oft nur anonymisierte Daten zur Verfügung, deswegen müssen die einzelnen Momente in (4-24) geschätzt werden. Es wird jedoch davon ausgegangen, dass die Varianz-Kovarianz-Matrix der Fehlervariablen mit den Daten mitgeliefert wurde.

Wegen Erwartungstreue

$$E[x_{itk}^a] = E[x_{itk} u_{itk}] = E[x_{itk}] \quad (4-25)$$

ist der Erwartungswert des wahren Regressors durch den Mittelwert der anonymisierten Daten schätzbar:

$$\widehat{\mu}_k = \bar{x}_k^a = \frac{1}{NT} \sum_i \sum_t x_{itk}^a \quad (4-26)$$

Für die einzelnen Varianzen lässt sich herleiten:

$$\text{Var}(x_{itk}^a) = \text{Var}(x_{itk} u_{it}) = \sigma_k^2 (\sigma_{u_k}^2 + 1) + \mu_k^2 \sigma_{u_k}^2 \quad (4-27)$$

Daraus folgt<sup>16</sup>

$$\hat{\sigma}_k^2 = \frac{s_{x_k^a}^2 - \bar{x}_k^{a2} \sigma_{u_k}^2}{\sigma_{u_k}^2 + 1} \quad (4-28)$$

Für die Kovarianz zwischen den Regressoren erhält man

$$\text{Cov}(x_k^a, x_l^a) = \sigma_{kl} (\sigma_{u_k u_l} + 1) + \sigma_{u_k u_l} \mu_k \mu_l \quad (4-29)$$

Es resultiert

$$\widehat{\sigma}_{kl} = \frac{s_{x_k^a x_l^a} - \sigma_{u_k u_l} \bar{x}_k \bar{x}_l}{\sigma_{u_k u_l} + 1} \quad (4-30)$$

Bei korrelierten Fehlervariablen werden die Kovarianzen zwischen Regressoren und Individualeffekten benötigt. Bis jetzt ist uns leider noch unklar, wie sie aus den anonymisierten Daten geschätzt werden könnten.

### Simulationsergebnisse

In diesem Abschnitt werden die Ergebnisse der Simulationsstudie für den Fall unkorrelierter Überlagerungsvariablen wiedergegeben.<sup>17</sup>

Das Simulationsdesign bleibt weitgehend wie im Kapitel 3.2 bei vaabMA. Die Regressoren wurden identisch und unabhängig verteilt (iid)<sup>18</sup>

$$X_1 \sim L(4, 35; 1, 75^2), \quad X_2 \sim L(3, 45; 1, 4^2)$$

Alle Überlagerungsvariablen sind iid verteilt und wurden aus einer Lognormalverteilung mit Erwartungswert 1 erzeugt. Neben der naiven Schätzung

<sup>16</sup> $s_{x^a}^2$  ist die empirische Varianz des anonymisierten Regressors:

$$s_{x^a}^2 = \frac{1}{NT} \sum_i \sum_t (x_{it}^a - \bar{x}^a)^2$$

<sup>17</sup>Es wurden ebenfalls Simulationen für korrelierte Fehler und bekannte Kovarianz zwischen Regressor und Individualeffekt durchgeführt. Die Ergebnisse decken sich mit den theoretischen Ausführungen. Auf ihre Darstellung wird hier jedoch aus Platzgründen verzichtet.

<sup>18</sup>Eine Simulationsstudie mit normalverteilten und autokorrelierten (AR(1)-Prozess) Regressoren hat zu ähnlichen Ergebnissen geführt, die hier deswegen nicht berichtet werden.



wird der korrigierte 'Within'-Schätzer untersucht. Die Korrektur wurde sowohl mit wahren als auch mit geschätzten Momenten für Originalregressoren durchgeführt. Da zwischen den beiden kein Unterschied gefunden werden konnte, werden hier nur die Ergebnisse mit geschätzten Werten berichtet.

Tabelle 4/2 zeigt die Ergebnisse. Der naive Schätzer weist eine Verzerrung auf, die mit größerer Überlagerung stark wächst. Die Korrektur reduziert diese Verzerrung auch bei hoher Anonymisierung ( $\sigma_{Fehler} = 0,2$ ) und verbessert wesentlich RMSE. Je höher die Überlagerung, desto schlechter fällt der relative Standardfehler (RELSE) aus. Dabei erleidet der korrigierte Schätzer einen größeren Effizienzverlust als die naive Schätzung.

Tabelle 4/2: 'Within'-Schätzung bei allgemeiner Überlagerung

$\sigma_{u_1} = \sigma_{u_2} = \sigma_v = 0,06$  :

	Ø Schätzer	Std.Abweichung	Ø Verzerrung	RMSE	RELSE
Original	1,001	0,011	0,001	0,011	0,971
	-2,500	0,013	0,000	0,013	1,004
Naiv	0,975	0,013	-0,025	0,028	0,968
	-2,438	0,018	0,062	0,065	0,841
Korrigiert	1,000	0,013	0,000	0,013	0,959
	-2,500	0,019	0,000	0,019	0,821

$\sigma_{u_1} = \sigma_{u_2} = \sigma_v = 0,1$  :

	Ø Schätzer	Std.Abweichung	Ø Verzerrung	RMSE	RELSE
Original	1,000	0,010	0,000	0,010	1,002
	-2,500	0,013	0,000	0,013	0,980
Naiv	0,934	0,015	-0,066	0,068	0,994
	-2,334	0,026	0,166	0,168	0,726
Korrigiert	1,000	0,017	0,001	0,017	0,951
	-2,500	0,028	0,000	0,028	0,707

$\sigma_{u_1} = \sigma_{u_2} = \sigma_v = 0,2$  :

	Ø Schätzer	Std.Abweichung	Ø Verzerrung	RMSE	RELSE
Original	0,999	0,010	-0,001	0,010	1,069
	-2,500	0,013	0,000	0,013	0,997
Naiv	0,775	0,023	-0,225	0,226	0,961
	-1,952	0,039	0,549	0,550	0,696
Korrigiert	0,999	0,030	-0,001	0,030	0,852
	-2,506	0,053	-0,006	0,054	0,610

## 4.2 Überlagerung mit einem konstanten Grundüberlagerungsfaktor

Eine spezielle Variante der multiplikativen Überlagerung wurde von Höhne (2004) vorgeschlagen. Sie hat zwei wesentliche Vorteile gegenüber dem allgemeinen Fall. Erstens zielt das Verfahren darauf ab, die Verhältnisse zwischen den Merkmalen nach der Anonymisierung zu erhalten (Ronning 2007). Zweitens wird aufgrund der anderen Modellierung der Fehlervariablen die Korrektur der Verzerrung auch ohne Kenntnis der Kovarianz zwischen Individualeffekten und Regressoren möglich. Die folgenden Ausführungen folgen Ronning (2007).

