
Methodenbericht

Modellbasierte Datenerganzung der Konjunkturstatistik im Produzierenden Bereich

Gerlinde Dinges
Fachbereich Methodik
Statistik Austria

Martin Haitzmann
Unternehmensstatistik
Statistik Austria

Abstract

Die im Jahr 1996 eingefuhrt EU-harmonisierte monatliche Konjunkturstatistik im Produzierenden Bereich stellt eine der zentralen Informationsquellen zur Beurteilung der konjunkturellen Entwicklung osterreichs und des gesamten Europaischen Wirtschafts- und Wahrungsraums dar. Ihre nationalen und internationalen Verwendungszwecke sind vielfaltig. Dem Anliegen der Wirtschaft entsprechend wird diese Statistik als Vollerhebung mit variablen Abschneidegrenzen (Schwellenwerte) unter Berucksichtigung eines standardisierten Reprasentanzkriteriums gefuhrt, wodurch Klein- und Kleinstunternehmen (rund 80% der Grundgesamtheit) von statistischen Verpflichtungen weitestgehend befreit werden. Um daruber hinaus auch den Bedurfnissen der Datennutzer zu entsprechen und EU-Vorgaben hinsichtlich vorgegebener Reprasentanzkriterien vollstandig erfullen zu konnen, erfolgt ab 2009 zusatzlich zur bisherigen Veroffentlichung primarstatistisch erfasster Ergebnisse der groeren meldepflichtigen Einheiten auch eine entsprechende Ergebnisdarstellung der Wirtschaftsleistung aller Unternehmen und Betriebe, die im Unternehmensregister der STATISTIK AUSTRIA als aktive Einheiten der Grundgesamtheit des Produzierenden Bereichs erfasst sind. Der folgende Beitrag ist vor allem der methodischen Konzeption des Modells zur Datenerganzung und dessen Umsetzung in der Praxis gewidmet, wobei Ergebnisse und Datenqualitat auch aus Sicht der Unternehmensstatistikexperten beleuchtet werden. Zur praktischen Veranschaulichung dienen ausgewahlte Wirtschaftszweige und Merkmale der Konjunkturstatistik 2009.

1. Einleitung

Die Konjunkturstatistik im Produzierenden Bereich (KJP) besteht in dieser Form seit 1996 und wird nach einem Konzept durchgefuhrt, das den EU-Vorgaben entspricht. Sie dient dazu, Informationen uber kurzfristig beobachtbare und messbare Phanomene im Bereich der Sachguterproduktion und im Bauwesen zu sammeln und zu verarbeiten. Die Ergebnisse der KJP bilden - bei Anwendung geeigneter statistischer Methoden - eine der wesentlichen Grundlagen fur die Beobachtung des Wirtschaftszyklus und liefern politischen Entscheidungstragern auf nationaler wie auch europaischer Ebene sowie fur die Europaische Zentralbank wichtige empirische Evidenz. Daruber hinaus bietet die KJP den Marktproduzenten selbst harmonisierte Informationen zum Verstandnis ihrer Markte und zum Vergleich ihrer Tatigkeit und Leistung mit Wettbewerbern desselben Wirtschaftszweigs auf nationaler und internationaler Ebene.

Dem Anliegen der Wirtschaft entsprechend wird die KJP als **Vollerhebung mit variablen Abschneidegrenzen** unter Berucksichtigung eines standardisierten Representanzkriteriums gefuhrt. Diese Erhebungsform wird auch als „**Konzentrationsstichprobe**“ bezeichnet. Auskunftspflicht besteht nur bei berschreitung normierter Schwellenwerte (grundsatzlich Beschaftigungsschwelle, bei Nichterreichung eines bestimmten Representanzgrades zusatzlich Auskunftspflicht unterhalb der Beschaftigungsschwelle bei berschreitung einer vorgegebenen Umsatzschwelle), wodurch Klein- und Kleinstunternehmen von statistischen Verpflichtungen weitestgehend ausgenommen sind.

Um den Bedurfnissen der Datennutzer zu entsprechen und EU-Vorgaben hinsichtlich vorgegebener **Representanzkriterien** vollstandig zu erfullen, erfolgt ab 2009 neben der bisherigen Veroffentlichung der primarstatistisch erfassten Ergebnisse fur die groeren Unternehmen und Betriebe auch eine entsprechende Ergebnisdarstellung fur die statistische Grundgesamtheit des Produzierenden Bereichs.

Da die vorgegebenen rechtlichen und daraus resultierenden konzeptiven Rahmenbedingungen keine konventionelle Hochrechnung zulassen,¹ wurde ein **modellbasierter Ansatz zur Datenerganzung** entwickelt, welcher den Informationsvorrat der Primarerhebung der KJP ausschopft und zusatzlich Informationen aus Verwaltungsquellen effizient nutzt.

Ein ahnlicher Weg zur Berechnung der Grundgesamtheit wird seit mehreren Jahren bereits in den jahrlichen Leistungs- und Strukturstatistiken² durch Einbindung von einerseits **primarstatistischen Ergebnissen**³ und andererseits **sekundarstatistischen Datenquellen**⁴ beschritten. Die ungleich hohere Komplexitat der KJP gestaltet die Erstellung eines plausiblen Datenkorpers fur die statistische Grundgesamtheit jedoch besonders schwierig, denn im Gegensatz zu den jahrlichen Leistungs- und Strukturstatistiken mussen in der KJP kurzfristige monatliche Ergebnisse vorliegen, und daruber hinaus ist zusatzlich zum Aktivitatsansatz (Ergebnisdarstellung nach der schwerpunktmaigen wirtschaftlichen Aktivitat der Unternehmen und Betriebe) eine Darstellung nach dem Guteransatz (Ergebnisdarstellung aller gleichartig klassifizierten Guter) erforderlich.

Im vorliegenden Methodenbericht wird die Konzeption des Modells zur Datenerganzung und dessen Umsetzung in der Praxis beschrieben.

2. Monatliche Primarerhebung der KJP

Bei der KJP⁵ handelt es sich um ein Instrument mit zwei unterschiedlichen Zielsetzungen.

¹Meldepflicht besteht nur bei berschreitung normierter Schwellenwerte, folglich liefert dieses Auswahlverfahren kein representatives Sample fur eine konventionelle Hochrechnung auf die Grundgesamtheit.

²Siehe auch Standarddokumentation zur Leistungs- und Strukturstatistik im Produzierenden Bereich auf der Homepage der STATISTIK AUSTRIA unter www.statistik.at > Dokumentationen > Produktion und Bauwesen.

³Primardaten werden eigens fur die betreffende Statistik direkt beim auskunftspflichtigen Unternehmen bzw. Betrieb erhoben.

⁴Sekundardaten konnen aus Register- oder Verwaltungsquellen gewonnen werden oder stehen aus anderen statistischen Erhebungen zur Verfugung.

⁵Eine ausfuhrliche Methodenbeschreibung findet sich in der Standarddokumentation zur KJP auf der Homepage der STATISTIK AUSTRIA unter www.statistik.at > Dokumentationen > Produktion und Bauwesen.

Es sollen einerseits die empirischen Voraussetzungen für Konjunkturanalysen für den Produzierenden Bereich geschaffen und andererseits das heimische Güteraufkommen dieses Wirtschaftsbereichs abgebildet werden. Somit werden mit der monatlichen KJP die Vorgaben der **EU-Verordnung über Konjunkturstatistiken**, Anhang A (Industrie) und Anhang B (Baugewerbe),⁶ und jene zur Erfassung der **nationalen Güterproduktion im Sinne der EU-Verordnung zur Einführung einer Gemeinschaftserhebung über die Produktion von Gütern** (PRODCOM)⁷ erfüllt. Grundsätzlich erstreckt sich die KJP auf alle Unternehmen, Betriebe (fachliche Einheiten), Arbeitsgemeinschaften (ARGEN) und Betriebe gewerblicher Art von Körperschaften des öffentlichen Rechts gemäß § 2 KStG 1988,⁸ die den **Abschnitten der ÖNACE 2008**⁹

- Bergbau und Gewinnung von Steinen und Erden (Abschnitt B),
- Herstellung von Waren (Abschnitt C),
- Energieversorgung (Abschnitt D),
- Wasserversorgung; Abwasser- und Abfallentsorgung und Beseitigung von Umweltverschmutzungen (Abschnitt E) sowie
- Bauwesen (Abschnitt F)

zuzuordnen waren und die Tätigkeit selbständig, regelmäßig und in der Absicht zur Erzielung eines Ertrags oder sonstigen wirtschaftlichen Vorteils ausüben.

Dabei besteht **Auskunftspflicht** für

- alle **Unternehmen** (Ein- und Mehrbetriebsunternehmen) und **Betriebe** (fachliche Einheiten) sowie Betriebe gewerblicher Art von Körperschaften des öffentlichen Rechts gemäß § 2 KStG 1988, die am 30. September des der Berichtsperiode vorangegangenen Kalenderjahres **20 und mehr Beschäftigte** hatten,
- alle **Arbeitsgemeinschaften** (ARGEN), unabhängig von der Beschäftigtenzahl, ab ihrer Gründung bis zu ihrer Auflösung sowie
- alle im Kalenderjahr der Berichtsperiode **neu gegründeten** bzw. durch Umstrukturierung **entstandenen statistischen Einheiten**, in Abhängigkeit von der Beschäftigtenzahl zum Zeitpunkt der Neugründung bzw. Umstrukturierung.

Darüber hinaus muss die Erhebungsmasse mindestens 90% des Gesamtumsatzes in jedem der Wirtschaftszweige gemäß den Abteilungen 05 bis 43 der ÖNACE 2008 enthalten (Repräsentanz bzw. Deckungsgrad). Wird dieses **Repräsentanzkriterium** mit Hilfe der voll erhobenen

⁶Verordnung (EG) Nr. 1165/98 über Konjunkturstatistiken, ABl. Nr. L 162 vom 05.06.1998 S. 1, in der Fassung der Verordnung (EG) Nr. 1178/2008.

⁷Verordnung (EWG) Nr. 3924/91 zur Einführung einer Gemeinschaftserhebung über die Produktion von Gütern, ABl. Nr. L 374 vom 31.12.1991, S. 1, zuletzt geändert durch Verordnung (EG) Nr. 1893/2006 zur Aufstellung der statistischen Systematik der Wirtschaftszweige NACE Revision 2 und zur Änderung der Verordnung (EWG) Nr. 3037/90 des Rates sowie einiger Verordnungen der EG über bestimmte Bereiche der Statistik, ABl. Nr. L 393 vom 30.12.2006 S. 1. bzw. vergleiche auch www.statistik.at > Klassifikationen > Klassifikationsdatenbank > Güter > CPA 2008.

