

Anpassung von Aktivitätenketten mittels wiederholter proportionaler Anpassung

Thomas Hettinger

Prof. Kay W. Axhausen

Betreut durch Nicolas Lefebvre

Semesterarbeit für den MSc Angewandte Mathematik

Juni 2007

Inhaltsverzeichnis

1	Einleitung	2
2	Der IPF-Algorithmus	2
2.1	Ein einführendes Beispiel	2
2.2	Der IPF-Algorithmus	3
3	Konvergenzbeweis	4
4	Anpassung von Aktivitätenketten	10
5	Mikrozensus 2000 und Mikrozensus 2005	15
6	Ausblick	18
	Bibliography	20

Abbildungsverzeichnis

1	Geradenscharen g_{α, ω_1} und h_{α, ω_2}	6
2	Visualisierung des IPF-Algorithmus	8

Tabellenverzeichnis

1	Ausgangsdaten aus Mikrozensus 2000 und 2005	16
2	Daten aus Mikrozensus nach Anwenden des IPF-Algorithmus	16
3	Vergleich Mikrozensus 2000 und 2005 sowie mittels IPF angepasste Daten . . .	17
4	Längen-Aktivitäten-Verteilung in der Schätzung \underline{n}^1	18

Semesterarbeit für den MSc Angewandte Mathematik

Anpassung von Aktivitätenketten mittels wiederholter proportionaler Anpassung

Thomas Hettinger
Kirchgasse 6
6318 Walchwil

Tel: 041 758 16 53

Fax:

hetthoma@student.ethz.ch

Juni 2007

Zusammenfassung

In dieser Semesterarbeit wird zuerst der *Iterative-Proportional-Fitting-Algorithmus* (IPF) beschrieben. Dabei handelt es sich um ein iteratives Verfahren zur Anpassung von Matrizen, so dass diese vorgegebene Randsummen-Bedingungen erfüllen. Des Weiteren wird ein anschaulicher Beweis der Konvergenz dieses Verfahrens gegeben.

In einem zweiten Teil wird folgendes Problem gestellt und ein Lösungsvorschlag beschrieben: Gegeben sei ein Satz von Aktivitätenketten-Häufigkeiten. Aus diesem soll ein zweiter Satz von Aktivitätenketten-Häufigkeiten so bestimmt werden, dass dieser Zweite gewisse vorgegebene Aktivitäten-Verteilungen und Kettenlängen-Verteilungen hat. Der vorgeschlagene Lösungsweg gliedert sich in drei Teilschritte: Der Bestimmung der Aktivitäten-Längen-Verteilung im ersten Datensatz, der Anpassung dieser Verteilung an die Aktivitäten- und Längen-Randverteilungen mittels IPF und der Bestimmung des zweiten Satzes von Ketten-Häufigkeiten aus dieser angepassten Aktivitäten-Längen-Verteilung.

Dieses Vorgehen wird noch auf die Mikrozensus 2000 und Mikrozensus 2005 - Daten angewendet, um zu überprüfen, ob es brauchbare Resultate liefert. Es stellt sich heraus, dass die Anpassung an die Aktivitäten-Längen-Verteilung bis auf kleine Abweichungen ihr Ziel erreicht.

Schlüsselwörter

Aktivitätenkette, IPF, Furness-Methode, Mikrozensus 2000, Mikrozensus 2005, Proportional Fitting

Bevorzugter Zitierstil

Hettinger, T. (2007) Anpassung von Aktivitätenketten mittels wiederholter proportionaler Anpassung, *Semesterarbeit für den MSc Angewandte Mathematik*, , Institut für Verkehrsplanung und Transportsysteme (IVT), ETH Zürich, Zürich.

1 Einleitung

Diese Semesterarbeit behandelt Iterative-Proportional-Fitting (IPF), ein iteratives Verfahren mit welchem mehrdimensionale Wahrscheinlichkeitsverteilungen an vorgegebene Randverteilungen angepasst werden können. Seine Anwendungen gehen weit über Verkehrsplanerische Themen hinaus, aber wir wollen uns in dieser Arbeit darauf beschränken.

Die ursprüngliche Prozedur geht zurück auf Bregman (1967), welcher die Lösung eines allgemeineren Optimierungsproblems beschrieb. In der Verkehrsplanung ist der IPF-Algorithmus auch bekannt unter dem Namen *Furness-Method* und beschreibt den Algorithmus zur Anpassung einer Start-Ziel-Fahrten-Matrix an vorgegebene Randsummen. Dabei wird die Ausgangsmatrix oftmals durch ein Gravitationsmodell erzeugt.

Beiwiese für die Konvergenz der des IPF-Algorithmus wurden unter anderem gegeben von Fienberg (1970), Evans (1970) oder, für eine dreidimensionale Erweiterung, von Kirby und Evans (1974). Die letzteren beiden Beweise sind eher analytischer Natur, während Fienberg den Algorithmus geometrisch interpretiert und so auch seine Konvergenz zeigt. Der Beweis der hier in Kapitel 3 wiedergegeben wird, folgt im Wesentlichen den Überlegungen von Fienberg.

Ob und wie der IPF-Algorithmus für die Anpassung von Aktivitätenketten gebraucht werden kann, wird im Kapitel 4 erörtert. Es stellt sich heraus, dass der IPF nicht allein zur Lösung des Problems führt. Jedoch wird ein Verfahren angegeben, in welchem die iterative proportionale Anpassung einen wesentlichen Teilschritt darstellt. Im letzten Kapitel wird dieses Verfahren zur Veranschaulichung noch auf Daten aus den Mikrozensus 2000 (ARE und BfS, 2001) und aus dem Mikrozensus 2005 (ARE und BfS, 2007) angewendet und mit Hilfe dieser Ergebnisse das Verfahren kurz bewertet.

2 Der IPF-Algorithmus

2.1 Ein einführendes Beispiel

Gegeben sei ein Verkehrsgebiet, welches in K Zonen unterteilt ist. Aus Konzepten der Fahrten-Erzeugung und Fahrten-Anziehung kann man schätzen wie viele Fahrten in jeder Zone starten (s_1, \dots, s_K) beziehungsweise enden (e_1, \dots, e_K) . Andererseits sei aus einer früheren Erhebung, aus einem Gravitationsmodell oder aus einer anderen Quelle eine Fahrtenmatrix N mit den Einträgen n_{ij} , $i, j = 1, \dots, K$ bekannt. Dabei bezeichnet n_{ij} die Anzahl Fahrten welche in Zone i starten und in Zone j enden.

Die Randsummen $\sum_j n_{ij}$, $i = 1, \dots, K$ geben an, wie viele Fahrten in der Erhebung in Zone

i starten. Diese Zahl sollte eigentlich mit den Ergebnissen der theoretischen Überlegungen übereinstimmen, dies ist aber nur in den wenigsten Fällen erfüllt. Genau wie die Zeilensummen sollten die Spaltensummen von N mit den theoretisch berechneten Werten für die Anzahl der in einer Zone endenden Fahrten übereinstimmen. Auch dies ist wohl nie exakt erreicht.

Da davon ausgegangen wird, dass die Genauigkeit der theoretischen Überlegungen besser ist als die der Quelle der Matrix N , versucht man nun die Matrix so zu verändern, dass sie die folgenden Randsummenbedingungen erfüllt:

$$\sum_j n_{ij} = s_i, \quad i = 1, \dots, K \quad (1)$$

$$\sum_i n_{ij} = e_j, \quad j = 1, \dots, K \quad (2)$$

Eine Möglichkeit, wie dies erreicht werden kann wird im folgenden Abschnitt beschrieben.

