

# 統計学

第 107 号

---

## 論文

国勢調査マイクロデータを用いたスワッピングの有効性の検証  
..... 伊藤 伸介・星野なおみ (1)

標本交代方式を採る統計調査の標本バイアス  
..... 山口 幸三 (17)

## 書評

吉田 忠著『近代オランダの確率論と統計学』（八朔社，2014年）  
..... 上藤 一郎 (33)

泉 弘志著『投下労働量計算と基本経済指標：新しい経済統計学の探求』  
（大月書店，2014年）  
..... 橋本 貴彦 (38)

## 海外統計事情

ロシア統計学会について  
..... イリーナ エリセーエワ・山口 秋義 (43)

## 本会記事

経済統計学会第58回（2014年度）全国研究大会 ..... (46)

---

2014年9月

経済統計学会

## 創刊のこ と ば

社会科学の研究と社会的実践における統計の役割が大きくなるにしたがって、統計にかんする問題は一段と複雑になってきた。ところが統計学の現状は、その解決にかならずしも十分であるとはいえない。われわれは統計理論を社会科学の基礎のうえにおくことによって、この課題にこたえることができると考える。このためには、われわれの研究に社会諸科学の成果をとりいれ、さらに統計の実際と密接に結びつけることが必要であろう。

このような考えから、われわれは、一昨年来経済統計研究会をつくり、共同研究を進めてきた。そしてこれを一層発展させるために本誌を発刊する。

本誌は、会員の研究成果とともに、研究に必要な内外統計関係の資料を収めるが同時に会員の討論と研究の場である。われわれは、統計関係者および広く社会科学研究者の理解と協力をえて、本誌をさらによりよいものとするを望むものである。

1955年4月

## 経 済 統 計 研 究 会

## 経 済 統 計 学 会 会 則

第1条 本会は経済統計学会（JSES : Japan Society of Economic Statistics）という。

第2条 本会の目的は次のとおりである。

1. 社会科学に基礎をおいた統計理論の研究
2. 統計の批判的研究
3. すべての国々の統計学界との交流
4. 共同研究体制の確立

第3条 本会は第2条に掲げる目的を達成するために次の事業を行う。

1. 研究会の開催
2. 機関誌『統計学』の発刊
3. 講習会の開催、講師の派遣、パンフレットの発行等、統計知識の普及に関する事業
4. 学会賞の授与
5. その他本会の目的を達成するために必要な事業

第4条 本会は第2条に掲げる目的に賛成した以下の会員をもって構成する。

- (1) 正会員
- (2) 院生会員
- (3) 団体会員
- 2 入会に際しては正会員2名の紹介を必要とし、理事会の承認を得なければならない。
- 3 会員は別に定める会費を納入しなければならない。

第5条 本会の会員は機関誌『統計学』等の配布を受け、本会が開催する研究大会等の学術会合に参加することができる。

- 2 前項にかかわらず、別に定める会員資格停止者については、それを適用しない。

第6条 本会に、理事若干名をおく。

- 2 理事から組織される理事会は、本会の運営にかかわる事項を審議・決定する。
- 3 全国会計を担当する全国会計担当理事1名をおく。
- 4 渉外を担当する渉外担当理事1名をおく。

第7条 本会に、本会を代表する会長1名をおく。

- 2 本会に、常任理事若干名をおく。
- 3 本会に、常任理事を代表する常任理事長を1名おく。
- 4 本会に、全国会計監査1名をおく。

第8条 本会に次の委員会をおく。各委員会に関する規程は別に定める。

1. 編集委員会
2. 全国プログラム委員会
3. 学会賞選考委員会
4. ホームページ管理運営委員会
5. 選挙管理委員会

第9条 本会は毎年研究大会および会員総会を開く。

第10条 本会の運営にかかわる重要事項の決定は、会員総会の承認を得なければならない。

第11条 本会の会計年度の起算日は、毎年4月1日とする。

- 2 機関誌の発行等に関する全国会計については、理事会が、全国会計監査の監査を受けて会員総会に報告し、その承認を受ける。

第12条 本会会則の改正、変更および財産の処分は、理事会の審議を経て会員総会の承認を受けなければならない。

付 則 1. 本会は、北海道、東北、関東、関西、九州に支部をおく。

2. 本会に研究部会を設置することができる。
3. 本会の事務所を東京都町田市相原4342法政大学日本統計研究所におく。

1953年10月9日（2010年9月16日一部改正[最新]）

# 国勢調査マイクロデータを用いた スワッピングの有効性の検証

伊藤伸介\*・星野なおみ\*\*

## 要旨

わが国ではこれまで、攪乱的手法を含む匿名化技法に関する実証的な研究が、諸外国と比較して非常に少なかった。そのため、マイクロデータに対する攪乱的手法の適用可能性を追究することによって、匿名データの作成において実用的な匿名化技法の範囲が拡大することが期待される。そこで、本稿では、攪乱的手法の1つであるスワッピングの適用可能性について検討を行うだけでなく、スワッピング済データにおける有用性と秘匿性の定量的な評価を行った。本分析結果によれば、ターゲット・スワッピングにおける秘匿性は、ランダム・スワッピングにおけるそれよりも全般的に高くなっている。このことは、有用性がある水準に設定された場合、ターゲット・スワッピングのほうが少ないスワッピング率でより高い秘匿性を確保することが可能なことを意味している。このように秘匿の観点から見ると、本分析の結果においては、ランダム・スワッピングよりもターゲット・スワッピングのほうがより有効な手法であると言える。

## キーワード

国勢調査, マイクロデータ, 匿名化技法, スワッピング

### 1. はじめに

諸外国では、様々な政府統計マイクロデータが提供されており、それによって主として社会経済の分野におけるマイクロレベルの実証研究に大きく寄与してきた。マイクロデータには個体情報が含まれていることから、マイクロデータの提供において個々人が特定化されるリスクを低減するためには、マイクロデータに対して法制度的あるいは技術的な匿名化措置を施すことが求められる。前者の法制度的な匿名化措置については、例えばアメリカセン

サス局の開示評価委員会（Disclosure Review Board）において匿名化措置に関するチェックリスト等を用いて政府統計マイクロデータの提供可能性を検討していることを指摘することができる。他方、後者の匿名化の技術的な手法は、原数値における区分を変更する等の加工を行う非攪乱的な（non-perturbative）手法と原数値にノイズを追加する等の加工を施す攪乱的な（perturbative）手法に類別される。非攪乱的な手法については、リコーディング（区分統合）、データの削除（レコード削除あるいは変数の削除）、トップ（ボトム）・コーディング（分布の上位あるいは下位における区分統合）が存在する。一方、攪乱的な手法については、ノイズの付加（加法ノイズ、乗法ノイズ）、スワッピング（レコード間の

\* 中央大学経済学部  
（独）統計センター非常勤研究員  
e-mail : ssitoh@tamacc.chuo-u.ac.jp

\*\*（独）統計センター  
e-mail : nsaitou2@nstat.go.jp

入れ替え), ラウンディング (丸め), ミクロアグリゲーション (変数値を層内の平均値等の代表値に置き換えること) 等の手法がある (Domingo-Ferrer and Torra, 2001a ; Willenborg and de Waal, 2001)<sup>1)</sup>。

諸外国では, 個票データに対して秘匿処理を施したマイクロデータ (以下「匿名化マイクロデータ」と呼称) を作成する上では, リコーディング, トップ (ボトム) ・コーディング等の非攪乱的な手法が用いられることが少なくない。その一方で, 匿名化マイクロデータの作成において攪乱的な手法が適用される場合もある。例えば, アメリカセンサス局は, 2000年のアメリカ人口センサスの一般公開用マイクロデータ (Public Use Microdata Sample ; PUMS) において, 加法ノイズやラウンディングを採用している (Zayatz, 2007)。また, イギリスでも, 2001年人口センサスの匿名化標本データ (Samples of Anonymised Records) において, PRAM (Post RAndomisation Method) が用いられている (De Kort and Wathan, 2009)。

