

証拠に基づく政策立案のための 潜在プロファイル分析と数法則発見法を 用いた社会実情のモデル化と可視化

Modeling and Visualization of Social Reality
Using Latent Profile Analysis and Number Law Discovery Methods
for Evidence-Based Policy Making

長瀬 永遠 (Towa Nagase)
u255013@st.pu-toyama.ac.jp

富山県立大学大学院 工学研究科 電子・情報工学専攻
情報基盤工学講座

N212, 09:30-10:00 Tuesday, February 13, 2024.

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

情報技術の発達により、社会における様々なデータを観測・収集することが可能に

→ 政策分野においても、証拠に基づく政策立案（Evidence Based Policy Making: EBPM）に注目が集まる

EBPM とは

政策における意思決定をデータに基づいて行うという考え方

課題

大規模なデータから有用な情報を取り出すデータマイニングの技術が必要

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

どんな情報を取り出すことが有効？

政策分野では、原因と結果の間に成り立つ関係性が重要

→ 複数の要因が複雑に影響しあうため、人間が把握するには限界がある

アプローチ

行政が持つ統計データを用いて分析を行い、データ間の関係性を数理モデルによって表す

数理モデルの例

$$(\text{データ A}) = 2.0 \cdot (\text{データ B}) + 1.0 \cdot (\text{データ C}) - 1.0 \cdot (\text{データ D})$$

地域経済分析システム (RESAS)

経済に関する項目を中心に自治体単位でデータを公開するオープンデータサイト

Table 1: RESAS から自動取得可能なデータ

データ項目	
農地平均取引価格	農業産出額
林地平均取引価格	海面漁獲物等販売金額
住宅用地平均取引価格	林産物販売金額
商業用地平均取引価格	林作業請負収入
マンション等平均取引価格	企業数
一人あたりの固定資産税	事業所数
一人当たりの地方税	就業者数
一人当たりの法人住民税	総人口
製造品出荷額	経営耕地面積
年間商品販売額	etc

自治体区分	総数
市	792
町	743
村	183
特別区	23
合計	1741

疑問

すべての市区町村を同じ尺度のサンプルとして扱うのは適切？

はじめに

統計データの特徴
と研究の概要

潜在的クラスター
ングと数法則発見

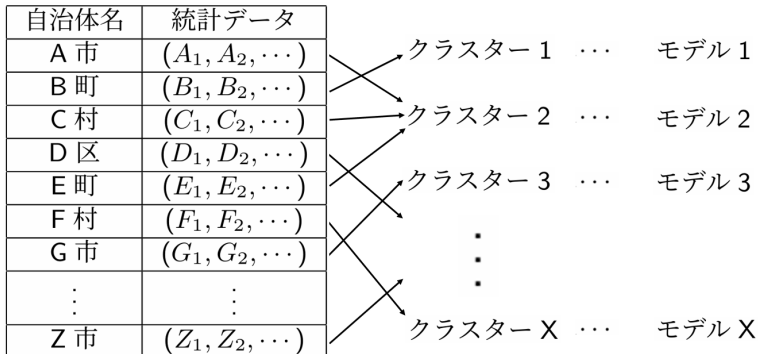
提案手法

数値実験並びに
考察

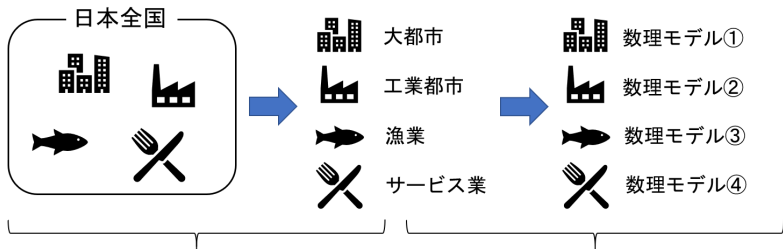
おわりに

目的

行政が持つ統計データを用いて自治体をクラスタリングし、その結果を考慮しながら統計データ間の関係性をモデル化する手法を提案する。



使用する手法の概要



潜在プロファイル分析 (LPA)

- データの潜在的な特徴を考慮
- ソフトクラスタリング

RF6.4法 (Rule extraction method from Fact 6.4)

- 表現力が高く複雑なデータに対応
- 式の可読性も比較的高い

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

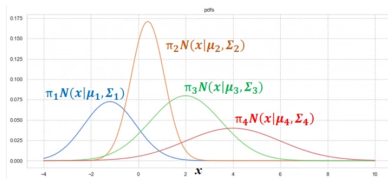
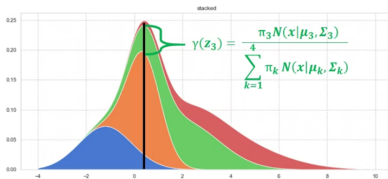
LPA