⁸Bundesgesetz vom 7. Juli 1988 über die Besteuerung des Einkommens von Körperschaften (Körperschaftsteuergesetz 1988 - KStG 1988), BGBl. Nr.401/1988, idgF.

⁹Vergl. Publikation der STATISTIK AUSTRIA „Systematik der Wirtschaftstätigkeiten - ÖNACE 2008; Band 1 und 2“ bzw. unter www.statistik.at > Klassifikationen > Klassifikationsdatenbank > Wirtschaftszweige > ÖNACE 2008.

Schicht nicht erreicht, so besteht Auskunftspflicht auch über statistische Einheiten mit weniger als 20 Beschäftigten, die in den zwölf Monaten, die dem Stichtag¹⁰ vorangegangen sind, in Summe einen Umsatz (exklusive Umsatzsteuer) von mindestens einer Million Euro erzielten (eingeschränkte Auskunftspflicht unterhalb der gesetzlichen Beschäftigungsschwelle).¹¹

3. Weshalb die KJP-Primärerhebung modellbasiert ergänzen?

3.1. Nichterfüllung der Repräsentanz

Das Bundesstatistikgesetz 2000 enthält die Verpflichtung, **Klein- und Kleinstunternehmen in höchstmöglicher Weise von der Auskunftspflicht auszunehmen**. Den Bestrebungen, den geforderten Abdeckungsgrad im Sinne der europäischen Vorgaben durch eine primärstatistische Erhebung zu erreichen, wird damit eine natürliche Grenze gesetzt.

Die nationale Durchführungsverordnung¹² sieht als **Repräsentanzkriterium die Erfassung von 90% des Gesamtumsatzes aller in diesem Zweig tätigen Einheiten je ÖNACE-2008-Abteilung** vor. Die EU-Verordnung über die Güterproduktion¹³ normiert den Deckungsgrad von 90% bezogen auf die Inlandsproduktion auf noch tieferer Gliederungsebene, nämlich je NACE-Rev.2-Klasse (entspricht der ÖNACE-2008-Klasse). Nach der nationalen Gesetzeslage dürfen Klein- und Kleinstunternehmen mit weniger als 20 Beschäftigten und einem Jahresumsatz von weniger als einer Million Euro selbst dann nicht zur Erfüllung des Repräsentanzkriteriums in die Erhebung einbezogen werden, wenn der branchenspezifische Gesamtumsatz der in die Erhebung einbezogenen statistischen Einheiten zur Erreichung des Deckungsgrades nicht ausreicht.

Wie aus den Abbildungen 1 und 2 ersichtlich, unterscheidet sich die Primärabdeckung der verschiedenen Wirtschaftsbereiche teils deutlich, was letztendlich auch die Vergleichbarkeit bzw. einheitliche Interpretation von Ergebnissen erschwert. Zusätzlich wird in einigen Wirtschaftsbereichen durch die primärstatistische Erhebung weder das Repräsentanzkriterium der nationalen Durchführungsverordnung (Abbildung 1) noch jenes der PRODCOM-Verordnung (Abbildung 2) erreicht.

Gemäß den europäischen Bestimmungen können jedoch die Mitgliedstaaten die erforderlichen Daten nach dem Grundsatz der verwaltungstechnischen Vereinfachung durch eine Kombination verschiedener Quellen (verbindliche Erhebungen, andere Quellen, die in Bezug auf Genauigkeit und Qualität zumindest gleichwertig sind, oder statistische Schätzverfahren) beschaffen. Dieser Weg wurde nun mit der Modellbasierten Datenergänzung im Produzierenden Bereich beschritten, mit dem Ziel, ein Maximum an Repräsentanz und damit einhergehend, ein Optimum an Qualität des Datenkörpers zu gewährleisten.

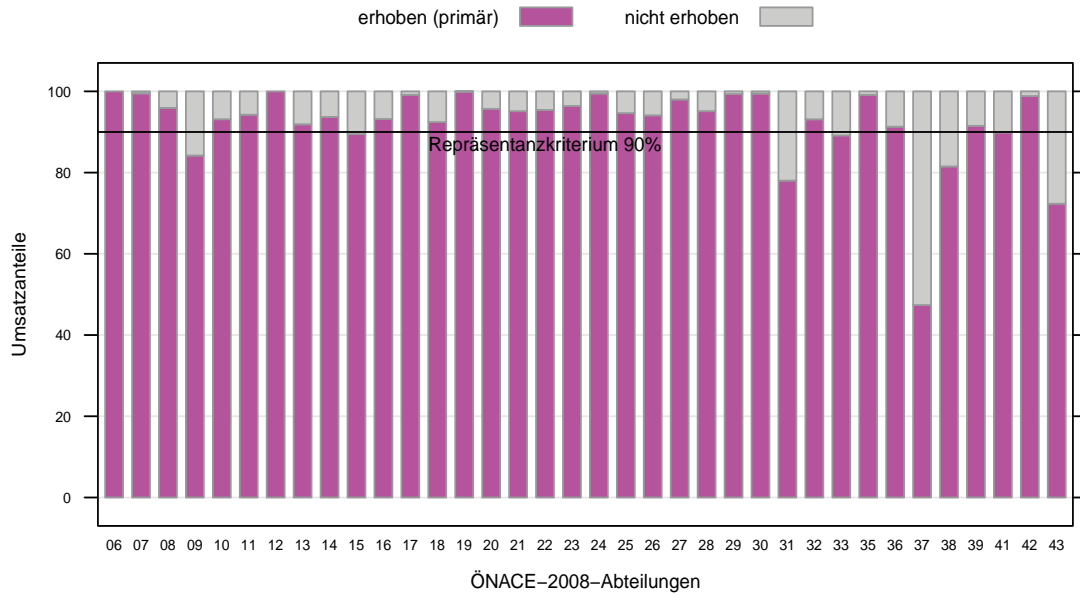
¹⁰Stichtag ist der 30. September des der Berichtsperiode vorangegangenen Kalenderjahres.

¹¹Mehr auf: www.statistik.at > Fragebögen > Unternehmen > Konjunkturerhebung im Produzierenden Bereich.

¹²BGBl. II Nr. 210/2003, zuletzt geändert durch BGBl. II Nr. 315/2007.

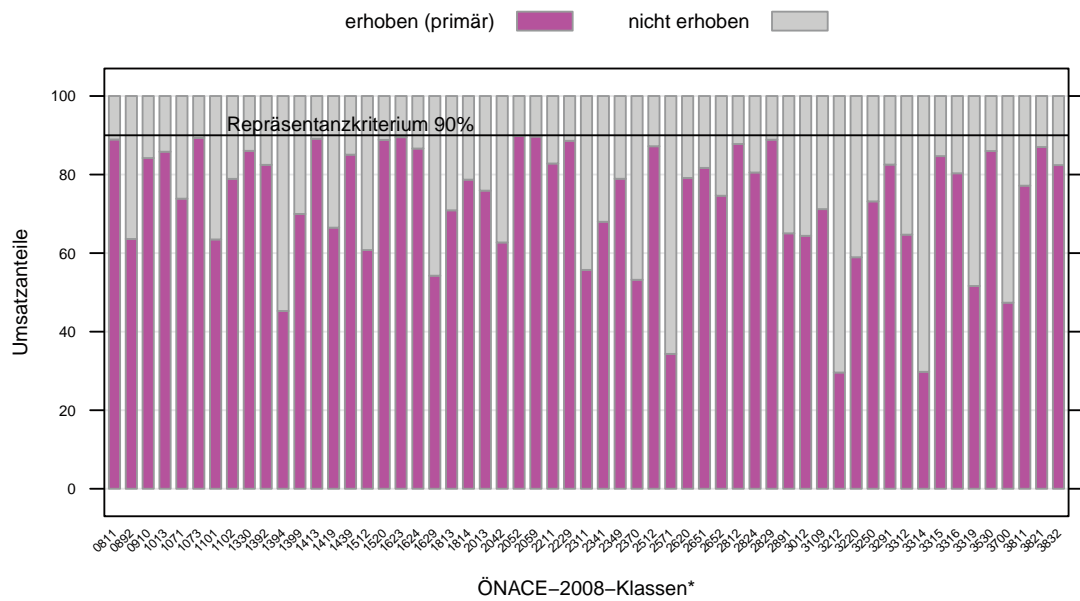
¹³Verordnung (EWG) Nr. 3924/91 des Rates vom 19. Dezember 1991 zur Einführung einer Gemeinschaftserhebung über die Produktion von Gütern, (ABl. Nr. L 196 vom 31. Dezember 1991, S.1), zuletzt geändert durch die VO (EG) Nr. 1893/2006.

Abbildung 1: Erfassung des Umsatzes im Produzierenden Bereich durch die KJP-Erhebung im 1. Quartal 2009 (Gemäß nationaler Durchführungsverordnung)



Quelle: KJP 2009 (1. Quartal)

Abbildung 2: Erfassung des Umsatzes im Produzierenden Bereich durch die KJP-Erhebung im 1. Quartal 2009 (Gemäß Europäischer Verordnung über die Güterproduktion)

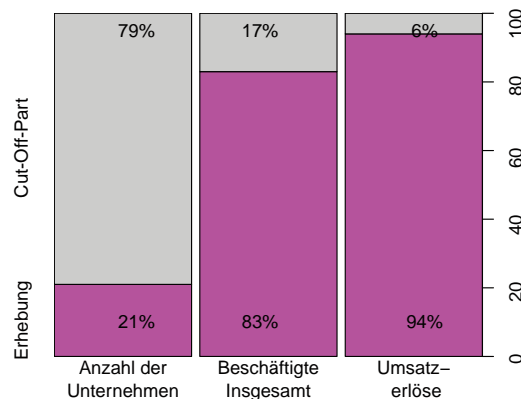


Quelle: KJP 2009 (1. Quartal). - *) Nur ÖNACE-2008-Klassen, die die Mindestrepräsentanz nicht erfüllen.