ところで, 諸外国では, ミクロデータに含まれる個体情報の露見リスク (disclosure risk) の低減 (露見制御, disclosure control) を図るために, ミクロデータおよび集計表の作成においてスワッピングを適用していることが知られている。アメリカセンサス局は, 1990年人口センサス以降, 集計表における秘匿処理として, 人口センサスの個票データにスワッピングを適用している (Federal Committee on Statistical Methodology, 1994 ; Gbur and Zelenak, 2004)。このスワッピングされた個票データに基づいて, PUMSおよび集計表が作成されている (Zayatz, 2007)。なお, イギリスにおいても, 人口センサスの個票データの作成において, レコードスワッピングが適用されている (Shlomo, 2007)。スワッピングの適用対象となるレコードは, 他のレコードと入れ替えられることから, 特定化の

リスクを回避することができることが主な理由だと考えられる。

一方, わが国における攪乱的手法に関する実証的な研究については, Takemura (2002) による人口動態調査死亡票の個票データを用いたスワッピングの研究, 伊藤他 (2008, 2009, 2010) による全国消費実態調査の個票データを用いたマイクロアグリゲーションの適用可能性に関する実証研究, さらには伊藤・村田 (2013) による家計調査の個票データを用いたマイクロアグリゲーションや加法ノイズの有効性の研究等があるが, 諸外国と比べると実証研究に関する蓄積は非常に少ないと思われる。マイクロデータに対する攪乱的手法の適用可能性を検証することによって, 匿名データの作成において実用的な匿名化技法の範囲が拡大することが期待されることから, わが国でも攪乱的手法についてはさらなる実証的な研究の必要性は高いと思われる。

現在, わが国では平成12年と17年の国勢調査の匿名データが提供されているが, 攪乱的手法としてスワッピングが初めて適用されている。将来的には, 小地域分析用の匿名データ等, 別のタイプの国勢調査の匿名データの要望が出てくる可能性があり, その予備的な研究として, 攪乱的手法の中でもスワッピングについてその方法的な可能性をさらに追究することは有用であると考えられる。

そこで, 本稿では, 匿名化技法としてのスワッピングに焦点を当て, わが国の政府統計マイクロデータに対するスワッピングの有効性について検討を試みる。本稿では, 最初に露見リスクの基本的な考え方とスワッピングの特徴を述べる。つぎに, スワッピングの有効性を評価するために, 匿名化技法を適用した場合の有用性 (data utility) と秘匿性 (data confidentiality) の定量的な評価方法および有用性と秘匿性の相対比較の方法についてのサーベイを行う。これらの議論を踏まえて, 本研究では, 政府統計マイクロデータを用いて

スワッピングの実験を行う。具体的には、スワッピングの対象となるレコードを探索した上で、該当するレコードに対してスワッピングを試行的に適用するだけでなく、スワッピングが施されたデータ（以下、「スワッピング済データ」と呼称）について有用性と秘匿性の定量的な評価を行うことによって、スワッピングの有効性の検証を試みる。

## 2. 露見リスクとスワッピング

露見リスクを議論する場合、主として、個体識別漏洩 (identification disclosure) に伴うリスクと予測漏洩 (prediction disclosure) によって発生するリスクに大別することができる (Duncan and Lambert, 1989; Skinner, 1992)。個体識別漏洩とは、マイクロデータに含まれるレコードからある個体が特定化されることによって、個体に関するセンシティブな情報が露見されることである。それに対して、予測漏洩とは、マイクロデータに含まれる個体が特定されなくても、その個体のセンシティブな属性に関しては狭い範囲で予測することが可能になることである (Skinner, 1992: p.23)。

以下では、個体識別を例に、露見リスクを議論することにしたい。マイクロデータの入手者 (侵入者, intruder) が、特定の個体に関する識別情報を含むファイル (識別ファイル) を持っていることを想定する。マイクロデータの入手者によって、①識別ファイルに含まれるレコードとマイクロデータ中のレコードとの間で、キー変数 (key variable) による1対1のマッチングが行われ、②そのマッチングされたレコードが特定の個体のものであることが突き止められた場合、個体識別が成立する (Müller *et al.*, 1995)。

もし、マイクロデータの入手者が、識別ファイルに相当する母集団に関する外部情報を持っていた場合、個体を特定するために外部情報とマイクロデータのマッチングを行うこと

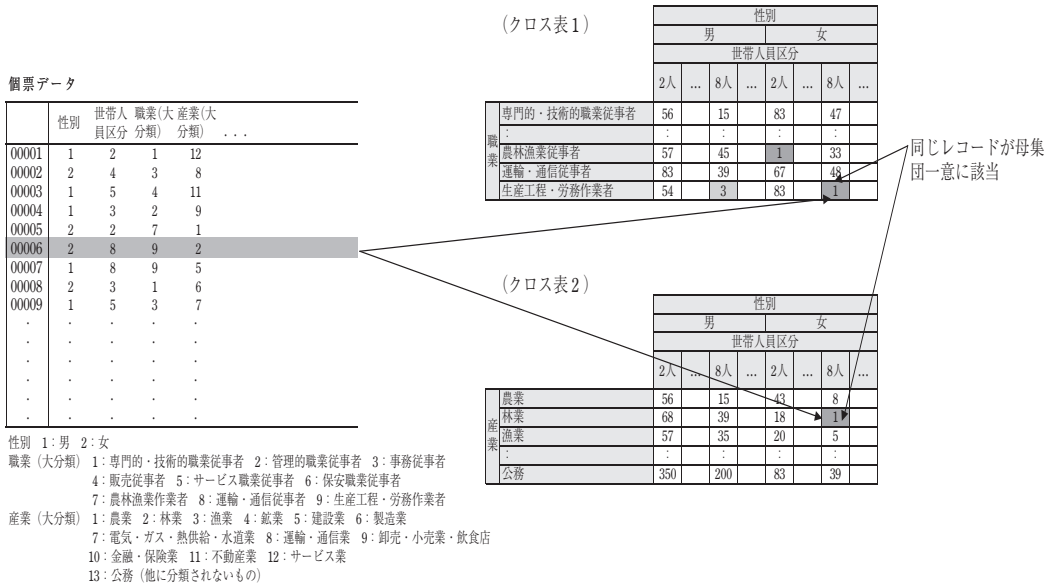
が考えられる。マイクロデータが母集団についてのレコードから構成され、母集団において属性の組み合わせがただ1つしか存在しない母集団一意 (population unique) に該当するレコードが含まれるのであれば、外部情報とのマッチングにおいて個体が特定化されるリスクが高まる。

それに対して、マイクロデータが標本に関するレコードから構成される場合、回答者の属性の組み合わせによって、一意となるレコードが存在したとしても、それは、標本一意 (sample unique: SU) であって、母集団一意とは異なる。一方で、標本一意の中で母集団一意に該当するレコードは、匿名化技法の適用対象となることが考えられる。

他方、「疫学的に特異であるために、本質的に (intrinsically) まれな属性群の組み合わせを有する」レコードは、「特殊な一意 (special uniques)」とみなされ (Elliot, 2001)、標本一意の中で母集団一意に該当するレコードの中でも個人が特定化される可能性が特に高くなる。特殊な一意とは、Elliot and Manning (2004) によれば、「K個のキー変数の集合において標本一意であるだけでなく、Kの部分集合であるk個 (のキー変数の集合) においても標本一意となること」であって、「少数のキー変数の組み合わせでも標本一意になるレコード」が特殊な一意に該当するとみなされる (Gross *et al.*, 2004)。図1は、特殊な一意の概略図を示したものである。議論を簡単にするために、図1に示される個票データは母集団を表すレコードを含んでいると仮定し、例えば一連番号00006のレコードは、性別については女、世帯人員区分に関しては8人世帯、職業 (大分類) については生産工程・労務作業、さらに産業 (大分類) に関しては林業の属性値を有しているとする。また、クロス表1は、性別、世帯人員区分と職業 (大分類) のクロス表であり、クロス表2は、性別、世帯人員区分と産業 (大分類) のクロス