Die Regressoren und die abhängige Variable werden multiplikativ wie folgt überlagert:

$$x_{itk}^a = x_{itk}(1 + \delta d_i + \varepsilon_{itk}) \quad (4-31)$$

$$y_{it}^a = y_{it}(1 + \delta d_i + \varepsilon_{ity}), \quad (4-32)$$

$$i = 1, \dots, N, \quad t = 1, \dots, T \text{ und } k = 1, \dots, K$$

$\delta$  ist ein Parameter, der den Abschlag oder Zuschlag zu Originalvariablen definiert. Die Variable  $d_i$  ist

$$d_i = \begin{cases} +1 & \text{mit Wahrscheinlichkeit } 0,5 \\ -1 & \text{mit Wahrscheinlichkeit } 0,5 \end{cases} \quad (4-33)$$

Ein  $i$ -tes Unternehmen erhält denselben Abschlag- oder Zuschlagparameter  $\delta$  für alle Merkmale und alle Zeitpunkte. Dadurch bleibt die Proportionalität zwischen zwei Merkmalen annähernd erhalten (Ronning 2007).<sup>19</sup> Über die gemeinsame Variable  $d_i$  sind die Fehlervariablen der Regressoren und der abhängigen Variablen miteinander korreliert.

Jede einzelne Beobachtung wird noch mit einem eigenen Zufallsfehler  $\varepsilon$  zusätzlich überlagert, wobei i.d.R.

$$\varepsilon_{itk} \sim N(0, \sigma_\varepsilon^2) \quad (4-34)$$

ist.

$d_i$ ,  $\varepsilon_{itk}$  und  $\varepsilon_{ity}$  sind miteinander und mit  $x_{itk}$  und  $y_{it}$  unkorreliert für alle  $i, t$  und  $k$ .

Weiter gilt

$$E[1 + \delta d_i + \varepsilon_{it}] = 1 \quad (4-35)$$

In der Arbeit von Ronning (2007) wird das Höhne-Verfahren ausführlich erläutert. Hier werden nur die wichtigsten Ergebnisse dargestellt.

---

<sup>19</sup>  $E[\frac{X}{Y}] \approx E[\frac{X^a}{Y^a}]$

Für den Grenzwert des naiven Schätzers ergibt sich (Ronning 2007, Formel C-13):

$$\begin{aligned}
\underset{N \rightarrow \infty}{plim} \widehat{\beta}_w^a &= \left( \underset{N \rightarrow \infty}{plim} \frac{1}{NT} (\mathbf{X}^a)' \mathbf{Q} \mathbf{X}^a \right)^{-1} \underset{N \rightarrow \infty}{plim} \frac{1}{NT} (\mathbf{X}^a)' \mathbf{Q} \mathbf{y}^a \\
&= (1 + \delta^2) \left( \begin{array}{cccc} (1 + \delta^2)\sigma_1^2 + \sigma_\varepsilon^2 (\sigma_1^2 + \mu_1^2) & (1 + \delta^2)\sigma_{12} & \dots & (1 + \delta^2)\sigma_{1K} \\ (1 + \delta^2)\sigma_{21} & (1 + \delta^2)\sigma_2^2 + \sigma_\varepsilon^2 (\sigma_2^2 + \mu_2^2) & \dots & (1 + \delta^2)\sigma_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ (1 + \delta^2)\sigma_{K1} & (1 + \delta^2)\sigma_{K2} & \dots & (1 + \delta^2)\sigma_K^2 + \sigma_\varepsilon^2 (\sigma_K^2 + \mu_K^2) \end{array} \right)^{-1} \times \\
&\quad \underbrace{\begin{pmatrix} \widehat{\sigma}_1^2 & \dots & \widehat{\sigma}_{1K} \\ \vdots & \ddots & \vdots \\ \widehat{\sigma}_{K1} & \dots & \widehat{\sigma}_K^2 \end{pmatrix}}_{=\Sigma_{\mathbf{x}\mathbf{x}}} \underbrace{\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_K \end{pmatrix}}_{=\beta} \quad (4-36) \\
&= (1 + \delta^2) \mathbf{E}^{-1} \Sigma_{\mathbf{x}\mathbf{x}} \beta
\end{aligned}$$

$\Sigma_{\mathbf{x}\mathbf{x}}$  ist die Kovarianzmatrix der Originalregressoren.

Der korrigierte 'Within'-Schätzer ist dann

$$\widehat{\beta}_w^{korr} = \frac{1}{(1 + \delta^2)} \widehat{\mathbf{S}}_{\mathbf{x}\mathbf{x}}^{-1} \widehat{\mathbf{E}} \widehat{\beta}_w^a, \quad (4-37)$$

Aufgrund der Modellierung der Fehlervariablen in (4-31) und (4-32) lassen sich im Gegensatz zur allgemeinen Überlagerung Individualeffekte bei der Durchführung der Within-Transformation eliminieren und man erhält eine Korrekturformel, in der alle Momente schätzbar sind.

### Schätzung der Momente

Da der Erwartungswert der anonymisierten Variable gleich 1 ist, wird  $\widehat{\mu}_x = \bar{x}^a$  verwendet.

