

匿名データを用いた統計教育 事例と直交変換による合成 データ生成法の紹介

東京情報大学

佐野夏樹

概要

1. 住宅・土地統計調査の匿名データを利用した基礎分析教育の紹介
 - 変数の種類と尺度
 - 変数間の関連性の尺度
 - クラメールの連関係数が計算出来ないケース
2. 全国消費実態調査の匿名データを利用した研究紹介
 - 直交変換による合成データ生成法
 - 全国消費実態調査の匿名データによるリスク評価と有用性評価

変数の種類と尺度

質的変数と量的変数

- 質的変数：離散的な値をとる変数, カテゴリカル変数とも呼ばれる.
- 量的変数：連続的な値をとる変数

質・量	尺度の種類	特徴
質的	名義尺度	性別や職業の様にカテゴリーの違いだけを表す
質的	順序尺度	優・良・可・不可のように順序に意味があるが, カテゴリー間の差は同じではない
量的	間隔尺度	温度の様に, 順序も間隔も意味があるが, 原点に意味はない
量的	比率尺度	長さや重さの様に, 間隔尺度であり, さらに原点に意味がある.

平成25年住宅・土地統計調査匿名データ における変数の例（1）

名義尺度の例：住宅以外の建物の種類4区分

会社等の寮	学校等の寮	旅館・宿泊所	その他の建物	対象外
90	51	56	461	340971

順序尺度の例：子の居住地6区分

一緒に住んでいる	徒歩5分程度の場所	片道15分未満の場所	片道1時間未満の場所	片道1時間以上の場所	子はいない	不詳	対象外
123630	7049	12576	23315	27818	53742	45183	48316

平成25年住宅・土地統計調査匿名データ における変数の例（2）

間隔尺度の例：建物の階数（うち一戸建・長屋建3区分）

1階建	2階建	3階建以上	対象外
41273	165291	6629	128436

比率尺度の例：住宅の延べ面積

回答者	不詳	対象外
286801	6379	48449

平均: 102.1868 m²

問題：変数の種類と尺度

平成25年住宅・土地統計調査匿名データの変数を

(1) 名義尺度 (2) 順序尺度 (3) 間隔尺度 (4) 比率尺度
のいずれかに分類せよ.

全変数の数 (157変数)

- 政府統計コード, 一連番号等のメタ変数関係 (8変数)
- 名義尺度 (118変数)
- 順序尺度 (13変数)
- 間隔尺度 (3変数)
- 比率尺度 (15変数)

代表的な関連性の尺度（1）

量的変数↔量的変数 ピアソン相関係数

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$-1 \leq r \leq 1$, \bar{x} : x の平均, \bar{y} : y の平均

質的変数↔質的変数 クラメールの連関係数

$$V = \sqrt{\frac{\chi_0^2}{n\{\min(a, b) - 1\}}}$$

$0 \leq V \leq 1$

χ_0^2 : カイ2乗値
 n : サンプル数
 a : 分割表の行数
 b : 分割表の列数

代表的な関連性の尺度 (2)

量的変数 ↔ 質的変数 相関比

質的変数 x	量的変数 y	平均
x_1	$y_{11}, y_{12}, \dots, y_{1n_1}$	\bar{y}_1
x_2	$y_{21}, y_{22}, \dots, y_{2n_2}$	\bar{y}_2
x_3	$y_{31}, y_{32}, \dots, y_{3n_3}$	\bar{y}_3
x_4	$y_{41}, y_{42}, \dots, y_{4n_4}$	\bar{y}_4

\bar{y} : 総平均

総平方和 S_T : $\sum_{i=1}^a \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2$

級間平方和 S_A : $\sum_{i=1}^a n_i (\bar{y}_{i\cdot} - \bar{y})^2$

相関比 : $\eta^2 = \frac{S_A}{S_T}$
 $0 \leq \eta^2 \leq 1$

平成25年住宅・土地統計調査匿名データ に対するRを使った計算例

相関係数

1ヶ月当たり家賃・間代（量的）と1ヶ月当たり共益費・管理費（量的）

```
> cor(yachin,kyoeki,use="complete.obs")  
[1] 0.215262
```

クラメールの連関係数

建物の階数_うち共同住宅・その他（9区分）（質的）とエレベーターの有無（質的）

```
> library(vcd)  
> assocstats(table(kaisu,elevator))$cramer  
[1] 0.8628889
```

