

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」  
令和3(2021)年11月18日(木)

# 公的統計マイクロデータの 利活用推進に資する Synthetic Dataの 作成方法について

高部 勲

立正大学データサイエンス学部

# Contents

---

1. 公的統計マイクロデータの概要
2. 疑似データの重要性と課題
3. Synthetic Dataに基づく方法
4. 実データに基づく試作
5. まとめと今後の課題

# **1. 公的統計ミクロデータの概要**

2. 疑似データの重要性と課題
3. Synthetic Dataに基づく方法
4. 実データに基づく試作
5. まとめと今後の課題

# 公的統計マイクロデータ（調査票情報）

- 国の行政機関が実施した統計調査の結果について、**調査対象の秘密の保護を図った上で、世帯単位や事業所単位といった集計する前の個票形式のデータ（マイクロデータ、調査票情報）**を提供
- マイクロデータ（調査票情報）を用いることで、研究者の方々は、より自由で多様な分析を行うことが可能となるため、新たな発見につながることを期待



## マイクロデータ（調査票情報）

Year	Industry	Company	Region	Product	Value	Unit	Count	Weight
2019	製造業	株式会社	関東	自動車	1000000	円	10000	1.0
2019	製造業	株式会社	中部	自動車	800000	円	8000	0.8
2019	製造業	株式会社	関西	自動車	1200000	円	12000	1.2
2019	製造業	株式会社	中国	自動車	600000	円	6000	0.6
2019	製造業	株式会社	四国	自動車	400000	円	4000	0.4
2019	製造業	株式会社	九州	自動車	500000	円	5000	0.5
2019	製造業	株式会社	北海道	自動車	300000	円	3000	0.3
2019	製造業	株式会社	東北	自動車	700000	円	7000	0.7
2019	製造業	株式会社	関東	自動車	1100000	円	11000	1.1
2019	製造業	株式会社	中部	自動車	900000	円	9000	0.9
2019	製造業	株式会社	関西	自動車	1300000	円	13000	1.3
2019	製造業	株式会社	中国	自動車	700000	円	7000	0.7
2019	製造業	株式会社	四国	自動車	500000	円	5000	0.5
2019	製造業	株式会社	九州	自動車	600000	円	6000	0.6
2019	製造業	株式会社	北海道	自動車	400000	円	4000	0.4
2019	製造業	株式会社	東北	自動車	800000	円	8000	0.8
2019	製造業	株式会社	関東	自動車	1200000	円	12000	1.2
2019	製造業	株式会社	中部	自動車	1000000	円	10000	1.0
2019	製造業	株式会社	関西	自動車	1400000	円	14000	1.4
2019	製造業	株式会社	中国	自動車	800000	円	8000	0.8
2019	製造業	株式会社	四国	自動車	600000	円	6000	0.6
2019	製造業	株式会社	九州	自動車	700000	円	7000	0.7
2019	製造業	株式会社	北海道	自動車	500000	円	5000	0.5
2019	製造業	株式会社	東北	自動車	900000	円	9000	0.9
2019	製造業	株式会社	関東	自動車	1300000	円	13000	1.3
2019	製造業	株式会社	中部	自動車	1100000	円	11000	1.1
2019	製造業	株式会社	関西	自動車	1500000	円	15000	1.5
2019	製造業	株式会社	中国	自動車	900000	円	9000	0.9
2019	製造業	株式会社	四国	自動車	700000	円	7000	0.7
2019	製造業	株式会社	九州	自動車	800000	円	8000	0.8
2019	製造業	株式会社	北海道	自動車	600000	円	6000	0.6
2019	製造業	株式会社	東北	自動車	1000000	円	10000	1.0
2019	製造業	株式会社	関東	自動車	1400000	円	14000	1.4
2019	製造業	株式会社	中部	自動車	1200000	円	12000	1.2
2019	製造業	株式会社	関西	自動車	1600000	円	16000	1.6
2019	製造業	株式会社	中国	自動車	1000000	円	10000	1.0
2019	製造業	株式会社	四国	自動車	800000	円	8000	0.8
2019	製造業	株式会社	九州	自動車	900000	円	9000	0.9
2019	製造業	株式会社	北海道	自動車	700000	円	7000	0.7
2019	製造業	株式会社	東北	自動車	1100000	円	11000	1.1
2019	製造業	株式会社	関東	自動車	1500000	円	15000	1.5
2019	製造業	株式会社	中部	自動車	1300000	円	13000	1.3
2019	製造業	株式会社	関西	自動車	1700000	円	17000	1.7
2019	製造業	株式会社	中国	自動車	1100000	円	11000	1.1
2019	製造業	株式会社	四国	自動車	900000	円	9000	0.9
2019	製造業	株式会社	九州	自動車	1000000	円	10000	1.0
2019	製造業	株式会社	北海道	自動車	800000	円	8000	0.8
2019	製造業	株式会社	東北	自動車	1200000	円	12000	1.2
2019	製造業	株式会社	関東	自動車	1600000	円	16000	1.6
2019	製造業	株式会社	中部	自動車	1400000	円	14000	1.4
2019	製造業	株式会社	関西	自動車	1800000	円	18000	1.8
2019	製造業	株式会社	中国	自動車	1200000	円	12000	1.2
2019	製造業	株式会社	四国	自動車	1000000	円	10000	1.0
2019	製造業	株式会社	九州	自動車	1100000	円	11000	1.1
2019	製造業	株式会社	北海道	自動車	900000	円	9000	0.9
2019	製造業	株式会社	東北	自動車	1300000	円	13000	1.3
2019	製造業	株式会社	関東	自動車	1700000	円	17000	1.7
2019	製造業	株式会社	中部	自動車	1500000	円	15000	1.5
2019	製造業	株式会社	関西	自動車	1900000	円	19000	1.9
2019	製造業	株式会社	中国	自動車	1300000	円	13000	1.3
2019	製造業	株式会社	四国	自動車	1100000	円	11000	1.1
2019	製造業	株式会社	九州	自動車	1200000	円	12000	1.2
2019	製造業	株式会社	北海道	自動車	1000000	円	10000	1.0