観測されたデータを混合分布モデルとして扱い、モデルの背後に存在する潜在変数を分析。その結果を考慮しながらデータをクラスタリングする手法である。

(混合ガウス分布)

$$N(\mathbf{x}_i | \pi, \mu, \Sigma) = \sum_{k=1}^K \pi_k N(\mathbf{x}_i | \mu_k, \Sigma_{rk}) \quad (1)$$

K : 潜在変数の数, π_k : 母集団における潜在変数の構成割合,
 $\mathbf{x}_i = (x_1, \dots, x_r)_i \quad i = 1, \dots, I$: 観測変数,
 $\mu_k = (\mu_{1k}, \dots, \mu_{rk})$: 観測変数の平均, Σ_{rk} : 各観測変数の共分散行列,



分かってない情報

- 各分布のパラメータ : $\pi_k, \mu_k, \Sigma_{rk}$
- 潜在変数 z_k の事後分布の期待値 = 各サンプルの存在確率 : $E_{z_{ik}}[z_{ik}]$
- クラスタ数 = 潜在変数の数 : K

(EM アルゴリズム)

■ E ステップ

その時点で最適と考えられるパラメータを式 (2) に代入することによって $E_{z_{ik}}[z_{ik}]$ の値を求める.

$$E_{z_{ik}}[z_{ik}] = \frac{\pi_k N_k(\mathbf{x}_i | \mu_k, \Sigma_{rk})}{\sum_k \pi_k N_k(\mathbf{x}_i | \mu_k, \Sigma_{rk})} \quad (2)$$

■ M ステップ

E ステップで求めた $E_{z_{ik}}[z_{ik}]$ を式 (3) の右辺に代入し, 式 (3) が最大となるようなパラメータを求める. 求めたパラメータをその時点で最適なパラメータとし, 再度 E ステップに戻る.

$$E_z[\ln p(\mathbf{x}, z | \pi, \mu, \sigma)] = \sum_i \sum_k E_{z_{ik}}[z_{ik}] (\ln \pi_k + \ln N_k(\mathbf{x}_i | \mu_k, \Sigma_{rk})) \quad (3)$$

はじめに

統計データの特徴
と研究の概要

潜在的クラスター
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

モデル選択

EM アルゴリズムでは、クラスター数 K を仮置きしているので、ベイズ情報量基準 (Bayesian information criterion: BIC) が最小となるクラスター数を最適とする。

$$BIC = -2L + K \ln n \quad (4)$$

L : M ステップで最大化した値, K : クラスター数, n : サンプル数

各サンプルにおける存在確率の算出

最適なモデルが決定した後、各サンプルにおける観測変数ベクトル \mathbf{y}_i を用いて存在確率を算出。

$$\pi_{k|\mathbf{x}_i} = \frac{\pi_k N_k(\mathbf{x}_i | \mu_k, \Sigma_{rk})}{\sum_k \pi_k N_k(\mathbf{x}_i | \mu_k, \Sigma_{rk})} \quad (5)$$

各変数各変数に対するデータ項目

目的変数 y : 自治体の財源

量的説明変数 x_1, x_2 : 自主財源, 依存財源

質的説明変数 q_1, q_2, q_3, q_4 : 大都市, 工業都市, 漁業, サービス業



大都市

$$if \ q_1 \quad y = v_0 + 5.5x_1^{1.5}x_2^{-0.5} + 0.1x_1^0x_2^{2.0}$$



工業都市

$$if \ q_2 \quad y = v_0 + 4.0x_1^{1.5}x_2^{-0.5} + 0.5x_1^0x_2^{2.0}$$



漁業

$$if \ q_3 \quad y = v_0 + 0.5x_1^{1.5}x_2^{-0.5} + 3.5x_1^0x_2^{2.0}$$



サービス業

$$if \ q_4 \quad y = v_0 + 0.7x_1^{1.5}x_2^{-0.5} + 2.0x_1^0x_2^{2.0}$$



RF6.4 法

目的変数 y と M 個の量的説明変数 x_m , K 個の質的説明変数 q_{kl} からなるデータセットを与え, 4 層パーセプトロンの学習を用いてパラメータ $c_{0d}, c_{gd}, c_{dkl}, v_{gm}$ を最適化することで質的条件付き多変量多項式を発見する.