Weiter wird noch der Schätzer für die Varianz der wahren Regressoren benötigt.

$$\begin{aligned}
Var(x_{itk}^a) &= E [x_{itk}^{a2}] - (E [x_{itk}^a])^2 \\
&= E [x_{itk}^2 (1 + \delta d_i + \varepsilon_{itk})^2] - (E [x_{itk} + x_{itk} \delta d_i + x_{itk} \varepsilon_{itk}])^2
\end{aligned}$$

Wegen (4-35), (4-34), Modellannahmen und  $E [d_i^2] = 1$  gelangt man zu

$$Var(x_{itk}^a) = \sigma_k^2 (1 + \delta^2 + \sigma_\varepsilon^2) + \mu_k^2 (\delta^2 + \sigma_\varepsilon^2) \quad (4-38)$$

Damit wird die wahre Varianz geschätzt mit

$$\widehat{\sigma}_k^2 = \frac{s_{k^a}^2 - (\delta^2 + \sigma_\varepsilon^2)\bar{x}_k^{a2}}{1 + \delta^2 + \sigma_\varepsilon^2} \quad (4-39)$$

Die geschätzte Kovarianz zwischen zwei Regressoren  $k$  und  $l$  wird wie folgt hergeleitet.

$$\begin{aligned} Cov(x_{itk}^a, x_{itl}^a) &= E[x_{itk}^a x_{itl}^a] - E[x_{itk}^a] E[x_{itl}^a] \\ &= E[x_{itk} x_{itl}] (1 + \delta^2) - E[x_{itk}] E[x_{itl}] \\ &= \sigma_{kl} (1 + \delta^2) + \delta^2 \mu_k^a \mu_l^a \end{aligned} \quad (4-40)$$

Unter Ausnutzung von (4-40) gelangt man zu dem Schätzer der Kovarianz:<sup>20</sup>

$$\widehat{\sigma}_{kl} = \frac{s_{k^a l^a} - \delta^2 \bar{x}_k^a \bar{x}_l^a}{(1 + \delta^2)} \quad (4-41)$$

### Simulationsergebnisse

In einer Simulationsstudie wurde die Auswirkung der Anonymisierung mit dem Höhne-Verfahren, nämlich der Effekt unterschiedlicher  $\delta$ -Parameter und Standardabweichungen der stochastischen Fehlervariable  $\varepsilon$ , auf die 'Within'-Schätzung überprüft und die Leistungsfähigkeit des Korrekturverfahrens getestet. Das Design der Simulation bleibt grundsätzlich wie im Kapitel 4.1. Die Überlagerungsvariablen wurden jedoch folgend den theoretischen Erläuterungen aus der Normalverteilung erzeugt. Trotz der Normalverteilung bleiben die Fehler bei gewählten Größen von  $\delta$  und  $\varepsilon$  im positiven Bereich.<sup>21</sup>

Die Ergebnisse sind in der Tabelle 4/3 präsentiert. Die naive Schätzung verhält sich, wie es zu erwarten war. Die Verzerrung steigt und die Effizienz sinkt mit größerer Anonymisierung. Dabei beeinträchtigt die Erhöhung von  $\sigma_\varepsilon$  stärker die Ergebnisse als die Wahl eines Zu- bzw. Abschlagsparameters  $\delta$ . Die korrigierte Schätzung gelingt gut unabhängig von der Höhe der Überlagerung, erleidet jedoch einen Effizienzverlust.

<sup>20</sup>  $s_{k^a l^a}$  ist die empirische Kovarianz der anonymisierten Regressoren.

<sup>21</sup> Die für die Anonymisierung empfohlenen Werte sind  $\delta = 0,11$  und  $\varepsilon = 0,03$ . Das würde der Standardabweichung der Fehlervariablen bei allgemeiner Überlagerung von ca. 0,11 entsprechen.

Tabelle 4/3: 'Within'-Schätzung bei Höhne-Überlagerung

 $\delta = 0,05, \sigma_{u_1} = \sigma_{u_2} = \sigma_v = 0,1 :$ 

	$\emptyset$ Schätzer	Std.Abweichung	$\emptyset$ Verzerrung	RMSE	RELSE
Original	1,001	0,010	0,001	0,010	0,997
	-2,501	-0,001	0,013	0,013	1,007
Naiv	0,933	0,016	-0,067	0,069	0,959
	-2,336	0,024	0,164	0,166	0,809
Korrigiert	1,000	0,018	0,000	0,018	0,901
	-2,501	0,026	-0,001	0,026	0,758

 $\delta = 0,11, \sigma_{u_1} = \sigma_{u_2} = \sigma_v = 0,03 :$ 

	$\emptyset$ Schätzer	Std.Abweichung	$\emptyset$ Verzerrung	RMSE	RELSE
Original	1,000	0,010	0,000	0,010	1,023
	-2,500	0,013	0,000	0,013	0,981
Naiv	0,994	0,011	-0,006	0,013	0,993
	-2,484	0,015	0,016	0,022	0,897
Korrigiert	1,000	0,011	0,000	0,011	0,987
	-2,500	0,015	0,000	0,015	0,893

 $\delta = 0,11, \sigma_{u_1} = \sigma_{u_2} = \sigma_v = 0,1 :$ 

	$\emptyset$ Schätzer	Std.Abweichung	$\emptyset$ Verzerrung	RMSE	RELSE
Original	1,000	0,010	0,000	0,010	0,986
	-2,500	0,013	-0,001	0,013	0,974
Naiv	0,933	0,016	-0,067	0,069	0,965
	-2,339	0,025	0,161	0,163	0,760
Korrigiert	1,000	0,017	0,000	0,017	0,924
	-2,503	0,028	-0,003	0,028	0,715

## 5 Zusammenfassung

Dieser Beitrag untersuchte den Effekt der Anonymisierungsverfahren *variablen-spezifische abstandsorientierte Mikroaggregation (vaabMA)* und *multiplikative stochastische Überlagerung* auf die 'Within'-Schätzung eines linearen Panelmodells.