2019	製造業	株式会社	東北	自動車	1400000	円	14000	1.4
2019	製造業	株式会社	関東	自動車	1800000	円	18000	1.8
2019	製造業	株式会社	中部	自動車	1600000	円	16000	1.6
2019	製造業	株式会社	関西	自動車	2000000	円	20000	2.0
2019	製造業	株式会社	中国	自動車	1400000	円	14000	1.4
2019	製造業	株式会社	四国	自動車	1200000	円	12000	1.2
2019	製造業	株式会社	九州	自動車	1300000	円	13000	1.3
2019	製造業	株式会社	北海道	自動車	1100000	円	11000	1.1
2019	製造業	株式会社	東北	自動車	1500000	円	15000	1.5
2019	製造業	株式会社	関東	自動車	1900000	円	19000	1.9
2019	製造業	株式会社	中部	自動車	1700000	円	17000	1.7
2019	製造業	株式会社	関西	自動車	2100000	円	21000	2.1
2019	製造業	株式会社	中国	自動車	1500000	円	15000	1.5
2019	製造業	株式会社	四国	自動車	1300000	円	13000	1.3
2019	製造業	株式会社	九州	自動車	1400000	円	14000	1.4
2019	製造業	株式会社	北海道	自動車	1200000	円	12000	1.2
2019	製造業	株式会社	東北	自動車	1600000	円	16000	1.6
2019	製造業	株式会社	関東	自動車	2000000	円	20000	2.0
2019	製造業	株式会社	中部	自動車	1800000	円	18000	1.8
2019	製造業	株式会社	関西	自動車	2200000	円	22000	2.2
2019	製造業	株式会社	中国	自動車	1600000	円	16000	1.6
2019	製造業	株式会社	四国	自動車	1400000	円	14000	1.4
2019	製造業	株式会社	九州	自動車	1500000	円	15000	1.5
2019	製造業	株式会社	北海道	自動車	1300000	円	13000	1.3
2019	製造業	株式会社	東北	自動車	1700000	円	17000	1.7
2019	製造業	株式会社	関東	自動車	2100000	円	21000	2.1
2019	製造業	株式会社	中部	自動車	1900000	円	19000	1.9
2019	製造業	株式会社	関西	自動車	2300000	円	23000	2.3
2019	製造業	株式会社	中国	自動車	1700000	円	17000	1.7
2019	製造業	株式会社	四国	自動車	1500000	円	15000	1.5
2019	製造業	株式会社	九州	自動車	1600000	円	16000	1.6
2019	製造業	株式会社	北海道	自動車	1400000	円	14000	1.4
2019	製造業	株式会社	東北	自動車	1800000	円	18000	1.8
2019	製造業	株式会社	関東	自動車	2200000	円	22000	2.2
2019	製造業	株式会社	中部	自動車	2000000	円	20000	2.0
2019	製造業	株式会社	関西	自動車	2400000	円	24000	2.4
2019	製造業	株式会社	中国	自動車	1800000	円	18000	1.8
2019	製造業	株式会社	四国	自動車	1600000	円	16000	1.6
2019	製造業	株式会社	九州	自動車	1700000	円	17000	1.7
2019	製造業	株式会社	北海道	自動車	1500000	円	15000	1.5
2019	製造業	株式会社	東北	自動車	1900000	円	19000	1.9
2019	製造業	株式会社	関東	自動車	2300000	円	23000	2.3
2019	製造業	株式会社	中部	自動車	2100000	円	21000	2.1
2019	製造業	株式会社	関西	自動車	2500000	円	25000	2.5
2019	製造業	株式会社	中国	自動車	1900000	円	19000	1.9
2019	製造業	株式会社	四国	自動車	1700000	円	17000	1.7
2019	製造業	株式会社	九州	自動車	1800000	円	18000	1.8
2019	製造業	株式会社	北海道	自動車	1600000	円	16000	1.6
2019	製造業	株式会社	東北	自動車	2000000	円	20000	2.0
2019	製造業	株式会社	関東	自動車	2400000	円	24000	2.4
2019	製造業	株式会社	中部	自動車	2200000	円	22000	2.2
2019	製造業	株式会社	関西	自動車	2600000	円	26000	2.6
2019	製造業	株式会社	中国	自動車	2000000	円	20000	2.0
2019	製造業	株式会社	四国	自動車	1800000	円	18000	1.8
2019	製造業	株式会社	九州	自動車	1900000	円	19000	1.9
2019	製造業	株式会社	北海道	自動車	1700000	円	17000	1.7
2019	製造業	株式会社	東北	自動車	2100000	円	21000	2.1
2019	製造業	株式会社	関東	自動車	2500000	円	25000	2.5
2019	製造業	株式会社	中部	自動車	2300000	円	23000	2.3
2019	製造業	株式会社	関西	自動車	2700000	円	27000	2.7
2019	製造業	株式会社	中国	自動車	2100000	円	21000	2.1
2019	製造業	株式会社	四国	自動車	1900000	円	19000	1.9
2019	製造業	株式会社	九州	自動車	2000000	円	20000	2.0
2019	製造業	株式会社	北海道	自動車	1800000	円	18000	1.8
2019	製造業	株式会社	東北	自動車	2200000	円	22000	2.2
2019	製造業	株式会社	関東	自動車	2600000	円	26000	2.6
2019	製造業	株式会社	中部	自動車	2400000	円	24000	2.4
2019	製造業	株式会社	関西	自動車	2800000	円	28000	2.8
2019	製造業	株式会社	中国	自動車	2200000	円	22000	2.2
2019	製造業	株式会社	四国	自動車	2000000	円	20000	2.0
2019	製造業	株式会社	九州	自動車	2100000	円	21000	2.1
2019	製造業	株式会社	北海道	自動車	1900000	円	19000	1.9
2019	製造業	株式会社	東北	自動車	2300000	円	23000	2.3
2019	製造業	株式会社	関東	自動車	2700000	円	27000	2.7
2019	製造業	株式会社	中部	自動車	2500000	円	25000	2.5
2019	製造業	株式会社	関西	自動車	2900000	円	29000	2.9
2019	製造業	株式会社	中国	自動車	2300000	円		