(質的条件付き多変量多項式)

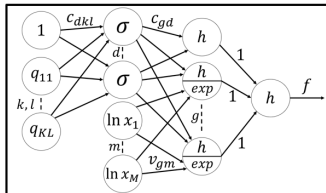
$$\text{if } \bigwedge_k \bigvee_{q_{kl} \in Q_k^i} q_{kl} \text{ then } y = v_0 + \sum_{g=1}^G v_g^i \prod_{m=1}^M x_m^{v_{gm}^i} + \epsilon, \quad i = 1, \dots, I \quad (6)$$

RF6.4 法の 4 層パーセプトロン

$$f(\mathbf{q}, \mathbf{x}; \phi) = v_0 + \sum_{g=1}^G v_g s_g,$$

$$v_0 = \sum_{d=1}^D c_{0d} \sigma_d, \quad v_g = \sum_{d=1}^D c_{gd} \sigma_d,$$

$$\sigma_d = \sigma \left(\sum_{k=1}^K \sum_{l=1}^{L_k} c_{dkl} q_{kl} \right), \quad s_g = \exp \left(\sum_{m=1}^M v_{gm} \ln x_m \right)$$



分かってない情報

- モデルのパラメータ : $\phi = (c_0, c_{gd}, c_{dkl}, v_{gm})$
- パーセプトロンの中間層の数 : G, D

最終的に求めたい数理モデル

$$f(\mathbf{q}, \mathbf{x}; \phi) = w_0 + \sum_{g=1}^G w_g s_g,$$

$$w_0 = \sum_{d=1}^D c_{0d} \sigma_d, \quad w_g = \sum_{d=1}^D c_{gd} \sigma_d, \quad (7)$$

$$\sigma_d = \sigma \left(\sum_{k=1}^K \sum_{l=1}^{L_k} c_{dkl} q_{kl} \right), \quad s_g = \exp \left(\sum_{m=1}^M v_{gm} \ln x_m \right)$$

下記の目的関数を BFGS 法によって最小化することでパラメータを求める

目的関数

$$E(\mathbf{q}, \mathbf{x}; \phi) = \frac{1}{2} \sum_{n=1}^N (f(\mathbf{q}, \mathbf{x}; \phi)^n - y^n)^2 \quad (8)$$

はじめに

統計データの特徴
と研究の概要潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

パーセプトロンの最適中間層数

パーセプトロンの学習では、中間層の数 G, D を仮置きしているの、BIC が最小となるときのモデルを採用。本研究では G, D ともに 1~5 まで変化させて学習した。

$$BIC(G, D) = \frac{N}{2} \ln \left(\frac{1}{N} \sum_{n=1}^N (f(\mathbf{q}^n, f(\mathbf{x}^n; \hat{\phi}_{G,D}) - y^n)^2 \right) + \frac{Z}{2} \ln N \quad (9)$$

N : サンプル数, Z : パラメータ数

ルールの復元

パーセプトロンの学習結果では、各サンプルに対するパラメータベクトルが得られるが、より一般的な数法則として記述するために、以下の処理を行う。

- 1 各サンプルに対する数法則の係数値ベクトル $\mathbf{c}^\mu = (c_0^\mu, \dots, c_J^\mu)$ を求める
- 2 k-means 法を用いて I 個のクラスターを求め、それぞれの重心ベクトル $\mathbf{a}^I = (a_0^I, \dots, c_J^I)$ を数法則の係数値とする

提案手法

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

自治体名	統計データ
A 市	(A_1, A_2, \dots, Y_A)
B 町	(B_1, B_2, \dots, Y_B)
C 区	(C_1, C_2, \dots, Y_C)
\vdots	\vdots
Z 村	(Z_1, Z_2, \dots, Y_Z)



所属確率 1	所属確率 2	...	所属確率 X
0.9	0.1	...	0
0.3	0	...	0.1
0.5	0.5	...	0
\vdots	\vdots	\vdots	\vdots
0	0	...	1.0



$$\tilde{y} = \frac{y - \text{mean}(y)}{\text{std}(y)} \quad \downarrow \quad \tilde{x} = \frac{x}{\max(x)}$$

目的変数	量的説明変数
\tilde{Y}_A	$(\tilde{A}_1, \tilde{A}_2, \dots)$
\tilde{Y}_B	$(\tilde{B}_1, \tilde{B}_2, \dots)$
\tilde{Y}_C	$(\tilde{C}_1, \tilde{C}_2, \dots)$
\vdots	\vdots
\tilde{Y}_Z	$(\tilde{Z}_1, \tilde{Z}_2, \dots)$