2.2 Der IPF-Algorithmus

Gegeben seien Randsummen $\mathbf{a} = (a_1, \dots, a_K)$ und $\mathbf{b} = (b_1, \dots, b_L)$ sowie eine $(K \times L)$ -Matrix N . Gesucht ist eine Matrix welche die Zeilensummen \mathbf{a} und die Spaltensummen \mathbf{b} hat und sich „möglichst wenig“ von der Matrix N unterscheidet. Was man genau unter dem Begriff „möglichst wenig“ versteht, wird anschliessend erläutert.

Folgendes Verfahren findet eine solche Matrix:

IPF-Algorithmus

1. setze

$$p_{ij}^0 := \frac{n_{ij}}{\sum_{i'j'} n_{i'j'}},$$

$$p_{i\cdot}^0 := \sum_j p_{ij}^0, \quad p_{\cdot j}^0 := \sum_i p_{ij}^0,$$

$$\alpha_i := \frac{a_i}{\sum_{i'} a_{i'}}, \quad \beta_j := \frac{b_j}{\sum_{j'} b_{j'}}$$

2. für $t = 1, 2, 3, \dots$ berechne

(a) $p_{ij}^{2t-1} := p_{ij}^{2t-2} \frac{\alpha_i}{p_{i\cdot}^{2t-2}}$
 $p_{i\cdot}^{2t-1}$ und $p_{\cdot j}^{2t-1}$ werden analog zu Punkt 1. definiert.

$$(b) \quad p_{ij}^{2t} = p_{ij}^{2t-1} \frac{\beta_j}{p_{\cdot j}^{2t-1}}$$

$p_{i\cdot}^{2t}$ und $p_{\cdot j}^{2t}$ werden analog zu Punkt 1. definiert.

3. für $i = 1, \dots, K$ und $j = 1, \dots, L$ berechne $\hat{n}_{ij} = p_{ij}^{2t} \cdot \sum_{i'=1}^K a_{i'}$
gib die Matrix \hat{N} mit den Einträgen \hat{n}_{ij} aus

Im ersten Schritt des Algorithmus werden alle Matrixeinträge mit

$$\frac{1}{\sum_{ij} n_{ij}}$$

multipliziert, so dass sie sich zu 1 aufsummieren. Diese Einträge p_{ij}^0 können auch als Wahrscheinlichkeit interpretiert werden, mit welcher man Zelle (i, j) wählt, wenn man zufällig eine der Zellen der Matrix auswählt. Genauso lassen sich die Zeilen- beziehungsweise die Spaltentotalen $p_{i\cdot}$ und $p_{\cdot j}$ verstehen: Als Wahrscheinlichkeit eine Zelle aus Zeile i oder aus Spalte j zu wählen. Sie werden auch Randverteilungen genannt. Die zu erreichenden Randverteilungen sind durch die α_i und die β_j gegeben.

Unter Punkt 2 (a) wird jeweils jede Zeile i mit einem Faktor $\frac{\alpha_i}{p_{i\cdot}^{2t-2}}$ multipliziert, was dazu führt dass die Zeilensummen genau die gewünschte Randverteilung $(\alpha_1, \dots, \alpha_K)$ erreichen. Unter Punkt 2 (b) wird dieselbe Anpassung Spaltenweise durchgeführt was zu einer Übereinstimmung der Spaltensummen mit $(\beta_1, \dots, \beta_L)$ führt. Allerdings ist nach jeder Anpassung der Spaltensummen die vorherige Anpassung der Zeilensummen und umgekehrt nach jeder Anpassung der Zeilensummen die vorherige Anpassung der Spaltensummen wieder verletzt.

Im folgenden Kapitel wird aber bewiesen dass dieses Verfahren gegen eine Matrix konvergiert, welche beide Randsummenbedingungen erfüllt.

3 Konvergenzbeweis

Der Konvergenzbeweis gliedert sich in zwei wesentliche Schritte: Als erstes wird der Beweis für den Fall einer 2×2 -Matrix geführt. Danach wird die Idee dieses Beweises auf den Fall einer $r \times c$ -Matrix ausgedehnt.

Eine 2×2 -Wahrscheinlichkeitsmatrix hat vier Einträge $P = \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix}$ welche sich zu 1 aufsummieren $p_{11} + p_{12} + p_{21} + p_{22} = 1$.

In Schritt 2 (a) des IPF-Algorithmus wird jede Zeile und in Schritt 2 (b) wird jede Spalte mit einem Faktor multipliziert. Daraus sieht man, dass nach jedem Schritt die *Cross-Product-Ratio* von P , $cpr(P) := \frac{p_{11}p_{22}}{p_{21}p_{12}}$, erhalten bleibt.

Wir betrachten nun die Menge $M_\alpha = \left\{ M \in R^{2 \times 2} \mid cpr(M) = \alpha, \sum_{ij} m_{ij} = 1 \right\}$ aller Wahrscheinlichkeitsmatrizen mit $cpr = \alpha$. Diese Menge umfasst alle Matrizen welche während des Algorithmus besucht werden, wenn mit einer Matrix gestartet wird, welche $cpr = \alpha$ hat.

Im folgenden wird eine Bijektion von M_α nach dem Einheitsquadrat $[0, 1] \times [0, 1]$ angegeben, um dann zu zeigen dass die Hintereinanderschaltung der beiden Schritte 2 (a) und 2 (b) eine Kontraktion im Einheitsquadrat darstellt, welche die gesuchte Lösung als Fixpunkt hat.

Für gegebene ω_1, ω_2 und α definieren wir die Geraden

$$g_{\alpha, \omega_1}(t) = \begin{pmatrix} (1-t)\omega_1 + t \frac{\omega_1}{\alpha(1-\omega_1) + \omega_1} \\ t \end{pmatrix}, t \in (0, 1)$$

und

$$h_{\alpha, \omega_2}(s) = \begin{pmatrix} 1-s \\ (1-s)(1-\omega_2) + s \frac{\alpha(1-\omega_2)}{\alpha(1-\omega_2) + \omega_2} \end{pmatrix}, s \in (0, 1)$$

Setzt man $\omega_1 = \frac{p_{11}}{p_{11} + p_{12}}$ und $\omega_2 = \frac{p_{11}}{p_{11} + p_{21}}$, so lässt sich jede Matrix aus M_α mit dem Schnittpunkt der beiden Geraden $g_{\alpha, \omega_1}(t)$ und $h_{\alpha, \omega_2}(s)$ identifizieren und so hat man die oben erwähnte Bijektion konstruiert. Dass es wirklich eine Bijektion ist, sieht man dadurch ein, dass zwei Geraden $g_{\alpha, \omega_1}(t)$ und $g_{\alpha, \omega'_1}(t)$ mit $\omega_1 \neq \omega'_1$ und $t \in (0, 1)$ keine gemeinsamen Punkte haben. Dasselbe gilt für Geraden $h_{\alpha, \omega_2}(s)$ und $h_{\alpha, \omega'_2}(s)$ mit $\omega_2 \neq \omega'_2$.