図1 本研究における「特殊な一意」のイメージ



注 本図では、性別、世帯人員区分と職業(大分類)と産業(大分類)をキー変数と仮定している。

表を示している。クロス表1において、母集団一意のセルが2つ存在するが、その中で一連番号00006のレコードは、性別が女、世帯人員区分が8人世帯、職業が生産工程・労務作業従事者であるセルに該当するだけでなく、クロス表2においても、性別が女、世帯人員区分が8人世帯で産業が林業であるセルに該当しているとする。図1において、性別、世帯人員区分と職業(大分類)と産業(大分類)をキー変数と仮定すると、一連番号00006のレコードは、性別、世帯人員区分と職業(大分類)というキー変数の組み合わせと性別、世帯人員区分と産業(大分類)という2つの組合せにおいて母集団一意であるということが出来る。さらに、いずれの場合も少数のキー変数の組み合わせであることから、一連番号00006のレコードは、母集団一意に該当するレコードの中でも、リスクが相対的に高いレコードということができ、特殊な一意となるレコードの可能性が高いことがわかる。

このように低次元のクロス表をもとに、特殊な一意に該当すると思われるレコードを探

索することが求められる。

こうした特殊な一意に該当するレコードに対して適用される匿名化技法が、スワッピングである。スワッピング(data swapping)とは、「マイクロデータに含まれるレコード同士で属性値を入れ替える」ことである(Willenborg and Waal, 2001 : p.126)。スワッピングの概略図については、図2で示している。図2では、個票データに対して地域が異なるレコード同士でスワッピングが行われている。具体的には、地域が「三大都市圏」、性別が「女」、年齢が「35~44歳」、雇用形態が「正規の職員・従業員」、週間就業時間が「35~48時間」となっているレコードを、地域が「三大都市圏以外」であるレコードに入れ替える。スワッピングのために使用するキー変数は、性別、年齢と雇用形態とする。図2を見ると、地域が「三大都市圏以外」で性別等のキー変数の値が同じレコードに入れ替えることによって、スワッピング済データにおいて作成された性別、年齢、雇用形態別のクロス表は、スワッピング前の個票データにおけるク

図2 スワッピングのイメージ

個票データ							スワッピング済データ						
番号	地域	性別	年齢	雇用形態	週間就業時間		番号	地域	性別	年齢	雇用形態	週間就業時間	
1	1	1	2	2	1		1	1	1	2	2	1	
2	1	2	4	1	2	入れ替え	2	2	2	4	1	3	
3	1	1	3	1	4		3	1	1	3	1	4	
4	1	1	5	3	1		4	1	1	5	3	1	
5	1	1	6	2	3		5	1	1	6	2	3	
6	1	1	4	3	2		6	1	1	4	3	2	
7	2	2	4	1	3		7	1	2	4	1	2	
8	2	1	5	1	4		8	2	1	5	1	4	
9	2	2	2	2	3		9	2	2	2	2	3	

性別 1:男 2:女

地域 1:三大都市圏 2:三大都市圏以外

年齢 1:15歳未満 2:15~24歳 3:25~34歳 4:35~44歳 5:45~54歳 6:55~64歳 7:65歳以上

雇用形態 1:正規の職員・従業員 2:パート・アルバイト 3:派遣・契約社員

週間就業時間 1:35時間未満 2:35~48時間 3:49~59時間 4:60時間以上

ロス表の数値と変わらないことが確認できる。

スワッピングは、特殊な一意のような露見リスクの高いレコードを対象に、特定のスワッピング率において適用されるが、スワッピングの適用によって基本的な属性間関係性が変わらないことが求められる。また、スワッピングは、特定化のリスクが特に高いと思われるレコードにターゲットを絞ってスワッピングを行うターゲット・スワッピング (targeted data swapping) と、無作為にスワッピングの対象となるレコードを選別した上でスワッピングを行うランダム・スワッピング (random data swapping) に大別される (Shlomo *et al.*, 2010)。

スワッピングの実際の適用においては、小地域レベルで特定の人口社会的属性群に基づいて一意性を有する世帯のレコードを対象に、別の地域における他の世帯との入れ替えが行われている。2000年アメリカ人口センサスの場合、スワッピングは、short formとlong formの2種類の調査票情報に適用され、特殊な一意の対象となるレコードを探索した上で、異なる地域に居住する世帯の組に対して、地域間におけるスワッピングが適用されてい

る。また、スワッピングの対象となる世帯の組については、最低限の人口社会的な属性に基づいた対応付けが行われている (Zayatz, 2007)。

### 3. ミクロデータにおける有用性と秘匿性の評価について

ミクロデータに対する匿名化技法の適用可能性を検証するために、匿名化ミクロデータにおける情報量損失 (information loss) の程度を表す有用性と秘匿処理に伴う個体情報の露見リスクの程度を表す秘匿性の定量的な評価に関する研究が行われてきた (Domigo-Ferrer and Torra, 2001b; Karr *et al.*, 2006; Shlomo, 2010)。したがって、スワッピング済データにおいても、こうした有用性と秘匿性の評価方法が適用されてきた (Shlomo *et al.*, 2010)。

#### 3-1 有用性の評価方法

ミクロデータの有用性の定量的な評価方法については、以下のような方法を指摘することができる (伊藤・村田, 2013)。第1は、平均や分散等の基本統計量、絶対距離の平均

値 (average absolute distance) 等を用いたクロス集計表における度数の比較 (Domigo-Ferrer and Torra, 2001b), クラメールのVといった関連性の指標 (Shlomo, 2010) 等を用いて, 個票データと匿名化マイクロデータの近似性の比較を行うことである。第2は, 個票データに対する匿名化マイクロデータの情報量損失を計測することである。具体的には, 量的属性に関しては, 属性値, 相関係数行列や分散共分散行列等を用いて, 平均平方誤差 (mean square error), 平均絶対誤差 (mean absolute error), および平均変化率 (mean variation) に基づく情報量損失を計測することが提案されている (Domigo-Ferrer and Torra, 2001a)。また, 質的属性については, エントロピーをもとに, 情報量損失を計測する方法が議論されている (De Waal and Wiltenborg, 1999)。なお, 有用性の評価方法については, 回帰分析における決定係数の比較や回帰係数の信頼区間に基づいた評価方法 (Karr *et al.*, 2006), さらに, 傾向スコア, クラスタ分析, 経験分布関数等を用いて有用性を定量的に評価する方法も提唱されている (Woo *et al.*, 2009 : pp.113-115)。

### 3-2 秘匿性の評価方法

秘匿性の定量的な評価方法は, ファイルレベルのリスク評価法 (file-level risk metrics) とレコードレベルのリスク評価法 (record-level risk metrics) に類別することが可能である (Elliot, 2001 : pp.80-84)。前者のファイルレベルのリスク評価法については, シナリオに基づいてキー変数を設定した上で, 母集団一意を計測することが指摘できる (Gross *et al.*, 2004)。母集団一意の評価指標に関しては, 母集団全体に占める母集団一意数の比率である母集団の一意性 (population uniqueness) や, 母集団一意かつ標本一意としての共通一意 (union uniques : UU) となるレコード数の標本一意 (SU) となるレコー

ド数に対する比率であるUUSU比率 (UUSU ratio) は, 母集団一意に関する主要な指標と考えることができる (Elliot, 2001)。

後者のレコードレベルのリスク評価法に関しては, 低次元のクロス表をもとに, 特殊な一意に該当すると思われるレコードを探索する特殊な一意の分析 (Special Uniques Analysis) (Elliot *et al.*, 2002) がある。このような特殊な一意のレコードが匿名化マイクロデータにおいてどの程度減少したのかを計測することも, 秘匿性の評価指標の1つと考えられる。さらに, 個票データと匿名化マイクロデータとのレコードリンケージ (record linkage) による評価研究 (Duncan *et al.*, 2011) もレコードレベルのリスク評価法の1つと思われる。これについては, わが国においても, 全国消費実態調査や家計調査のマイクロデータを用いて, レコードリンケージに基づく秘匿性の評価を行った研究がある (伊藤他, 2009 ; 伊藤他, 2010 ; 伊藤・村田, 2013)。