質的説明変数 1	質的説明変数 2	...	質的説明変数 X
[1, 0, ..., 0]	[0, 0, ..., 1]	...	[0, 0, ..., 1]
[0, 0, ..., 0]	[0, 0, ..., 1]	...	[0, 0, ..., 1]
[0, 1, ..., 0]	[0, 1, ..., 0]	...	[0, 0, ..., 1]
\vdots	\vdots	\vdots	\vdots
[0, 0, ..., 1]	[0, 0, ..., 1]	...	[1, 0, ..., 0]

RF6.4法

$$Y_1 = v_{1,0} + \sum_{g=1}^G v_{1-g} \prod_{m=1}^M x_m^{v_{gm}} \quad Y_2 = v_{2,0} + \sum_{g=1}^G v_{2-g} \prod_{m=1}^M x_m^{v_{gm}} \quad \dots \quad Y_X = v_{X,0} + \sum_{g=1}^G v_{X-g} \prod_{m=1}^M x_m^{v_{gm}}$$

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

各自治体が所属するクラスターと対応する数理モデルを視覚的に提示するため、GIS を用いたシステムを作成

GIS とは

地理情報システム（Geographic Information System）の略称. 地図上にデータを表示し、可視化や重ね合わせなどを行うことが可能

GIS を用いる狙い

- 分析結果全体を俯瞰的に把握
- クラスターの分布と地理的特徴の関係を可視化

動画

提案システムのながれを動画でお見せします.

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

使用したデータ

項目：11 種類の量的データ

サンプル数：データ欠損のない 650 の自治体におけるオープンデータ

対象年：2020 年

参照元：地域経済分析システム-RESAS-

Table 2: 数値実験に用いたデータ

データ項目	単位	データ項目	単位
1 人あたりの固定資産税	千円/人	総人口	人
1 人あたりの地方税	千円/人	住宅用地平均取引価格	円 / m^2
1 人当たりの法人住民税	千円/人	商業用地平均取引価格	円 / m^2
経営耕地面積	畝 / 経営体	農地平均取引価格	円 / m^2
製造品出荷額	万円	林地平均取引価格	円 / m^2
年間商品販売額	百万円		

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

目的

以下の 3 つの事柄について、検証するためそれぞれ数値実験を行った

- LPA は自治体の特徴をどのように捉えているか
- 自治体のクラスタリング結果を考慮することが数理モデルの精度向上に有効か
- クラスタリング手法として LPA を用いることは適切か

実験 1

総務省から発行されている市町村類型の内容と本研究で得られたクラスタリング結果を比較

実験 2

クラスタリングを用いない手法と提案手法それぞれにおいて得られた数理モデルの精度を比較

実験 3

クラスタリングに k-means を用いた手法と提案手法それぞれにおいて得られた数理モデルの精度を比較

市町村類型とは

「人口」と「産業構造」の二軸から市町村を分類したもの。

類型の数

- 政令指定都市，特別区，中核市，特例市：それぞれ 1 類型
- 一般市：16 類型
- 町村：15 類型

Table 3: 市の類型

産業構造		Ⅱ次,Ⅲ次90%以上		Ⅱ次,Ⅲ次90%未満	
		Ⅲ次65%以上		Ⅲ次55%以上	
		Ⅲ次65%以上	Ⅲ次65%未満	Ⅲ次55%以上	Ⅲ次55%未満
人口	0以上～ 50,000未満	I -3	I -2	I -1	I -0
	50,000以上～ 100,000未満	Ⅱ -3	Ⅱ -2	Ⅱ -1	Ⅱ -0
	100,000以上～ 150,000未満	Ⅲ -3	Ⅲ -2	Ⅲ -1	Ⅲ -0
	150,000以上～	Ⅳ -3	Ⅳ -2	Ⅳ -1	Ⅳ -0

Table 4: 町村の類型

産業構造		Ⅱ次,Ⅲ次80%以上		Ⅱ次、Ⅲ次80%未満
		Ⅲ次60%以上	Ⅲ次60%未満	
人口	0以上～ 5,000未満	I -2	I -1	I -0
	5,000以上～ 10,000未満	Ⅱ -2	Ⅱ -1	Ⅱ -0
	10,000以上～ 15,000未満	Ⅲ -2	Ⅲ -1	Ⅲ -0
	15,000以上～ 20,000未満	Ⅳ -2	Ⅳ -1	Ⅳ -0
	20,000以上～	V -2	V -1	V -0