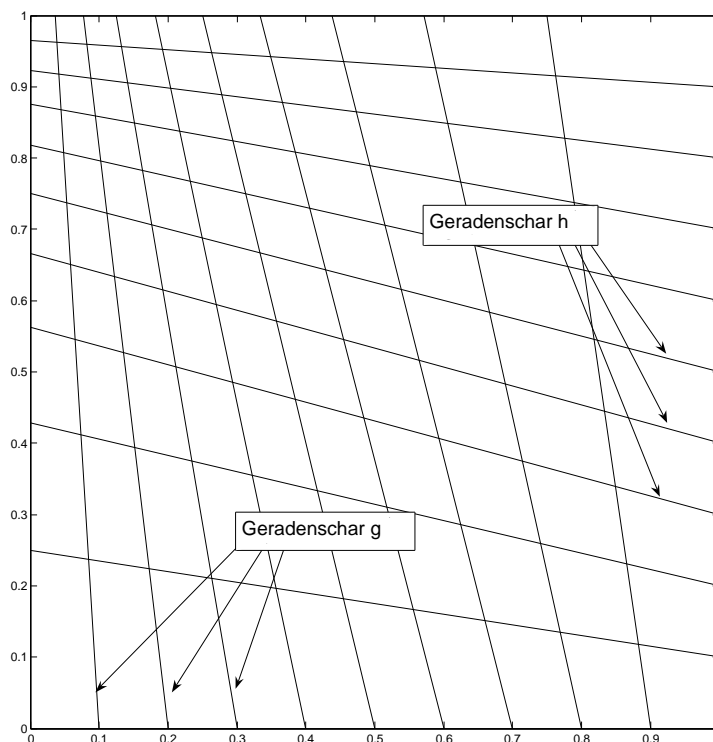
Da Schritt 2 (a) im IPF-Algorithmus den Wert von ω_1 unverändert lässt, entspricht dieser Schritt einer Translation entlang der Geraden $g_{\alpha, \omega_1}(t)$. Analog sieht man leicht ein, dass Schritt 2 (b) einer Translation entlang der Geraden $h_{\alpha, \omega_2}(s)$ entspricht.

Aus den Gleichungen

$$p_{11} + p_{12} + p_{21} + p_{22} = 1 \tag{3}$$

und

$$\frac{p_{11}p_{22}}{p_{21}p_{12}} = \alpha \tag{4}$$

Abbildung 1: Geradenscharen g_{α, ω_1} und h_{α, ω_2} für $\alpha = 3$ und $\omega_1, \omega_2 = 0, 0.1, \dots, 1$ 

Quelle: nach Fienberg (1970), Seite 912

folgt, dass

$$p_{11} = \frac{\alpha p_{21}(1 - p_{21} - p_{22})}{\alpha p_{21} + p_{22}} \quad (5)$$

$$p_{12} = \frac{p_{22}(1 - p_{21} - p_{22})}{\alpha p_{21} + p_{22}} \quad (6)$$

$$p_{11} = \frac{\alpha p_{12}(1 - p_{12} - p_{22})}{\alpha p_{12} + p_{22}} \quad (7)$$

$$p_{21} = \frac{p_{22}(1 - p_{12} - p_{22})}{\alpha p_{12} + p_{22}} \quad (8)$$

Setzt man nun $t = p_{21} + p_{22}$ und $s = 1 - p_{11} - p_{21} = p_{12} + p_{22}$ so sieht man, dass

$$\begin{aligned}
 g_{\alpha, \omega_1}(t) &= \begin{pmatrix} (p_{11} + p_{12}) \frac{p_{11}}{p_{11} + p_{12}} + (p_{21} + p_{22}) \frac{\frac{p_{11}}{p_{11} + p_{12}}}{\alpha \frac{p_{12}}{p_{11} + p_{12}} + \frac{p_{11}}{p_{11} + p_{12}}} \\ p_{21} + p_{22} \end{pmatrix} \\
 &= \begin{pmatrix} p_{11} + (p_{21} + p_{22}) \frac{p_{11}}{\alpha p_{12} + p_{11}} \\ p_{21} + p_{22} \end{pmatrix} \\
 &\stackrel{(5),(6)}{=} \begin{pmatrix} p_{11} + (p_{21} + p_{22}) \frac{\frac{\alpha p_{21}(1-p_{21}-p_{22})}{\alpha p_{21} + p_{22}}}{\frac{\alpha p_{22}(1-p_{21}-p_{22})}{\alpha p_{21} + p_{22}} + \frac{\alpha p_{21}(1-p_{21}-p_{22})}{\alpha p_{21} + p_{22}}} \\ p_{21} + p_{22} \end{pmatrix} \\
 &= \begin{pmatrix} p_{11} + (p_{21} + p_{22}) \frac{p_{21}}{p_{21} + p_{22}} \\ p_{21} + p_{22} \end{pmatrix} \\
 &= \begin{pmatrix} p_{11} + p_{21} \\ p_{21} + p_{22} \end{pmatrix}
 \end{aligned}$$

und analog berechnet sich mit Hilfe von (7) und (8), dass

$$h_{\alpha, \omega_2}(s) = \begin{pmatrix} p_{11} + p_{21} \\ p_{21} + p_{22} \end{pmatrix}$$

Damit haben wir auch die x - und die y -Koordinate des Punktes gefunden, welcher mit der Matrix $P = \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix}$ identifiziert wird:

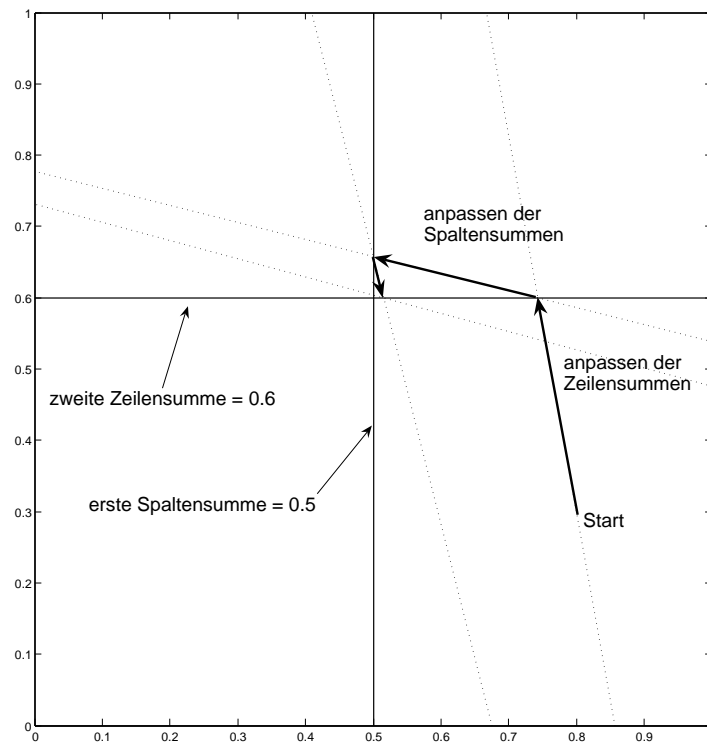
$$P \hat{=} \begin{pmatrix} p_{11} + p_{21} \\ p_{21} + p_{22} \end{pmatrix} = \begin{pmatrix} \text{erste Spaltensumme von } P \\ \text{zweite Zeilensumme von } P \end{pmatrix}$$

Insbesondere wurde auch gezeigt, dass alle Matrizen mit den gleichen Zeilensummen dieselbe y -Koordinate haben, beziehungsweise alle Matrizen mit der gleichen Spaltensumme dieselbe x -Koordinate haben.

Schritt 2 (a) des IPF Algorithmus geht also entlang der Gerade $g_{\alpha, \omega_1}(t)$ bis die entsprechende Zeilensumme beziehungsweise y -Koordinate erreicht ist und Schritt 2(b) folgt dann der Geraden $h_{\alpha, \omega_2}(s)$ bis die entsprechende Spaltensumme beziehungsweise x -Koordinate erreicht ist.