3.2. Konzeptionelle Vorgaben der Primärerhebung der KJP

Das Erhebungskonzept der KJP wird auch als **Cut-Off-Census**¹⁴ oder „Konzentrationsstichprobe“ bezeichnet, da weniger (ge-)wichtige Information nicht erfasst wird. Der „bewusste Daten-Cut-Off“ liefert Informationen, die keinem repräsentativen Sample im Sinne der Stichprobentheorie entsprechen. In der Wirtschaftstatistik sind Konzentrationserhebungen gängige Praxis. Grundidee dieser Methode ist es, die hohe Merkmalskonzentration wirtschaftsstatistischer Daten zu nutzen und anhand möglichst weniger erhobener Einheiten eine möglichst hohe primärstatistische Abdeckung jener Merkmale zu erreichen, die im Mittelpunkt des Interesses stehen. Abbildung 3 zeigt, dass im 1. Berichtsquartal 2009 nur Daten für 21% der rund 58.000 Unternehmen erhoben werden mussten, um bereits 83% der Gesamtbeschäftigten und 94% der Umsatzerlöse primärstatistisch zu erfassen. Dabei ist jedoch zu bedenken, dass, auch wenn ein bestimmtes vorgegebenes Repräsentanzkriterium erreicht werden kann, diese bewusste Beschränkung auf die „wesentlichen Elemente“ die Annahme eines „vernachlässigbaren Beitrags“ des Cut-Off-Parts impliziert. Für die Berechnung der Konjunkturindikatoren zum Beispiel ist diese Beschränkung auf die „wesentlichen Elemente“ per definitionem ausreichend, manche Analysen, für detaillierte Konjunkturprognosen etwa, erfordern aber vollständige Daten über die statistische Grundgesamtheit.

Abbildung 3: Merkmalskonzentration in der KJP-Grundgesamtheit
Primärabdeckung ausgewählter Unternehmensmerkmale



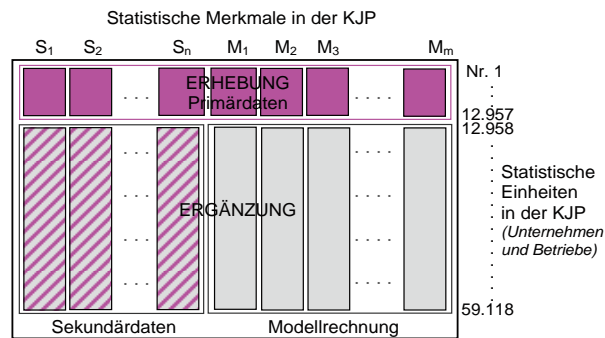
Quelle: KJP 2009 (1. Quartal).

4. Modellbasierte Datenergänzung (MDE)

Durch die **Modellbasierte Datenergänzung (MDE)** wird zu Beginn eines Berichtsjahres (Berichtsmonat Jänner) der primärstatistisch erhobene Datenkörper - wie in Abbildung 4 skizziert - um den ursprünglichen „Cut-Off-Part“ ergänzt. Die Struktur der Grundgesamtheit kann durch die MDE als Kombination aus Primärdaten, Sekundärdaten und Modellrechnung ohne grundsätzlichen Informationsverlust dargestellt werden.

¹⁴Ein *Cut-Off-Census* ist die Vollerhebung aller Einheiten, die eine bestimmte Schwelle überschreiten. Ein *Cut-Off-Sample* ist eine Stichprobenziehung nur aus den Einheiten, die eine bestimmte Schwelle überschreiten.

Abbildung 4: Datenkörper der KJP-Grundgesamtheit



Quelle: KJP 2009 (Berichtsmonat Jänner).

Die Datenerstellung erfolgt grundsätzlich auf **Mikroebene**. Für jede statistische Einheit der KJP-Grundgesamtheit steht somit ein vollständiger Datensatz für weitere Ergebnisdarstellungen zur Verfügung. Bei den zu ergänzenden statistischen Einheiten handelt es sich in der Regel um **Einbetriebsunternehmen**, vereinzelt können aber auch Mehrbetriebsunternehmen (und deren Betriebe) in die Schätzbasis fallen (im Jänner 2009 scheinen fünf Fälle auf). Diese Ausnahmefälle werden bei der MDE jedoch als „vernachlässigbar“¹⁵ betrachtet, weshalb bei der weiteren Schätzmodellbeschreibung nicht ausdrücklich auf diese Unterscheidung hingewiesen wird. Grundsätzlich werden in weiterer Folge also Unternehmen und Betriebe gleichgesetzt.

4.1. Schritte zur Mikrodaterstellung

Zu Beginn eines Berichtsjahres (Berichtsmonat Jänner) wird die Grundgesamtheit N (59.118 Unternehmen und Betriebe im Berichtsjahr 2009) ermittelt. Eine statistische Einheit ist genau dann der Grundgesamtheit der KJP zugeordnet, wenn sie gemäß Unternehmensregister der STATISTIK AUSTRIA (UR)¹⁶ im Berichtszeitraum wirtschaftlich aktiv ist und eine Haupttätigkeit im Sinne einer Unterklasse der ÖNACE 2008 ausübt, welche dem Produzierenden Bereich zuzurechnen ist. Die Schätzbasis umfasst im Berichtsjahr 2009 rund 46.000 nicht meldepflichtige Einheiten. Auch einige Meldeausfälle ohne historische Struktur fallen in die Schätzbasis.¹⁷ Für die zu ergänzenden statistischen Einheiten $n^{(erg)}$ können sogenannte **Eckdaten** (wirtschaftliche Tätigkeit, Anzahl der unselbständig Beschäftigten, monatliche Umsatzerlöse) unter Einbindung von Verwaltungsquellen erstellt werden. Andere Merkmale (wie

¹⁵Das bedeutet, dass beim selten auftretenden Fall von Mehrbetriebsunternehmen in der Schätzbasis, der Unternehmensumsatz und die Beschäftigten, die aus Verwaltungsquellen für das Berichtsmonat verfügbar sind, über Registerinformationen anteilmäßig auf das Unternehmen und seine Betriebe aufgeteilt werden. Die so erzeugten Dummy-Einheiten werden im Schätzmodell wie Einbetriebsunternehmen behandelt. Nach dem Schätzvorgang werden die Unternehmensmerkmale über die Betriebe ermittelt und die relevanten Merkmalsblöcke der unterschiedlichen statistischen Einheiten erstellt.

¹⁶Das UR erfasst alle Einheiten mit mindestens einem unselbständig Beschäftigten oder mehr als 10.000 Euro Jahresumsatz sowie Einheiten des Staates und Non-Profit Organisationen.

¹⁷Fälle von Unit-Non-Response werden in der KJP über Vorperiodenmeldungen anhand der „historischen Struktur“ des Meldeausfalls imputiert. Für Unit-Non-Response ohne Meldung aus Vorperioden fehlt diese Substitutionsgrundlage, weshalb diese Einheiten durch die MDE in der Schätzbasis mitberücksichtigt werden.

Arbeitsvolumen, Arbeitskosten, Produktion eines Unternehmens) werden ber diese Eckdaten, basierend auf aktuellen primarstatistischen Zusammenhangen, mittels geeigneter Schatzmodelle erzeugt. **Beschaftigtenbezogene Merkmale** (*Lohne, Arbeitsstunden* usw.) werden dabei ber ein robustes lineares Regressionsmodell berechnet, wahrend fur **umsatzbezogene Merkmale** (*abgesetzte Produktion, technische Gesamtproduktion, ...*) ein kombinierter Ansatz angewendet wird, dessen Ziel es ist, fur die statistischen Einheiten der verschiedenen Wirtschaftszweige eine moglichst plausible Verteilung der Produktionsarten und deren Darstellung auf Guterebene (OPRODCOM-Gliederung) zu erhalten. Beim Guteransatz erfolgt die Darstellung der Merkmalsgruppen nach der statistischen Einheit *Betrieb*. Beim Aktivitatsansatz werden Merkmalsgruppen nach den statistischen Einheiten *Unternehmen und Betrieb* in Verbindung mit der schwerpunktmaigen wirtschaftlichen Aktivitat derselben dargestellt. Vorrangiges Ziel der MDE ist es, Mikrodaten zu erzeugen, die eine beliebige Ergebnisdarstellung fur die **Grundgesamtheit der KJP auf Aktivitatsebene** und auf **Guterebene** erlauben.

4.2. Schatzmodell auf Aktivitatsebene

Die Schatzung erfolgt auf Aktivitatsebene ber die schwerpunktmaige wirtschaftliche Aktivitat einer statistischen Einheit. Konzeptionell wurde hier das fur die Leistungs- und Strukturstatistik entwickelte LS-Schatzmodell¹⁸ implementiert, wobei jedoch aufgrund abweichender Merkmalskataloge und der zum Zeitpunkt der Schatzung unterschiedlichen Verfugbarkeit von Verwaltungsdaten entsprechende Modelladaptierungen erforderlich waren. Im Wesentlichen basieren die berlegungen zum KJP-Grundmodell jedoch auf Erfahrungen und Erkenntnissen, die im Zuge der Entwicklung des LS-Schatzmodells gewonnen wurden. Grundidee und Vorgehensweise der MDE auf Aktivitatsebene werden in diesem Abschnitt erlautert.

Erstellung der Eckdaten

ber die **Verknupfung der Unternehmen des UR** mit den vorliegenden Umsatzsteuervoranmeldungen (UVA) der **Finanzbehorde** sowie den Beschaftigtenmeldungen des **Hauptverbands der sterreichischen Sozialversicherungstrager (HV)** erfolgt fur jede einzelne der $n^{(erg)}$ Einheiten der Schatzbasis zuerst die Erstellung der zugehorigen Eckdaten. Beginnend mit der Zuweisung der jeweiligen unternehmens- bzw. betriebsspezifischen Merkmale, wie etwa Wirtschaftsaktivitat oder Rechtsform, die aus dem UR fur jede statistische Einheit zur Verfugung stehen, erfolgt anschlieend die Erweiterung dieser Eckdaten durch die Einbindung der unselbstandig Beschaftigten des HV, gegliedert nach Geschlecht und Qualifikationen (Arbeiter, Angestellte, Lehrlinge), und durch die ubernahme der monatlichen UVA-Meldungen eines Unternehmens. Bei fehlender Beschaftigtenmeldung vom HV wird aufgrund des hohen Verknupfungsgrades im UR davon ausgegangen, dass das Unternehmen keine unselbstandig Beschaftigten hat. Fehlende UVA-Meldungen¹⁹ werden unter Berucksichtigung

¹⁸Ein ausfuhrlicher *Methodenbericht zur Modellbasierten Datenerganzung in der Leistungs- und Strukturstatistik* erscheint im 1. Quartal 2010. Eine kurze Methodenbeschreibung findet sich auch in der Standarddokumentation zur Leistungs- und Strukturstatistik im Produzierenden Bereich auf der Homepage der STATISTIK AUSTRIA unter www.statistik.at > Dokumentationen > Produktion und Bauwesen.

¹⁹Keine fristgerechte Erstattung der UVA-Meldung bei den Finanzbehorden; keine UVA-Meldepflicht fur Einheiten mit Jahresumsatz <100.000 €, etc.

der individuellen Unternehmensentwicklung und der zugehörigen Branchenentwicklung bzw. auch anhand historischer Informationen imputiert. Für die weiteren Berechnungen gilt:

Die KJP-Grundgesamtheit ist unterteilt in L Branchen B_1, \dots, B_L mit $B = \bigcup_{j=1}^L B_j, B_j \cap B_r = \emptyset$.