# 公的統計マイクロデータのオンサイト利用

**公的統計の  
マイクロデータを使って  
分析してみませんか？**

公的統計データの新たな利用方法  
**オンサイト**のご案内

公表されている集計結果からは  
分からないことも…

マイクロデータを使用した  
分析によって…

新たな視点  
自由な分析  
新しい発見

公的統計の調査結果は、公表されている集計表に加え、  
マイクロデータを使用することにより、研究者の方々は  
更に多様な分析を行うことができます。

総務省統計局  
独立行政法人 統計センター

## オンサイトとは？

オンサイトとは、情報セキュリティが確保された環境で、  
許可を受けた研究者がマイクロデータを用いて、独自の集計・分析を行うことができる専用室です。  
現在、オンサイト施設及び利用可能な統計調査の拡充に向けて、  
関係府省、関係機関等と順次調整しています。



入退室管理や  
監視カメラを備えた  
オンサイト室



(マイクロデータ利用ポータルサイト“miripo”から引用  
<https://www.e-stat.go.jp/microdata/>)

# 公的統計マイクロデータ研究コンソーシアム

公的統計マイクロデータ研究コンソーシアム：  
ウェブサイト

<http://jmodc.org/>

公的統計マイクロデータ  
研究コンソーシアム

トップページ  
コンソーシアム概要  
オンサイトネットワーク  
活動予定・報告  
お問い合わせ

お問い合わせ

お問い合わせはメールにて承ります。下記メールアドレス宛に以下の内容をご記載の上、お送りください。

- お名前
- ご所属
- お問い合わせ内容

公的統計マイクロデータの  
利活用推進に向けて

「公的統計マイクロデータ研究コンソーシアム」は、我が国における公的統計マイクロデータの研究利用（二次利用）を促進するために、学官産の関係機関が一体となり、取り組むことを目的として立ち上げるものです。

個人公開募集中

2021.5.11  
公的統計マイクロデータを利用した研究にご関心のある研究者の皆様へ、  
オンラインのアンケートを実施します。  
(回答期限:2021年6月11日(金) 所要時間:15分程度 無記名可)

重要なお知らせ

- [●アンケートのお願い\(依頼\)](#)
- [●アンケート \(Googleドキュメント\)](#)

1. 公的統計マイクロデータの概要

 **2. 疑似データの重要性**

3. Synthetic Dataに基づく方法

4. 実データに基づく試作

5. まとめと今後の課題

# 疑似データの重要性

## ① プログラムテスト用データとしての活用

- 公的統計ミクロデータの利用申請を行う前にデータの利活用のイメージを把握（※オンサイト利用）
- 各種のアプリケーションの開発・テスト
- 匿名加工情報等の作成におけるニーズの事前把握 など

## ② データを利用した教育への活用

- 経済・社会の計量モデルは、個人・企業等を単位としているものが多い（重回帰分析など）
  - ⇒ 集計データではないレコード単位のデータの重要性
  - ⇒ 元データと結果がそれほど変わらないのが望ましい
- 複雑な手続きを必要とせず自由に利用できるデータの重要性