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

分析対象

年代：2020 年
 クラスター数：15 クラスター
 範囲：九州地方
 サンプル数：112 の自治体

分析方法

15 個のクラスターそれぞれに所属する自治体の類型を調査し、クラスターごとの内訳を分析

Table 5: 数値実験に用いたデータ

データ項目	単位	データ項目	単位
1 人あたりの固定資産税	千円/人	総人口	人
1 人あたりの地方税	千円/人	住宅用地平均取引価格	円 / m^2
1 人当たりの法人住民税	千円/人	商業用地平均取引価格	円 / m^2
経営耕地面積	畝 / 経営体	農地平均取引価格	円 / m^2
製造品出荷額	万円	林地平均取引価格	円 / m^2
年間商品販売額	百万円		

実験 1 の結果と考察

21/26

Figure 1: クラスター 1 における市の
の類型割合

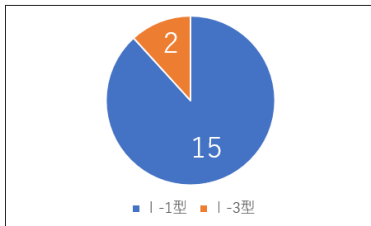
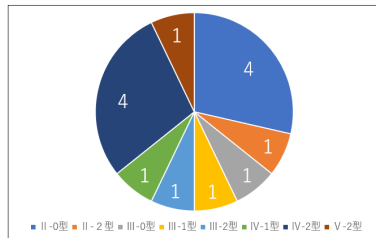


Figure 2: クラスター 1 における町村の
の類型割合



結果

- 市においては全体的に人口が近い自治体が同じクラスターとなった
- 町村は比較的そうでない場合が多かった

考察

- クラスタリングに用いたデータが産業面で不十分であった
- 町村には人口とは別の潜在的共通要素が存在する

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

検証内容

自治体のクラスタリング結果を考慮することが数理モデルの精度向上に有効か

比較対象

- クラスタリングを用いない手法 (RF5 法)
- 提案手法 (LPA + RF6.4 法)

データ

データ : 11 項目のオープンデータ

サンプル : データ欠損のない 650 の自治体

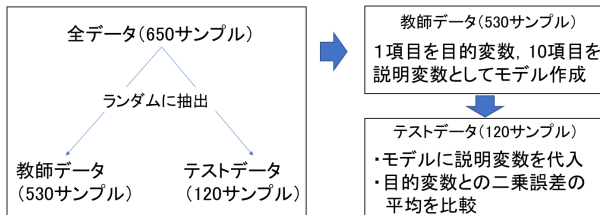
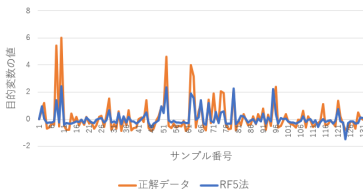


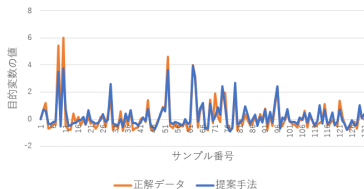
Table 6: 最小の BIC と二乗誤差の平均値

手法	BIC	二乗誤差
RF5 法	-234.7011	0.5543
提案手法	-254.8814	0.2262

RF5法と正解データ



提案手法と正解データ



考察

- BIC, 二乗誤差ともに提案手法が優れていることから自治体のクラスターを考慮することは数理モデルの精度向上に有効である
- 提案手法においても正解データに対して誤差が大きいため, 改善の余地が残される

検証内容

クラスタリング手法として LPA を用いることは適切か

比較対象

- クラスタリングに k-means を用いた手法 (k-means + RF6.4 法)
- 提案手法 (LPA + RF6.4 法)

データ

データ : 11 項目のオープンデータ

サンプル : データ欠損のない 650 の自治体

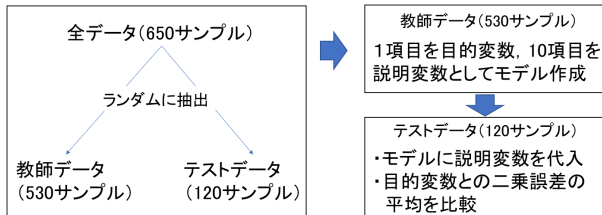


Table 7: 最小の BIC と二乗誤差の平均値

手法	BIC	二乗誤差
k-means+RF6.4 法	-244.8024	0.3929
提案手法	-254.8814	0.2262

はじめに

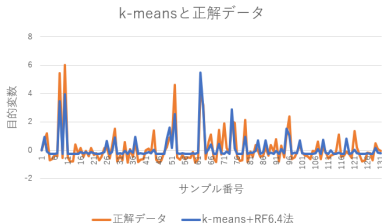
統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに



考察

- BIC, 二乗誤差ともに提案手法が優れていることからクラスタリング手法に LPA を用いることは適切である
- グラフから k-means は小さい値への適合に難がある可能性が示唆された

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

まとめ

行政が持つ統計データを用いて、自治体をクラスタリングし、データの関係を数理モデルで表す手法を提案した。提案手法の結果に対して、クラスタリングの傾向とモデルの精度に関する検証を行った。

- 既存の市町村分類と比較して、人口の観点では概ね合致したが、産業構造では異なる点が見られた
- クラスタリングを用いない場合と比較してモデルの精度が向上した
-
- モデルの精度自体には改善の余地がある

今後の展望

- 項目数、サンプル数ともにより大規模なデータを用いた議論が必要
- 専門的な知見をもって検証と考察することが必要
- データの関係性が既知な対象への適用も検討すべき

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

以下, 差分

類団の表

26/26

はじめに

統計データの特徴
と研究の概要

潜在的クラスター
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

1s- I 1	15
総数 s- I 3	2
31m- II 0	4
m- II 2	1
m- III 0	1
m- III 1	1
m- III 2	1
m- IV 1	1
m- IV 2	4
m- V 2	1

2s- III 3	1
総数 中核市	1
2	

3大都市	1
総数	
1	

5s- III 1	3
総数 s- III 3	2
6s- IV 1	1

3大都市	1
総数	
1	

6s- I 0	1
総数 s- I 1	6

14s- I 3	3
m- II 2	1
m- IV 0	1
m- IV 1	1
m- IV 2	1

7s- I 0	2
総数 s- I 1	4
23s- I 2	1

s- I 3	1
s- II 1	8
s- II 2	3
s- II 3	2
m- V 2	2

8s- I -1	1
総数 s- I 2	1

10s- I 3	1
s- II 2	1
s- II 3	4
m- V 2	2

9s- II 3	1
総数 m- II 2	1
3m- V 2	1

10s- II 1	1
総数 s- II 2	1
5s- II 3	1
m- IV 2	1
m- V 1	1

11政令指定都市	1
総数 大都市	1
3中核市	1

12中核市	1
総数	
1	

はじめに

統計データの特徴
と研究の概要

潜在的クラスタリ
ングと数法則発見

提案手法

数値実験並びに
考察

おわりに

$$y = 1.0504 + 0.0915h_1 - 0.4932h_2 - 1.5345h_3$$

$$y = 1.5630 - 0.2676h_1 + 0.4047h_2 + 0.2488h_3$$

$$y = 2.6819 - 0.1679h_1 - 0.1266h_2 - 1.3990h_3$$

$$y = 3.7813 - 0.6816h_1 + 1.0757h_2 + 0.8165h_3$$

$$y = 0.9802 - 0.1457h_1 + 0.1915h_2 + 0.0174h_3$$

$$h_1 = x_1^{-0.3317} x_2^{0.8776} x_3^{0.4363} x_4^{0.2339} x_5^{0.8738} x_6^{0.5663} x_7^{-0.1736} x_8^{0.6950} x_9^{1.1073} x_{10}^{0.2303}$$

$$h_2 = x_1^{-0.0301} x_2^{-1.3162} x_3^{0.4784} x_4^{0.1623} x_5^{0.0419} x_6^{-0.2382} x_7^{-0.1257} x_8^{-0.0275} x_9^{-0.0340} x_{10}^{-0.0340}$$

$$h_3 = x_1^{0.3044} x_2^{-0.0357} x_3^{0.1758} x_4^{0.5694} x_5^{1.2188} x_6^{1.2825} x_7^{-0.5933} x_8^{-0.7382} x_9^{0.3629} x_{10}^{0.4825}$$

(10)

(目的変数)

住宅地取引平均価格

(説明変数)

固定資産税，地方税，住民税，経営耕地面，製造品出荷額，商業地取引平均価格，年間商品販売額，農地取引平均価格，林地取引平均価格，総人口

3.1 因果探索によるデータ間の関係性

26/26

因果探索とは、観測データを用いて、データ群の因果グラフ（複数の観測データにおいて、それぞれの値がお互いに及ぼしあっている影響の度合いを構造的に示したもの）を導出するための教師なし学習である。