### 3-3 有用性と秘匿性の比較分析の方法

近年では, 各種の匿名化マイクロデータにおける有用性と秘匿性の比較・検証が行われている。有用性と秘匿性の比較分析を行うための主な方法としては, ①総合指標による評価, ②R-Uマップ (R-U confidentiality map ; Rはrisk, Uはutilityの略) の作成がある。

前者の総合指標による評価については, Domingo-Ferrer等が, 情報量損失とリスクに関するスコアをもとに総合指標を作成した上で, 有用性と秘匿性に関する相対評価を行っている (Domigo-Ferrer and Torra, 2001b)。具体的には, 様々な匿名化マイクロデータを対象に, 相関係数行列の平均平方誤差等を用いて情報量損失のスコアを計測するだけでなく, レコードリンケージに基づいてリスクに関するスコアの計算を行っている。スコアに基づいて有用性と秘匿性に関する定量的な総合指標を作成していることから, 匿名化技法の有



効性について定量的に評価することが容易であるが、スコアの計算方法や総合指標の算定式の設定によって、評価結果が変わることも考えられる。

後者のR-Uマップに関しては、Duncan等が、有用性と秘匿性について相対比較を行うために、R-Uマップの作成を提唱している（Duncan *et al.*, 2001）。R-Uマップによって、有用性と秘匿性がトレードオフの関係にあることが視覚的に把握できることから、R-Uマップでの位置を確認した上で、R-Uマップ上で有用性と秘匿性の相対的な変化の程度を明示することによって、各種の匿名化技法を比較・検討することが可能である。その一方で、R-Uマップにおいて有用性と秘匿性に関する許容可能な水準（閾値）を設定しない場合、R-Uマップ上で、有用性と秘匿性の両面から最適な匿名化技法を選ぶのは困難である。わが国では、全国消費実態調査や家計調査の個票データを例に、R-Uマップの試行的な作成が行われている（伊藤他，2010；伊藤・村田，2013）。

#### 4. 国勢調査のマイクロデータに対するスワッピングの方法

本節では、わが国の国勢調査のマイクロデータを用いて行ったスワッピングに関する研究の概要を述べる。本研究で国勢調査を使用する理由は、本研究の成果が、将来国勢調査の小地域分析用マイクロデータの作成を検討する上で基礎資料として寄与しうると考えたからである。なお、本研究では、平成17年国勢調査の個票データにおける特定の地域（以下「地域A」と呼称）の記録をもとに個人単位で抽出した約100,000レコードを使用する。

本研究では、(1)スワッピングの対象となるレコードを探索するために、スワッピングの対象レコードの中で相対的にリスクの高いレコードをスコアに基づいて選び出し、(2)リスクの高いレコードに対してスワッピングを適

用する。

スワッピングの対象となるレコードの探索にあたっては、最初にキー変数を用いて、母集団一意の計測を行った。母集団一意に該当するレコードは、露見リスクの可能性があると考えられるために、スワッピングの適用対象となりうるからである。本研究で使用するキー変数については、外観識別性等を考慮した結果、つぎの11個の変数が選ばれた。

- ・世帯主との続き柄（13区分）
- ・男女の別（2区分）
- ・年齢5歳階級（25区分）
- ・配偶関係（5区分）
- ・国籍（13区分）
- ・労働力状態（9区分）
- ・従業上の地位（8区分）
- ・産業大分類（19区分）
- ・職業大分類（10区分）
- ・住居の種類（9区分）
- ・住居の建て方（4区分）+建物の階数（30区分）(建物の階数については共同住宅のみ)

この11変数をキー変数として母集団一意を計測した結果、母集団一意に該当するレコードは32,064レコードとなった。これらのレコードがスワッピングの対象となるレコードとして設定される。

つぎに、本研究は、スワッピングの対象レコードの中で相対的にリスクの高いレコードを選び出すために、母集団一意の対象レコードについて、キー変数のすべての組み合わせでクロス集計を行い、ある特定のレコードが母集団一意に該当した回数をレコードごとに計測し、その計測結果をもとにスコアを算定した。例えば、10個のクロス表で母集団一意に該当するのであれば、10点のスコアが算出される。このようなスコアの算出を行う理由は、スコアが高いレコードについては、相対的にリスクがより大きなレコードとすることができ、特殊な一意に該当するレコードの可能性が高くなると考えられるからである。

本研究において、キー変数11変数のすべての組み合わせ（全部で2,047通り）についてスコアを計算した結果、スコアの最大値は1,518、最小値は2となった。また、スコアの平均値と中央値はそれぞれ、260と192となっている。

最後に、スワッピングの対象レコードを選んだ上で、スワッピングが実行される。本研究では、地域Aのレコードから住居の建て方が空欄であるレコードを削除した上でスワッピングを適用する。また、本研究においては、(1)ターゲット・スワッピングと(2)ランダム・スワッピングの2種類のスワッピングを行う。ターゲット・スワッピングの場合、スコアの高い上位 $p\%$  ( $p=1, 2, 3, 4, 5, 8, 10, 15, 20$ )に該当するレコードをスワッピングの対象レコードとした。一方、ランダム・スワッピングについては、母集団一意に該当するレコードの中から、 $p\%$ にしたがってランダムに選んだレコードをスワッピング対象レコードとした。なお、本実験では、対象レコードに対して入れ替えの候補となるレコードについては、地域Aとは異なる地域（以下「地域B」と呼称）から作成したドナーファイル（約50,000レコード）から探索する。

ところで、スワッピングの対象となるレコードは、特殊な一意として出現する可能性が高いことから、スワッピングの対象レコードとキー変数の値が完全に一致するレコードがドナーファイルで見つかる可能性は低いと考えられる。そこで、本実験では、スワッピングの対象レコードに対して、ドナーファイルに含まれるレコードとの距離を計測し、ドナーファイルの中で最も距離が小さいレコードとスワッピングを行った。具体的には、以下の手順に従っている。

最初に、 $i$  ( $i=1, \dots, m$ ) および  $j$  ( $j=1, \dots, n$ ) を、それぞれスワッピング対象レコードの番号およびドナーファイルのレコード番号とする ( $m$ と $n$ は、それぞれスワッピング対象レ

コードの数およびドナーファイルのレコード数)。また、 $k$  ( $j=1, \dots, 11$ ) をキー変数の番号とする。このとき、 $i$  番目のレコードにおけるキー変数  $k$  の分類区分の数値を  $Cs_{ki}$ 、また、 $j$  番目のドナーファイルのレコードにおけるキー変数  $k$  の分類区分の数値を  $Cd_{kj}$  とすれば、キー変数  $k$  に関する  $i$  と  $j$  の質的属性値間の距離 (distance for categorical variables)  $Sd_{kij}$  は次の(1)式のように定義できる (Domingo-Ferrer and Torra, 2001a : pp.105-106)。

$$Sd_{kij} = |Cs_{ki} - Cd_{kj}| \quad (1)$$

なお、年齢および住居の建て方の「共同住宅」以外の場合、 $|Cs_{ki} - Cd_{kj}| > 0$  であれば、 $Sd_{kij} = 1$  とする。

次に、質的属性値間の距離をスコア化するために、 $k$  番目のキー変数における分類区分数  $C_k$  で  $Sd_{kij}$  を除することによって、 $k$  番目のキー変数におけるスコアである  $Score_{kij}$  が(2)式によって算出される。すなわち、

$$Score_{kij} = \frac{1}{C_k} \cdot Sd_{kij} \quad (2)$$

さらに、各キー変数のスコアを合計することで、 $i$  番目と  $j$  番目のレコード間の距離について、全てのキー変数を総合した指標  $D_{ij}$  が(3)式によって計算される。

$$D_{ij} = \sum_k Score_{kij} \quad (3)$$

最後に、スワッピングの対象レコードとドナーファイルとの間の距離計測型リンケージを行い (Domingo-Ferrer and Torra, 2001a : Takemura, 1999)、ドナーファイルの中でこの距離が最も小さいレコードを、スワッピング対象レコードと置き換える<sup>2)</sup>。

