

Median Polish 法を用いたデータ解析の特性に関する考察

2013SE089 小林優

指導教員：松田真一

1 はじめに

Median Polish 法はジョン・チューキー氏によって提唱された方法であり，多因子の様々な要因を重要性を調べるためのデータの分析方法とされている．二元配置分散分析と Median Polish 法の二つの方法を用いて，はずれ値を含むデータの比較を行い，Median Polish 法の性質について明らかにしていく．

2 事例研究について

研究事例をいくつか行ったがここでは紙面の都合上1つだけ紹介する．Web[1]のデータを用いて，47都道府県の県別人口と県庁所在地人口の都道府県人口に対する割合から都道府県の平均地価を求めた．ここでは県庁所在地の人口÷都道府県の人口を県庁所在地人口率という言葉で定義する．47都道府県別人口は600万人以上，600万人-200万人，200万人以下の大，中，小で区分けし，県庁所在地人口率は40%以上，40%-20%，20%以下の大，中，小で区分けした．そして図1のように各区分に入る9都道府県を取り出した．表1は縦は人口，横は県庁所在地人口率に対し，それぞれの都道府県の平均地価をあらわしている．

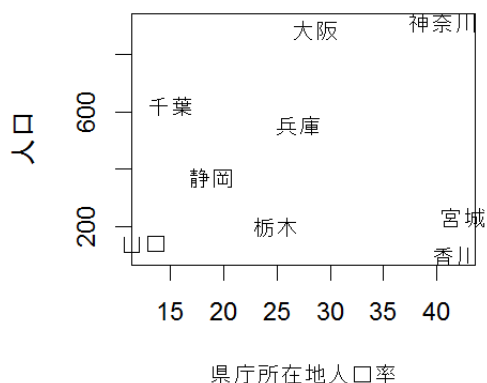


図1 人口と県庁所在地人口率の関係

表1 平均地価

	大	中	小
大	23.91	23.60	10.71
中	8.15	13.50	8.76
小	4.77	4.15	3.47

3 二元配置分散分析について

二元配置分散分析とは，二つの因子 A と B が存在し個々のデータのもっている情報を縮約して1つの値にした統計値への影響を調査する手法である．ここでは繰り返しのないデータを使用する．一般的な手法であるため一般式などは省略する．事例研究の平均地価を二元配置分散分析を用いて得られた残差が表2である．

表2 二元配置分散分析を用いて得られた残差

	大	中	小
大	3.4511	1.6678	-5.1189
中	-3.0389	0.8378	2.2011
小	-0.4122	-2.5056	2.9178

これから県庁所在地人口率が小，人口が大である千葉県が少し外れ値になっていることがわかる．

4 Median Polish 法について

Median Polish 法は竹内 [2] に基づき事例研究のデータで説明する．R では medpolish 関数で計算が可能である．

表3 Median Polish 法の計算手順その1

23.6	0.31	0.00	-12.89
8.76	-0.61	4.74	0.00
4.15	0.62	0.00	-0.68
	0.31	0.00	-0.68
0.00	0.00	0.00	-12.21
0.68	-0.92	4.74	0.68
0.00	0.31	0.00	0.00

表4 Median Polish 法の計算手順その2

0.00	0.00	0.00	-12.21
0.68	-1.6	4.06	0.00
0.00	0.31	0.00	0.00
	0.00	0.00	0.00
0.00	0.00	0.00	-12.21
0.00	-1.6	4.06	0.00
0.00	0.31	0.00	0.00

1. 表1の各行の中央値（行中央値）を求め，表3左上に記入する．

2. 行ごとの各データから行中央値を引いた残差を求め、表3右上に記入する。
3. 手順2で求めた残差について列ごとの中央値(列効果)を求め、表3右中央に記入する。
4. 手順3で求めた残差から列効果を引き、表3の右下に記入する。
5. 手順4で求めた残差について行ごとの中央値(行効果)を求め、表3の左下に記入する。このとき表3の左下のすべての数値が0にならない場合、表3の右下の数値を元のデータとして用いて、この手順1~手順5すべての数値が0になるまで繰り返す。
(ここでは0になっていないので繰り返す。)
6. 表3と表4の列効果成分同士、行効果成分同士を足して新たな列効果と行効果を得る。このとき新たに得た行効果の中央値を総中央値とする。
7. 行効果から総中央値を引いて新たな行効果を得る。
この事例研究のデータでは総中央値が9.44、行効果は上から14.16, 0, -5.29, 列効果は0.31, 0, -0.68となる。