Zurückgreifend auf die Ergebnisse von Matthias Schmid (2007) konnte gezeigt werden, dass die vaabMA die 'Within'-Schätzung nicht beeinflusst und der naive Schätzer konsistent bleibt. Wenn die Daten mittels vaabMA anonymisiert werden, könnten damit in den empirischen Analysen naive Schätzer verwendet werden.

Wird hingegen der Paneldatensatz multiplikativ stochastisch überlagert, führt die Verwendung der naiven Schätzer zu verzerrten Ergebnissen. Es ist aber möglich, einen korrigierten 'Within'-Schätzers herzuleiten. Im Fall der

allgemeinen Überlagerung mit korrelierten Fehlervariablen ist die Korrektur der Verzerrung nur dann möglich, wenn die Kovarianz zwischen Regressor und Individualeffekt bekannt ist. Die Korrektur ist ohne Kenntnis dieser Kovarianz jedoch bei allgemeiner Überlagerung mit unkorrelierten Fehlern und bei Anonymisierung nach dem Hühne-Verfahren durchführbar.

In der vorliegenden Arbeit wurde bei der Untersuchung der multiplikativen stochastischen Überlagerung von den identisch und unabhängig verteilten Regressoren ausgegangen. Da ein Paneldatensatz Informationen über die Untersuchungseinheiten zu mehreren Zeitpunkten anbietet, wird die zeitliche Korrelation eine Rolle spielen. Ein weiterer Arbeitsschritt wäre die Aufhebung der iid-Annahme.

Eine wichtige Frage ist, wie die beschriebenen Anonymisierungsverfahren Analysen mit empirischen Daten beeinflussen. Der nächste Schritt wird sich mit der Untersuchung der Auswirkungen von variablenspezifischer abstandsorientierter Mikroaggregation und multiplikativer stochastischer Überlagerung auf die 'Within'-Schätzung unter Heranziehen der dem Projekt zur Verfügung stehenden empirischen Daten befassen.

## Literatur

- [1] Biørn, E. (1996) Panel Data with Measurement Errors. In: Mátyás, L., Sevestre, P. *The Econometrics of Panel Data*, 236-279.
- [2] Greene, W. H. (2000) *Econometric Analysis*, 4th ed. Prentice Hall International.
- [3] Höhne, J. (2004) Varianten der Zufallsüberlagerung. Arbeitspapier des Projekts "Faktische Anonymisierung wirtschaftsstatistischer Einzeldaten".
- [4] Hsiao, C. (2003) *Analysis of Panel Data*, 2nd ed. Cambridge: Cambridge University Press.
- [5] Iturria, S. J., Carroll, R. J. and Firth, D. (1999) Regression and Estimating Functions in the Presence of Multiplicative Measurement Error. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 61, 547 – 561.
- [6] Lechner, S., Pohlmeier, W.(2003) Schätzung ökonometrischer Modelle auf der Grundlage anonymisierten Daten. In: Gnoss, R. und Ronning, G. (Hrsg.), *Anonymisierung wirtschaftsstatistischer Einzeldaten*, Bd.42 von *Forum der Bundesstatistik*, Wiesbaden, 115-137.
- [7] Lechner, S.(2007) The Multiplicative Simulation-Extrapolation Approach. Center for Quantitative Methods and Survey Research. University of Konstanz. Working Paper.
- [8] Mundlak, Y.(1978) On the Pooling of Time Series and Cross-Section Data. *Econometrica*, 46, 69-84.
- [9] Ronning, G., Sturm, R., Hoehne, J., Lenz, R., Rosemann, M., Scheffler, M., und Vorgrimler, D. (2005) *Handbuch zur Anonymisierung wirtschaftsstatistischer Mikrodaten*. Statistisches Bundesamt: Statistik und Wissenschaft, Band 4.
- [10] Ronning, G. (2007) IAW-Diskussionspapiere. Discussion Paper 30. Stochastische Überlagerung mit Hilfe der Mischungsverteilung. April 2007.
- [11] Rosemann, M. (2006) Auswirkungen datenverändernder Anonymisierungsverfahren auf die Analyse von Mikrodaten. *IAW-Forschungsberichte*, 66.
- [12] Schmid, M., Schneeweiß, H. (2005) The Effect of Microaggregation Procedures on the Estimation of Linear Models: A Simulation Study. In: *Econometrics of Anonymized Micro Data*, ed. by W. Pohlmeier, G. Ronning, and J. Wagner. *Jahrbücher für Nationalökonomie und Statistik*, 225, No. 5, Stuttgart: Lucius & Lucius, 529-543.

- [13] Schmid, M. (2006a) Estimation of a Linear Model under Microaggregation by Individual Ranking. Allgemeines Statistisches Archiv, 90, 419-438.
- [14] Schmid, M. (2006b) Estimation of a Linear Regression with Microaggregated Data. München : Dr. Hut.
- [15] Stefanski, L. (1985) The Effects of Measurement Error on Parameter Estimation. Biometrika, 72, 583-592.

## A Variablenspezifische abstandsorientierte Mikroaggregation

Die empirische 'Within'-Kovarianz  $S_{x_k^a x_l^a}^w$  kann umgeschrieben werden als

$$\begin{aligned}
S_{x_k^a x_l^a}^w &= \frac{1}{N} \sum_i \frac{1}{T} \sum_t (x_{itk}^a - \bar{x}_{ik}^a) (x_{itl}^a - \bar{x}_{il}^a) \\
&= \frac{1}{N} \sum_i \left( \frac{1}{T} \sum_t x_{itk}^a x_{itl}^a - \bar{x}_{il}^a \frac{1}{T} \sum_t x_{itk}^a - \bar{x}_{ik}^a \frac{1}{T} \sum_t x_{itl}^a + \frac{1}{T} \sum_t \bar{x}_{ik}^a \bar{x}_{il}^a \right) \\
&= \frac{1}{N} \sum_i \left( \frac{1}{T} \sum_t x_{itk}^a x_{itl}^a - \bar{x}_{ik}^a \bar{x}_{il}^a \right) \\
&= \frac{1}{N} \sum_i \left( \frac{1}{T} \sum_t x_{itk}^a x_{itl}^a - \frac{1}{T} \sum_t x_{itk}^a \frac{1}{T} \sum_s x_{isl}^a \right) \\
&= \frac{1}{T} \sum_t \left( \frac{1}{N} \sum_i x_{itk}^a x_{itl}^a \right) - \frac{1}{T^2} \sum_t \sum_s \left( \frac{1}{N} \sum_i x_{itk}^a x_{isl}^a \right) \quad (\text{A-42})
\end{aligned}$$