LiNGAM

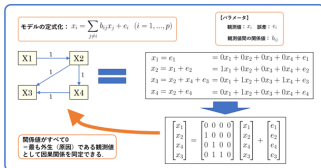


図15 LiNGAMのアルゴリズム

LiNGAMの制約

1. 内生変数と外生変数をつなぐ関数は線形。
2. 外生変数の分布は非ガウス連続分布。
3. 因果グラフは非巡回
4. 外線変数同士は互いに独立。

内生・・・観測済み 外生・・・未観測

Direct-LiNGAM

Direct-LiNGAMのアプローチ

・・・回帰分析を用いる手法

- ・ 内生変数群から2変数を取り出しそれらの変数間に成り立つ因果関係を同定することを繰り返して因果グラフの始点を探索。
- ・ その変数を内生変数群から除外し、残った変数のみで内生変数群を再形成。

適当な2変数

$$\begin{cases} x_1 = e_1 \\ x_2 = b_{21}x_1 + e_2 \end{cases}$$

目的変数を x_1 、説明変数を x_2 としたときの回帰残差(r_2).

$$r_2 = \left\{ 1 - \frac{b_{21} \text{cov}(x_1, x_2)}{\text{var}(x_2)} \right\} e_1 - \frac{b_{21} \text{var}(x_1)}{\text{var}(x_2)} e_2$$

ダルモア・スキットビッチの定理

2つの確率変数 y_1, y_2 が互いに独立な確率変数 $s_i (i = 1, \dots, q)$ を用いて下記のように表されるとき、 y_1, y_2 が独立なら、 $\alpha_j, \beta_j \neq 0$ となるような変数 s_j はガウス分布に従う。

$$y_1 = \sum_{i=1}^q \alpha_i s_i \quad y_2 = \sum_{i=1}^q \beta_i s_i$$

3.2 DEA による効率値と入力・出力改善値の導出

26/26

DEA とは、ある分野における組織の集合において、対象の組織の業績を評価するために生み出されたノンパラメトリックなアプローチである。組織とは、その活動においていくつかの種類の入力（投入）をいくつかの出力（産出）に変換することに携わる生産体（Decision Making Unit: DMU）を指す。

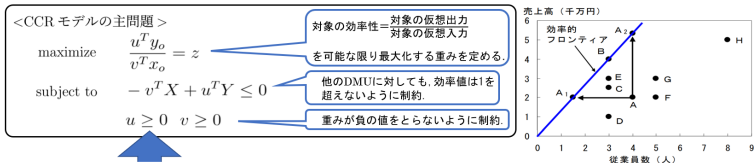


図16 1入力・1出力のDMU群

双対問題

<入力指向モデル>

$$\begin{aligned} &\text{minimize} && w = \theta \\ &\text{subject to} && Y\lambda \geq y_o \\ &&& -X\lambda + x_o\theta \geq 0 \\ &&& \lambda \geq 0 \end{aligned}$$

<出力指向モデル>

$$\begin{aligned} &\text{maximize} && w = \eta \\ &\text{subject to} && X\mu \leq x_o \\ &&& -Y\mu + y_o\eta \leq 0 \\ &&& \mu \geq 0 \end{aligned}$$

<入力改善案>

$$\hat{x}_i = \sum_{k=1}^K x_{ik}\lambda_k \quad i = 1, 2, \dots, m$$

<出力改善案>

$$\hat{y}_j = \sum_{k=1}^K y_{jk}\mu_k \quad j = 1, 2, \dots, n$$

それぞれのDMUに対して各入力をどれだけ減少、各出力をどれだけ増加させれば評価値が1になるかが算出できる。また、その際に参考としたDMUもわかる。

LPA

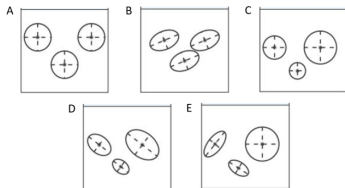
観測されたデータを混合分布モデルとして扱い、モデルの背後に存在する潜在変数を分析。その結果を考慮しながらデータをクラスタリングする方法である。

(混合ガウス分布)

$$N(\mathbf{x}_i | \pi, \mu, \Sigma) = \sum_{k=1}^K \pi_k N(\mathbf{x}_i | \mu_k, \Sigma_{rk}) \quad (11)$$