# 疑似データの重要性（教育への活用）

## 数理・データサイエンス・AI リテラシーレベルの教育方法

▶ 「導入」「基礎」「心得」「選択」のそれぞれの分類ごとに、推奨される具体的な教育方法を以下のとおりまとめた。

導入	<p><b>1. 社会におけるデータ・AI利活用</b></p> <ul style="list-style-type: none"> <li>データ・AI利活用事例を紹介した動画（MOOC等）を使った<b>反転学習</b>を取り入れ、講義ではデータ・AI活用領域の広がりや、技術概要の解説を行うことが望ましい</li> <li>学生がデータ・AI利活用事例を調査し発表する<b>グループワーク</b>等を行い、一方通行で事例を話すだけの講義にしないことが望ましい</li> </ul>	<p>教育方法（例）※</p> <p>1, 2, 3, 4</p> <p>1, 4</p>
基礎	<p><b>2. データリテラシー</b></p> <ul style="list-style-type: none"> <li>各大学・高専の特徴に応じて<b>適切なテーマ</b>を設定し、<b>実データ</b>（あるいは模擬データ）を用いた講義を行うことが望ましい</li> <li>実際に手を動かしてデータを可視化する等、学生自身がデータ利活用プロセスの一部を<b>体験</b>できることが望ましい</li> <li>必要に応じてデータハンドリングスキルを埋めるためのフォローアップ講義（<b>補講</b>等）を準備することが望ましい</li> </ul>	<p>教育方法（例）※</p> <p>1, 2, 3, 4</p> <p>1, 4</p>
心得	<p><b>3. データ・AI利活用における留意事項</b></p> <ul style="list-style-type: none"> <li>身近で起こったデータ・AI活用における負の事例を通して、データ駆動型社会のリスクを<b>自分ごと</b>として考えさせることが望ましい（必要に応じてMOOC等の活用も検討する）</li> <li>データ・AIが引き起こす課題について<b>グループディスカッション</b>等を行い、一方通行で事例を話すだけの講義にしないことが望ましい</li> </ul>	<p>教育方法（例）※</p> <p>1, 2, 3, 4</p> <p>1, 4</p>
選択	<p><b>4. オプション</b></p> <ul style="list-style-type: none"> <li>本内容は<b>オプション</b>扱いとし、大学・高専の特徴に応じて学修内容を選択する</li> <li>各大学・高専の特徴に応じて<b>適切なテーマ</b>を設定し、<b>実データ</b>（あるいは模擬データ）を用いた講義を行うことが望ましい</li> <li>学生が希望すれば本内容を受講できるようにしておくことが望ましい（<b>大学間連携</b>等）</li> </ul>	<p>教育方法（例）※</p> <p>1, 2, 3, 4</p>

※上記の「教育方法」欄の手法・形式 1～4 については次頁以降を参照

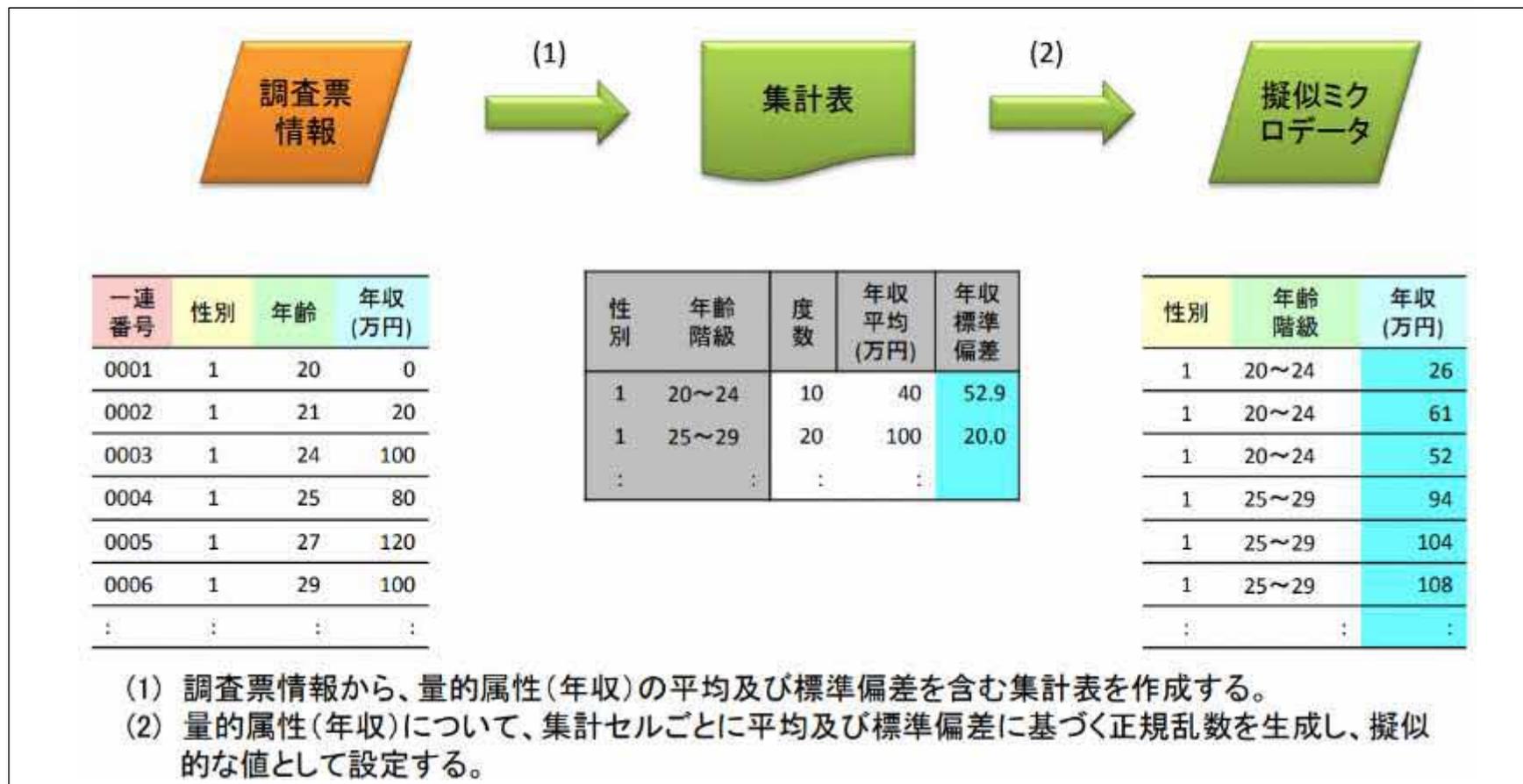
19

「数理・データサイエンス・AI（リテラシーレベル）モデルカリキュラム」  
（2020年4月、数理・データサイエンス教育強化拠点コンソーシアム） から引用

# 一般用マイクロデータ

(「一般用マイクロデータ(仮称)の作成及び利活用について」  
(2015年度統計関連学会連合大会)から引用)