- $j = 1, \dots, L$... laufender Branchenindex
- $i = 1, \dots, N$... Index der statistischen Einheiten
- n_j ... Anzahl der statistischen Einheiten einer Branche j
- $N = \sum_{j=1}^L n_j$... Anzahl der statistischen Einheiten insgesamt ($N = n^{(erh)} + n^{(erg)}$)
- $UVA_{i,m}$... UVA-Meldung einer statistischen Einheit i im Berichtsmonat m
- $M_i := \{(m-t) : t \in \{0, 1, \dots, 11\}, UVA_{i,(m-t)} \text{ liegt vor} \}$
- $J_j := \{i : \text{statistische Einheit } i \text{ ist in Branche } j \text{ klassifiziert}\}$
- $J_j^{(erh)} := \{i \in J_j : \text{Einbetriebsunternehmen in Erhebung}\}$
- $J_j^{(erg)} := \{i \in J_j : \text{statistische Einheiten in Ergänzung}\}$

Nur Unternehmen, für die im gleitenden Jahresverlauf die letzten zwölf UVA-Meldungen vollständig vorliegen, gehen in die Berechnung von Formel (1) ein. Dem branchenspezifischen mittleren Umsatzanteil des Berichtsmonats m entspricht:

$$\bar{u}_{j,m} = \frac{\sum_{i \in J_j^*} \frac{UVA_{i,m}}{\sum_{t=0}^{11} UVA_{i,(m-t)}}}{\sum_{i \in J_j^*} 1}, \quad J_j^* := J_j^{(erg)} \setminus \{i : |M_i| < 12\}. \quad (1)$$

Einheiten mit fehlender Meldung für den Berichtsmonat, jedoch mindestens sechs vorliegenden UVA-Meldungen im gleitenden Jahresverlauf, werden mittels eines Umsatzsubstituts berücksichtigt. Dem imputierten Umsatz einer Unternehmenseinheit i für den Berichtsmonat m entspricht:

$$\hat{U}_{i^\circ,m} = \frac{\sum_{(m-t) \in M_{i^\circ}} UVA_{i^\circ,(m-t)}}{11} \bar{u}_{j^\circ,m}, \quad \begin{aligned} i^\circ &\in \{i \in J_j^{(erg)} : |M_i| \geq 6, UVA_{i,m} \text{ fehlend}\}, \\ j^\circ &\in \{j : i^\circ \in J_j^{(erg)}\}. \end{aligned} \quad (2)$$

Mit $I_{i^\circ,(m-t)} = \begin{cases} 1 & \text{wenn } UVA_{i^\circ,(m-t)} \text{ vorliegt,} \\ 0 & \text{sonst.} \end{cases}$

Für Einheiten ohne verwendbare UVA-Meldungen wird der erforderliche Monatsumsatz anhand des für das Unternehmen letzt verfügbaren Jahreswertes (Ergebnisse der Umsatzsteuererklärung, der Leistungs- und Strukturstatistik bzw. KJP)²⁰ und der monatlichen UVA-Branchenentwicklung imputiert (8% der ergänzten Unternehmensumsätze).

²⁰Fortgeschrieben über branchenspezifische Umsatzentwicklungen im Zeitverlauf.

Nach Erstellung und Vervollstandigung der Eckdaten erfolgt die modellbasierte Berechnung der restlichen beschaftigten- und umsatzbezogenen Merkmale (Lohne, Arbeitsstunden, Eigenproduktion usw.).

Auswahl der Modellbasis

Um den zu schatzenden $n^{(erg)}$ nicht erhobenen Einheiten moglichst ahnliche erhobene Einheiten aus $n^{(erh)}$ zugrunde zu legen, erfolgt eine iterative Auswahl der zur Schatzung erforderlichen **Modellbasis** innerhalb von Umsatzklassen und Wirtschaftstatigkeit - „bottom up“, beginnend mit der tiefsten Gliederungsebene.²¹

Fur die Klassenauswahl der Modellbasis gilt:

$$J_{j\alpha} := \{i \in J_j^{(erh)} : U_i < \tilde{U}_{j\alpha}, \alpha = \min(\{\alpha' \in (1, 7] : (\alpha'/10)n_j^{(erh)} \geq 20\})\}.$$

U_i ... Umsatz einer statistischen Einheit i

$\tilde{U}_{j\alpha}$... α -Dezil der Umsatzverteilung einer Branche j der Erhebung

$n_{j\alpha}$... Anzahl der statistischen Einheiten einer Branche j der Modellbasis

Bei den $n_{j\alpha}$ Einheiten der Modellbasis handelt es sich um jene erhobenen Einheiten²² einer Branche j , deren Monatsumsatz unter dem α -Dezil $\tilde{U}_{j\alpha}$ der primarstatistischen Umsatzverteilung der betreffenden Branche liegt. $\tilde{U}_{j\alpha}$ entspricht dabei dem kleinsten α -Dezil bei dem mindestens 20 Beobachtungen in der Modellbasis aufscheinen. Weisen jedoch auch die kleinsten 70% der Umsatzverteilung keine ausreichende Primarbesetzung auf, erfolgt die Auswahl auf einer ubergeordneten Gliederungsebene (ONACE-Aggregat).²³

Im Janner des Berichtsjahres 2009 gelangten auf diese Weise von rund 12.000 erhobenen Unternehmen nur die kleinsten 24% iterativ in die Modellbasis, um als Strukturspender fur die zu schatzenden Einheiten der 300 Wirtschaftsbereiche zu dienen. Gemessen am Eckwert Beschaftigte konnte fur fast 90% der Schatzbasis im Basismonat Janner 2009 die Parameterberechnung bereits auf Ebene der jeweiligen ONACE-(Unter-)Klasse (4- bzw. 5-Steller) erfolgen.

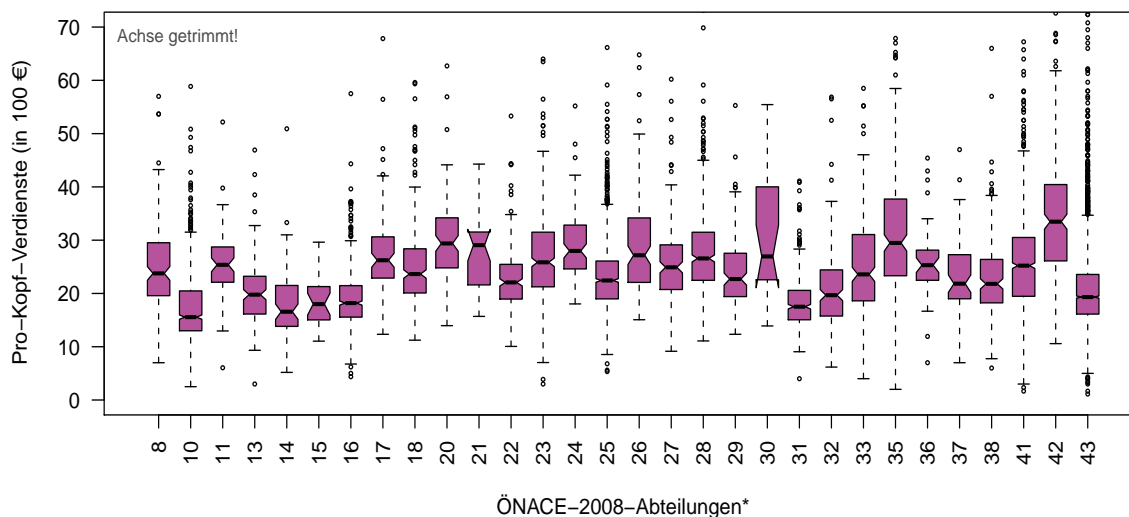
Schon am einfachen Beispiel der durchschnittlichen *Bruttoverdienste pro unselbstandig Beschaftigten* in Abbildung 5 ist ersichtlich, dass sich die durch *Box-Plots* visualisierten Verteilungen der Pro-Kopf-Verdienste der ONACE-2008-Abteilungen deutlich unterscheiden, und die Berucksichtigung der branchenspezifischen Strukturunterschiede (also der Wirtschaftstatigkeit einer Einheit) unbedingt erforderlich ist.

²¹Entspricht bei der Wirtschaftstatigkeit den rund 300 ONACE-2008-Unterklassen des Produzierenden Bereichs.

²²Ausschlielich *Einbetriebsunternehmen* dienen auf Aktivitatsebene als Strukturspender.

²³Im Gegensatz zur Aktivitatsebene wird bei der spateren Erstellung des Guteransatzes nicht auf ubergeordnete Branchen gewechselt.

Abbildung 5: Branchenspezifische Strukturunterschiede in der KJP-Erhebung
Verteilung der Pro-Kopf-Verdienste der Unternehmen der ÖNACE-2008-Abteilungen



Quelle: KJP 2009 (Berichtsmonat Jänner). - *) Ohne ÖNACE-2008-Abteilung 06, 07, 09, 12, 19, 39.

Aufgrund der relativ starken Fluktuation an der Meldeschwelle²⁴ vom Zeitpunkt der Stichprobenziehung bis zum Zeitpunkt der Durchführung der MDE kann innerhalb einer Branche den größeren zu schätzenden Einheiten eine vergleichbare Basisstruktur zugrunde gelegt werden. Je näher die zu schätzenden Einheiten an den Meldeschwellen liegen, umso größer ist einerseits ihr Gewicht, aber umso ähnlicher sind sie andererseits auch den Strukturspendern in der Modellbasis. Mit wachsendem Abstand zur Modellbasis verringert sich zugleich auch das Gewicht der zu schätzenden Einheiten und somit deren Einfluss auf die Merkmalssummen eines Wirtschaftsbereichs. Diesen Eigenschaften der zu schätzenden Einheiten in Kombination mit der grundsätzlich hohen primärstatistischen Merkmalsabdeckung der KJP sollten den Erwartungen nach qualitativ hochwertigen Ergebnissen pro futuro sowohl hinsichtlich zeitlicher Vergleichbarkeit als auch hinsichtlich der **Kohärenz mit anderen wirtschaftsstatistischen Projekten** (insbesondere der Leistungs- und Strukturstatistik) entsprechen - selbst wenn unterstellt wird, dass innerhalb eines Wirtschaftsbereichs **unternehmensgrößenabhängige Strukturunterschiede** vorliegen. Als problematisch muss aber eine Konstellation aus geringer Primärbesetzung und starken **branchenspezifischen Strukturunterschieden** betrachtet werden, die bei sehr detaillierter Ergebnisdarstellung auftreten kann. Hier sind aufgrund *gemischter Strukturen*²⁵ auch mit der vorgenommenen Modellbasisabgrenzung deutlich verzerrte Schätzer zu erwarten, und wenn zusätzlich noch eine schwache Primärabdeckung²⁶

²⁴Die „Fluktuation an den Meldeschwellen“ bezieht sich auf jene meldepflichtigen und nichtmeldepflichtigen Einheiten, die zum Zeitpunkt der MDE (t+90 Tage) in den jeweils anderen Bereich (unter bzw. über die Meldeschwellen) fallen.

