

Stataによる多重代入

株式会社 ライトストーン 2020年8月 修正2020年10月





- 第2回では欠損値が存在するときの対応方法の一般論について解説した
- 今回はStataを使って欠損値がモノトーンパターンの場合と、 モノトーンパターンに従わないときの、代表的な対応方法に ついて説明する





- ◆ 欠損メカニズムがモノトーンパターンの時
 - ◆ mi impute monotone:多変量正規分布とは無関係。モノトーンパターンに従う場合にだけ利用する代入コマンド
- ◆ データ拡大アルゴリズム(Data Augmentation)
 - ◆ mi impute mvn:欠損メカニズムがモノトーンパターンでない時に利用するベイズ統計による代入コマンド。変数の正規性を利用する
- キーワード: データ拡大アルゴリズム、Average RVI、Largest FMI、MCMCの収束、WLF



1.欠損メカニムズがモノトーンパターンの時

- 欠損値の発生メカニズムがモノトーンパターンの時は比較的、簡単な計算方法で多重代入が実行できる
- もう少し具体的に言えば、複数の変数に欠損値が存在しても、一変数ごとにモデルを利用した計算を繰り返していけばよい
- 今、欠損のある変数Xがp個あり、欠損値を含まない完全データの変数Zが存在するものとする

*モノトーンパターンに関する解説は第2回にもあり





次に示す方法で理論値x*を求める(Rubin 1987)

$$x_1^* \sim f_1(x_1|z)$$
 $x_2^* \sim f_2(x_2|x_1^*, z)$
...
 $x_p^* \sim f_p(x_p|x_1^*, x_2^*, \dots, x_{p-1}^*, z)$

- f()は欠損値の代入に利用する関数
- ポイントは欠損の多い変数の代入値の計算に、すでに計算した代入値を利用するところ





操作:第1回と同じサンプルデータを利用して、モノトーンパターンの確認と、多重代入、そしてロジスティックモデルの推定を行う

- . use https://www.stata-press.com/data/r16/mheart5s0, clear
- . mi des
- . misstable nested

1. age(12) -> bmi(28)

ageとbmiに欠損値が存在し、ageの12個の欠損値はbmiの28個の欠損値にネストしている





bmiには回帰モデルを利用し、ageにはpmm(predictive mean matching)を利用した代入をおこなう。pmmの詳細は第4回で解説する

.mi impute monotone (reg) bmi (pmm,knn(3)) age = attack smokes hsgrad female ,add(10)

Conditional models:

age: pmm age_attack smokes hsgrad female , knn(3)
bmi: regress bmi age attack smokes hsgrad female

ageの12個の欠損値に代入値を作成し、それを利用してbmiの代入値を作成している。代入回数は10回

例題:モノトーンパターン



- . mi est,dots:logit attack smokes age bmi hsgrad female
- オプションdotsは計算プロセスを示すもの。計算に長い時間を 要する場合に利用する

attack	Coef.	Std. Err.	t	P> t	[95% Conf.	Interval]
smokes	1.168286	.357026	3.27	0.001	.4685075	1.868065
age	.0279243	.0172083	1.62	0.106	0060082	.0618569
bmi	.1039346	.0513625	2.02	0.045	.0023066	.2055627
hsgrad	.1415342	.40294	0.35	0.725	6482223	.9312907
female	0824037	.4128752	-0.20	0.842	8916261	.7268186
_cons	-4.991369	1.743229	-2.86	0.005	-8.425476	-1.557262

- オッズ比を表示する場合はeformオプションを次の位置に入力する
- . mi estimate, eform: logit attack smokes age bmi hsgrad female



- 欠損値の割合が高いほど、推定 値の分散は大きくなる
- RVIIはrelative variance increase。
 この値が小さいほど、欠損値が分散に与える影響が小さいことを示す
- したがって、もし、完全データで同じ推定を実行すれば、RVIはゼロになる

```
Imputations
                           10
Number of obs
                          154
Average RVI
                       0.1139
Largest FMI
                       0.2762
DF:
       min
                       128.19
                 = 113,489.58
       avg
                 = 530,022.84
       max
F( 5, 2710.4)
                         3.09
Prob > F
                       0.0087
```





- MI推定値の変動(分散)は代入した1セットのデータ内変動と、複数の代入セット間の変動により構成される
- つまり、サンプルサイズと代入回数に左右される
- サンプルサイズが変更できない場合は代入回数を増やして、Average RVIを小さくする





- Largest Fraction Missing Information (欠損している情報割合の最大値)
- 次の計算式で代入回数を設定する

$$M \ge 100 \times \text{FMI}$$

 いま、0.2762×100=27.62なので、代入回数は28回 以上に変更するべき





- Stataは推定コマンドによって、小/大標本の仮定を自動的に使い分ける
- DFはRVIの逆数に関連する統計量なので、RVI が小さい時は自動的に大きくなる
- 係数に関するF検定は、すべての係数に関する欠 損情報の割合は一定であるという仮定を、便宜 的においている