- 現行の制度上、マイクロデータから直接的にレコード単位の疑似データを作成・提供することはできない
- 中間的な集計表を公表し、そこから疑似データを作成



⇒ 変数間(離散・連続)の関係(回帰モデル)  
を考慮した疑似データの作成が課題

1. 公的統計マイクロデータの概要

2. 疑似データの重要性と課題

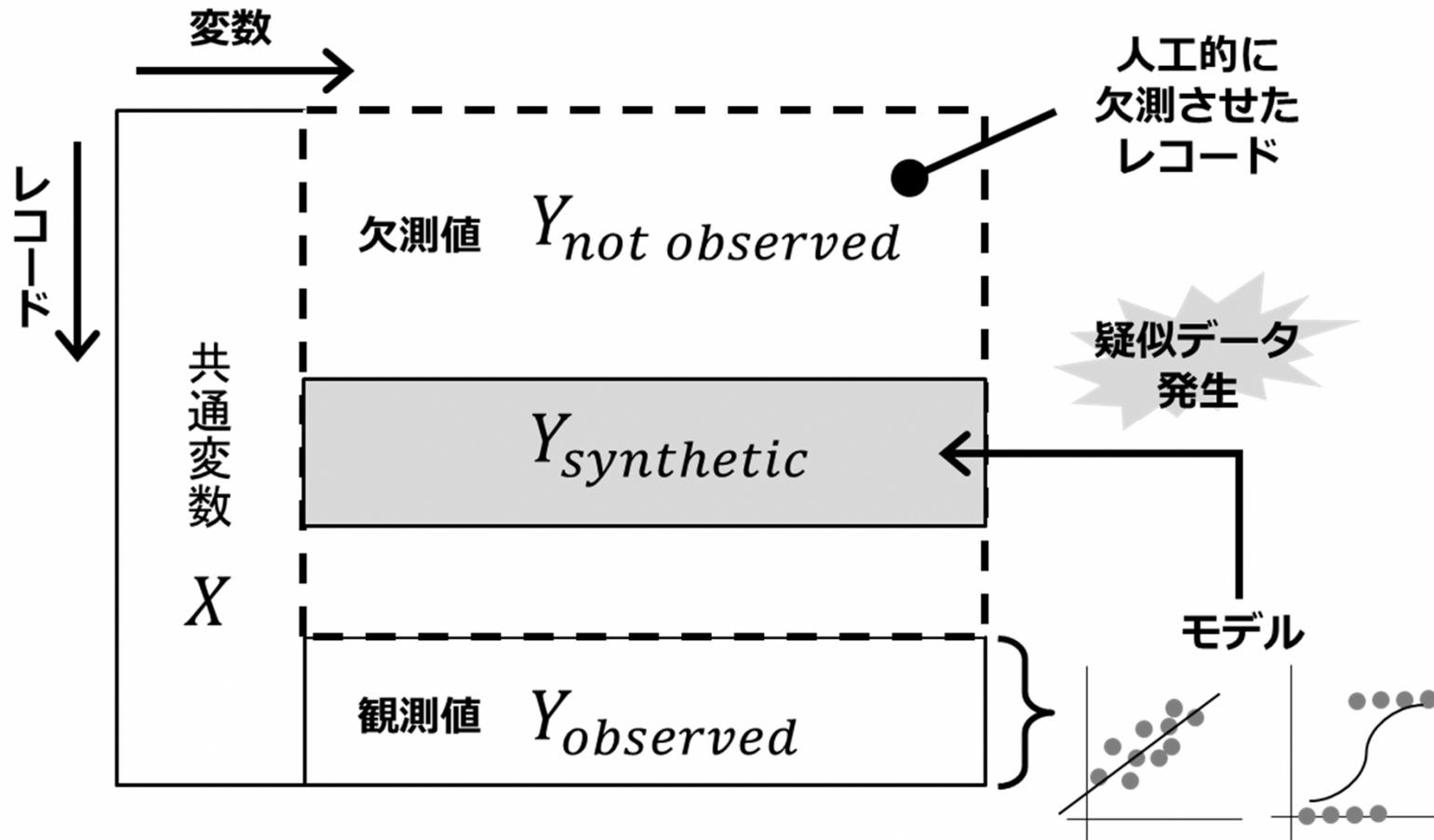
 **3. Synthetic Dataに基づく方法**

4. 実データに基づく試作

5. まとめと今後の課題

# Synthetic Data

- 一部のレコードを人工的に欠測させ回帰モデル、ロジットモデル等により逐次的に疑似データを作成  
⇒変数間の関係を保持したデータの作成が可能  
(Rのパッケージも提供されている (“synthpop”))



# 疑似データの作成方法

- ただし現行制度上、マイクロデータから直接的にレコード単位の疑似データを作成・提供できない  
⇒ 現行の制度に合わせた形でSynthetic Dataを作成するための工夫が必要
- 本研究では、事前に集計・推定した統計表とモデルの結果を基に、Synthetic Dataの方法に基づいて疑似データを作成する方法について検討  
⇒ どのような集計結果、モデルの結果を公表すれば、秘匿性を確保した上で、元データの構造を保持した疑似的なデータが作成できるかを検討

