

# 全国消費実態調査の匿名データから SASによる新擬似マイクロデータの作成

周防節雄 (公財)統計情報研究開発センター

高橋行雄 BioStat研究所(株)

宮内亨 (独)統計センター

2017/11/17

平成29年度共同研究集会

「官民オープンデータ利活用の動向及び人材育成の取組」(29-共研-5010)

@ 統計数理研究所

## はじめに

### 目的

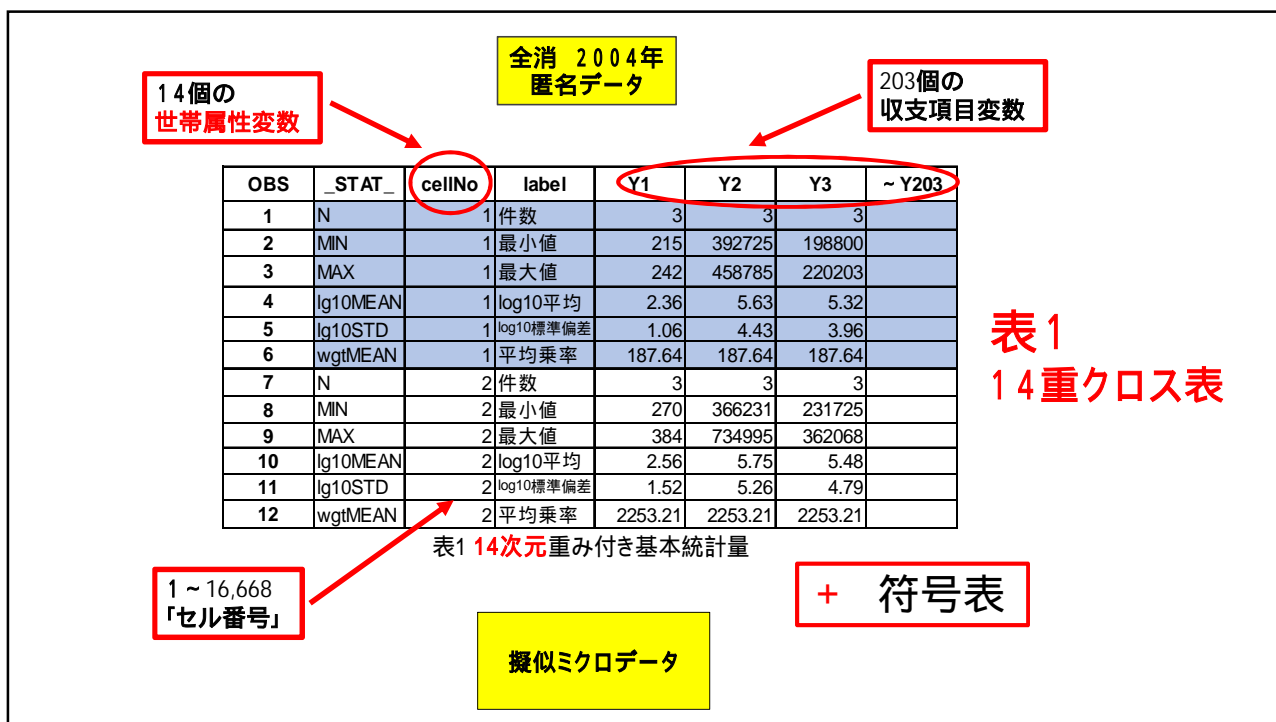
全消(2004年)の匿名データ → SAS擬似マイクロデータ

### 方法

匿名データ → **基本統計表** → 擬似マイクロデータ

### 結果

作成したデータの品質？



## 秘匿性のために

出現頻度1のセル

コピーを2つ追加 → 出現頻度3にする

出現頻度2のセル

コピーを一つずつ追加 → 出現頻度4にする

コピーしたデータには  
乱数を使って若干の誤差変動を与えた

## 2. 擬似マイクロデータの構造

変数名	項目名	変数リスト(付録1)					
	<b>世帯事項</b>	Y38	支出総額	Y94	家庭用耐久財	Y150	(特掲)イ
Year	調査年	Y39	実支出	Y95	家事用耐久財	Y151	その他の
No	レコード連番号	Y40	消費支出	Y96	冷暖房用器具	Y152	諸雑費
X01	大都市圏の別	Y41	食料	Y97	一般家具	Y153	理美容
X02	世帯区分	Y42	穀類	Y98	室内装備・装飾品	Y154	理美容
X03	世帯人員	Y43	米	Y99	寝具類	Y155	身の回
X04	就業人員	Y44	パン	Y100	家事雑貨	Y156	たばこ
X05	住居の構造	Y45	めん類	Y101	家事用消耗品	Y157	その他の
X06	住居の建て方	Y46	他の穀類	Y102	家事サービス	Y158	こづかい
X07	住居の所有関係	Y47	魚介類	Y103	被服及び履物	Y159	交際費
X08	世帯主の性別	Y48	生鮮魚介	Y104	和服	Y160	交際費
X09	世帯主の年齢	Y49	塩干魚介	Y105	洋服	Y161	交際費
X10	企業区分・従業者規模	Y50	魚肉練製品	Y106	男子用洋服	Y162	交際費
X11	家族分類	Y51	他の魚介加工品	Y107	婦人用洋服	Y163	交際費
X12	未就学児の有無	Y52	肉類	Y108	子供用洋服	Y164	他の物
X13	学校に通う世帯員の有無	Y53	生鮮肉	Y109	シャツ・セーター類	Y165	贈与金
X14	65歳以上の世帯員数	Y54	加工肉	Y110	男子用シャツ・セーター類	Y166	他の交
Weight	集計用乗率	Y55	乳卵類	Y111	婦人用シャツ・セーター類	Y167	仕送り金