- regressコマンドの場合は小標本を前提とする
- regressコマンドを実行すると、他の推定コマンドの場合とは異なり、欠損値が存在しない場合の自由度 Complete DFを追加表示する
- これは完全にデータが揃っている場合の、理想的な状態 の自由度である
- データが完全に揃っていても自由度が小さい時に注意する





- 欠損値の発生メカニズムがモノトーンパターンに 従っていない場合はどう対応する?
- 第2回ではEMアルゴリズムを紹介したが、解析的 に条件付き期待値を求める必要があった
- 一般的にはベイズ統計を利用したデータ拡大アル ゴリムを利用することが多い
- ベイズ統計では事前分布と尤度、そしてMCMCを 利用して事後分布の最頻値を求める
- Tanner and Wong (1987)の提案したベイズ統計に よるデータ拡大アルゴリズムは、EMアルゴリズム と似たテクニックを利用して変数の分布を求める

データ拡大アルゴリズム



- Tanner and Wong (1987)
- 完全データなxは観測データyと欠損したzにより構成されるものと考える

$$x^T = (y^T, z^T)$$

• 完全データxによるパラメータ θ の事後分布を $\pi(\theta|x)$ とすると、それは観測データyを使って、

$$\pi(\theta|y) = \int \pi(\theta|x) f(z|y) dz = \int \pi(\theta|y,z) f(z|y) dz$$

ただし、

$$f(z|y) = \int f(z|y,\theta)\pi(\theta|y)d\theta$$

モンテカルロ積分によ る近似計算





$$f(z|y) = \int f(z|y,\theta)\pi(\theta|y)d\theta$$

Step1.π(θ|y)に従う乱数θiを作成 Step2.f(z|y,θi)に従う乱数ziを作成

このようにして作成したziを利用して、 $\pi(\theta|y)$ に対する次の近似式を計算する

$$\frac{1}{m}\sum_{i=1}^{m}\pi(\theta|y,z_i)$$

*データ拡大アルゴリズムは欠損値ziを何回も作成するところがポイント





- Stataの多重代入コマンドはmi impute mvn
- データ拡大アルゴリズムはベイズ統計のシミュレーション手法あるMCMCを利用する
- 欠損値のパターンはモノトーンに限定しない。どんな欠損パターンでも利用可能
- ただし、欠損している変数は連続変数で、正規分布に従 うことを仮定する
- MCMCを利用しているので、代入実行後に診断を行う必要がある





- □ 住宅価格の決定要因を探すヘドニックモデルは単純な重 回帰モデルであるが、築年数(age)と課税(tax)にいくつ かの欠損値がある
- 2つの欠損値の発生メカニズムはMARで、モノトーンパターンにはなっていない
- □ したがって、2つの変数の分布をきるだけ正規分布に近づけてMIを実行する
- ロ サンプルデータはmhouses1993.dta

データ



- . use https://www.stata-press.com/data/r16/mhouses1993,clear
- . des
- . misstable nested
- . misstable pattern

Missing-value patterns
 (1 means complete)

	Pattern
Percent	1 2
56%	1 1 ネストしていない
35	1 0
7	0 0
2	0 1
100%	

Variables are (1) tax (2) age

正規性



- □ 2つの変数の対数をとって、正規化する
 - . gen lnage = ln(age)
 - . gen Intax = In(tax)

- □ 多重代入の設定
 - . mi set mlong
 - . mi register imputed Inage Intax
 - . mi register regular price sqft nfeatures ne custom corner
- □ 最終的に推定するヘドニックモデルではInageとIntaxは元の変数に戻すことに注意





- データ拡大アルゴリズムによる多重代入の実行
 - . mi impute mvn Inage Intax = price sqft nfeatures ne custom corner, add(20)

Performing EM optimization:

note: 8 observations omitted from EM estimation because of all imputation variables missing

observed log likelihood = 112.1464 at iteration 48

代入値の初期値を求めるEMアルゴリズムの繰り返 し計算回数





□ EMアルゴリズムの計算回数48がデフォルトの稼働検査期間(burn-in)100よりも小さいことが必要

Prior: uniform

Iterations = 2000

burn-in = 100

between = 100

	Observations per m				
Variable	Complete	Incomplete	Imputed	Total	
lnage lntax	68 107	49 10	49 10	117 117	

MCMCの収束

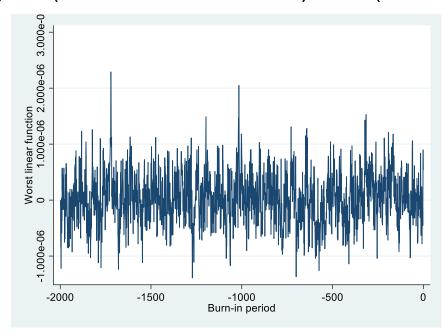


- 稼働検査期間が条件を満たしていれば、MCMCによるパラメータの分布はうまく収束できている可能性が高い
- Schafer (1997)は複数のモデルパラメータの定常状態への収束状況を判定するためのworst linear functionを提案した
 - . mi impute mvn Inage Intax = price sqft nfeatures ne /// custom corner, mcmconly burnin(2000) rseed(23) /// savewlf(wlf, replace)
 - . save mhouses1993out,replace