²⁵Gemischte Strukturen sind dann anzunehmen, wenn etwa aufgrund geringer Primärbesetzung die Parameterschätzung einer ÖNACE-2008- Unterklasse (5-Steller) auf dem übergeordneten 4-Steller (ÖNACE-2008-Klasse) erfolgen muss, und die anderen 5-Steller dieser ÖNACE-2008-Klasse bspw. ein deutlich höheres Lohnniveau aufweisen als jener 5-Steller für den die Parameterschätzung erfolgen soll.

²⁶Zu unterscheiden ist die Problematik einer schwachen Primärabdeckung (Anteil der MDE an Merkmalssumme ist hoch) von jener der geringen Primärbesetzung (zu wenig Strukturspender für die branchenspezifische

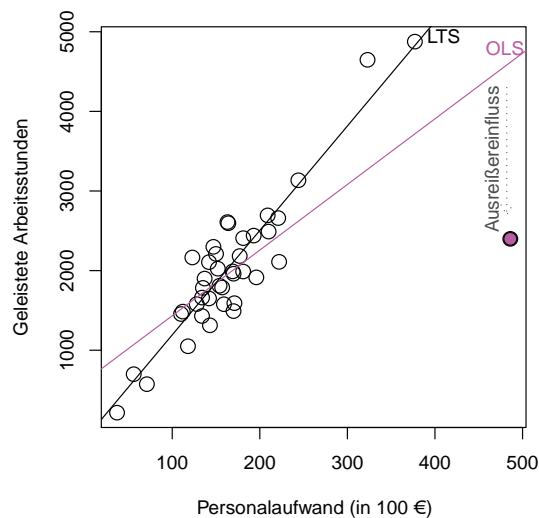
vorliegt, der Anteil der MDE also hoch ist, dann kann ein gewisser systematischer Effekt in den Ergebnissen kaum vermieden werden. Im Einzelfall muss deshalb zusatzliches qualifiziertes Expertenrating²⁷ zur Gewahrleistung der erforderlichen Datenqualitat beitragen.

Parameterschatzung

Wirtschaftsstatistische Zusammenhange von Mikrodaten konnen haufig schon durch ein einfaches lineares Modell gut beschrieben werden. Da wirtschaftsstatistische Daten jedoch stets mit **Ausreißern**²⁸ behaftet sind, ist eine robuste Methode zur Gewahrleistung einer stabilen Qualitat der Modellanpassung unerlasslich.

Abbildung 6 veranschaulicht dies an einem einfachen Beispiel fur eine der 300 ONACE-2008-Unterklassen der Modellbasis. Es ist ersichtlich, dass im klassischen linearen Regressionsmodell (OLS-Regression)²⁹ schon eine einzige Beobachtung ausreicht, um vollig unsinnige Schatzer zu erhalten - die Losung wird quasi zum Ausreißer hingezogen. Die LTS-Regression³⁰ hingegen ist robust, reagiert also nicht so empfindlich auf Ausreißer und passt sich gut an die Datenmehrheit an.

Abbildung 6: Ausreißereinfluss im Regressionsmodell
(Auszug aus Modellbasis)



Quelle: KJP 2009 (Berichtsmonat Janner).

Basierend auf der iterativ bestimmten Modellbasis werden bei der MDE deshalb alle be-

Modellbildung).

²⁷Beim sogenannten „Expertenrating“ handelte es sich um Einschatzungen und Modifizierungen einzelner Werte und Parameter durch Fachexperten der Wirtschaftsstatistik aufgrund ihrer Kenntnisse branchenspezifischer Eigenheiten und Kriterien, die beim Vorliegen bekannter systematischer Abweichungsmechanismen vorgenommen werden.

²⁸Als Ausreißer gelten im vorliegenden Fall jene Beobachtungen, die nicht dem (linearen) Muster der Datenmehrheit folgen.

²⁹Ordinary Least Squares Regression.

³⁰Least Trimmed Squares Regression.

schäftigtenbezogenen Hauptmerkmale branchenspezifisch mit Hilfe eines **robusten linearen Regressionsmodells** berechnet. Um den Einfluss additiver Effekte bei der Mikrodatenerstellung zu vermeiden, wird grundsätzlich ein lineares Modell ohne Interzept angewandt.

Für das Modell $y_i = x_{i1}\beta_{j1} + \dots + x_{ip}\beta_{jp} + \epsilon_i$, mit ϵ_i als Störterm, und $i \in J_{j\alpha}$ ist die Least Trimmed Squares (LTS) Regression folgendermaßen definiert (*Rousseeuw, 1984*):³¹

$$\hat{\beta}_{n_{j\alpha}}^{LTS} = \underset{\hat{\beta}_j}{\operatorname{argmin}} \sum_{k=1}^{h_j} (e^2)_{k:n_{j\alpha}} \quad \dots \text{ist die Zielfunktion der LTS Optimierung.} \quad (3)$$

Die Summe der h_j kleinsten quadratischen Residuen sind dabei zu minimieren, wobei die quadrierten Residuen $e_i = (y_i - \hat{y}_i)$ in geordneter Form wie folgt definiert sind:

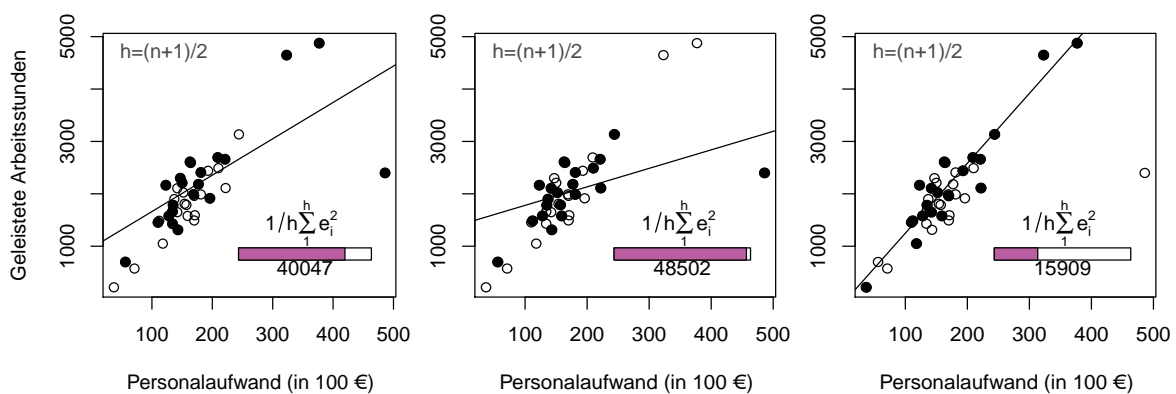
$$(e^2)_{1:n_{j\alpha}} \leq (e^2)_{2:n_{j\alpha}} \leq \dots \leq (e^2)_{k:n_{j\alpha}} \leq \dots \leq (e^2)_{n_{j\alpha}:n_{j\alpha}}. \quad (4)$$

Der Umfang h_j ist ein Subset aus den gegebenen $n_{j\alpha}$ Beobachtungen der Branche j , wobei h_j definiert ist als

$$h_j \in \left[\frac{n_{j\alpha} + 1}{2}, n_{j\alpha} \right]. \quad (5)$$

Die unter (3) bis (5) formal beschriebene LTS-Technik kann vereinfacht als iteratives Verfahren erklärt werden wie es in Abbildung 7 demonstriert wird. Wird der Umfang des Subsets beispielsweise auf $h = (n + 1)/2$ gesetzt, so werden wiederholt zufällig 21 Beobachtungen (dunkle Punkte) aus den $n = 40$ Beobachtungen ausgewählt und basierend darauf die OLS-Regression und die Summe der h kleinsten quadratischen Residuen berechnet. Ziel ist es, jene Teilmenge h aus den n Beobachtungen zu finden, bei der die Summe der h kleinsten quadratischen Residuen minimiert wird. Dem LTS-Schätzer entspricht dann die OLS-Anpassung an diese h Punkte.

Abbildung 7: Demo der LTS-Technik - zufällig gewählte Subsets (Greedy Algorithmus)



Quelle: KJP 2009 (Berichtsmonat Jänner).

³¹*P.J. Rousseeuw (1984). "Least median of squares regression". Journal of the American Statistical Association.*

Diese Technik ist extrem rechenintensiv, da theoretisch der „ganze Raum“ abgesucht werden musste, um das Minimum zu finden. Deshalb werden spezielle Algorithmen eingesetzt, die die Losung (fur groere Datensatze) approximieren. Die beschriebene LTS-Prozedur kann uber SAS/IML (Interactive Matrix Language) aufgerufen werden, wobei standardmaig der FAST-LTS Algorithmus³²) implementiert ist.

Mit der LTS-Methode werden stets $n - h$ Beobachtungen ausgeschlossen, wodurch relevante Information verloren gehen kann. Mit wachsendem Subset h gewinnt der Schatzer zwar mehr Information aus den Daten, verliert aber an Robustheit und vice versa. Wird $h = 0.75n$ gewahlt, so gilt dies als guter Kompromiss zwischen Robustheit und Effizienz. Es ist jedoch kaum zu erwarten, dass in der Modellbasis alle Wirtschaftsbereiche und Merkmalszusammenhange bis zu 25% ausreißerbehaftet sind. Andererseits ist die tatsachliche Ausreißeranzahl im Vorhinein nicht bekannt, d.h. auch nicht exakt uber h bestimmbar. Deshalb wird fur die MDE eine **gewichtete LS-Regression** angewandt, die in einem zusatzlichen Schritt - basierend auf dem besten LTS-Subset und der zugehorigen LTS-Funktion - nur noch Beobachtungen mit groen Residuen ausschliet und basierend auf dem neu gesetzten Subset h' den finalen LTS-Schatzer berechnet (im Regelfall gilt: $h' > h$). Fur die Beispieldaten in Abbildung (7) gilt $h' = 39$, d.h. alle bis auf eine der Beobachtungen werden in die finale Berechnung aufgenommen.