相関比

1ヶ月当たり家賃・間代（量的）と建築の時期1-4区分（質的）

```
> summary(lm(rent.fee~as.factor(year)))$r.squared  
[1] 0.08834161
```

問題：クラメール連関係数

以下の3つのパターンの行列に対してクラメール連関係数を計算出来ない理由を計算式とデータの側面から考察し、対策を考えよ。

(A)

	V11				
V9	1	2	3	4	5
	1	0	0	0	0
	2	0	0	0	0
	3	0	0	0	0
	4	0	0	0	0

(B)

	V39
V10	1
	125609
1	0
2	0
V	0

(C)

	V11				
V10	1	2	3	4	5
	108102	106753	100192	24847	1077
1	0	0	0	0	0
2	0	0	0	0	0
V	0	0	0	0	0

パターン (A)

$$\chi^2 = \sum_{i=1}^l \sum_{j=1}^k \frac{(f_{ij} - \hat{f}_{ij})^2}{\hat{f}_{ij}}$$

f_{ij} : (i,j)セルの観測度数

\hat{f}_{ij} : (i,j)セルの期待度数

l : 分割表の行数

k : 分割表の列数

	V11				
V9	1	2	3	4	5
1	0	0	0	0	0
2	0	0	0	0	0
3	0	0	0	0	0
4	0	0	0	0	0

計算出来ない理由

特定の行(列)が全て0だと、その行(列)確率が0になり、期待度数が0となるセルが生じるため。



欠損値NAも含めた分割表

	V11					
V9	1	2	3	4	5	<NA>
1	0	0	0	0	0	90
2	0	0	0	0	0	51
3	0	0	0	0	0	56
4	0	0	0	0	0	461
<NA>	108102	106753	100192	24847	1077	0

2つの質問項目の対象が異なり、対象外がNAとして処理されていると、このパターンになる。

V9 : 住宅以外の建物の種類4区分 (対象 : 住宅以外の建物に居住する世帯)

V11 : 建物の構造5区分 (対象 : 住宅)

パターン (B)

$$V = \sqrt{\frac{\chi_0^2}{n\{\min(a, b) - 1\}}}$$

$0 \leq V \leq 1$

	V39
V10	1
	125609
1	0
2	0
V	0

計算出来ない理由

分割表の行（列）数が1だと
クラメール連関係数の平方
根中の分母が0となるため

欠損値NAも含めた分割表

	V39	
V10	1	<NA>
	125609	215887
1	0	60
2	0	70
V	0	3

「無し」に該当するカテゴリがなく、「無し」と
対象外が区別されず、NAとして処理されている
場合、このパターンになる。

V10：住宅以外の建物の所有の関係2区分

V39：高齢者等のための設備状況_手すりの有無

パターン (C)

V11						
V10	1	2	3	4	5	
	108102	106753	100192	24847	1077	
1	0	0	0	0	0	
2	0	0	0	0	0	
V	0	0	0	0	0	



欠損値NAも含めた分割表

V11							
V10	1	2	3	4	5	<NA>	
	108102	106753	100192	24847	1077	525	
1	0	0	0	0	0	60	
2	0	0	0	0	0	70	
V	0	0	0	0	0	3	

計算出来ない理由

特定の行（列）が全て0だと、その行（列）確率が0になり、期待度数が0となるセルが生じるため。

2つの質問項目の対象が異なり、片方の質問は、対象外をNAとして処理し、もう片方の質問項目は、対象外を空白として処理している場合、このパターンになる。

V10：住宅以外の建物の所有の関係2区分
(対象：住宅以外の建物に居住する世帯)

V11：建物の構造5区分（対象：住宅）

基礎分析教育のまとめ

- 住宅・土地統計調査の匿名データを用いて、名義尺度、順位尺度、間隔尺度、比率尺度のデータとは、具体的にどのようなデータか、教育
- 変数の尺度によって、基礎分析においても、異なる分析手法を適用する必要があることを教育
- 単に、分析パッケージ利用するだけでなく、NaN等の計算不可の出力結果から、分析指標の意味やデータの性質について理解することの重要性を教育



直交変換による合成データ生成法

統計的開示制御の方法

分類	手法	概要
非攪乱手法	グローバルリコーディング	いくつかのカテゴリを結合し, 新しいカテゴリーを作成する.
	局所抑制	特定に結びつく危険な組み合わせは, 開示せず, 抑制する.
攪乱手法	ノイズ追加	連続値のデータにノイズを付加し, 異なるデータに攪乱する.
	PRAM	カテゴリカルデータを確率的に攪乱する.
	マイクロアグリゲーション	連続値のデータをクラスタリングし, 同じグループに所属するクラスタは, クラスタの平均値で置き換える.
	データシャッフリング	連続値のデータに対して, 原データの周辺分布と同じ周辺分布になることを保証しながら, 異なるサンプル間のデータを取り替える.
合成データの生成	多重代入法	原データを多重代入法によって補間したデータで置き換える.