1. 公的統計マイクロデータの概要

2. 疑似データの重要性と課題

3. Synthetic Dataに基づく方法

 **4. 実データに基づく試作**

5. まとめと今後の課題

# 分析に用いたデータ

---

## ➤ 商用データ(平成24年2月分)

- ある県のデータを使用
- サイズ：7,558レコード
- 連続変数：
  - ・ 従業員数 (対数変換)
  - ・ 売上高 (対数変換)
- 離散変数：
  - ・ 産業 (建設業、製造業、小売業、それ以外)
  - ・ 地域 (3地域)
  - ・ 経営組織 (株式会社、有限会社)
  - ・ 開設年 (~1984年, 1985年~1994年, 1995年~)
  - ・ 資本金 (300万円以上500万円未満, 500万円以上1000万円未満, 1000万円以上2000万円未満, 2000万円以上)

# データ作成の流れ

## ● 周辺分布（地域(3)×産業(4)）の集計

	建設業	製造業	卸・小売業	左記以外	合計
地域1	1065	748	645	512	2970
地域2	935	508	831	827	3101
地域3	389	357	466	275	1487
合計	2389	1613	1942	1614	7558

- 上記の周辺分布に合うような7,558レコードを生成
- 元データから事前に推定したモデルで変数を逐次的に推測
  - (1) 地域・産業から経営組織を推測【2項ロジット】
  - (2) 地域・産業・経営組織から開設年を推測【順序ロジット】
  - (3) 地域・産業・経営組織・開設年から資本金を推測【順序ロジット】
  - (4) 地域・産業・経営組織・開設年・資本金から従業者数を推測【重回帰】
  - (5) 地域・産業・経営組織・開設年・資本金・従業者数から売上高を推測【重回帰】

※(4)・(5)：重回帰の残差から中央値(MD)、中央絶対偏差(MAD)を算出

⇒  $N(MD, MAD^2)$  に従う正規乱数を付与（0で切断）

# モデルの推定結果

## ➤ 売上高を推測するモデルの推定結果

	Estimate	Std. Error	z value	Pr(> z )
定数項	2.948865	0.029566	99.737	< 2e-16
地域：				
地域 1 【ベースライン】				
地域 2	0.105774	0.019574	5.404	6.73E-08
地域 3	0.007828	0.024043	0.326	0.745
産業：				
建設業 【ベースライン】				
製造業	-0.103277	0.024815	-4.162	3.19E-05
小売業	0.683146	0.023473	29.104	< 2e-16
その他	-0.147661	0.024546	-6.016	1.87E-09
経営組織：				
株式会社 【ベースライン】				
有限会社	-0.152575	0.027627	-5.523	3.45E-08
開設年：				
～1984年 【ベースライン】				
1985年～1994年	0.085169	0.017506	4.865	1.17E-06
1995年～	0.01336	0.015361	0.87	0.384
資本金：				
300万円以上500万円未満 【ベースライン】				
500万円以上1000万円未満	0.350539	0.029154	12.024	< 2e-16
1000万円以上2000万円未満	0.14328	0.018607	7.7	1.53E-14
2000万円以上	0.083222	0.01993	4.176	3.00E-05
従業員数（対数）	0.912558	0.009206	99.129	< 2e-16

中央値(MD) : -0.02701386

中央絶対偏差(MAD) : 0.8777895

# 試行データの作成結果

## ➤ 経営組織

【レコード数】	株式会社	有限会社	合計
元データ	4381	3177	7558
疑似データ	4416	3142	7558

【構成比】	株式会社	有限会社	合計
元データ	0.58	0.42	1.00
疑似データ	0.58	0.42	1.00

## ➤ 開設年

【レコード数】	～1984年	1985年 ～1994年	1995年～	合計
元データ	3018	2656	1884	7558
疑似データ	3092	2563	1903	7558

【構成比】	～1984年	1985年 ～1994年	1995年～	合計
元データ	0.40	0.35	0.25	1.00
疑似データ	0.41	0.34	0.25	1.00

# 試行データの作成結果

## ➤ 資本金

【レコード数】	300万円 ～1000万円	500万円 ～1000万円	1000万円 ～2000万円	2000万円～	合計
元データ	1896	1170	2864	1628	7558
疑似データ	1873	1169	2920	1596	7558

【構成比】	300万円 ～1000万円	500万円 ～1000万円	1000万円 ～2000万円	2000万円～	合計
元データ	0.25	0.15	0.38	0.22	1.00
疑似データ	0.25	0.15	0.39	0.21	1.00

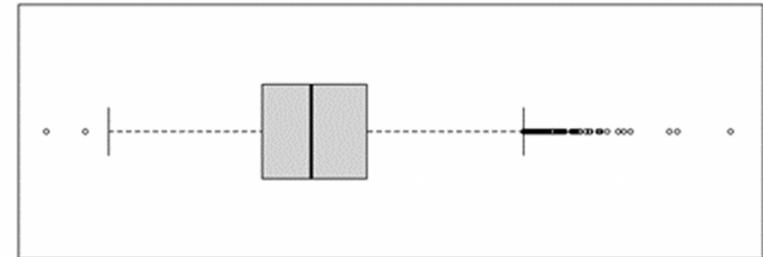
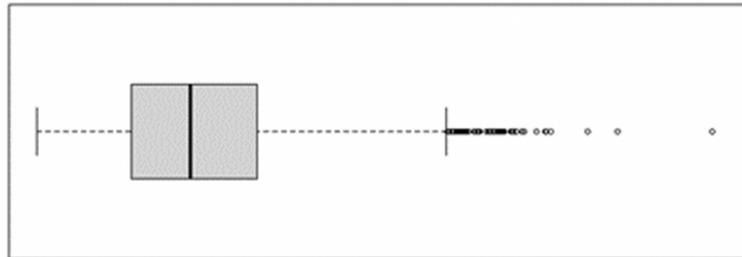
# 試行データの作成結果