Mit Hilfe der berechneten branchenspezifischen Regressionsparameter erfolgt, gebunden an die Eckdaten der 46.000 zu schatzenden statistischen Einheiten, die Berechnung von Merkmalen wie *Bruttolohne*, *Teilzeitbeschaftigte*, *geleistete Arbeitsstunden* usw.

4.3. Schatzmodell auf Guterebene

Fur die Erganzung des primarstatistischen Datenkorpers sind auch die Produktionsarten (wie *abgesetzte Produktion* oder *durchgefuhrte Lohnarbeit*, um zwei Beispiele zu nennen) einer statistischen Einheit zu berechnen. Zusatzlich ist die **Aufgliederung der Produktion nach PRODCOM-Positionen** (genauer: 8-Steller des nationalen Guterverzeichnisses OPRODCOM, da das nationale Guterverzeichnis OPRODCOM uber die Erfordernisse der europaischen PRODCOM-Liste hinausgeht) erforderlich, um eine entsprechende Ergebnisdarstellung im Rahmen der Guterproduktion (Guteransatz) zu ermoglichen.

Die Schatzung erfolgt auch hier wiederum aktivitatsbezogen und unter Einbindung der Eckdaten aus Verwaltungsquellen. Dabei kann aufgrund des hohen statistischen Zusammenhanges das Produktionsvolumen einer statistischen Einheit uber dessen Unternehmensumsatz bestimmt werden.

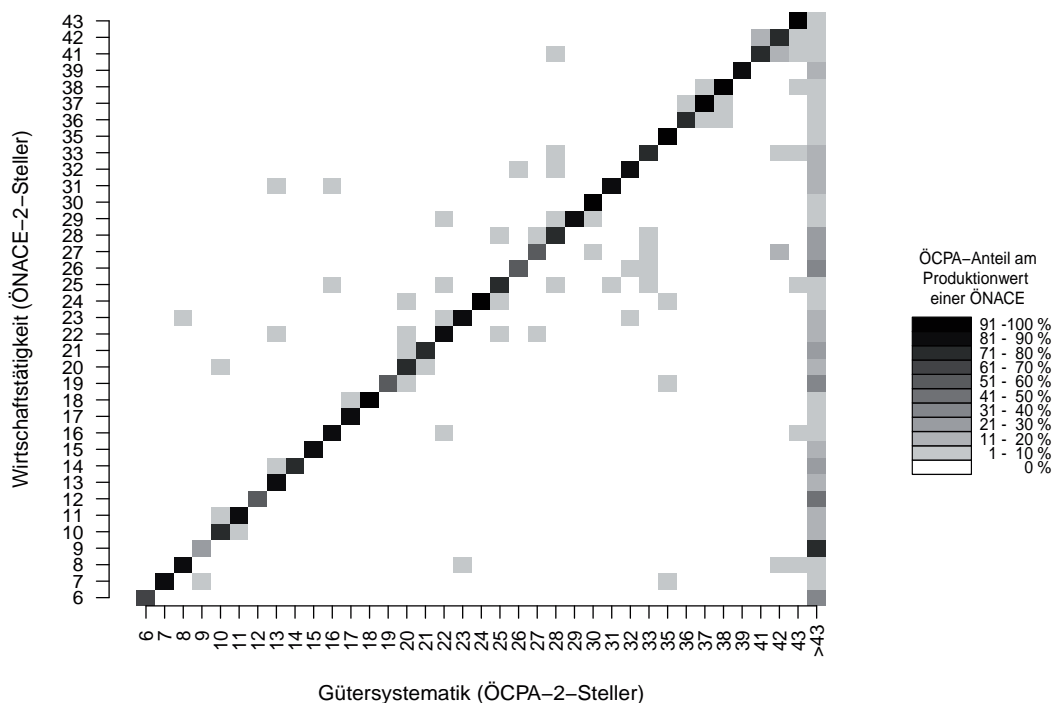
Daruber hinaus ist auch eine plausible OPRODCOM-Zuordnung uber die wirtschaftliche Aktivitat einer statistischen Einheit moglich. Die Matrixdarstellung nach der Systematik der Wirtschaftsaktivitaten (ONACE 2008) und der statistischen Guterklassifikation in Verbindung mit den Wirtschaftszweigen (OCPA 2008)³³ weist generell eine hohe Konzentration des

³²Rousseeuw, P.J. and Van Driessen, K.: “Computing LTS Regression for Large Data Sets”, Springer, Netherlands 2006.

³³Vgl. Anhang zur Verordnung (EG) Nr. 451/2008 des Europaischen Parlaments und des Rates vom 23. April 2008 zur Schaffung einer neuen Guterklassifikation in Verbindung mit den Wirtschaftszweigen (CPA) und zur Aufhebung der Verordnung (EWG) Nr. 3696/93 des Rates.

relevanten Produktionsvolumens an den Hauptachsen auf. Abbildung 8 veranschaulicht die Gliederung des Produktionsvolumens nach Wirtschaftsbereichen und Gütergruppen. An der Intensität der Farbe ist die Konzentration der Produktion eines Wirtschaftsbereichs in den verschiedenen Gütergruppen ersichtlich. In der ersten Zeile der Makematrix ist beispielsweise die Konzentration der *abgesetzten Produktion* für die ÖNACE-2008-Abteilung 43 (*Vorbereitende Baustellenarbeiten, Bauinstallationen und sonstiges Ausbaugewerbe*) veranschaulicht. Für die ÖNACE-2008-Abteilung 43 liegt die Konzentration der charakteristischen Produktion (ÖCPA-2008-Abteilung 43) bei rund 92%, die restlichen 8% fallen auf *Produktbegleitende Umsätze* (ÖCPA-2008-Abteilung > 43).

Abbildung 8: Makematrix nach ÖNACE-2008-Abteilungen und ÖCPA-2008-Abteilungen
Konzentration der abgesetzten Produktion in KJP-Erhebung



Quelle: KJP 2009 (Berichtsmonat Jänner).

Schätzung der Parameter zur Gütergliederung

Als Ausgangsbasis für die Aufgliederung der Produktionswerte nach Gütern dient eine Matrix, in der den ~ 300 Wirtschaftsbereichen (ÖNACE-2008-Unterklassen) plausible ÖPRODCOM-Codes zugeordnet sind. Dies kann wie zuvor auf Aktivitätsebene beschrieben, grundsätzlich durch empirisch beobachtbare Zuordnungen innerhalb einer Branche erfolgen.

Nur in einigen nicht primarstatistisch erfassten Wirtschaftsbereichen wird ein OPRODCOM-Set unter Einbindung der Klassifikationsdatenbank der STATISTIK AUSTRIA zugewiesen.³⁴ Fur alle Gutercodes eines Wirtschaftsbereichs werden empirische Gewichte wie etwa *Besetzungsgewicht, mittlerer Mengenanteil, mittleres Codegewicht* usw., berechnet.

Fur die $j = 1, \dots, 300$ Wirtschaftsbereiche erfolgt die Berechnung der Parameter zur Gutergliederung des Produktionswertes basierend auf den jeweils $n_{j\alpha}$ erhobenen Einheiten der Modellbasis. Es gilt:

$c = 1, \dots, T$... laufender Gutercodesindex
 $w_{i,c}$... Produktionswert eines Gutercodes c einer statistischen Einheit i
 $w_i = \sum_{c=1}^T w_{i,c}$... Produktionswert einer statistischen Einheit i

$$g_{j,c} = \frac{1}{n_{j\alpha}} \sum_{i \in J_{j\alpha}} I_{i,c}, \quad \dots \text{Besetzungsgewicht eines Gutercodes } c \text{ einer Branche } j. \quad (6)$$

Mit $I_{i,c} = \begin{cases} 1 & \text{wenn Einheit } i \text{ den Code } c \text{ aufweist,} \\ 0 & \text{sonst.} \end{cases}$

Fur den branchenspezifischen mittleren Anteil des Gutercodes c am Produktionsvolumen w gilt:

$$\bar{w}_{j,c} = \frac{\sum_{i \in J_{j\alpha}} \frac{w_{i,c}}{w_i}}{\sum_{i \in J_{j\alpha}} I_{i,c}} \quad (7)$$

Der in (6) und (7) beschriebenen Parameterberechnung wird wieder eine iterativ gewahlte Modellbasis kleinerer primarstatistisch erhobener Einheiten zugrunde gelegt, wodurch sich auch die zu berucksichtigenden OPRODCOM-Codes von mehr als 3.200 unterschiedlichen Codes der Primarerhebung der KJP auf etwa 1.200 Gutercodes der Modellbasis reduzieren. Die Vorgehensweise bei der Modellbasisauswahl entspricht im Wesentlichen der bereits auf Aktivitatsebene beschriebenen Methode (vgl. dazu den Abschnitt Auswahl der Modellbasis), wobei jedoch die Berechnungen beim Guterausatz ausschlielich auf Ebene der ONACE-2008-Unterklasse erfolgen.³⁵

Bei nicht ausreichender primarstatistischer Besetzung wird beim Guterausatz (im Gegensatz zum Aktivitatsansatz) fur die Parameterberechnung nicht auf eine ubergeordnete Wirtschaftsebene gewechselt, sondern es werden innerhalb der betreffenden ONACE-2008-Unterklasse schrittweise auch jene groen erhobenen Einheiten (groe Einbetriebsunternehmen und Betriebe von Mehrbetriebsunternehmen)³⁶ in die Modellbasis aufgenommen, die beim Aktivitatsansatz nicht berucksichtigt werden. Die Einbindung beschrankt sich in diesen Fallen jedoch auf

³⁴Automatische ubernahme jener Gutercodes der Klassifikationsdatenbank, deren Codierung auf Ebene der 4- und 5-Steller mit der betreffenden ONACE 2008-(Unter)Klasse ubereinstimmen, und nachtragliche uberprufung der erfolgten Zuordnung (Hinweis: der symmetrische Aufbau von CPA-2008 und ONACE-2008 ist in einigen wenigen Fallen nicht eindeutig und muss entsprechend berucksichtigt werden (etwa beim ONACE-3-Steller Bau von Gebauden (412)).

³⁵Uber ein hoheres ONACE-Aggregat kann in der Regel keine plausible OPRODCOM-Zuweisung erfolgen.

³⁶Betriebe eines Mehrbetriebsunternehmens flieen in die Berechnung der Modellgewichte nur als eine zusammengefasste Einheit ein.

Gütercodes, die dem Kriterium des „wirtschaftlichen Ursprungs“ folgen (deren Codierung auf Ebene der 4- und 5-Steller mit der betreffenden ÖNACE-2008-(Unter)Klasse übereinstimmt). Gütercodes, die in der primärstatistischen Datenmasse nur vereinzelt auftreten, werden grundsätzlich nicht auf die zu ergänzenden Einheiten übertragen, um eine mögliche Überzeichnung dieser Codes zu vermeiden (außer es handelt sich um Hauptcodes, und für die Branche liegen keine weiteren relevanten Codes vor).