合成データの生成法

1. 完全合成データ：原データは無く、合成データのみが公開される。
2. 部分的合成データ：機密性の高い変数やサンプル（レコード）のみ合成データで置き換えられ、公開される。
3. ハイブリッドデータ：原データと合成データを組み合わせて生成されるデータ。組み合わせる割合によって、原データ寄り、もしくは合成データ寄りのデータになる。

主成分分析

直交変換による多変量データの情報縮約手法

$$\mathbf{C} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T = \sum_{i=1}^p \lambda_i \mathbf{u}_i \mathbf{u}_i^T$$

\mathbf{C} : 共分散行列 ($p \times p$)

$\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$: 対角行列

$\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p)$: 直交行列

主成分得点：新しい基底上の座標

$f_i = \mathbf{u}_i^T \mathbf{z}, i = 1, \dots, p$: 第 i 主成分

$\mathbf{z} = \mathbf{x} - \bar{\mathbf{x}}$: 中心化サンプル

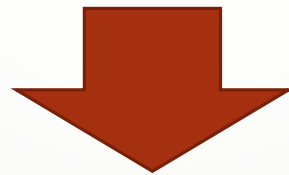
\mathbf{x} : 原データのベクトル

$\bar{\mathbf{x}}$: 平均ベクトル

直交変換 (OT) による合成データ生成法

q 次元空間 ($q < m$) への射影

$$\mathbf{f}^{(q)} = (\mathbf{u}_1^T \mathbf{z}, \mathbf{u}_2^T \mathbf{z}, \dots, \mathbf{u}_q^T \mathbf{z})^T = \mathbf{U}_{(q)}^T \mathbf{z}, \quad \mathbf{U}_{(q)} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_q)$$



原データの空間での再構成

$$\hat{\mathbf{x}}^{(q)} = \bar{\mathbf{x}} + \mathbf{U}_{(q)} \mathbf{f}^{(q)} = \bar{\mathbf{x}} + \mathbf{U}_{(q)} \mathbf{U}_{(q)}^T \mathbf{z}$$

提案手法による合成データ の情報損失指標

$$ILO T = \frac{\sum_{i=q+1}^p \lambda_i^R}{p}$$

λ_i^R : 相関係数行列 R の第 i 固有値

$$R = \sum_{i=1}^p \lambda_i^R \mathbf{u}_i^R \mathbf{u}_i^{R^T}$$

Cf. 回帰による合成データ生成

i 番目の変数を目的変数, 残りの変数を説明変数とした回帰モデル

$$\mathbf{x}_i = \mathbf{X}^{-i} \boldsymbol{\beta}^{-i} + \boldsymbol{\varepsilon}_i, i = 1, \dots, p, \boldsymbol{\varepsilon}_i \sim N(0, \sigma_i^2)$$

$$\mathbf{X}^{-i} = (\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_p)$$

撓乱項無し of データ生成 (DRI)