## ➤ 連続変数の分布

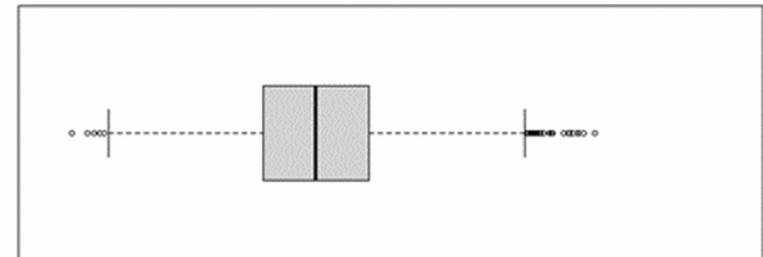
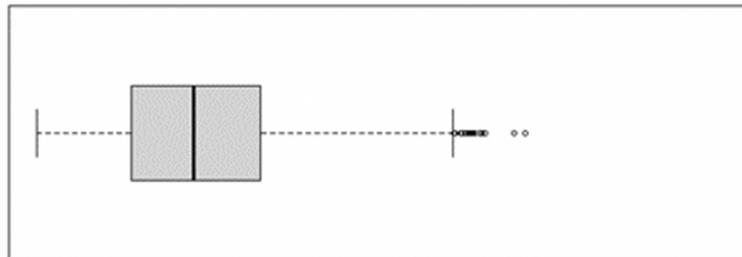
従業員数

売上高

元データ



疑似  
データ



# 試行データの作成結果

## ➤ 相関係数

元データ	(1)地域	(2)産業	(3)経営組織	(4)開設年	(5)資本金	(6)従業員数	(7)売上高
(1)地域		0.09	-0.03	-0.04	0.01	-0.01	0.02
(2)産業			-0.02	-0.01	-0.04	-0.11	-0.04
(3)経営組織				0.27	-0.71	-0.40	-0.42
(4)開設年					-0.42	-0.25	-0.23
(5)資本金						0.48	0.49
(6)従業員数							0.79
(7)売上高							

疑似データ	(1)地域	(2)産業	(3)経営組織	(4)開設年	(5)資本金	(6)従業員数	(7)売上高
(1)地域		0.09	-0.03	-0.04	0.02	0.01	0.03
(2)産業			-0.03	0.01	-0.05	-0.12	-0.05
(3)経営組織				0.27	-0.72	-0.43	-0.44
(4)開設年					-0.41	-0.28	-0.25
(5)資本金						0.51	0.53
(6)従業員数							0.80
(7)売上高							

# 試行データの作成結果

## 重回帰モデル：売上高～地域＋産業＋経営組織＋開設年＋従業員数

	元データ				疑似データ			
	Estimate	Std. Error	z value	Pr(> z )	Estimate	Std. Error	z value	Pr(> z )
定数項	2.948865	0.029566	99.737	< 2e-16	2.945073	0.027631	106.584	< 2e-16
地域： 地域1【ベースライン】								
地域2	0.105774	0.019574	5.404	6.73E-08	0.124272	0.017619	7.053	1.90E-12
地域3	0.007828	0.024043	0.326	0.745	-0.014865	0.021664	-0.686	0.493
産業： 建設業【ベースライン】								
製造業	-0.103277	0.024815	-4.162	3.19E-05	-0.126383	0.022369	-5.65	1.66E-08
小売業	0.683146	0.023473	29.104	< 2e-16	0.673735	0.02099	32.098	< 2e-16
その他	-0.147661	0.024546	-6.016	1.87E-09	-0.166643	0.022165	-7.518	6.19E-14
経営組織：株式会社【ベースライン】								
有限会社	-0.152575	0.027627	-5.523	3.45E-08	-0.140898	0.023638	-5.961	2.63E-09
開設年：～1984年【ベースライン】								
1985年～1994年	0.085169	0.017506	4.865	1.17E-06	0.092696	0.015629	5.931	3.14E-09
1995年～	0.01336	0.015361	0.87	0.384	0.008486	0.013555	0.626	0.531
資本金：300万円以上500万円未満【ベースライン】								
500万円以上1000万円未満	0.350539	0.029154	12.024	< 2e-16	0.372735	0.02591	14.386	< 2e-16
1000万円以上2000万円未満	0.14328	0.018607	7.7	1.53E-14	0.161312	0.016809	9.597	< 2e-16
2000万円以上	0.083222	0.01993	4.176	3.00E-05	0.079515	0.016976	4.684	2.86E-06
従業員数（対数）	0.912558	0.009206	99.129	< 2e-16	0.908968	0.0092	98.798	< 2e-16

# 変数の順序を変更した場合の影響

- 先に連続変数を推測するよう、変数の順序を変更

- (1) 地域・産業から従業者数を推測【重回帰】
- (2) 地域・産業・従業者数から売上高を推測【重回帰】
- (3) 地域・産業・従業者数・売上高から経営組織を推測【2項ロジット】
- (4) 地域・産業・従業者数・売上高・経営組織から開設年を推測  
【順序ロジット】
- (5) 地域・産業・従業者数・売上高・経営組織・開設年から  
資本金を推測【順序ロジット】

※(1)・(2)：重回帰の残差から中央値(MD)、中央絶対偏差(MAD)を算出  
⇒  $N(MD, MAD^2)$  に従う正規乱数を付与 (0で切断)