In Abbildung 2 sieht man, dass die Geraden $g_{\alpha, \omega_1}(t)$ mit der Horizontalen und die Geraden $h_{\alpha, \omega_2}(s)$ mit der Vertikalen jeweils einen Winkel grösser als 45° einschliessen. Um diese Aussage für die Geraden $g_{\alpha, \omega_1}(t)$ zu beweisen, reicht es zu zeigen dass der Betrag der Ableitung

Abbildung 2: Visualisierung des IPF-Algorithmus mit Start-Matrix $\begin{pmatrix} 0.6 & 0.1 \\ 0.2 & 0.1 \end{pmatrix}$ ($\alpha = 3$),
Ziel-Zeilensummen (0.4, 0.6) und Ziel-Spaltensummen (0.5, 0.5)



Quelle: nach Fienberg (1970), Seite 913

der ersten Komponente nach t für alle ω_1 kleiner als 1 ist:

$$\begin{aligned} \left| \frac{\partial g_{\alpha, \omega_1}(t)_1}{\partial t} \right| &\leq 1 \\ \left| -\omega_1 + \frac{\omega_1}{\alpha(1-\omega_1) + \omega_1} \right| &\leq 1 \\ \left(-\omega_1 + \frac{\omega_1}{\alpha(1-\omega_1) + \omega_1} \right)^2 &\leq 1 \\ \omega_1^2 - 2 \frac{\omega_1^2}{\alpha(1-\omega_1) + \omega_1} + \frac{\omega_1^2}{\alpha^2(1-\omega_1)^2 + \alpha\omega_1(1-\omega_1) + \omega_1^2} &\leq 1 \\ \omega_1^2(\alpha(1-\omega_1) + \omega_1)^2 - 2\omega_1^2(\alpha(1-\omega_1) + \omega_1) + \omega_1^2 &\leq (\alpha(1-\omega_1) + \omega_1)^2 \end{aligned}$$

$$\begin{aligned}\omega_1^2(\alpha(1 - \omega_1) + \omega_1 - 1)^2 &\leq (\alpha(1 - \omega_1) + \omega_1)^2 \\ \omega_1^2(1 - \omega_1)^2(\alpha - 1)^2 &\leq \alpha^2(1 - \omega_1)^2 + 2\alpha(1 - \omega_1)\omega_1 + \omega_1^2\end{aligned}$$

Die letzte Ungleichung ist sicher richtig für $\alpha \leq 2$, denn dann ist $(\alpha - 1)^2 \leq 1$ und somit ist die linke Seite kleiner als der dritte Summand der rechten Seite, weil $\omega_1 \in (0, 1)$ und $\alpha \geq 0$. Und da die anderen beiden Summanden auch ≥ 0 sind, ist die Ungleichung erfüllt. Für $\alpha \geq 2$ ist die linke Seite sicher kleiner als der erste Summand der rechten Seite und somit ist die Aussage für die Geraden $g_{\alpha, \omega_1}(t)$ bewiesen. Für die Geraden $h_{\alpha, \omega_2}(t)$ lässt sie sich ganz analog beweisen.

Weil diese Winkel nun alle grösser als 45° sind, ist man jeweils nach jedem Schritt näher am Ziel als vorher

Punkt 2 des Algorithmus als Ganzes definiert deshalb eine Kontraktion ϕ auf $[0, 1] \times [0, 1]$ welche genau die Matrix als Fixpunkt hat, welche die gewünschten Randsummenbedingungen erfüllt und die gleiche cross-product-ratio hat wie die Ausgangsmatrix.

Damit ist der Konvergenzbeweis des IPF-Algorithmus für eine 2×2 -Matrix erbracht.

Eine wichtige Konstante der Matrix, welche während des ganzen Algorithmus erhalten bleibt, ist die cross-product-ratio α . Für die Verallgemeinerung von einer (2×2) -Matrix auf eine $(r \times c)$ -Matrix muss zuerst das Pendant zur cross-product-ratio gefunden werden. Die Matrix M habe die Einträge m_{ij} , ($i = 1, \dots, r$; $j = 1, \dots, c$). Man kann hier nun wieder die cross-product-ratio bezüglich Zeilen i, k und Spalten j, l bestimmen: $cpr_{ijkl} := \frac{m_{ij}m_{kl}}{m_{il}m_{jk}}$. Alle diese Werte bleiben während des IPF-Algorithmus ebenfalls erhalten. Sie sind allerdings nicht unabhängig voneinander. Beispielsweise kann aus $cpr_{1122} = \frac{m_{11}m_{22}}{m_{12}m_{21}}$ und $cpr_{1223} = \frac{m_{12}m_{23}}{m_{13}m_{22}}$ cpr_{1123} einfach durch $cpr_{1123} = \frac{m_{11}m_{23}}{m_{21}m_{13}} = \frac{m_{11}m_{22}}{m_{21}m_{12}} \frac{m_{12}m_{23}}{m_{22}m_{13}} = cpr_{1122} \cdot cpr_{1223}$ bestimmt werden. Man sieht, dass alle cross-product-ratios eindeutig festgelegt sind durch die $cpr_{i(i+1)j(j+1)}$, ($i = 1, \dots, r - 1$; $j = 1, \dots, c - 1$).

Man definiert

$$M_\alpha = \left\{ M \in R \mid cpr_{i(i+1)j(j+1)} = \alpha_{ij}, (i = 1, \dots, r - 1; j = 1, \dots, c - 1), \sum_{ij} m_{ij} = 1 \right\}$$

wobei α hier eine $(r - 1) \times (c - 1)$ -Matrix mit Einträgen α_{ij} bezeichnet. Es wird nun wieder eine Bijektion von M_α nach einem Teilraum des Einheitsquaders $[0, 1]^2$ angegeben: Die Matrix $M \in M_\alpha$ wird identifiziert mit dem Punkt $\mathbf{x} = (\sum_j m_{1j}, \dots, \sum_j m_{(r-1)j}, \sum_i m_{i1}, \dots, \sum_i m_{i(c-1)})$ das heisst die ersten $r - 1$ Komponenten von \mathbf{x} entsprechen den ersten $r - 1$ Zeilensummen von M und die weiteren $c - 1$ Komponenten von \mathbf{x} entsprechen den ersten $c - 1$ Spaltensummen von M .

Wiederum verändern sich während eines Schritts 2 (a) die Werte

$$\omega_j^1 := \frac{m_{1,j}}{\sum_{l=1}^c m_{1,l}} \quad j = 1, \dots, c-1$$

nicht. Alle Matrizen welche die gleichen Werte ω_j^1 $j = 1, \dots, c-1$ haben liegen auf einer $(r-1)$ -dimensionalen Hyperebene

$$\epsilon_{\omega_1^1, \dots, \omega_{c-1}^1, \alpha}(t_1, \dots, t_{r-1}) = \begin{pmatrix} t_1 \\ \vdots \\ t_{r-1} \\ \sum_{i=1}^r m_{i1}(t_1, \dots, t_{r-1}) \\ \vdots \\ \sum_{i=1}^r m_{i(c-1)}(t_1, \dots, t_{r-1}) \end{pmatrix}$$

mit

$$m_{ij}(t_1, \dots, t_{r-1}) = \frac{t_i \omega_j^1 \prod_{k=1}^{i-1} \prod_{l=1}^{j-1} \alpha_{kl}}{\sum_{j'=1}^c \omega_{j'}^1 \prod_{k=1}^{i-1} \prod_{l=1}^{j'-1} \alpha_{kl}},$$

(wobei $t_r = 1 - \sum_{i=1}^{r-1} t_i$ und $\omega_c^1 = 1 - \sum_{j=1}^{c-1} \omega_j^1$)

auf der man sich also während des Schrittes 2(a) bewegt. Analog zum (2×2) -Fall schneiden sich diese Ebenen nicht im Einheitsquader und schliessen mit den „horizontalen“ Ebenen $(const_1, \dots, const_{r-1}, s_1, \dots, s_{c-1})^\top$ einen Winkel grösser als 45° ein.