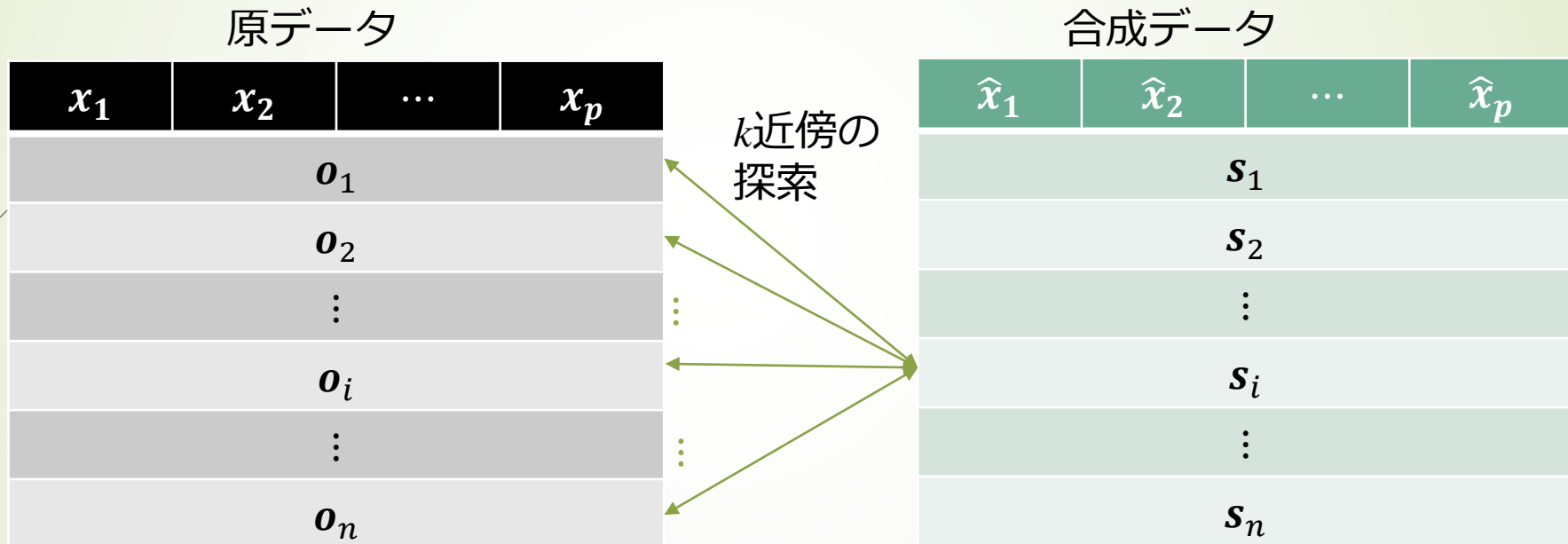
$$\hat{\mathbf{x}}_i = \mathbf{X}^{-i} \hat{\boldsymbol{\beta}}^{-i}, i = 1, \dots, p$$

撓乱項付き of データ生成 (SRI)

$$\hat{\mathbf{x}}_i = \mathbf{X}^{-i} \hat{\boldsymbol{\beta}}^{-i} + \mathbf{e}_i, i = 1, \dots, p, \mathbf{e}_i \sim N(0, \hat{\sigma}_i^2)$$

レコードリンケージによるリスク評価

原データの*i*番目のレコード o_i から生成した合成データを s_i とする



侵入者が s_i に対して、原データの中から k 近傍 $knn(s_i)$ を探した時に、 o_i が含まれていれば、特定されたと考え、特定されたレコード割合をリスク評価値とする。

$$Risk(k) = \frac{\sum_{i=1}^n I[o_i \in knn(s_i)]}{n}$$

有用性評価指標

平均絶対誤差により, 有用性 (情報損失) を評価する

$$\text{MAE} = \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p |o_{ij} - s_{ij}|$$

利用データ

- ▶ 平成16年全国消費実態調査の匿名データ
- ▶ 2人世帯データ
- ▶ サンプル数：43861
- ▶ 利用する項目（変数）：現金で購入された食料品82項目.
- ▶ ただし, 左の肉類, 生鮮肉, 加工肉の様に, 他の項目を集計した項目は, 冗長なので除く

削除する集計項目の例

符合	項目
1.3	肉類
1.3.1	生鮮肉
220	牛肉
221	豚肉
222	鶏肉
22X	合いびき肉
224	他の生鮮肉
1.3.2	加工肉
225	ハム・ソーセージ
229	他の加工肉

全国消費実態調査の匿名データを用いた 有用性とリスクの評価結果

1. 有用性評価

手法	DRI	SRI	OT(0.01)	OT(0.05)	OT(0.10)	OT(0.20)
MAE	585.114	904.207	64.791	182.304	274.295	382.404

OTの括弧内の数値はILOTの値

2. リスク評価

手法	DRI	SRI	OT(0.01)	OT(0.05)	OT(0.10)	OT(0.20)
$k=1$	0.000	0.000	0.996	0.930	0.694	0.271
$k=3$	0.001	0.000	0.997	0.949	0.775	0.377
$k=5$	0.002	0.000	0.997	0.955	0.801	0.425

結論

- 直交変換による合成データ生成法および、生成データの有用性（情報損失）評価指標を提案.
- 提案手法のリスクおよび有用性を全国消費実態調査の匿名データを利用して評価.
- 提案手法は、有用性評価指標に対応した主成分数によって、リスクおよび有用性をコントロールすることができ、回帰による合成データ生成よりも、高い有用性を持つデータを生成することもできる.



ご静聴ありがとうございました