K : 潜在変数の数, π_k : 母集団における潜在変数の構成割合,
 $\mathbf{x}_i = (x_1, \dots, x_r)_i \quad i = 1, \dots, I$: 観測変数,
 $\mu_k = (\mu_{1k}, \dots, \mu_{rk})$: 観測変数の平均, Σ_{rk} : 各観測変数の共分散行列,

$$\begin{array}{ccc} \text{A} \begin{bmatrix} \sigma_1^2 & & \\ 0 & \sigma_2^2 & \\ \vdots & \vdots & \ddots \\ 0 & 0 & \dots & \sigma_r^2 \end{bmatrix} & \text{B} \begin{bmatrix} \sigma_{21}^2 & & \\ \sigma_{21} & \sigma_{22}^2 & \\ \vdots & \vdots & \ddots \\ \sigma_{r1} & \sigma_{r2} & \dots & \sigma_r^2 \end{bmatrix} & \text{C} \begin{bmatrix} \sigma_{1k}^2 & & \\ 0 & \sigma_{2k}^2 & \\ \vdots & \vdots & \ddots \\ 0 & 0 & \dots & \sigma_{rk}^2 \end{bmatrix} \\ \text{D} \begin{bmatrix} \sigma_{1k}^2 & & \\ \sigma_{21} & \sigma_{2k}^2 & \\ \vdots & \vdots & \ddots \\ \sigma_{r1} & \sigma_{r2} & \dots & \sigma_{rk}^2 \end{bmatrix} & \text{E} \begin{bmatrix} \sigma_{1k}^2 & & \\ \sigma_{21k} & \sigma_{2k}^2 & \\ \vdots & \vdots & \ddots \\ \sigma_{r1k} & \sigma_{r2k} & \dots & \sigma_{rk}^2 \end{bmatrix} & \end{array}$$



パラメータ推定

観測変数 \mathbf{y} と潜在変数 z の同時確率の尤度を目的関数とし、その最大値を求める。値の最大化には EM アルゴリズムを用いる。

(EM アルゴリズム)

■ E ステップ

その時点で最適と考えられるパラメータを式 (2) に代入することによって $E_{z_{ik}}[z_{ik}]$ の値を求める。

$$E_{z_{ik}}[z_{ik}] = \frac{\pi_k f_k(\mathbf{y}_i | \mu_k, \Sigma_{rk})}{\sum_k \pi_k f_k(\mathbf{y}_i | \mu_k, \Sigma_{rk})} \quad (12)$$

■ M ステップ

E ステップで求めた $E_{z_{ik}}[z_{ik}]$ を式 (3) の右辺に代入し、式 (3) が最大となるようなパラメータを求める。求めたパラメータをその時点で最適なパラメータとし、再度 E ステップに戻る。

$$E_z[\ln p(\mathbf{y}, z | \pi, \mu, \sigma)] = \sum_i \sum_k E_{z_{ik}}[z_{ik}] (\ln \pi_k + \ln f_k(\mathbf{y}_i | \mu_k, \Sigma_{rk})) \quad (13)$$

パーセプトロンの学習方法

RF6.4 法の学習には、準ニュートン法を基本枠組みとする学習アルゴリズムであるBPQ (Back propagation based on Partial Quasi-Newton) 法を用いる。

(最小化する関数)

$$E(\mathbf{q}, \mathbf{x}; \phi) = \frac{1}{2} \sum_{n=1}^N (f(\mathbf{q}, \mathbf{x}; \phi)^n - y^n)^2 \quad (14)$$

Step 1: パラメータの初期化

ϕ_1 を初期化. $\mathbf{H}_1 = \mathbf{I}, b = 1$ とおく. \mathbf{I} は ϕ と次元の等しい単位行列

Step 2: 探索方向 $\Delta\phi_s$ の計算

探索方向 $\Delta\phi_s = -\mathbf{H}_s \nabla E(\mathbf{q}, \mathbf{x}; \phi_s)$ を計算. $\nabla E(\mathbf{q}, \mathbf{x}; \phi_s)$ の全要素が 0.0001 未満の場合, 反復を終了.

Step 3: 最適探索幅 α_s の計算

$E(\mathbf{x}; \phi_b + \alpha_b \Delta\phi_b)$ を最小にする最適探索幅 α_b を計算.

Step 4: 結合重み ϕ_s の更新

$\phi_{b+1} = \phi_b + \alpha_s \Delta\phi_b$ に従って結合重み ϕ_b を更新.

Step 5: 二次微分の逆行列の近似値 \mathbf{H} の更新

$b \equiv 0 \pmod{Z}$ (Z : 全パラメータ数) のとき, $\mathbf{H}_{b+1} = \mathbf{I}$ とし, それ以外のとき, BFGS 公式で \mathbf{H}_{b+1} を更新. $b \leftarrow b + 1$ として, Step2 の戻る.