## 5. スワッピングにおける有用性と秘匿性の評価

本研究では、スワッピング済データにおいて有用性と秘匿性の評価に関する定量的な評価を行った。第1に、有用性の評価について

は、Shlomo *et al.* (2010) に基づいて、絶対距離の平均値を用いて評価を行う<sup>3)</sup>。具体的には、絶対距離の平均値による有用性の評価指標DU (data utility) に関しては、個票データとスワッピング済データの両方についてクロス表を作成した上で、個票データを用いて作成したクロス表におけるセルの度数  $T^O(c)$  とスワッピング済データを用いて作成したクロス表におけるセルの度数  $T^S(c)$  の差の絶対値を集計表におけるセルの数  $n_T$  で除することによって求められる。すなわち、

$$DU = \frac{\sum_c |T^S(c) - T^O(c)|}{n_T} \quad (4)$$

他方、本研究では、秘匿性の評価指標DR (disclosure risk) として、個票データにおけるクロス表の中で度数1であるセルの数  $\sum_c I(T^O(c)=1)$  に対するスワッピング済データにおけるクロス表の中で度数1であるセルの数  $\sum_c I(T^O(c)=1, T^S(c)=1)$  の比率が用いられた。

$$DR = \frac{\sum_c I(T^O(c)=1, T^S(c)=1)}{\sum_c I(T^O(c)=1)} \quad (5)$$

この秘匿性の評価指標DRによって、スワッピングを行った場合に、個票データにおいて度数1だったセルのどの程度が度数0あるいは度数2以上に置き換えられたかがわかることから、スワッピングの効果を定量的に評価することが可能になっている<sup>4)</sup>。

先述のように、スワッピングは、特殊な一意となる可能性の高いレコードを対象に適用されることから、低次元のクロス表においてその効果を計測することが望ましい。したがって、本研究では、キー変数の中から3変数を選んだ場合のすべての組み合わせについてクロス表を作成した上で、有用性の評価を試みた<sup>5)</sup>。表1-1は、一例として、①年齢 (5歳階級) × 性別 × 国籍、②年齢 (5歳階級) × 世帯主との続き柄 × 労働力状態における有用性の評価指標DUの結果を示したものである。また、③キー変数における3変数のすべての

表1-1 有用性の評価指標に関する試算結果

スワッピング率とスワッピングの種類	年齢×性別×国籍	年齢×世帯主の続き柄×労働力状態	3変数のすべての組み合わせに関する平均値
ターゲット・スワッピング			
1%	0.9785	0.2790	0.7830
2%	1.5569	0.4855	1.3234
3%	2.0492	0.6475	1.7503
4%	2.3754	0.8253	2.1656
5%	2.6769	0.9668	2.5370
8%	3.3692	1.3354	3.6276
10%	3.7108	1.5385	4.2739
15%	4.5108	1.9385	5.8221
20%	5.1938	2.5347	7.9918
ランダム・スワッピング			
1%	0.2554	0.1149	0.2582
2%	0.3815	0.2072	0.4502
3%	0.4738	0.2735	0.6104
4%	0.5908	0.3344	0.7833
5%	0.7569	0.3870	0.9610
8%	1.1662	0.5983	1.5289
10%	1.4738	0.7268	1.9086
15%	2.2185	1.0393	2.9229
20%	3.3200	1.5856	4.8096

組み合わせにおける有用性の平均値についても示している。年齢、性別と国籍のクロス表については、年齢、世帯主との続き柄と労働力状態におけるクロス表と比較して、情報量損失が大きいことがわかる。その要因として、国籍については日本人以外の分類区分に該当するレコードは相対的に少なく、クロス表において度数が0になるセルが数多く存在するため、スワッピング率を上げた場合、情報量損失がより大きくなることが考えられる。その一方で、表1-1のいずれの結果でも、スワッピング率を上げるにつれて、有用性の程度が低くなることが確認される。また、ランダム・スワッピングのほうが、ターゲット・スワッピングと比較して、全般的に有用性が高いことがわかる<sup>6)</sup>。

一方、表1-2では、上記の①～③の3つのクロス表における秘匿性の評価指標DRの結果の一部も示されている。表1-2を見ると、年齢、性別と国籍のクロス表については、年齢、世帯主との続き柄と労働力状態における

クロス表と比較して、スワッピングを行った場合の秘匿性の程度がより大きくなっていることが確認できる。有用性の検証結果と同様、国籍における分布特性が秘匿性の評価結果に影響を及ぼしていることが推察される。また、スワッピング率を上げるにつれて、秘匿性の評価指標の数値が相対的に小さくなっていることから、秘匿性の程度が高くなることが確認される。また、ターゲット・スワッピングのほうが、ランダム・スワッピングと比較して、全般的に秘匿性が高くなっていることがわかる。

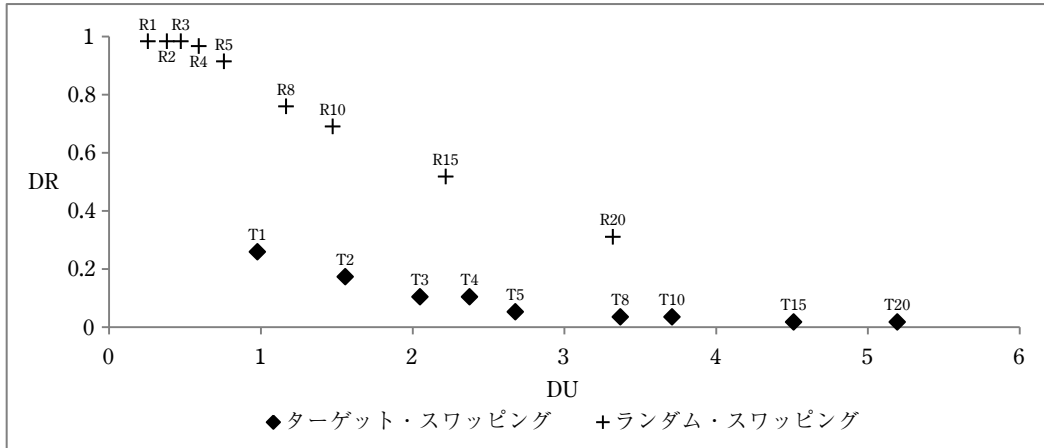
つぎに、本研究では、有用性と秘匿性の評価指標をもとに、R-Uマップを作成し、有用性と秘匿性の相対比較を試みた。R-Uマップで使用する有用性と秘匿性の評価指標に関しては、キー変数の中のあらゆる3変数の組み合わせについて計算された評価指標の平均値がそれぞれ用いられている。図3は、表1-1と表1-2をもとに作成したR-Uマップの結果を示したものである。年齢、性別と国籍のク