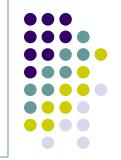
WLF



- wlfファイルに保存した関数値が定常状態になっていることを確認する
 - . use wlf, clear
 - . tsset iter
 - . tsline wlf, ytitle(Worst linear function) xtitle(Burn-in period)

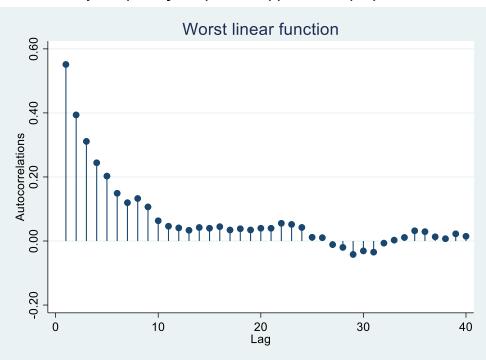






● WLF値の自己相関を確認する。自己相関の変化にトレンドがなく、短いラグでゼロに収束することが望ましい

 ac wlf, title(Worst linear function) ytitle(Autocorrelations) ciopts(astyle(none)) note("")







- データセットをもとに戻し、MI推定を実行する。対 数変換して登録した2つの変数を登録しなおす
 - . use mhouses1993out,clear
 - *MIIによるモデル推定
 - . mi passive: gen newage = exp(lnage)
 - . mi passive: gen newtax = exp(Intax)
 - . mi estimate:reg price newage newtax sqft nfeatures ne





● 任意の欠損パターンのデータで体重代入を実行した ところ、ほとんどの変数の有意性は失われた

price	Coef.	Std. Err.	t	P> t	[95% Conf.	Interval]
newage	4416874	.7900943	-0.56	0.582	-2.084179	1.200804
newtax	.6095571	.1479177	4.12	0.000	.3118694	.9072448
sqft	.3096159	.0916235	3.38	0.002	.1237405	.4954912
nfeatures	8.272231	14.22707	0.58	0.562	-20.01092	36.55538
ne	2.771427	35.54388	0.08	0.938	-67.71154	73.25439
_cons	51.71564	67.63651	0.76	0.447	-82.86647	186.2978





- 欠損値の発生メカニズムがMARで、モノトーンパターンである時と、無い時の多重代入の実行方法を紹介した
- モノトーンパターンでない時、変数には正規性を仮定した。必要に応じて正規化して、データ拡大アルゴリズム (DA)を利用する
- DAを利用したら、MCMCの収束判定を行う。問題があるときは、モデルを再考する
- WLFは複数のパラメータの収束状況を単一の尺度で判断するための統計量である





- ALLISON, P.D. (2002). Missing Data. Sage Publications
- STATA MULTIPLE IMPUTATION REFERENCE MANUAL RELEASE 16

研究者向けの統計解析ソフトウェア 512 評価版のお申込み



技術サポートの

対象となります



機能制限なし!30日間無料で使える評価版でまずはお試しください! https://www.lightstone.co.jp/stata/evaluate.html

※学生の方への評価版の提供はございません。大学の講義で使用するなど単位取得に関わる形でのご利用となる場合は学生版のご購入をご検討ください。

ライセンスは サブスクリプションがおすすめ!

常にStataの 最新バージョンが

毎年の経費として Stataを導入 できます

マルチユーザライセンス(サブスクリプション) 2ユーザ以上

サブスクリプションライセンスを選ぶメリット

利用者の増減に 柔軟に対応可能

初期導入費を 抑えたい場合にも おすすめ

シングルライセンス (サブスクリプション)

職場のPC・ノートPC・自宅のPC等、個人所有 のPCにインストール可能です。



- ボリュームライセンス

利用できます

シングルライセンスを複数人 でまとめてご購入いただく際 のボリュームディスカウント商 品です。



同時起動ライセンス

ご利用環境により下記3通りの運用方法からご選択いただけます。







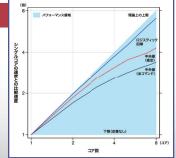
Stata/MPの演算能力

データ処理時間の削減ならMP

Stata MPはPCの持つマルチコアの特性を活かして、処理を分散・並列する機能を備え ます。およそ85%以上のコマンドで処理速度が向上し、コア数に応じた計算時間の短縮 が期待できます。

コマンドごとの処理速度の向上の度合いについては以下のページや資料をご覧ください。

コア数	全てのコマンド	推定コマンド	ロジスティック回帰
2	1.7倍	1.8倍	1.9倍
4	2.6倍	3.1倍	3.8倍
8	3.3倍	4.2倍	6.8倍





弊社Webページ内の『Stata/MP』

https://www.lightstone.co.jp/stata/statamp.html

株式会社ライトストーン

〒101-0031

東京都千代田区東神田 2-5-12 龍角散ビル 7F

TEL: 03-3864-5211 FAX: 03-3865-0050 e-mail: sales@lightstone.co.jp WEB: https://www.lightstone.co.jp/