	<b>収支項目(単位千円)</b>	Y56	牛乳	Y112	子供用シャツ・セーター類	Y168	(再掲)教
Y1	年間収入	Y57	乳製品	Y113	下着類	Y169	(再掲)教

## 収支項目の203変数の関係

親子関係の変数 (186個)

足し上げ構造

親の変数 = 子の全変数の合計

付録3

再掲・特掲関係の変数(16個)

再掲・特掲の変数

大小関係

それにぶら下がる変数

無関係の変数

1個

「年間収入」

表1 (14重クロスの基本統計量)

+

変数間の密接な関係



擬似マイクロデータ

擬似マイクロデータ作成の  
作業手順

(付録4を参照)

# 作成された擬似マイクロデータの 問題点

JMPによる作成作業で行った

「ゼロ円収支項目の出現割合調整」  
「年間収入3階区分別主要21項目間の相関係数行列」？  
組み込まなかった

足し上げ関係は全て補正した  
再掲・特掲関係の一部の変数に再調整が必要

## 検証と改善案

匿名データ (47,797 observations)



我々の擬似マイクロデータ (73,519 observations)

最初の段階で

出現頻度1と2のデータをコピーした結果

秘匿性のために  
必要なのか？

## 検証と改善案

擬似マイクロデータと匿名データについて  
203変数  
平均値、標準偏差、最小値、最大値  
(集計用乗率:不使用/使用) **付録5**

### 平均値

集計用乗率なし:匿名データの平均値にかなり近い  
集計用乗率あり:両者の差が広がっている

### 箱ひげ図

擬似マイクロデータに外れ値多い

**付録6**

## 検証と改善案

外れ値の値と出現割合を匿名データにほぼ見合うようにして、新擬似マイクロデータの変数を発生させれば散らばり具合をかなりコントロールできる



### 提案

公開する「14次元重み付き基本統計量」(表1)の他に、公表の許容範囲内で、**セル毎の外れ値に関する情報**も表形式にまとめて公表し、擬似マイクロデータ作成時に利用できるようにする。

?

## まとめ

他の年次の全消匿名データも使って、  
複数回次の擬似マイクロデータの作成。

利用者にとって、どのような家族分類が最適？

最終的に擬似マイクロデータの作成方法が確定した後、  
作成マニュアルを整備・公開

## 謝辞

本研究の推進に全国消費実態調査の匿名データを使用した  
が、その際、(株)SASジャパンインスティテュートから  
使用料の援助を受けた。また、匿名データの利用に際しては、  
独立行政法人統計センターから便宜を図っていただいた。  
ここに記して謝意を述べる。

ご静聴、  
ありがとうございました。