Die empirische Gesamtkovarianz zu festen Zeitpunkten  $t$  und  $s$  ist

$$S_{x_{tk}^a x_{sl}^a} = \frac{1}{N} \sum_i x_{itk}^a x_{isl}^a - \bar{x}_{tk}^a \bar{x}_{sl}^a, \quad (\text{A-43})$$

wobei  $t = s$  sein kann und  $\bar{x}_{tk}^a = \frac{1}{N} \sum_i x_{itk}^a$  bzw.  $\bar{x}_{tl}^a = \frac{1}{N} \sum_i x_{itl}^a$  sind.  
Einsetzen von (A-43) in (A-42) resultiert in

$$S_{x_k^a x_l^a}^w = \frac{1}{T} \sum_t (S_{x_{tk}^a x_{tl}^a} + \bar{x}_{tk}^a \bar{x}_{tl}^a) - \frac{1}{T^2} \sum_t \sum_s (S_{x_{tk}^a x_{sl}^a} + \bar{x}_{tk}^a \bar{x}_{sl}^a) \quad (\text{A-44})$$

Auf analoge Weise lassen sich  $S_{x_k^a}^w$  und  $S_{x_k^a y^a}^w$  herleiten.



## B Multiplikative stochastische Überlagerung

### B.1 'Within'-Varianz eines Regressors

Die 'Within'-Varianz des  $k$ -ten anonymisierten Regressors ist

$$\begin{aligned}
 S_{x_k}^{w^2} &= \frac{1}{NT} \sum_i \sum_t (x_{itk}^a - \bar{x}_{ik}^a)^2 \\
 &= \frac{1}{NT} \sum_i \sum_t x_{itk}^{a^2} - \frac{1}{N} \sum_i \left( 2 \frac{1}{T} \sum_t x_{itk}^a \bar{x}_{ik}^a - \frac{1}{T} \sum_t \bar{x}_{ik}^{a^2} \right) \\
 &= \frac{1}{NT} \sum_i \sum_t x_{itk}^{a^2} - \frac{1}{N} \sum_i \bar{x}_{ik}^{a^2} \tag{B-45}
 \end{aligned}$$

mit

$$\bar{x}_{ik}^a = \frac{1}{T} \sum_t x_{itk}^a$$

Wegen Umformung in (B-45) und unter Verwendung der Regeln für das Rechnen mit Wahrscheinlichkeitsgrenzwerten<sup>22</sup> ergibt sich

$$\begin{aligned}
 & \underset{N \rightarrow \infty}{plim} \frac{1}{N} \sum_i \frac{1}{T} \sum_t (x_{itk}^a - \bar{x}_{ik}^a)^2 \\
 &= \underset{N \rightarrow \infty}{plim} \frac{1}{N} \sum_i \frac{1}{T} \sum_t x_{itk}^{a^2} - \underset{N \rightarrow \infty}{plim} \frac{1}{N} \sum_i \bar{x}_{ik}^{a^2} \tag{B-46}
 \end{aligned}$$

Da der Stichprobenmittelwert gegen das theoretische Moment in Wahrscheinlichkeit konvergiert, müssen die Erwartungswerte der beiden Ausdrücke auf der rechten Seite von (B-46) ausgerechnet werden.

$$E \left[ \frac{1}{T} \sum_t x_{itk}^{a^2} \right] = \text{var}(x_{itk}^a) + (E[x_{itk}^a])^2 \tag{B-47}$$

$$E[\bar{x}_{ik}^{a^2}] = \frac{\text{var}(x_{itk}^a)}{T} + (E[x_{itk}^a])^2 \tag{B-48}$$

Unter Ausnutzung der Annahmen

$$x_{itk} \sim iid(\mu_k, \sigma_k^2) \quad \text{und} \quad u_{itk} \sim iid(1, \sigma_{u_k}^2)$$

erhält man weiter

$$\underset{N \rightarrow \infty}{plim} \frac{1}{N} \sum_i \frac{1}{T} \sum_t (x_{itk}^a - \bar{x}_{ik}^a)^2 = \left( 1 - \frac{1}{T} \right) \text{var}(x_{itk}^a) \tag{B-49}$$

---

<sup>22</sup>siehe z.B. Greene 2000, S.113

$$\begin{aligned}
&= \left(1 - \frac{1}{T}\right) \left(E \left[ (x_{itk} u_{itk})^2 \right] - (E [x_{itk} u_{itk}])^2\right) \\
&= \left(1 - \frac{1}{T}\right) (\sigma_k^2 + (\sigma_k^2 + \mu_k^2) \sigma_{u_k}^2), \tag{B-50}
\end{aligned}$$

wobei  $\sigma_k^2$  und  $\mu_k^2$  die Varianz und der Erwartungswert des  $k$ -ten Regressors sind.  $\sigma_{u_k}^2$  bezeichnet die Varianz der Fehlervariable des  $k$ -ten Regressors.