表1-2 秘匿性の評価指標に関する試算結果

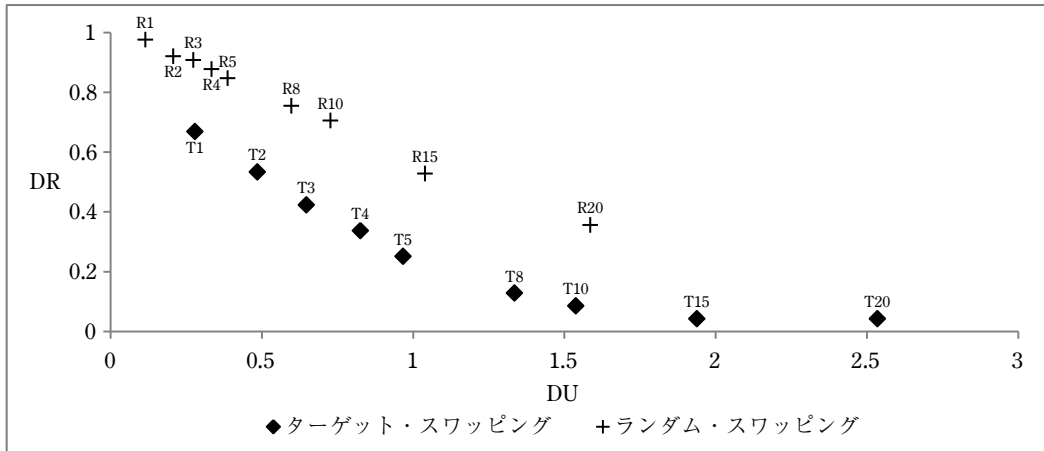
スワッピング率とスワッピングの種類	年齢×性別×国籍	年齢×世帯主の続き柄×労働力状態	3変数のすべての組み合わせに関する平均値
ターゲット・スワッピング			
1%	0.2586	0.6687	0.4493
2%	0.1724	0.5337	0.2859
3%	0.1034	0.4233	0.2010
4%	0.1034	0.3374	0.1561
5%	0.0517	0.2515	0.1138
8%	0.0345	0.1288	0.0704
10%	0.0345	0.0859	0.0577
15%	0.0172	0.0429	0.0448
20%	0.0172	0.0429	0.0422
ランダム・スワッピング			
1%	0.9828	0.9755	0.9644
2%	0.9828	0.9202	0.9341
3%	0.9828	0.9080	0.9070
4%	0.9655	0.8773	0.8767
5%	0.9138	0.8466	0.8418
8%	0.7586	0.7546	0.7314
10%	0.6897	0.7055	0.6706
15%	0.5172	0.5276	0.4830
20%	0.3103	0.3558	0.3191

図3 R-Uマップの結果

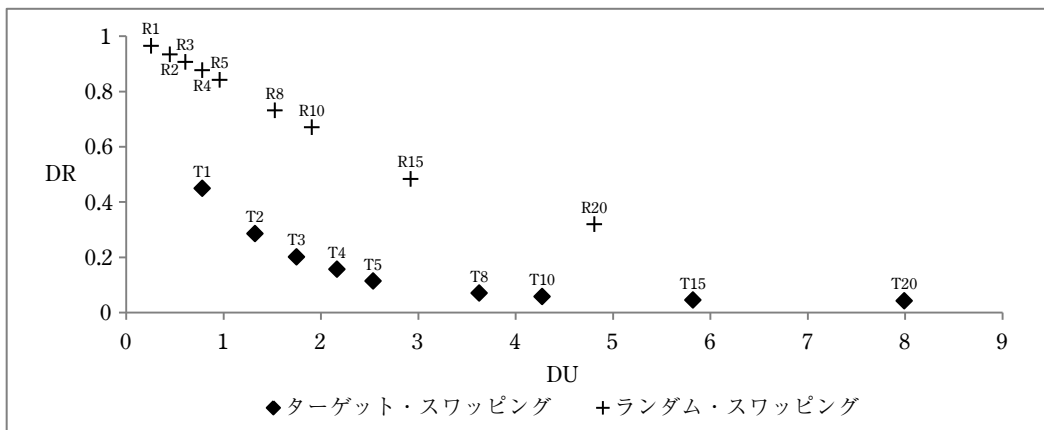
年齢×性別×国籍



年齢×世帯主との続き柄×労働力状態



キー変数における3変数のすべての組み合わせの平均値



注 Rp (pはスワッピング率) についてはランダム・スワッピング, Tp (pはスワッピング率) についてはターゲット・スワッピングを表す。



ロス表に関するR-Uマップを見ると、ターゲット・スワッピングにおいてスワッピング率を1%とした場合、あらゆるランダム・スワッピングよりも秘匿性が高くなることが確認できる。一方、有用性については、ターゲット・スワッピングにおいてスワッピング率が8%に設定された場合、ランダム・スワッピングにおいてスワッピング率を20%にした場合と比較しても、その有用性は低くなっている。こうしたターゲット・スワッピングとランダム・スワッピングにみられる傾向は、年齢、世帯主との続き柄と労働力状態におけるR-Uマップにおいても基本的には変わらない。

さらに、キー変数における3変数のすべての組み合わせの平均値に関するR-Uマップについても見ていくことにしたい。一例として2%のスワッピング率に着目すると、ターゲット・スワッピングを適用した場合、あらゆるランダム・スワッピングよりも秘匿性が高くなることが確認できる。一方、2%のスワッピング率において、ターゲット・スワッピングを適用すると、8%のスワッピング率でランダム・スワッピングを行った場合よりも有用性が高いことがわかる。このことは、有用性の指標がある水準に設定されたとき、ターゲット・スワッピングのほうがより小さなスワッピング率で効率的に秘匿性を高めることが可能なことを意味している。このように秘匿の観点を考慮した場合には、本分析結果から、ランダム・スワッピングよりもターゲット・スワッピングのほうがより有効な手法であると言える。

## 6. おわりに

わが国において政府統計マイクロデータの利

用を促進させるための1つの方向は、より広範な匿名化マイクロデータの作成・提供であるが、そのためには、マイクロデータに対する匿名化技法についての適用可能性の検討が必要である。そこで、本稿では、匿名化技法としてのスワッピングに焦点を当て、スワッピングの有効性について検証を試みた。本研究では、匿名データ作成のための実用性の観点も踏まえ、「特殊な一意」となるレコードの探索方法、スワッピングを行うための質的属性におけるリンケージ技法、クロス表を用いた秘匿性と有用性の評価方法について議論した。本分析結果に関しては、秘匿の観点からは、ランダム・スワッピングよりもターゲット・スワッピングのほうがより有効な手法であることが実証的に明らかになった。一方、本分析ではランダム・スワッピングにおける有用性は、ターゲット・スワッピングのそれよりも高いことが確認されることから、匿名化マイクロデータの作成においては、有用性と秘匿性のバランスを図ることが求められる。

スワッピングは、政府統計マイクロデータの作成のための有力な攪乱の手法の1つであり、諸外国で実用化もなされてきたにも関わらず、わが国における実証研究はこれまで非常に少なかった。本研究は、わが国の国勢調査のマイクロデータを用いてスワッピングの有効性に関する実証分析を行った初めての研究であって、わが国における政府統計の匿名化マイクロデータの作成において、スワッピングの適用可能性を検討する上で有益な研究成果であると考えている。今後、わが国でスワッピングを含む匿名化技法の実証研究がより一層進展することによって、わが国における政府統計の二次的利用のさらなる促進が図られることを期待したい。

## 付記

本稿の作成に当たり、総務省統計局および(独)統計センターの関係各位に大変お世話になった。記して謝意を表したい。また、本稿の旧稿の一部については、Privacy in Statistical Databases 2012 (2012年9月26日～9月28日、於イタリア、パレルモ)等で報告を行ったが、Robert McCaa名誉教授(ミネソタ大学)をはじめとして、多くの方々から貴重なコメントをいただいた。ここに記して感謝の意を表したい。なお、本稿の内容は筆者の個人的見解を示すものであり、(独)統計センターの見解を示すものではないことを申し述べておく。

## 注

- 1) ミクロデータに対する匿名化技法としての攪乱の手法に関する議論は、少なくとも1970年代に遡ることができ、スワッピングの可能性等が議論されてきた(Dalenius and Reiss, 1978)。
- 2) 距離を計算した際に、ドナーファイルの中で最も距離が小さいレコードが複数存在する場合もある。その場合には、最小の距離を有する複数のレコードの中からランダムに1つのレコードを選んでいる。
- 3) 本実験では、 $m \times n$ のクロス表における関連性の尺度であるクラメールのVを用いた有用性の検証も行っている。クラメールのVを用いた有用性の評価指標は、以下の(F1)式で与えられている(Shlomo *et al.*, 2010)。

$$\text{有用性の評価指標} = \frac{CV(T^S) - CV(T^O)}{CV(T^O)} \times 100 \quad (\text{F1})$$