Parameterzuweisung und Gütergliederung

Mit der im Abschnitt zuvor beschriebenen Parameterschätzung erhält man für jeden Wirtschaftsbereich ein Gütercode-Set mit zugehörigen Gewichten. In Abbildung (9) ist am Beispiel der ÖNACE-2008-Unterklasse 10711 (*Herstellung von Schwarz- und Weißbackwaren*) eines dieser 300 Gütercode-Sets mit einigen Parametern veranschaulicht.

Abbildung 9: Ausgewählte Parameter zur Gütergliederung am Beispiel der ÖNACE 10711 (*Herstellung von Schwarz- und Weißbackwaren* (Auszug aus Modellbasis))

ÖPROCDOM 2008		Gewicht in Modellbasis				Anzahl der Betriebe in Ergänzung
Bezeichnung	CODE	Eh	g (in %)	HC (in %)	\bar{w}	
Schwarzbrot	1071110001	kg	99	14	0,16275	1.148
Weißbrot	1071110002	kg	83	6	0,11716	965
Spezialbrot	1071110003	kg	83	4	0,09771	965
Weißgebäck	1071110005	kg	91	64	0,29394	1.065
Brot, Gebäck ... für diätische Verwendung	1071110006	kg	9	1	0,09310	100
Anderes frische Brot ohne Zusatz von ...	1071110007	kg	13	1	0,09561	150
Weichwaren (ohne Dauerbackwaren), gesüßt	1071120001	kg	4	3	0,10340	466
Feingebäck (ohne Dauerbackwaren), gesüßt	1071120002	kg	89	6	0,12784	1.032
Sonstiger Einzelhandel, a. n. g.	4700002000		83	0	0,17156	965
Dienstleistungen von Restaurants, ..., Cafés	5610101000		37	0	0,14143	433
Q: KJP 2009 (Berichtsmonat Januar)						

Für die zu ergänzenden statistischen Einheiten eines Wirtschaftsbereiches wird das über den Unternehmensumsatz berechnete Produktionsvolumen \hat{w} nun anhand dieser branchenspezifischen Gewichte nach ÖPRODCOM-Positionen aufgegliedert. Die ÖPRODCOM-Zuweisung erfolgt dabei über die jeweilige schwerpunktmäßige wirtschaftliche Tätigkeit einer statistischen Einheit und unter Berücksichtigung des Besetzungsgewichtes $g_{j,c}$ anhand eines probabilistischen Zuordnungsverfahrens.³⁷

Von den 1.179 zu schätzenden statistischen Einheiten der in der Abbildung 9 beispielhaft angeführten ÖNACE-2008-Unterklasse 10711 erhalten also 99% den Code 1071110001, 83% den Code 1071110002 usw.

Der Anteil eines Gütercodes c am über den Unternehmensumsatz geschätzten Produktionswert \hat{w} hängt letztendlich davon ab, welche Code-Kombination der statistischen Einheit per Zufallsauswahl innerhalb der Branche zugewiesen wurde. Jede statistische Einheit hat aber zumindest einen Hauptcode (HC) und $k \geq 0$ Nebencodes.

³⁷Liegen jedoch aus Vorperioden (früheren Erhebungen) brauchbare historische Informationen über die Gütergliederung einer statistischen Einheit vor, so wird diese Struktur auf die Einheit übertragen.

Es gilt:

$\hat{w}_{i,c}$... geschatzter Produktionswert eines Codes c einer stat. Einheit i
 $\hat{w}_i = \sum_{c=1}^T \hat{w}_{i,c}$... geschatzter Produktionswert einer statistischen Einheit i
 $C_{j\alpha} := \{c \in C_j : g_{j,c} > 0\}$... Menge der Gutercodes c in Branche j der Modellbasis

C_j umfasst alle primarstatistisch auftretenden Gutercodes einer Branche j , und $C_{j\alpha}$ alle zulassigen Codes der Modellbasis. Es gilt: $\hat{C}_i \subset C_{j\alpha} \subset C_j$, mit \hat{C}_i als die einer statistischen Einheit zugeordneten Gutercodemenge.

Die Funktion zur Verteilung der Gutercodes c einer Branche j auf die statistische Einheit i ist definiert als:

$$\hat{w}_{i,c} = \hat{w}_i \frac{\bar{w}_{j,c} I_{i,c}}{\sum_{c' \in C_{j\alpha}} \bar{w}_{j,c'} I_{i,c'}}, \quad (8)$$

mit $I_{i,c} = \begin{cases} 1 & \text{wenn Code } c \text{ Einheit } i \text{ zugewiesen,} \\ 0 & \text{sonst.} \end{cases}$

Fur die in der Abbildung 9 beispielhaft angefuhrte ONACE-Unterklasse 10711 gilt etwa $|C_{j=25}| = 42$ und $|C_{j=25 \alpha=2}| = 10$, d.h. es sind nur noch zehn aus ursprunglich 42 moglichen Gutercodes der Primarerhebung in der Modellbasis.

Von der Moglichkeit, die in der Erhebungsmasse vorkommende *Kombination von Guterkategorien* gema ihrer Auftrittswahrscheinlichkeit auf die Schatzbasis zu ubertragen wurde abgesehen, da diese fur die vorgesehenen statistischen Zwecke nicht unbedingt erforderlich ist, und ein multivariater Ansatz - aufgrund der Modellvorgaben (Beschrankungen durch Modellbasis, Restriktionen bei Gutercodewubertragung, Spenderstruktur kann nicht wie bei Hot-Deck einfach ubertragen werden, usw.) - zu nicht abschatzbaren Effekten fuhren kann. Die Ubertragung der univariaten Auftrittswahrscheinlichkeit wurde als ausreichend erachtet.

5. Unterjahrige Fortschreibung

Zu Beginn eines Berichtsjahres (Berichtsmonat Janner) erfolgt im Zuge der modellbasierten Erganzung des primarstatistischen Datenkorpers die Uberprufung einflussreicher Falle³⁸ der Schatzbasis durch Experten der Fachstatistik - insbesondere hinsichtlich der ubernommenen Verwaltungsdaten an die alle ubrigen Merkmale gebunden sind.

Unterjahrig (Berichtsmonat Februar bis Dezember) werden die mit dem Berichtsjahresmonat Janner erzeugten Daten in einem automatisierten Prozess fortgeschrieben.³⁹ Grundsatzlich wieder - wie im Abschnitt zuvor beschrieben - durch die Einbindung von Verwaltungsquellen und durch die Berucksichtigung von Branchenentwicklungen. Die Beschaftigtenstruktur einer statistischen Einheit kann dabei uber Registerverknupfungen zum HV, die Umsatzhohe

³⁸Als einflussreiche Falle gelten bei der MDE jene Einheiten, die von der Definition der Klein- und Kleinstunternehmen abweichen oder aus bestimmten Kriterien auffallig sind (bspw. durch einen hohen Anteil an der Merkmalssumme eines Gutercodes oder durch eine starke Umsatzanderung im Vergleich zur Vorjahresstatistik).

³⁹Fur unterjahrig neu auftretende Einheiten wird ein Vormonats-Dummydatensatz angelegt (mittels branchenspezifischer Spenderstruktur aus Unternehmen der Vormonatserganzung), inaktive Einheiten werden geloscht.

durch UVA-Meldungen der Finanzbehörde monatlich aktualisiert werden. Zur Berücksichtigung arbeitstätiger Konstellationen und saisonaler Muster werden zudem über die unterjährige Entwicklung der Primärdaten branchenspezifische (mediane) Änderungsraten berechnet. Die Vorgehensweise entspricht im Wesentlichen der LOCF-Methode⁴⁰, welche in der KJP für Meldeausfälle mit historischen Informationen entwickelt wurde. Da es sich bei den ergänzten Daten aber nicht um zufällige Meldeausfälle, sondern um eine systematisch abgeschnittene Masse handelt, muss den zu schätzenden Branchengewichten eine andere Basis als den Meldeausfällen zugrunde gelegt werden (vgl. dazu den Abschnitt Auswahl der Modellbasis). Zudem wird im Gegensatz zur LOCF-Substitution die monatliche Umsatzänderung nicht basierend auf der Primärerhebung der KJP gewonnen, sondern kann, wie zuvor beschrieben, über die unterjährige UVA-Entwicklung der zu ergänzenden Einheiten selbst berücksichtigt werden.⁴¹

Es gilt:

$v = 1, \dots, r$... Index der statistischen Merkmale

x ... Wert eines statistischen Merkmals mit Index v^\bullet , mit $v^\bullet \in V^\bullet$

y ... Wert eines statistischen Merkmals mit Index v° , mit $v^\circ \in V^\circ$

$V^\bullet = \{v : \text{aus Verwaltungsquellen verfügbar}\}$

$V^\circ = \{v : \text{nicht aus Verwaltungsquellen verfügbar}\}$

$$\hat{y}_{i^*,m} = \left[\hat{y}_{i^*,(m-1)} \left(\frac{x_{i^*,m}}{x_{i^*,(m-1)}} \right) \right]^{I_x^*} \left[\left(\underset{i \in J_{j\alpha}^\circ}{\text{median}} \left(\frac{y_{i,m}}{x_{i,m}} \right) \right) x_{i^*,m} \right]^{1-I_x^*} \left[\underset{i \in J_{j\alpha}^\circ}{\text{median}} \left(\frac{\frac{y_{i,m}}{x_{i,m}}}{\frac{y_{i,(m-1)}}{x_{i,(m-1)}}} \right) \right]^{I_x I_x^*} \quad (9)$$

Mit $i^* \in \{i : i \in J_j^{(erg)}\}$ und $J_{j\alpha}^\circ := J_{j\alpha} \setminus \{i : x_{i,(m-t)} = 0, t \in \{0, 1\}\}$

$I_x = 1$ wenn Quelle für x_y der HV, $I_x = 0$ sonst.

$I_x^* = 1$ wenn $x_{i,(m-1)}^y > 0$, $I_x^* = 0$ sonst.