## B.2 'Within'-Kovarianz zweier Regressoren

Die 'Within'-Kovarianz zwischen zwei anonymisierten Regressoren  $k$  und  $l$  ist definiert als

$$\begin{aligned}
S_{x_k x_l}^w &= \frac{1}{NT} \sum_i \sum_t (x_{itk}^a - \bar{x}_{ik}^a) (x_{itl}^a - \bar{x}_{il}^a) \\
&= \frac{1}{NT} \sum_i \sum_t x_{itk}^a x_{itl}^a - \frac{1}{N} \sum_i \left( \bar{x}_{il}^a \frac{1}{T} \sum_t x_{itk}^a - \bar{x}_{ik}^a \frac{1}{T} \sum_t x_{itl}^a + \frac{1}{T} \sum_t \bar{x}_{ik}^a \bar{x}_{il}^a \right) \\
&= \frac{1}{NT} \sum_i \sum_t x_{itk}^a x_{itl}^a - \frac{1}{N} \sum_i \bar{x}_{ik}^a \bar{x}_{il}^a \tag{B-51}
\end{aligned}$$

Damit ist der folgende Grenzwert zu bestimmen:

$$\begin{aligned}
&\underset{N \rightarrow \infty}{plim} \frac{1}{N} \sum_i \frac{1}{T} \sum_t (x_{itk}^a - \bar{x}_{ik}^a) (x_{itl}^a - \bar{x}_{il}^a) \\
&= \underset{N \rightarrow \infty}{plim} \frac{1}{N} \sum_i \frac{1}{T} \sum_t x_{itk}^a x_{itl}^a - \underset{N \rightarrow \infty}{plim} \frac{1}{N} \sum_i \bar{x}_{ik}^a \bar{x}_{il}^a \tag{B-52}
\end{aligned}$$

Für die Erwartungswerte von (B-52) lässt sich schreiben

$$\begin{aligned}
E \left[ \frac{1}{T} \sum_t x_{itk}^a x_{itl}^a \right] &= cov(x_{itk}^a, x_{itl}^a) + E[x_{itk}^a] E[x_{itl}^a], \tag{B-53} \\
E[\bar{x}_{ik}^a \bar{x}_{il}^a] &= cov(\bar{x}_{ik}^a, \bar{x}_{il}^a) + E[\bar{x}_{ik}^a] E[\bar{x}_{il}^a]
\end{aligned}$$

$$\begin{aligned}
&= cov \left( \frac{1}{T} \sum_t x_{itk}^a, \frac{1}{T} \sum_s x_{its}^a \right) + E \left[ \frac{1}{T} \sum_t x_{itk}^a \right] E \left[ \frac{1}{T} \sum_t x_{itl}^a \right] \\
&= \frac{1}{T^2} \sum_t \sum_s cov(x_{itk}^a, x_{its}^a) + E \left[ \frac{1}{T} \sum_t x_{itk}^a \right] E \left[ \frac{1}{T} \sum_t x_{itl}^a \right]
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{T^2} \left( \sum_t \text{cov}(x_{itk}^a, x_{itl}^a) + \sum_t \sum_{\substack{s \\ t \neq s}} \text{cov}(x_{itk}^a, x_{isl}^a) \right) + E \left[ \frac{1}{T} \sum_t x_{itk}^a \right] E \left[ \frac{1}{T} \sum_t x_{itl}^a \right] \\
&\stackrel{iid}{=} \frac{1}{T^2} \sum_t \text{cov}(x_{itk}^a, x_{itl}^a) + E[x_{itk}^a] E[x_{itl}^a] \\
&= \frac{\text{cov}(x_{itk}^a, x_{itl}^a)}{T} + E[x_{itk}^a] E[x_{itl}^a] \tag{B-54}
\end{aligned}$$

(B-54) wurde unter der Annahme hergeleitet, dass die Variablen über alle  $i$  und  $t$  identisch und unabhängig verteilt sind.

Abziehen (B-54) von (B-53) führt zu

$$\begin{aligned}
\text{plim}_{N \rightarrow \infty} \frac{1}{N} \sum_i \frac{1}{T} \sum_t (x_{itk}^a - \bar{x}_{ik}^a) (x_{itl}^a - \bar{x}_{il}^a) &= \left(1 - \frac{1}{T}\right) \text{cov}(x_{itk}^a, x_{itl}^a) \\
&= \left(1 - \frac{1}{T}\right) (E[x_{itk} u_{itk} x_{itl} u_{itl}] - E[x_{itk} u_{itk}] E[x_{itl} u_{itl}]) \\
&= \left(1 - \frac{1}{T}\right) (\sigma_{kl} + (\sigma_{kl} + \mu_k \mu_l) \sigma_{u_k u_l}) \tag{B-55}
\end{aligned}$$

mit  $\sigma_{kl}$  als Kovarianz der Originalregressoren  $k$  und  $l$  und  $\sigma_{u_k u_l}$  Kovarianz der Fehlervariablen von  $X_k$  und  $X_l$ .

### B.3 'Within'-Kovarianz eines Regressors und der abhängigen Variable

Des Weiteren wird der Grenzwert von  $S_{x_k y}^w$  benötigt.

Die Differenzen der einzelnen Beobachtungen von den Mittelwerten über die Zeit kann man schreiben als

$$x_{itk}^a - \bar{x}_{ik}^a = x_{itk} u_{itk} - \bar{x}_{ik} u_{ik},$$

$$y_{it}^a - \bar{y}_i^a = \alpha_i (v_{it} - \bar{v}_i) + \sum_{l=1}^K \beta_l (x_{itl} v_{it} - \bar{x}_{il} \bar{v}_i) + (\epsilon_{it} v_{it} - \bar{\epsilon}_i \bar{v}_i)$$

mit  $\bar{x}_{ik} = \frac{1}{T} \sum_t x_{itk} u_{itk}$ ,  $\bar{v}_i = \frac{1}{T} \sum_t v_{it}$ ,  $\bar{x}_{il} \bar{v}_i = \frac{1}{T} \sum_t x_{itl} v_{it}$  und  $\bar{\epsilon}_i \bar{v}_i = \frac{1}{T} \sum_t \epsilon_{it} v_{it}$ .