Analog zu obigem Abschnitt liegen die Matrizen, für welche die Werte

$$\omega_i^2 := \frac{m_{i,1}}{\sum_{k=1}^r m_{k,1}} \quad i = 1, \dots, r-1$$

gleich sind, auf einer $(c-1)$ -dimensionalen Hyperebene, auf der man sich während des Schrittes 2 (b) bewegt. Und wieder schneiden sich diese Ebenen nicht im Einheitsquader und schliessen mit den „vertikalen“ Ebenen einen Winkel grösser als 45° ein.

Daraus folgt wieder dass der Punkt 2 als ganzes im IPF-Algorithmus eine Kontraktion im Einheitsquader darstellt, deren Fixpunkt gerade derjenigen Matrix entspricht, welche die gewünschten Zeilen- und Spaltensummen hat.

4 Anpassung von Aktivitätsketten

Oft sind aus früheren Datensätzen Aktivitätsketten und deren Verteilungen (bzw Anzahlen in einer Erhebung) bekannt. Daraus lässt sich auf einfache Weise die (Rand-)Verteilung der Ak-

tivitäten und (Rand-)Verteilung der Aktivitätenkettenlängen bestimmen. Von einer aktuelleren Erhebung seien nun nur die entsprechenden Randverteilungen der Aktivitätenketten bekannt. Die Aufgabe besteht nun darin, die früheren vollständigen Daten an die postulierten Randverteilungen der neueren Erhebungen anzupassen, also eine neue Verteilung der Aktivitätenketten zu bestimmen.

Im Folgenden bezeichnen

- T_k , $k = 1, \dots, K$ die verschiedenen Aktivitätenkettentypen,
- L_i , $i = 1, \dots, I$ die verschiedenen Aktivitätenkettenlängen, sowie
- A_j , $j = 1, \dots, J$ die verschiedenen Aktivitäten.

Beispiel

Es seien 6 Aktivitätenketten gegeben $T_1 = \text{h-w-h}$, $T_2 = \text{h-e-h}$, $T_3 = \text{h-w-h-w-h}$, $T_4 = \text{h-e-h-e-h}$, $T_5 = \text{h-w-h-e-h}$, $T_6 = \text{h-w-e-w-h}$. Dabei kommen die Aktivitätenkettenlängen $L_1 = 3$ und $L_2 = 5$ sowie die Aktivitäten $A_1 = \text{h}$ (home), $A_2 = \text{w}$ (work) und $A_3 = \text{e}$ (education) vor.

An Daten seien gegeben

- n_k^0 , $k = 1, \dots, K$ die Anzahl der Aktivitätenketten vom Typ T_k (z.B. aus einer früheren Arbeit) und
- die postulierten Randsummen der Aktivitäten $p_j^{(A)}$, $j = 1, \dots, J$ und der Aktivitätenkettenlängen $p_i^{(L)}$, $i = 1, \dots, I$.

Beispiel (Fortsetzung)

Anzahl Aktivitätenketten: $\underline{n}^0 = (n_1^0, \dots, n_6^0) = (150, 50, 200, 30, 40, 10)$

postulierte Randverteilungen der Aktivitäten: $(p_1^{(A)}, p_2^{(A)}, p_3^{(A)}) = (700, 200, 300) =$ geschätzte zukünftige Anzahl Aktivitäten von den Typen (h,w,e)

postulierte Randverteilungen der Kettenlängen: $(p_1^{(L)}, p_2^{(L)}) = (420, 780) =$ geschätzte zukünftige Anzahl Aktivitäten innerhalb von Ketten der Länge 3 beziehungsweise 5. Die Anzahl der Ketten der Länge 3 (bzw 5) lässt sich durch Division durch die Anzahl Aktivitäten pro Kette bestimmen. Das heisst man postuliert hier 140 Ketten der Länge 3 und 156 Ketten der Länge 5.

Ziel ist es, aus den Daten eine zukünftige Verteilung von Aktivitätenketten n_k^1 , $k = 1, \dots, K$ zu bestimmen. Das Vorgehen gliedert sich nun wie folgt:

1. Aus der Verteilung der Aktivitätenketten n_k^0 wird die Längen-Aktivitäten-Matrix B^0 berechnet, deren Einträge b_{ij}^0 die totale Anzahl der Aktivitäten vom Typ A_j in Ketten der Länge L_i bezeichnen.
2. Auf die Matrix B^0 wird der IPF-Algorithmus mit Randsummenbedingungen $\underline{p}_i^{(L)}$ und $\underline{p}_j^{(A)}$ angewandt. B^1 sei das Ergebnis.
3. Aus den Einträgen in B^1 wird durch eine Art Umkehrung von Schritt 1 die geschätzte Verteilung der Aktivitätenketten n_k^1 berechnet.

Der erste Schritt kann folgendermassen ausgeführt werden: Die Matrix B^0 wird als Vektor interpretiert, der entsteht, wenn man die Matrix zeilenweise liest: $B^0 \hat{=} (b_{11}^0, \dots, b_{1,J}^0, \dots, b_{I,1}^0, \dots, b_{I,J}^0)^\top =: \underline{b}^0$. Dann ist

$$\underline{b}^0 = M \underline{n}^0$$

wobei M in der zu b_{ij}^0 gehörenden Zeile und der zu n_k^0 gehörenden Spalte die Anzahl Aktivitäten vom Typ A_j in einer Kette vom Typ T_k bezeichnet.

Beispiel (Fortsetzung)

$$M = \begin{matrix} & \begin{matrix} \text{h-w-h} \\ \text{h-e-h} \\ \text{h-w-h-w-h} \\ \text{h-e-h-e-h} \\ \text{h-w-h-e-h} \\ \text{h-w-e-w-h} \end{matrix} \\ \begin{pmatrix} 2 & 2 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 3 & 3 & 2 \\ 0 & 0 & 2 & 0 & 1 & 2 \\ 0 & 0 & 0 & 2 & 1 & 1 \end{pmatrix} & , \underline{n}^0 = \begin{pmatrix} 150 \\ 50 \\ 200 \\ 30 \\ 40 \\ 10 \end{pmatrix} \Rightarrow \underline{b}^0 = \begin{pmatrix} 400 \\ 150 \\ 50 \\ 830 \\ 460 \\ 110 \end{pmatrix} \end{matrix}$$

$$\Rightarrow B^0 = \begin{pmatrix} 400 & 150 & 50 \\ 830 & 460 & 110 \end{pmatrix}$$

Der zweite Schritt entspricht genau dem in den vorherigen Kapiteln ausführlich besprochenen IPF-Algorithmus mit Matrix B^0 und Randsummenbedingungen $\underline{p}^{(L)}$ und $\underline{p}^{(A)}$.