ここで

$CV(T^O)$ : 個票データを用いて作成したクロス表におけるクラメールのV

$CV(T^S)$ : スワッピング済データを用いて作成したクロス表におけるクラメールのV

(F1)式は、クラメールのVを用いた個票データに対するスワッピング済データの情報量損失を表したものであり、(F1)式における有用性の評価指標が大きいほど、情報量損失が大きくなることから、有用性は低いとみなすことができる。

- 4) 個票データにおけるクロス表の中で度数1であるセルが、スワッピング済データにおけるクロス表において度数1のセルとして同じ位置に存在していたとしても、その度数1に該当するレコードにスワッピングが適用されている可能性はある。しかしながら、本実験では、そのようなスワッピング済のレコードについては追跡することができなかった。なお、原データにおけるクロス表の中で度数1であるセルが、ある特定のスワッピング率(例えばスワッピング率が1%)でスワッピングを施すことによって度数0に置き換えられたものの、より高いスワッピング率(例えばスワッピング率が2%)が適用された場合においては、そのセルが再び度数1に置換されることもある。こうした場合には、より高いスワッピング率(例えばスワッピング率が2%)においてセルが度数1であったとしても、それに該当するレコードについては、スワッピングの処理がなされたものとみなしている。
- 5) 本研究では、2変数のすべての組み合わせについてもクロス表を作成し、有用性の評価の比較を行っているが、スワッピング率を変えた場合の情報量損失の変化がより明確に捉えられることから、本稿では、3変数のクロス表をもとに有用性の検証を行っている(これについては秘匿性の検証の場合も同様)。
- 6) 2変数のすべての組み合わせにおけるクロス表をもとに有用性を検証する場合、本研究では、クラメールのVによる指標と絶対距離の平均値による有用性の評価の比較をしている。有用性の評価指標として、クラメールのVを用いた指標の場合、スワッピング率を上げるにつれて、結果数値の動きが傾向的に示されない場合がある。具体的には、国籍と年齢のクロス表の場合、スワッピング率が上がっても、有用性の評価指標が、傾向的に大きくならないことが分かる。これに関しても、国籍において日本人以外の分類区分に該当するレコードが少ないために、クロス表において度数0

となるセルが多くなっており、このことが、クラメールのVにおける指標の結果に影響を及ぼしていると思われる。

#### 参考文献

- [ 1 ] Dalenius, T and Reiss, S.P. (1978) “Data-Swapping: A Technique for Disclosure Control (Extended Abstract)”, in Proceedings of the Section on Survey Research Methods, American Statistical Association, Washington, D.C., pp.191-194.
- [ 2 ] De Kort, S., and Wathan, J. (2009) “Guide to Imputation and Perturbation in the Samples of Anonymised Records”.  
<http://www.ccsr.ac.uk/sars/resources/imputation.doc>. 【2014年7月19日アクセス】
- [ 3 ] De Waal, T. and Willenborg, L. (1999) “Information Loss through Global Recoding and Local Suppression”, *Netherlands Official Statistics (special issue on SDC)*, Vol. 14, pp.17-20.
- [ 4 ] Domingo-Ferrer, J. and Torra, V. (2001a) “Disclosure Control Methods and Information Loss for Microdata”, Doyle *et al.* (eds.) *Confidentiality, Disclosure and Data Access: Theory and Practical Applications for Statistical Agencies*, Elsevier Science, Amsterdam, pp.91-110.
- [ 5 ] Domingo-Ferrer, J. and Torra, V. (2001b) “A Quantitative Comparison of Disclosure Control Methods for Microdata”, Doyle *et al.* (eds.) *Confidentiality, Disclosure, and Data Access: Theory and Practical Application for Statistical Agencies*, Elsevier Science, Amsterdam, pp.111-133.
- [ 6 ] Duncan, G. and Lambert, D. (1989) “The Risk of Disclosure for Microdata” *Journal of Business and Economic Statistics*, Vol. 7, pp.207-217.
- [ 7 ] Duncan, G.T., Keller-McNulty, S. and Stokes, S.L. (2001) “Disclosure Risk vs. Data Utility: the R-U Confidentiality Map” *Technical Report 121*, US National Institute of Statistical Sciences, Durham, North Carolina.
- [ 8 ] Duncan, G.T., Elliot, M., Salazar-González, J. (2011) *Statistical Confidentiality*, Springer, New York.
- [ 9 ] Elliot, M. (2001) “Disclosure Risk Assessment”, Doyle *et al.* (eds.) *Confidentiality, Disclosure, and Data Access: Theory and Practical Application for Statistical Agencies*, Elsevier Science, Amsterdam, pp.75-90.
- [10] Elliot, M.J., Manning, A.M., Ford, R.W. (2002) “A Computational Algorithm for Handling The Special Uniques Problem”, *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, Vol. 10, No. 5, pp.493-509.
- [11] Elliot, M.J. and Manning, A. (2004) “The Methodology used for the 2001 SARs Special Uniques Analysis”, Paper Presented to An Open Meeting on the Samples of Anonymised Records from the 2001 Census, CCSR.  
<http://www.ccsr.ac.uk/sars/events/2004-09-30/Elliot.pdf>. 【2014年7月19日アクセス】
- [12] Federal Committee on Statistical Methodology (1994) *Statistical Policy Working Paper 22: Report on Statistical Disclosure Limitation Methodology*, U.S. Office of Management and Budget, Office of Information and Regulatory Affairs, Washington, D.C..
- [13] Gbur, P.M., Zelenak, M.F. (2004) “Statistical Methodology for the Census 2000 Public Use Microdata Samples”, in Proceedings of the Section on Survey Research Methods, American Statistical Association, pp.3555-3562.
- [14] Gross, B., Guiblin, P., Merrett, K. (2004) “Risk Assessment of the Individual Sample of Anonymised Records (SAR) from the 2001 Census”.  
<http://www.ccsr.ac.uk/sars/guide/2001/Gross2.pdf>. 【2014年7月19日アクセス】
- [15] 伊藤伸介・磯部祥子・秋山裕美 (2008) 「匿名化技法としてのマイクロアグリゲーションの有効性に関する研究—全国消費実態調査を例に—」, 『製表技術参考資料』 No. 10, 33～66頁
- [16] 伊藤伸介・磯部祥子・秋山裕美 (2009) 「秘匿性の評価方法に関する実証研究—全国消費実態調査のマイクロアグリゲートデータを用いて—」, 『製表技術参考資料』 No. 11, 1～35頁

- [17] 伊藤伸介 (2010) 「マイクロデータにおける秘匿性の評価方法に関する一考察」, 明海大学『経済学論集』第22巻第2号, 1~17頁
- [18] 伊藤伸介・高野正博・秋山裕美・後藤武彦 (2010) 「マイクロデータにおける有用性と秘匿性の定量的な評価に関する研究」, 『製表技術参考資料』No. 14, 1~40頁
- [19] 伊藤伸介・村田磨理子 (2013) 「家計調査マイクロデータを用いた攪乱の手法の有効性に関する研究」『製表技術参考資料』No. 22, 1~26頁
- [20] Karr, A.F., Kohnen, C.N., Oganian, A., Reiter, J.P., Sanil, A.P. (2006) “A Framework for Evaluating the Utility of Data Altered to Protect Confidentiality”, *The American Statistician*, Vol. 60, No. 3, pp.1-9.
- [21] Müller, W., Blien, U., Wirth, H. (1995) “Identification Risks of Micro Data: Evidence from Experimental Studies”, *Sociological Methods and Research*, Vol. 24, No. 2, pp.131-157.
- [22] Shlomo, N. (2007) “Statistical Disclosure Control Methods for Census Frequency Tables”, *S3RI Methodology Working Papers M07/04*, pp.1-40.  
<http://eprints.soton.ac.uk/44610/1/44610-01.pdf>. 【2014年7月19日アクセス】
- [23] Shlomo, N. (2010) “Releasing Microdata: Disclosure Risk Estimation, Data Masking and Assessing Utility”, *The Journal of Privacy and Confidentiality*, Vol. 2, No. 1, pp.73-91.
- [24] Shlomo, N., Tudor, C., Groom, P. (2010) “Data Swapping for Protecting Census Tables”, Domingo-Ferrer, J. and Magkos, E. (eds) *Privacy in Statistical Databases UNESCO Chair in Data Privacy International Conference, PSD 2010 Corfu, Greece, September, 2010 Proceedings*, Springer, pp.41-51.
- [25] Skinner, C.J. (1992) “On Identification Disclosure and Prediction Disclosure for Microdata”, *Statistica Neerlandica*, Vol. 46, No. 1, pp.21-32.
- [26] Takemura, A. (1999) “Local Recoding by Maximum Weight Matching for Disclosure Control of Microdata sets”, *ITME Discussion Paper*, No. 11, Faculty of Economics, Univ. of Tokyo.
- [27] Takemura, A. (2002) “Local Recoding and Record Swapping by Maximum Weight Matching for Disclosure Control of Microdata Sets”, *Journal of Official Statistics*, Vol. 18, No. 2, pp.275-289.
- [28] Willenborg, L. and de Waal, T. (2001) *Elements of Statistical Disclosure Control*, Springer, New York.
- [29] Woo, M., Reiter, J.P., Oganian, A., Karr, A.F. (2009) “Global Measures of Data Utility for Microdata Masked for Disclosure Limitation”, *The Journal of Privacy and Confidentiality*, Vol. 1, No. 1, pp.111-124.
- [30] Zayatz, L. (2007) “Disclosure Avoidance Practices and Research at the U.S. Census Bureau: An Update”, *Journal of Official Statistics*, Vol. 23, No. 2, pp.253-265.