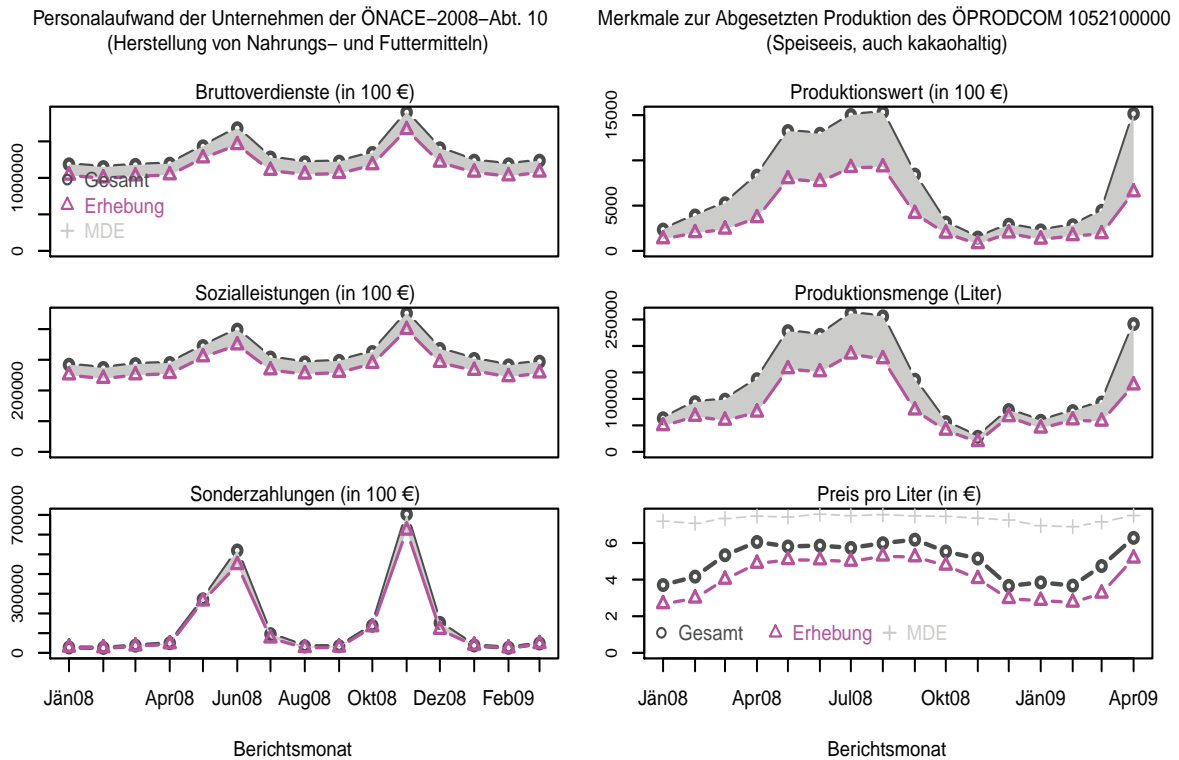
In Abbildung 10 ist am Beispiel ausgewählter beschäftigtenbezogener Merkmale der ÖNACE-2008-Abteilung 10 (*Herstellung von Nahrungs- und Futtermitteln*) die Entwicklung der ergänzten Datenmasse auf Aktivitätsebene ersichtlich. Hervorzuheben wäre etwa der Aspekt, dass Sonderzahlungen im Modell Großteils zu den Zeitpunkten Juni und November Berücksichtigung finden (Urlaubsgeld, Weihnachtsrenumeration).

Abbildung 10 veranschaulicht zudem an der unterjährigen Entwicklung der ÖPRODCOM-Position 1052100000 (*Speiseeis, auch kakaohaltig*) die Datenergänzung auf Güterebene von Jänner 2008 bis April 2009. Hier ist nicht nur bei Produktionswert und -menge ein sehr ausgeprägter saisonaler Effekt ersichtlich, sondern auch beim Literpreis. Ursache der Preisschwankungen ist, dass die industrielle Eiserzeugung ganzjährig mit sehr niedrigen Literpreisen erfolgt, und die Vielzahl an Konditoren ihr Speiseeis vorwiegend saisonal produzieren - zu deutlich höheren Preisen. Dadurch steigt der aus verschiedenen Wirtschaftsbereichen „zusammengesetzte“ Literpreis des betreffenden Gütercodes saisonal stark an, obwohl die Preise innerhalb der jeweiligen Wirtschaftsbereiche relativ konstant bleiben.

⁴⁰Last Observation Carried Forward

⁴¹Einige Merkmale (z.B. Sonderzahlungen) können nicht auf diese Weise fortgeschrieben werden, hier erfolgt die monatliche Berechnung über branchenspezifische mittlere Anteilswerte.

Abbildung 10: Beispiel fur Fortschreibung auf Aktivitatsebene und Guterebene



Quelle: KJP (01/2008 bis 04/2009).

6. Datenqualitat

Die im vorliegenden Beitrag beschriebene Methode der KJP-Datenerganzung illustriert die Vielzahl an Faktoren, welche bei der Beurteilung der Datenqualitat zu beruckichtigen sind. Die Beachtung verschiedenster Qualitatsfaktoren ist zum einen in den konzeptionellen Vorgaben begrundet, die keine konventionelle Hochrechnung zulassen, und zum anderen auf die Komplexitat des zu erstellenden Datenkorpers selbst zuruckzufuhren. Die Darstellung der KJP-Grundgesamtheit erfordert neben geeigneter mathematischer Verfahren auch das Einbeziehen von Informationen aus verschiedenen primar- und sekundarstatistischen Datenquellen. Neben der Schlusselrolle, die dabei dem Unternehmensregister der STATISTIK AUSTRIA zukommt, sind es vor allem Verfugbarkeit und konzeptionelle Eignung der externen Verwaltungsquellen, die in die **Gesamtqualitat der Daten** einflieen. Umfangreiche Voranalysen und Testberechnungen zeigten, dass die sekundarstatistisch erganzten Eckdaten grundsatzlich als „voll erhoben“ betrachtet werden konnen und die Gute der sonstigen beschaftigten- und umsatzbezogenen Merkmale davon abhangt, inwieweit ausreichende primarstatistische Informationen zur Modellbildung vorliegen. Uberall dort, wo das eigentliche Ziel einer MDE umgesetzt werden kann - namlich die modellbasierte Erganzung eines vorwiegend primarstatistisch erfassten Datenkorpers - kann grundsatzlich von hoher Datenqualitat ausgegangen werden. Diese Annahme ist einerseits darin begrundet, dass fur die Datenerganzung ausreichend primarstatistische Informationen zur Modellentwicklung vorliegen, und andererseits, dass selbst in Wirtschaftsbereichen mit deutlicher Strukturabweichung mogliche Modelleffek-

te mit wachsender Primärabdeckung gegen Null tendieren.

Während die MDE auf Aktivitätsebene für die Mehrzahl von Hauptaggregaten auch ohne ressourcenintensive Arbeitsschritte zufriedenstellende Ergebnisse erbringt, stellt sich die Ergänzung im Rahmen der Güterproduktion ungleich schwieriger dar. Beim Güteransatz sind der methodischen Vorgehensweise auch „natürliche Grenzen“ gesetzt. Güter, die von Betrieben in der Spendermasse nicht hergestellt werden, können auch nicht imputiert werden. Es fehlen somit auch bei einer MDE grundsätzlich jene Güter, die ausschließlich von kleinen Einheiten unterhalb der Cut-Off-Grenze produziert werden. Darüber hinaus zeigten in primärstatistisch schwach besetzten, nicht homogenen Wirtschafts-(teil-)bereichen und in Bereichen mit deutlicher Strukturabweichung die im Vorfeld durchgeführten Analysen, dass die Ergänzung nicht ausschließlich anhand automationsunterstützter Verfahren erfolgen kann. Trotz aller Informationen, die bereits in dieses System der modellbasierten Datenergänzung auf Mikrodatenebene eingehen und automatisch verarbeitet werden, ist in Einzelfällen qualifiziertes Expertenrating unerlässlich.

7. Resümee und Ausblick

In Österreich führte die politische Intention, Klein- und Kleinstunternehmen in höchstmöglicher Weise von der Auskunftspflicht auszunehmen, zur Situation, dass in einigen Wirtschaftsbereichen weder das Repräsentanzkriterium der nationalen Durchführungsverordnung, noch jenes der PRODCOM-Verordnung ausreichend umgesetzt werden konnte. Um ein Maximum an Repräsentanz und damit einhergehend, ein Optimum an Effizienz und Qualität des Datenkörpers zu gewährleisten, können jedoch gemäß den europäischen Bestimmungen Mitgliedstaaten die erforderlichen Daten nach dem Grundsatz der verwaltungstechnischen Vereinfachung durch eine Kombination verschiedener Quellen erstellen, sofern die übermittelten Daten die Struktur der Grundgesamtheit der statistischen Einheiten widerspiegeln. Dies ist mit der Entwicklung des Konzepts der modellbasierten Datenergänzung zur Darstellung der Grundgesamtheit der KJP seit dem Referenzjahr 2009 sowohl nach dem Aktivitäts- als auch nach dem Güteransatz in vollem Umfang gewährleistet.

Beginnend mit den endgültigen Ergebnissen des Berichtsjahres 2008 werden zusätzlich zum primärstatistischen Datenkörper auch die Ergebnisse der KJP Grundgesamtheit veröffentlicht. Weitere Entwicklungen, wie die Berücksichtigung weiterer sekundärstatistischer Quellen zur Verstärkung der Eckdaten und der an sie gebundenen Modellrechnung werden im Sinne einer ständigen Qualitätsverbesserung angestrebt. Derzeit besteht die Möglichkeit, im Zuge des **MEETS-Project 2010**⁴² eine künftige Einbindung von Lohnzetteldaten zu prüfen. Weiters sind methodische Arbeiten geplant, um zu prüfen, ob es durch den Einsatz statistisch mathematischer Methoden möglich ist, dem Nutzer zur Abschätzung der Datengenauigkeit geeignete Kennzahlen zur Verfügung zu stellen.

⁴² *Feasibility study of implementing wage tax data in structural business statistics.* MEETS (Modernisation of European Enterprise and Trade Statistics) projects on modernisation of Business Statistics - Use of administrative data.

Stand: Dezember, 2009

Projekt-Team

- Mag. Johann Hameseder (Direktion Unternehmen / Stv. Direktor)
- ADir. RR. Leopold Milota (Direktion Unternehmen / Projektleiter KJP)
- Mag. Martin Haitzmann (Direktion Unternehmen / Analyse)
- Dr. Mag. Martin Hirsch (Direktion Unternehmen / Analyse)
- Mag. Markus Fröhlich (Methodik / LOCF-Substitution)
- Mag. Gerlinde Dinges (Methodik / Modellbasierte Datenergänzung)

Statistik Austria, Guglgasse 13, A-1110 Wien.

Dieser Methodenbericht ist online abrufbar unter

[http : //www.statistik.at/web_de/downloads/methodik/kjp.pdf](http://www.statistik.at/web_de/downloads/methodik/kjp.pdf)