Beispiel (Fortsetzung)

$$\underline{p}^{(L)} = \begin{pmatrix} 420 \\ 780 \end{pmatrix}, \underline{p}^{(A)} = \begin{pmatrix} 700 \\ 200 \\ 300 \end{pmatrix} \xrightarrow{\text{IPF}} B^1 = \begin{pmatrix} 257.3 & 56.5 & 106.2 \\ 442.7 & 143.5 & 193.8 \end{pmatrix}$$

Im dritten Schritt muss nur das Gleichungssystem

$$M \cdot \underline{n}^1 = \underline{b}^1 \quad (9)$$

gelöst werden. Dies ist allerdings nicht einfach, wie man folgendermassen einsehen kann: Das Gleichungssystem XY lässt sich aufteilen in I unabhängige Teil-Gleichungssysteme $M_i \underline{n}_i^1 = \underline{b}_i^1$. Ob diese Lösungen besitzen oder nicht hängt insbesondere davon ab, wieviele Aktivitätsketten der Länge L_i berücksichtigt werden, das heisst wieviele Spalten M_i hat. Es ist auch möglich, dass das Gleichungssystem lösbar ist, aber keine sinnvollen Lösungen hat, das heisst keine Lösungen welche nur Einträge ≥ 0 hat. Andererseits gibt es auch solche Gleichungssysteme, welche mehrere oder sogar unendlich viele sinnvolle Lösungen haben. Dann stellt sich die Frage, welche dieser Lösungen man nehmen soll.

Man kann folgendermassen vorgehen:

- Hat das i -te Teil-Gleichungssystem sinnvolle Lösungen, das heisst Lösungen deren Einträge alle nichtnegativ sind? Oder anders: Gibt es ein $\underline{n}_i^1 \geq \underline{0}$ so dass $M_i \cdot \underline{n}_i^1 = \underline{b}_i^1$? Existiert genau ein solches \underline{n}_i^1 so ist man fertig, gibt es aber mehrere so soll unter ihnen eines ausgewählt werden. Nach welchen Kriterien diese Auswahl geschehen soll, kann man unterschiedlicher Meinung sein. Der Autor schlägt folgendes vor: Man soll dasjenige \underline{n}_i^1 auswählen, welches zu einem (skalierten) $\lambda \cdot \underline{n}_i^0$ den geringsten euklidischen Abstand hat. \underline{n}_i^0 soll skaliert werden, damit die Summe der Einträge in $\lambda \cdot \underline{b}_i^0 = M_i \cdot \lambda \cdot \underline{n}_i^0$ gleich der Summe der Einträge in $\underline{b}_i^1 = M_i \cdot \underline{n}_i^1$ ist, welche wiederum gleich $p_i^{(L)}$ ist. Damit ist λ als $\frac{i\text{-te Zeilensumme in } B_1}{i\text{-te Zeilensumme in } B_0} = \frac{p_i^{(L)}}{p_i^{(L)0}}$ zu wählen.

Im Falle von mehreren sinnvollen Lösungen ist also das folgende Optimierungsproblem zu lösen:

$$\min_{\substack{M_i \cdot \underline{n}_i^1 = \underline{b}_i^1 \\ \underline{n}_i^1 \geq \mathbf{0}}} \left\| \underline{n}_i^1 - \lambda \cdot \underline{n}_i^0 \right\|_2, \quad \text{wobei } \lambda = \frac{p_i^{(L)}}{p_i^{(L)0}} \quad (A)$$

- Hat das i -te Teil-Gleichungssystem keine sinnvollen Lösungen, so muss eine Approximation einer sinnvollen Lösung gefunden werden. Wiederum kann eigentlich frei entschieden werden, nach welchen Kriterien die Güte einer Approximation bewertet werden soll. Auch hier hat sich der Autor für eine Minimierung des euklidischen Abstands entschieden, was zu folgendem Optimierungsproblem führt:

$$\min_{\underline{n}_i^1 \geq \mathbf{0}} \left\| M_i \cdot \underline{n}_i^1 - \underline{b}_i^1 \right\|_2 \quad (B)$$

Die Entscheidung, zu welchem der beiden Fälle ein Teil-Gleichungssystem gehört, kann ge-

troffen werden, indem zuerst das Problem (B) gelöst wird. Hat dieses eine Lösung mit Ziel-funktionswert Null, das heisst mit $M_i \cdot \underline{n}_i^1 = \underline{b}_i^1$, so gibt es eine exakte Lösung und es muss noch (A) gelöst werden.

Die beiden hier beschriebenen Optimierungsprobleme gehören zur selben Kategorie von Problemen, nämlich denjenigen mit quadratischer Zielfunktion und linearen Bedingungen. Zur Lösung solcher Aufgaben stellen gängige Programme zur Lösung mathematischer Probleme (beispielsweise MATLAB) effiziente Algorithmen bereit.

Beispiel (Fortsetzung)

Das erste Teil-Gleichungssystem behandelt die Aktivitätsketten der Länge 3:

$$M_1 \cdot \begin{pmatrix} n_1^1 \\ n_2^1 \end{pmatrix} = \underline{b}_1^1$$

wobei $M_1 = \begin{pmatrix} 2 & 2 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}$ und $\underline{b}_1^1 = \begin{pmatrix} 257.3 \\ 56.5 \\ 106.2 \end{pmatrix}$.

Dieses Gleichungssystem hat keine Lösung da es überbestimmt ist. Es gilt also das Optimierungsproblem

$$\min_{\underline{n}_1^1 \geq 0} \|M_1 \cdot \underline{n}_1^1 - \underline{b}_1^1\|_2$$

zu lösen. Es hat die Lösung $\underline{n}_1^1 = (41.3, 91.1)^\top$.

Das zweite Teil-Gleichungssystem behandelt die Aktivitätsketten der Länge 5:

$$M_2 \cdot \begin{pmatrix} n_3^1 \\ \vdots \\ n_6^1 \end{pmatrix} = \underline{b}_2^1$$

wobei $M_2 = \begin{pmatrix} 3 & 3 & 3 & 2 \\ 2 & 0 & 1 & 2 \\ 0 & 2 & 1 & 1 \end{pmatrix}$ und $\underline{b}_2^1 = \begin{pmatrix} 442.7 \\ 143.5 \\ 193.8 \end{pmatrix}$.

Dieses Gleichungssystem hat unendlich viele Lösungen. Wir suchen also die Lösung des Optimierungsproblems

$$\min_{\substack{M_2 \cdot \underline{n}_2^1 = \underline{b}_2^1 \\ \underline{n}_2^1 \geq 0}} \|\underline{n}_2^1 - \tilde{\underline{n}}_2^0\|_2, \quad \text{wobei } \tilde{\underline{n}}_2^0 = \frac{p_i^{(L)}}{p_i^{(L)0}} \cdot \underline{n}_2^0 = \frac{780}{1400} \cdot \underline{n}_2^0,$$

welche gegeben ist durch $\underline{n}_2^1 = (38.6, 76.4, 15.7, 25.3)^\top$

5 Mikrozensus 2000 und Mikrozensus 2005

Bisher wurde nur anhand von theoretischen Überlegungen und an Beispielen erklärt, wie Aktivitätenketten mit Hilfe des IPF-Algorithmus an gewisse Randbedingungen angepasst werden können. Dass dieses Vorgehen auch in der Wirklichkeit gebraucht werden kann wird in diesem Abschnitt aufgezeigt. Die Ausgangsdaten stammen aus dem Mikrozensus 2000 (ARE und BfS, 2001) und beinhalten im wesentlichen eine Liste von Aktivitätenketten und deren Häufigkeiten. Aus dem Mikrozensus 2005 (ARE und BfS, 2007) sind im wesentlichen dieselben Daten vorhanden, das heisst, eigentlich kennt man die Verteilung der Aktivitätenketten im Mikrozensus 2005 bereits. Hier sollen aber diese 2005-Daten nur dazu gebraucht werden, um die Randverteilungen zu bestimmen, an welche die 2000-Daten im IPF-Algorithmus angepasst werden.