事例研究の平均地価をMedian Polish法を用いて得られた結果、二元配置分散分析以上に千葉県が大幅に外れ値を取ることが分かった。

5 シミュレーション

事例研究からMedian Polish法が外れ値を含むデータに対して有効的な手段ではないかと考察されたが、実際にどのくらいの大きさの外れ値に有効な手段であるかをモンテカルロ法でシミュレーションを行った。

6 プログラムについて

小, 中, 大の効果(今回は0, 1, 2で固定)をA,Bの要素として、それぞれに与え、3行3列の9分割したデータを作成し、平均0, 標準偏差0.1の正規乱数を誤差として加えた。その行列にさらに外れ値を乗せて、二元配置分散分析とMedian Polish法を同時に100回ずつ行い、その平均を求めるプログラムを製作した。外れ値を与える位置は3行3列の行列の角である(1,1)成分, 辺である(1,2)成分, 真ん中である(2,2)成分の3種類に分けた。外れ値は0.1から2.0まで0.1ずつ増やした。

7 考察

(1,1)成分, (1,2)成分, (2,2)成分の3つのパターンともMedian Polish法のB中が0をとり、基準をあわせるために二元配置分散分析のA中が0となるように基準化した。その結果の内、(1,1)成分と(2,2)成分について、横軸に外れ値の大きさ、縦軸に推定された効果の大きさを示したグラフが図2,3である。

二元配置分散分析は線形的にプロットされていくパターンか0をプロットし続けるパターン、Median Polish法は全体的に疎らにプロットされていくパターンが多いことがわ

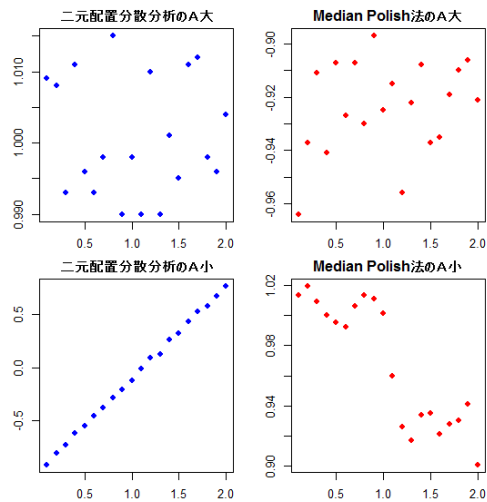


図2 (1,1) 成分の結果

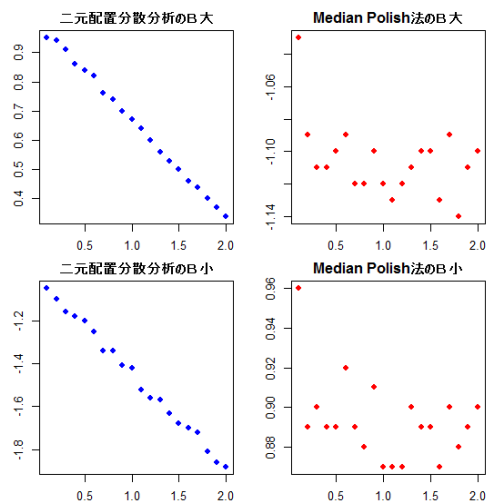


図3 (2,2) 成分の結果

かった。しかし例外も存在する。図2の二元配置分散分析A大は、まだらにプロットされている。これはA中が0となるように基準化し、A大にも影響が及んだためだと考察された。また、図2のMedian Polish法A小も例外であり、x軸の1付近で点の動き方が変化している。これはMedian Polish法は中央値を用いており、差が1でA中と中央値が入れ替わるためだと考察された。

8 おわりに

シミュレーションを行った事で、例外も存在するが基本的には外れ値を含むデータに対しMedian Polish法は疎らに分散される事がわかった。

参考文献

- [1] 観光庁:『統計情報・白書』, <http://www.mlit.go.jp/kankocho/siryou/toukei/ranking.html>, 2014.(2016/6 閲覧)
- [2] 竹内啓:『統計学辞典』, 東洋経済新報社, 1990.