Anschließend ist zu berechnen

$$\begin{aligned}
& plim_{N \rightarrow \infty} \frac{1}{N} \sum_i \frac{1}{T} \sum_t (x_{itk}^a - \bar{x}_{ik}^a) (y_{it}^a - \bar{y}_i^a) \\
&= plim_{N \rightarrow \infty} \frac{1}{N} \sum_i \frac{1}{T} \sum_t [\alpha_i (x_{itk} u_{itk} - \bar{x} u_{ik}) (v_{it} - \bar{v}_i)] \\
&+ plim_{N \rightarrow \infty} \frac{1}{N} \sum_i \frac{1}{T} \sum_t \left[ (x_{itk} u_{itk} - \bar{x} u_{ik}) \sum_l \beta_l (x_{itl} v_{it} - \bar{x}_l \bar{v}_i) \right] \\
&+ plim_{N \rightarrow \infty} \frac{1}{N} \sum_i \frac{1}{T} \sum_t [(x_{itk} u_{itk} - \bar{x} u_{ik}) (\epsilon_{it} v_{it} - \bar{\epsilon} \bar{v}_i)] \tag{B-56}
\end{aligned}$$

Analog zur Herleitung der 'Within'-Kovarianz zwischen zwei Regressoren ergibt sich

$$\begin{aligned}
plim_{N \rightarrow \infty} \frac{1}{N} \sum_i \frac{1}{T} \sum_t [(\alpha_i x_{itk} u_{itk} - \bar{\alpha}_i \bar{x} u_{ik}) (v_{it} - \bar{v}_i)] &= \left(1 - \frac{1}{T}\right) cov(\alpha_i x_{itk} u_{itk}, v_{it}) \\
&= \left(1 - \frac{1}{T}\right) (E[\alpha_i x_{itk} u_{itk} v_{it}] - E[\alpha_i x_{itk} u_{itk}] E[v_{it}]) \\
&= \left(1 - \frac{1}{T}\right) (\sigma_{\alpha k} + \mu_\alpha \mu_k) \sigma_{u_k v} \tag{B-57}
\end{aligned}$$

$\sigma_{\alpha k}$  ist die Kovarianz zwischen Individualeffekt und Regressor  $k$ .  $\mu_\alpha$  bezeichnet den Erwartungswert der Individualeffekte.

$$\begin{aligned}
& plim_{N \rightarrow \infty} \frac{1}{N} \sum_i \frac{1}{T} \sum_t \left[ (x_{itk} u_{itk} - \bar{x} u_{ik}) \sum_l \beta_l (x_{itl} v_{it} - \bar{x}_l \bar{v}_i) \right] \\
&= \left(1 - \frac{1}{T}\right) cov\left(x_{itk} u_{itk}, \sum_l \beta_l x_{itl} v_{it}\right) \\
&= \left(1 - \frac{1}{T}\right) \left( E\left[ x_{itk} u_{itk} \sum_l \beta_l x_{itl} v_{it} \right] - E[x_{itk} u_{itk}] E\left[ \sum_l \beta_l x_{itl} v_{it} \right] \right) \\
&= \left(1 - \frac{1}{T}\right) \left( \sum_l \beta_l \sigma_{kl} + \sum_l \beta_l (\sigma_{kl} + \mu_k \mu_l) \sigma_{u_k v} \right) \tag{B-58}
\end{aligned}$$

Anschließend erhält man

$$plim_{N \rightarrow \infty} \frac{1}{N} \sum_i \frac{1}{T} \sum_t [(x_{itk} u_{itk} - \bar{x} u_{ik}) (\epsilon_{it} v_{it} - \bar{\epsilon} \bar{v}_i)]$$

$$= \left(1 - \frac{1}{T}\right) \text{cov}(x_{itk}u_{itk}, \epsilon_{it}v_{it}) = 0 \quad (\text{B-59})$$

Hier wird ausgenutzt, dass der Störterm  $\epsilon_{it}$ , der Regressor  $x_{itk}$  und die Fehlervariablen  $u_{itk}$  sowie  $v_{it}$  stochastisch unabhängig sind und  $E[\epsilon_{it}] = 0$  gilt.

Die in (B-57), (B-58) und (B-59) hergeleiteten Ergebnisse führen zu dem Grenzwert<sup>23</sup>

$$\begin{aligned} & \underset{N \rightarrow \infty}{\text{plim}} \frac{1}{N} \sum_i \frac{1}{T} \sum_t (x_{itk}^a - \bar{x}_{ik}^a) (y_{it}^a - \bar{y}_i^a) \\ &= \left(1 - \frac{1}{T}\right) \left( \sum_l \beta_l \sigma_{kl} + \left( \sigma_{\alpha k} + \mu_\alpha \mu_k + \sum_l \beta_l (\sigma_{kl} + \mu_k \mu_l) \right) \sigma_{u_k v} \right) \end{aligned} \quad (\text{B-60})$$

---

<sup>23</sup>Wird ein Panelmodell mit einem Absolutglied definiert

$$y_{it} = c + \alpha_i + \sum_l x_{itl} \beta_l + \epsilon_{it}$$

mit  $\mu_\alpha = 0$ ,  
so erhält man statt (B-60)

$$\left(1 - \frac{1}{T}\right) \left( \sum_l \beta_l \sigma_{kl} + \left( \sigma_{\alpha k} + c \mu_k + \sum_l \beta_l (\sigma_{kl} + \mu_k \mu_l) \right) \sigma_{u_k v} \right)$$

# IAW-Diskussionspapiere

Bisher erschienen:

- Nr. 1 (September 2001)  
Das Einstiegsgeld – eine zielgruppenorientierte negative Einkommensteuer: Konzeption, Umsetzung und eine erste Zwischenbilanz nach 15 Monaten in Baden-Württemberg  
*Sabine Dann / Andrea Kirchmann / Alexander Spermann / Jürgen Volkert*
- Nr. 2 (Dezember 2001)  
Die Einkommensteuerreform 1990 als natürliches Experiment. Methodische und konzeptionelle Aspekte zur Schätzung der Elastizität des zu versteuernden Einkommens  
*Peter Gottfried / Hannes Schellhorn*
- Nr. 3 (Januar 2001)  
Gut betreut in den Arbeitsmarkt? Eine mikroökonomische Evaluation der Mannheimer Arbeitsvermittlungsagentur  
*Jürgen Jergler / Christian Pohnke / Alexander Spermann*
- Nr. 4 (Dezember 2001)  
Das IAW-Einkommenspanel und das Mikrosimulationsmodell SIMST  
*Peter Gottfried / Hannes Schellhorn*
- Nr. 5 (April 2002)  
A Microeconometric Characterisation of Household Consumption Using Quantile Regression  
*Niels Schulze / Gerd Ronning*
- Nr. 6 (April 2002)  
Determinanten des Überlebens von Neugründungen in der baden-württembergischen Industrie – eine empirische Survivalanalyse mit amtlichen Betriebsdaten  
*Harald Strotmann*
- Nr. 7 (November 2002)  
Die Baulandausweisungsumlage als ökonomisches Steuerungsinstrument einer nachhaltigkeitsorientierten Flächenpolitik  
*Raimund Krumm*
- Nr. 8 (März 2003)  
Making Work Pay: U.S. American Models for a German Context?  
*Laura Chadwick, Jürgen Volkert*

# IAW-Diskussionspapiere

- Nr. 9 (Juni 2003)  
Erste Ergebnisse von vergleichenden Untersuchungen mit anonymisierten und nicht anonymisierten Einzeldaten am Beispiel der Kostenstrukturerhebung und der Umsatzsteuerstatistik  
*Martin Rosemann*
- Nr. 10 (August 2003)  
Randomized Response and the Binary Probit Model  
*Gerd Ronning*
- Nr. 11 (August 2003)  
Creating Firms for a New Century: Determinants of Firm Creation around 1900  
*Joerg Baten*
- Nr. 12 (September 2003)  
Das fiskalische BLAU-Konzept zur Begrenzung des Siedlungsflächenwachstums  
*Raimund Krumm*
- Nr. 13 (Dezember 2003)  
Generelle Nichtdiskontierung als Bedingung für eine nachhaltige Entwicklung?  
*Stefan Bayer*
- Nr. 14 (Februar 2003)  
Die Elastizität des zu versteuernden Einkommens. Messung und erste Ergebnisse zur empirischen Evidenz für die Bundesrepublik Deutschland.  
*Peter Gottfried / Hannes Schellhorn*
- Nr. 15 (Februar 2004)  
Empirical Evidence on the Effects of Marginal Tax Rates on Income – The German Case  
*Peter Gottfried / Hannes Schellhorn*
- Nr. 16 (Juli 2004)  
Shadow Economies around the World: What do we really know?  
*Friedrich Schneider*
- Nr. 17 (August 2004)  
Firm Foundations in the Knowledge Intensive Business Service Sector. Results from a Comparative Empirical Study in Three German Regions  
*Andreas Koch / Thomas Stahlecker*

# IAW-Diskussionspapiere

- Nr. 18 (Januar 2005)  
The impact of functional integration and spatial proximity on the post-entry performance of knowledge intensive business service firms  
*Andreas Koch / Harald Strotmann*
- Nr. 19 (März 2005)  
Legislative Malapportionment and the Politicization of Germany's Intergovernmental Transfer System  
*Hans Pitlik / Friedrich Schneider / Harald Strotmann*
- Nr. 20 (April 2005)  
Implementation ökonomischer Steuerungsansätze in die Raumplanung  
*Raimund Krumm*
- Nr. 21 (Juli 2005)  
Determinants of Innovative Activity in Newly Founded Knowledge Intensive Business Service Firms  
*Andreas Koch / Harald Strotmann*
- Nr. 22 (Dezember 2005)  
Impact of Opening Clauses on Bargained Wages  
*Wolf Dieter Heinbach*
- Nr. 23 (Januar 2006)  
Hat die Einführung von Gewinnbeteiligungsmodellen kurzfristige positive Produktivitätswirkungen? – Ergebnisse eines Propensity-Score-Matching-Ansatzes  
*Harald Strotmann*
- Nr. 24 (März 2006)  
Who Goes East? The Impact of Enlargement on the Pattern of German FDI  
*Claudia M. Buch / Jörn Kleinert*
- Nr. 25 (Mai 2006)  
Estimation of the Probit Model from Anonymized Micro Data  
*Gerd Ronning / Martin Rosemann*
- Nr. 26 (Oktober 2006)  
Bargained Wages in Decentralized Wage-Setting Regimes  
*Wolf Dieter Heinbach*



# IAW-Diskussionspapiere

- Nr. 27 (Januar 2007)  
A Capability Approach for Official German Poverty and Wealth Reports:  
Conceptual Background and First Empirical Results  
*Christian Arndt / Jürgen Volkert*
- Nr. 28 (Februar 2007)  
Typisierung der Tarifvertragslandschaft – Eine Clusteranalyse der tarifvertraglichen  
Öffnungsklauseln  
*Wolf Dieter Heinbach / Stefanie Schröpfer*
- Nr. 29 (März 2007)  
International Bank Portfolios: Short- and Long-Run Responses to the Business Cycles  
*Sven Blank / Claudia M. Buch*
- Nr. 30 (April 2007)  
Stochastische Überlagerungen mit Hilfe der Mischungsverteilung  
*Gerd Ronning*
- Nr. 31 (Mai 2007)  
Openness and Growth: The Long Shadow of the Berlin Wall  
*Claudia M. Buch / Farid Toubal*
- Nr. 32 (Mai 2007)  
International Banking and the Allocation of Risk  
*Claudia M. Buch / Gayle DeLong / Katja Neugebauer*
- Nr. 33 (Juli 2007)  
Multinational Firms and New Protectionisms  
*Claudia M. Buch / Jörn Kleinert*
- Nr. 34 (November 2007)  
Within-Schätzung bei anonymisierten Paneldaten  
*Elena Biewen*