Man versucht also aus der Verteilung der Aktivitätenketten 2000 und den Randsummen 2005 die Verteilung der Aktivitätenketten 2005 zu bestimmen. Schlussendlich kann man die Daten aus dem Mikrozensus 2005 auch dazu verwenden, die Resultate der Anpassung in diesem Fall zu überprüfen und zu bewerten.

Es sei noch erwähnt, dass hier nur Aktivitätenketten die in beiden Mikrozensen vorkommen betrachtet werden. Eine Verallgemeinerung auf weitere Aktivitätenketten stellt aber kein Problem dar, wenn man einfach die Ketten mit den entsprechenden Häufigkeiten (auch mit Häufigkeit 0) in die Daten einfügt.

Unter Berücksichtigung dass im Mikrozensus 2000 etwa 4 mal mehr Aktivitäten registriert wurden als im Mikrozensus 2005, lässt sich doch festhalten, dass sich zwischen 2000 und 2005 eine Verschiebung zu Gunsten der Aktivitätenketten mit Längen 5,7 und 9 und auf Kosten der Ketten mit Längen 4,6 und 10 ergeben hat. Bei den Verteilungen der Aktivitäten sind diese Veränderungen praktisch nicht spürbar (siehe Tabelle 1).

In Tabelle 2 sehen wir die Aktivitäten-Längen-Matrix nachdem der IPF-Algorithmus auf die 2000-er Daten mit den Randsummen der 2005-er Daten angewendet wurde.

Im letzten Schritt werden nun Häufigkeiten von Aktivitätenketten gesucht, welche zusammengefasst die Aktivitäten-Längen-Verteilung in Tabelle 2 ergeben. Dies wird, wie bereits im Kapitel 4 beschrieben als eine Art Umkehrung des ersten Schritts gemacht, in welchem aus den Häufigkeiten der verschiedenen Aktivitätenketten die Aktivitäten-Längen-Verteilungs-Matrix berechnet wurde. Die Ergebnisse werden in Tabelle 3 dargestellt. Dabei wird auch gleich ein Vergleich mit den Ausgangsdaten und den wahren 2005-er Daten, die ja - wie erwähnt - eigentlich bereits zu Beginn vorhanden waren, möglich.

Dieser Vergleich fällt zwiespältig aus: Zum einen fällt sofort auf, dass die Häufigkeiten der Aktivitätenketten oftmals in der Grössenordnung der wahren Häufigkeiten 2005 liegen. Es gibt

Tabelle 1: Häufigkeiten der Aktivitäten in Ketten der jeweiligen Längen im Mikrozensus 2000 und entsprechende Randsummen aus dem Mikrozensus 2005
(h = home, w = work, e = education, l = leisure, s = shopping)

MZ 2000	e	h	l	s	w	Σ		MZ 2005	
Länge 3	5843	95356	13009	10868	17958	143034	34.6%	35103	33.8%
Länge 4	5899	56588	22060	14380	14249	113176	27.4%	8536	8.2%
Länge 5	3078	34547	14192	8901	17807	78525	19.0%	36395	35.1%
Länge 6	2443	17549	12120	6199	11735	50046	12.1%	9132	8.8%
Länge 7	822	8653	4974	2458	5080	21987	5.3%	12558	12.1%
Länge 8	108	1628	1355	509	1424	5024	1.2%	1128	1.1%
Länge 9	10	885	443	184	296	1818	0.4%	882	0.9%
Länge 10	0	80	74	37	9	200	0.1%	20	0.0%
Σ	18203	215286	68227	43536	68558	413810		103754	
	4.4%	52.0%	16.5%	10.5%	16.6%				
MZ 2005	3988	59878	17443	10470	11975	103754			
	3.8%	57.7%	16.8%	10.1%	11.5%				

Quelle: Mikrozensus 2000 und Mikrozensus 2005

Tabelle 2: Häufigkeiten der Aktivitäten in Ketten der jeweiligen Längen nach der Anpassung mittels IPF-Algorithmus an die Mikrozensus 2005-Daten

	e	h	l	s	w	Σ
Länge 3	1286	25635	3092	2497	2594	35103
Länge 4	409	4789	1650	1040	648	8536
Länge 5	1373	18826	6837	4146	5213	36395
Länge 6	436	3828	2337	1156	1375	9132
Länge 7	455	5857	2976	1422	1847	12558
Länge 8	24	446	328	119	210	1128
Länge 9	5	488	216	87	88	882
Länge 10	0	9	7	3	1	20
Σ	3988	59878	17443	10470	11975	103754

Quelle: eigene Berechnungen auf Datenbasis der Mikrozensus 2000 und 2005

aber andererseits auch Ketten, bei denen die Häufigkeiten sehr unterschiedlich sind. Beim Vergleich mit den Daten aus Häufigkeiten aus dem Mikrozensus 2000 muss beachtet werden, dass in diesem etwa vier mal mehr Aktivitätenketten beobachtet wurden als im Mikrozensus 2000.

Anstatt unser nicht gerade einfaches Verfahren durchzuführen, hätte man auch die Aktivitätenketten-

Tabelle 3: Ausgangsdaten aus Mikrozensus 2000 und Mikrozensus 2005 sowie die mittels IPF Algorithmus und nachträglichem Zurückrechnen erhaltenen Daten (aus Platzgründen werden nur die ersten 10 und die letzten 10 von Total 252 Aktivitätenkettentypen angegeben). Zum besseren Vergleich der Datensätze sind in Klammern die Anteile an der jeweiligen Gesamtzahl der Aktivitätenketten angegeben.

Typ	MZ 2000 \underline{n}^0		2005 $\underline{n}^{\text{wahr}}$		Schätzung 2005 \underline{n}^1	
heh	5843	(5.6278%)	583	(2.3631%)	2074	(8.0551%)
hlh	13009	(12.5076%)	5262	(21.3287%)	3880	(15.0683%)
hsh	10868	(10.4491%)	3139	(12.7234%)	3285	(12.7597%)
hwh	17958	(17.2658%)	2717	(11.0129%)	3382	(13.1344%)
helh	4618	(4.4400%)	105	(0.4256%)	231	(0.8954%)
hesh	780	(0.7499%)	39	(0.1581%)	78	(0.3029%)
hewh	101	(0.0971%)	6	(0.0243%)	0	(0.0000%)
hleh	89	(0.0856%)	11	(0.0446%)	231	(0.8954%)
hlsh	2171	(2.0873%)	422	(1.7105%)	408	(1.5866%)
hlwh	1131	(1.0874%)	61	(0.2473%)	290	(1.1277%)
⋮	⋮	⋮	⋮	⋮	⋮	⋮
hwhlhwhlh	11	(0.0106%)	2	(0.0081%)	9	(0.0339%)
hwhwhshlh	19	(0.0183%)	4	(0.0162%)	10	(0.0378%)
hwhwhwhlh	5	(0.0048%)	6	(0.0243%)	4	(0.0145%)
hwhwhwhwh	10	(0.0096%)	6	(0.0243%)	4	(0.0157%)
hwlwhlhlh	35	(0.0337%)	2	(0.0081%)	2	(0.0096%)
hwlwhshlh	5	(0.0048%)	1	(0.0041%)	0	(0.0000%)
hwlwlshlh	6	(0.0058%)	2	(0.0081%)	0	(0.0000%)
hwlwlwlwh	10	(0.0096%)	1	(0.0041%)	0	(0.0000%)
hslhslhlh	17	(0.0163%)	1	(0.0041%)	2	(0.0068%)
hwlwhwhlsh	3	(0.0029%)	1	(0.0041%)	0	(0.0012%)