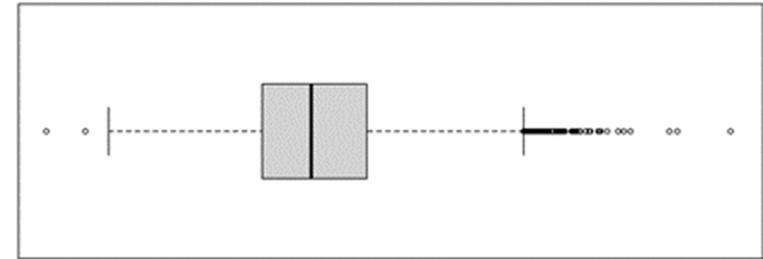
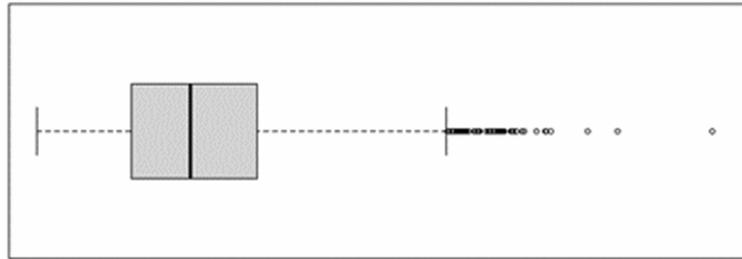
- 離散変数の分布、相関係数、売上高を予測する重回帰モデルについては、変数の順序変更の影響は、それほど大きくないことを確認

# 変数の順序を変更した場合の影響

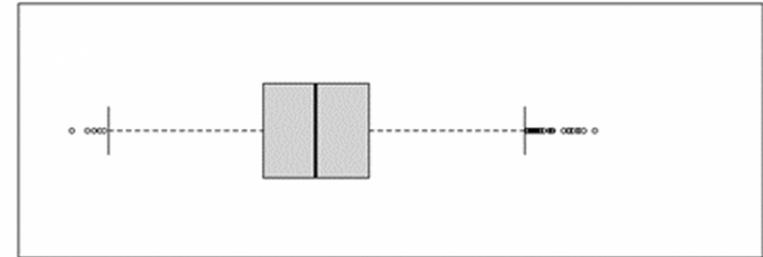
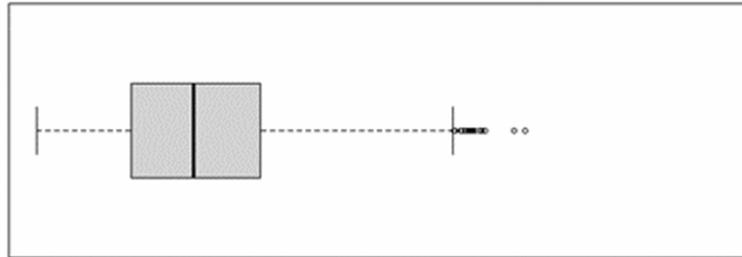
従業員数

売上高

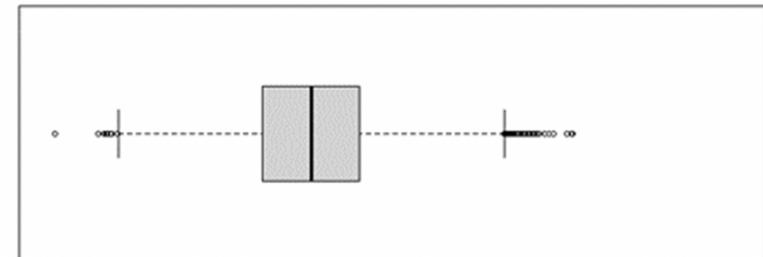
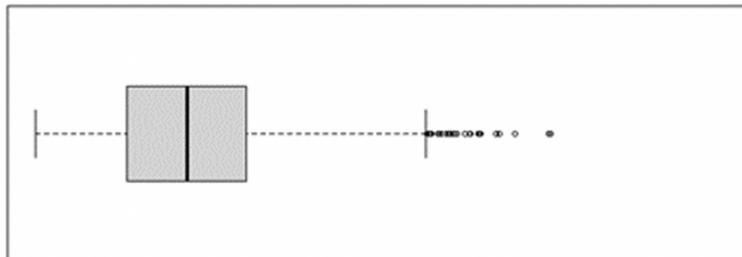
元データ



疑似  
データ



疑似  
データ  
(順序変更)



1. 公的統計マイクロデータの概要
2. 疑似データの重要性と課題
3. Synthetic Dataに基づく方法
4. 実データに基づく試作

## 5. まとめと今後の課題

# まとめと今後の課題

---

- 商用データを用いた試作では、周辺分布の統計表、回帰モデルの結果、残差の分布に関する情報を事前に作成（公開）することにより、回帰モデルの構造をある程度保持したデータの作成ができた
- ただし、一部の変数では相関係数の係数の符号が逆になった
- 極端に大きな値の発生方法（発生させる必要があるか？、トップコーディング）
- 資本金など、特定の値に集中している変数では、ロジットモデル・回帰モデルの2段階推定が必要

# まとめと今後の課題

- 周辺分布に合わせるための操作（比推定、繰返し比例補正など）の検討⇒必要な集計表の公開
  - 変数の順序の影響を、より精緻に分析
  - モデルの推測能力の向上、モデル推定の繰返し
  - ロジットモデル・回帰モデルの2段階推定  
（資本金など、特定の値に集中している変数）
  - 他の変数変換(Box-Coxなど)の検討⇒情報公開
  - 元データに一致あるいは非常に近いデータが生じた場合の対応（データの削除、ノイズ付与）
  - 差分プライバシーの考慮（適切なノイズの付与）
- ⇒世帯データを含む公的統計ミクロデータに対象を拡大して、さらに分析・検討を行っていく予定