# Effectiveness of Data Swapping Based on the Microdata from Population Census

Shinsuke ITO<sup>\*</sup>, Naomi HOSHINO<sup>\*\*</sup>

## Summary

Only a limited number of empirical studies on disclosure limitation methods including perturbation, disclosure risk and information loss have been conducted in Japan so far. More extensive research on perturbative methods could help expand their use in the creation of anonymized official microdata in Japan. This paper examines the potential of data swapping as a perturbative method for the anonymization of individual data from Japanese official statistics, and empirically determines data utility and data confidentiality for the swapped data. The results show an overall higher data confidentiality for targeted data swapping than for random data swapping, and therefore indicate that for a specific level of data utility, targeted data swapping achieves higher data confidentiality than random data swapping even at lower swapping rates. This suggests targeted data swapping is the more effective method to achieve data confidentiality.

## Key Words

Population Census, Microdata, Disclosure Limitation Methods, Data Swapping

---

<sup>\*</sup> Faculty of Economics, Chuo University  
(Visiting Fellow of National Statistics Center)

<sup>\*\*</sup> National Statistics Center



編集委員会からのお知らせ  
機関誌『統計学』の編集・発行について

編集委員会

1. 常時、投稿を受け付けます。
2. 次号以降の発行予定日は、  
第108号：2015年3月31日、第109号：2015年9月30日です。
3. 投稿に際しては、「投稿規程」、「執筆要綱」、「査読要領」などをご熟読願います。
4. 原稿は編集委員長（下記メールアドレス）宛にお送り願います。
5. 原稿はPDF形式のファイルとして提出して下さい。また、紙媒体での提出も旧規程に準拠して受け付けます。紙媒体の送付先は編集委員長宛をお願いいたします。
6. 原則としてすべての投稿原稿が査読の対象となります。
7. 通常、査読から発刊までに要する期間は、査読が順調に進んだ場合でも2ヶ月間程を要します。投稿にあたっては十分に留意して下さい。

編集委員会、投稿応募についての問い合わせは、  
下記メールアドレス宛に連絡下さい。  
また、編集委員長へのメールアドレスも下記になります。

[editorial@jsest.jp](mailto:editorial@jsest.jp)

編集委員長 岡部純一（横浜国立大学）

副委員長 長澤克重（立命館大学）

編集委員

栗原由紀子（弘前大学）

橋本貴彦（立命館大学）

山田 満（関東支部所属）

[注記] 2013年度より編集体制の見直しとして、第一次査読を従来のように支部選出委員が担当するのではなく、編集委員会全体で担当するように方針を変更しています。『統計学』の定期刊行にも力点をおく所存です。常時、投稿を受け付けていますので、できるかぎり早期のご投稿をお願いいたします。108号（2015年3月31日発行予定）への掲載を想定すると、A：「論文」・「研究ノート」の場合、2015年1月初旬、B：その他の場合、2015年1月末を目途に、それまでにご投稿いただく必要があります。

以上

編集後記

ご投稿いただいたすべての執筆者のみなさん、査読に関わってくださった会員のみなさんに心より御礼申し上げます。今回は書評や海外統計事情の執筆依頼にもご快諾いただきました。そうした掲載記事について、会員のみなさんから編集委員会にご提案ご推薦いただければ、紙面活性化にもつながりありがたいです。よろしく願います。

（岡部純一 記）

[訂正] 『統計学』第106号（2014年3月）p.40の「2013年度関西支部例会」5月19日(土)【報告者】  
(1) 桂政昭（誤）について、(1) 桂昭政（正）に訂正します。失礼いたしました。

## 執筆者紹介 (掲載順)

伊藤伸介	(中央大学経済学部)
星野なおみ	((独)統計センター)
山口幸三	(総務省統計研修所)
橋本貴彦	(立命館大学経済学部)
上藤一郎	(静岡大学人文社会科学部)
イリーナ・エリセーエワ	(ロシア統計学会会長)
山口秋義	(九州国際大学経済学部)

## 支部名

## 事務局

北海道	004-0042	札幌市厚別区大谷地西 2-3-1 北星学園大学経済学部 (011-891-2731)	古谷次郎
東北	986-8580	石巻市南境新水戸 1 石巻専修大学経営学部 (0225-22-7711)	深川通寛
関東	192-0393	八王子市東中野 742-1 中央大学経済学部 (042-674-3424)	芳賀寛
関西	525-8577	草津市野路東 1-1-1 立命館大学経営学部 (077-561-4631)	田中力
九州	870-1192	大分市大字且野原 700 大分大学経済学部 (097-554-7706)	西村善博

## 編集委員

岡部純一 (関東) [長]	長澤克重 (関西) [副]
山田満 (関東)	橋本貴彦 (関西)
栗原由紀子 (関東)	

## 統計学 No.107

---

2014年9月30日 発行	発行所	経済統計学会 〒194-0298 東京都町田市相原町4342 法政大学日本統計研究所内 TEL 042(783)2325 FAX 042(783)2332 <a href="http://www.jses.t.jp/">http://www.jses.t.jp/</a>
	発行人	代表者 菊地進
	発売所	音羽リスマチック株式会社 〒112-0013 東京都文京区音羽1-6-9 TEL/FAX 03(3945)3227 E-mail: <a href="mailto:otorisu@jupiter.ocn.ne.jp">otorisu@jupiter.ocn.ne.jp</a> 代表者 遠藤誠

---

# STATISTICS

---

No. 107

2014 September

---

## Articles

- Effectiveness of Data Swapping Based on the Microdata from Population Census  
..... Shinsuke ITO and Naomi HOSHINO (1)
- Estimation Bias in Statistical Survey applying the Sample Rotation System  
..... Kozo YAMAGUCHI (17)

## Book Reviews

- Tadashi YOSHIDA, *On the Progress of Probability Theory and Statistics in the Netherlands*,  
Hassakusha, 2014  
..... Ichiro UWAFUJI (33)
- Hiroshi IZUMI, *A Measurement of Embodied Labor and Basic Economic Indicators*,  
Ohtsuki Syoten, 2014  
..... Takahiko HASHIMOTO (38)

## Foreign Statistical Affairs

- Russian Association of Statisticians  
..... Irina ELISEEVA and Akiyoshi YAMAGUCHI (43)

## Activities of the Society

- The 58<sup>th</sup> Session of the Society of Economic Statistics ..... (46)

---

JAPAN SOCIETY OF ECONOMIC STATISTICS

---