Quelle: Mikrozensus 2000, 2005 und eigene Berechnungen

Häufigkeiten aus dem Mikrozensus 2000 einfach skalieren können, so dass deren Summe übereinstimmt mit der Summe der Aktivitätenketten-Häufigkeiten aus dem Mikrozensus 2005. \tilde{n}^0 bezeichne diese skalierten Häufigkeiten. Ein Vergleich der beiden Schätzungen \tilde{n}^0 und \underline{n}^1 (mittels einer Summe von normierten Abständen zu den wahren Werten aus dem Mikrozensus 2005 $\underline{n}^{\text{wahr}}$) zeigt Folgendes:

$$\sum_{k=1}^K \left| \frac{\tilde{n}_k^0 - \underline{n}_k^{\text{wahr}}}{\underline{n}_k^{\text{wahr}}} \right| = 522.5$$

$$\sum_{k=1}^K \left| \frac{\underline{n}_k^1 - \underline{n}_k^{\text{wahr}}}{\underline{n}_k^{\text{wahr}}} \right| = 944.1$$

Tabelle 4: Längen-Aktivitäten-Verteilung in der Schätzung \underline{n}^1

	e	h	l	s	w	Σ
Länge 3	2074	25242	3880	3285	3382	37863
Länge 4	618	4580	1858	1248	856	9160
Länge 5	1373	18830	6839	4147	5216	36405
Länge 6	435	3824	2338	1156	1373	9126
Länge 7	457	5859	2972	1417	1846	12551
Länge 8	23	441	324	118	206	1112
Länge 9	4	493	218	88	88	891
Länge 10	0	8	8	4	0	20
Σ	4987	59277	18437	11463	12967	107128

Quelle: eigene Berechnungen auf Datenbasis der Mikrozensen 2000 und 2005

Das heisst, betrachtet man nur dieses Mass, so wäre die simple Skalierung der Ausgangsdaten eine bessere Annäherung an die Aktivitätsketten-Häufigkeiten als unser kompliziertes Verfahren. Unser Vorgehen hatte aber noch ein weiteres Hauptziel, nämlich die Annäherung an die Längen- und die Aktivitäten-Verteilung, welche mittels Skalierung nicht erreicht werden kann.

Betrachtet man die aus der Schätzung der Aktivitätsketten entstandene Längen-Aktivitäten-Verteilung (siehe Tabelle 4), so fällt sofort auf, dass sie nicht mit der Längen-Aktivitäten-Verteilung nach dem IPF-Algorithmus (Tabelle 2) übereinstimmt. Die Unterschiede haben ihre Ursprünge zum einen darin, dass das Gleichungssystem $M * \underline{n}^1 = \underline{b}^1$ nicht exakt gelöst werden kann und zum andern in Rundungsfehlern in der Bestimmung der Aktivitätsketten-Häufigkeiten. Die Unterschiede sind aber nur von geringem Ausmass, weshalb die Schlussfolgerung gezogen werden kann, dass das in Kapitel 4 beschriebene Vorgehen bezüglich der Längen-Aktivitäten-Verteilung der Aktivitätsketten brauchbare Resultate liefert. Auch die Anpassung an die vorgegebenen Randverteilungen (Längen- und Aktivitäten-Verteilung aus dem Mikrozensus 2005) wird bis auf geringe Abweichungen erreicht.

6 Ausblick

Die in dieser Arbeit präsentierte Beschreibung des Problems zur Anpassung von Aktivitätsketten lässt noch einige Fragen offen. Eine der wichtigsten davon ist, wie die Anpassung von Aktivitätsketten durchzuführen ist, wenn in den Ausgangsdaten nicht dieselben Aktivitätsketten gegeben sind wie man sie in den Schlussdaten sucht. Das Problem wiegt besonders schwer, falls für die gesuchten Daten ganze Gruppen von Aktivitätsketten (zum Beispiel mehrere Ket-

ten der Länge l_{neu} , welche in den Ausgangsdaten nicht vorkommt) hinzugenommen werden. Dann ist zum einen nämlich die Konvergenz des IPF-Algorithmus nicht mehr gewährleistet, da die Ausgangsmatrix viele Nullen enthält, und zum andern wird im Fall von Konvergenz der Anteil der Ketten mit Länge l_{neu} gleich Null sein, da diese Ketten in den Ausgangsdaten nicht vorkommen.

Unterschiedliche Aktivitätenketten in Anfangs- und Enddaten führen auch zu Problemen bei der Optimierung falls ein Teilgleichungssystem aus (9) mehrere Lösungen hat. Es kann dann in (A) nicht einfach an die Ausgangsdaten angepasst werden und es muss ein anderer Weg gefunden werden, eine optimale aus vielen richtigen Lösungen auszuwählen.

Des weiteren kann auch darüber diskutiert werden ob der minimale euklidische Abstand ein gutes Kriterium für gute Anpassung ist, oder ob sich noch bessere finden lassen.

Im Bezug auf die Mikrozensus 2000 und Mikrozensus 2005-Daten stellt sich die Frage, ob das gleiche Vorgehen, wobei aber nur Aktivitäten *home*, *work* und *education* (hwe) betrachtet werden, zu besseren Ergebnissen führen würde. Diese drei Aktivitäten stellen eine besonders gute Datenbasis dar, da sie einfach zu erkennen, in den meisten Tagesabläufen fest eingebunden und in Umfragen deshalb mit grösserer Genauigkeit festgehalten sind. Andererseits werden beispielsweise leisure-Aktivitäten oftmals nicht erkannt oder dieselbe Aktivität von verschiedenen Personen einmal als Aktivität gewertet und ein anderes mal nicht. Man würde also erwarten dass sich bei einer Beschränkung auf hwe-Ketten die Genauigkeit der Vorhersage, aufgrund grösserer Genauigkeit der Ausgangsdaten, verbessert.

Literatur

ARE und BfS (2001) Mobilität in der Schweiz, Ergebnisse des Mikrozensus 2000 zum Verkehrsverhalten, *Schlussbericht*, Bundesamt für Raumentwicklung, Bundesamt für Statistik, Bern und Neuenburg, <http://www.portal-stat.admin.ch/mz05/>.

ARE und BfS (2007) Mobilität in der Schweiz, Ergebnisse des Mikrozensus 2005 zum Verkehrsverhalten, *Schlussbericht*, Bundesamt für Raumentwicklung, Bundesamt für Statistik, Bern und Neuchâtel, <http://www.portal-stat.admin.ch/mz05/>.

Bregman, L. (1967) The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming, *U.S.S.R. Computational Math. Mathematical Phys.*, **7** (3) 200–217.

Evans, A. W. (1970) Some properties of trip distribution methods, *Transportation Research*, **4** (1) 19–36.

Fienberg, S. E. (1970) An iterative procedure for estimation in contingency tables, *Annals of Mathematical Statistics*, **41** (3) 907–917.

Kirby, H. R. und S. P. Evans (1974) A three-dimensional furness procedure for calibrating gravity models, *Transportation Research*, **8** (2) 105–122.