



都市形成史研究の基礎と応用

-第二回-

2019年度スタートアップゼミ

2019/06/20

小林 里瑛 (M2)

図：1866年（慶應二年）頃の道後温泉本館周辺 道後温泉本館調査委員会編『道後温泉本館の歴史』，松山市，(1994)

今日の内容

• 都市形成史研究の基礎②

- 史料の種類（復習）
- 史料批判の方法と概要

• EMアルゴリズムの導入

今日のメイン

- 隠れ変数を含むロジットモデルのパラメータ推定方法
- EMアルゴリズムの定式化
- 実際の計算結果

モデルのパラメータ推定の一例

今日の目標

- ① 史料批判の重要性を知る
- ② EMアルゴリズムの導入方法を知る

THE FULL VIEW OF YUNOMACHI OF
DOGO FROM THE AIR, DOGO.
景全町の湯後道るたし瞰俯りよ上機（勝名後道）

前回の復習



具体的にどんな史料があるのか

都市形成史研究で
主に使う史料

沈黙史料

建築、土木建造物
そのもの（もしくは
遺構）

発言史料

無形

有形

伝承史料

グラフィック史料

キャラクタ史料

風俗、習慣、行事、地名、字名、音声など

地図、絵図、絵画、図面、写真、映像など

文書、記録、文学作品、紀行、日記、帳簿、議事録、新聞など

建築史家である玉井哲雄の史料分類を参考に作成

史料批判について

史料批判の概要

- 史料批判とは何か

- ある史料の価値（史料の有効性と信頼性）を見極める作業のこと
- 歴史研究におけるアプローチの基礎

- 一次史料、二次史料

- 「その時」「その場で」「その人が」叙述 = 「一次史料」（国立国会図書館の定義）
- 作成された年、作成された場所、作成した人が特定可能かどうか
 - 例：伊佐庭如矢日誌 = 日誌そのものが編纂され、関連文書の書誌が掲載されているので「二次史料」、掲載された書誌の原本が「一次史料」になる

形式批判と内容批判：キャラクタ史料の場合

・形式批判＝史料の真贋の検証

・検証事項の例

- ・ オリジナルな史料か（一次史料か否か）
- ・ 史料そのものの偽造や改変の可能性はないか
- ・ 史料の来歴（作成経緯や時代背景，作成者の社会的位置付けなど）

・内容批判＝史料内容の真贋の検証

・検証事項の例

- ・ 故意か過失かに関わらず、史料の内容に虚偽／錯誤の可能性はないか
- ・ 他の史料と比較して内容に矛盾がないか
 - ・ 「記録」が目的の地図や台帳⇔日記や書簡「なぜ書かれたのか」
- ・ 例：旧土地台帳、旧土地台帳付属図面
 - ・ 税務署→法務局へ移管され、公証において「正当な文書」とされていることから、公的な課税証明という範囲で有効かつ信頼性がある，と判断

形式批判と内容批判：キャラクタ史料の場合

・形式批判＝史料の真贋の検証

・検証事項の例

- ・ オリジナルな史料か（一次史料か否か）
- ・ 史料そのものの偽造や改変の可能性はないか
- ・ 史料の来歴（作成経緯や時代背景，作成者の社会的位置付けなど）

・内容批判＝史料内容の真贋の検証

・検証事項の例

- ・ 故意か過失かに関わらず、史料の内容に虚偽／錯誤の可能性はないか
- ・ 他の史料と比較して内容に矛盾がないか
 - ・ 「記録」が目的の地図や台帳⇔日記や書簡「なぜ書かれたのか」
- ・ 例：旧土地台帳、旧土地台帳付属図面

史料に限らず研究で使用する実データの有効性と信頼性には常に細心の注意を

EMアルゴリズムの導入

はじめに

- 対象のモデル構造に応じて，推定手法と推定アルゴリズムを検討するのが良い 最尤推定でいい？推定アルゴリズムは準-Newton法でいい？
- 最小化（最大化）問題に対して最適化関数optimを用いた

MNLモデルのパラメータ推定

対数尤度関数の最大化

```
76 ##### 対数尤度関数 fr の最大化#####  
77  
78 ##パラメータ値の最適化  
79 res <- optim(b0,fr, method = "Nelder-Mead", hessian = TRUE, control=list(fnscale=-1))  
80
```

最適化関数optim (Rがあらかじめ用意してくれている関数)

```
optim(par, fn, gr = NULL, method = "~~~~", lower, upper, control = list(), hessian)
```

b0というパラメータを初期値として，frを最大化するようにパラメータを動かしながら反復して探索。

その時の探索方法は"~~~~"。

ヘッセ行列を返すように指示して，最小化ではなく最大化。

という命令をしている。

対数尤度を最大化したときの結果のもろもろを変数resに入れている。

図：2019年度スタートアップゼミ4/24より抜粋

EMアルゴリズムとは

目的：隠れ(潜在)変数を持つモデルの最尤解を見つける

モデルが同時確率分布 $P(x, s|\theta)$ を持つ場合，対数尤度関数は

$$\ln P(x|\theta) = \ln \sum_s P(x, s|\theta)$$

x : 観測変数
 s : 隠れ変数
 θ : パラメータ

となり一般的に最大化問題を解くことが困難になる

そこで

基本方針

- 現在パラメータ θ^{old} を用いて s の事後分布 $P(s|x, \theta^{old})$ を計算
- 同時確率分布の対数尤度 $\log P(x, s|\theta)$ の期待値を計算
Expectation
- 期待値を最大にする θ^{new} を推定する
Expectation Maximization

EMアルゴリズムとは

1. E-step : θ^{old} を使って隠れ変数 s の事後分布 $P(s|x, \theta^{old})$ を計算
2. E-step : 対数尤度 $\ln P(x, s|\theta)$ の期待値(Q関数)を計算

$$Q(\theta, \theta^{old}) = \sum_s P(s|x, \theta^{old}) \ln P(x, s|\theta)$$

3. M-step : Q関数を θ について最大化し新しいパラメータ θ^{new} を決定

$$\theta^{new} = \arg \max_{\theta} Q(\theta, \theta^{old})$$

4. 収束判定 : 対数尤度関数 $\log P(x, s|\theta^{new})$ を計算

収束条件を満たしている : 終了

収束条件を満たしていない : $\theta^{old} \leftarrow \theta^{new}$ として θ を更新し2へ戻る

ロジットモデルへの適用例

- 隠れ変数の例

- 選択主体（歩行者，避難する人etc)の行動規範には異質性(heterogeneity)がある→潜在クラスモデル(latent class model)を導入すればいいのでは

- 潜在クラスモデルとは：ある説明変数 βx_n で特徴付けられる主体 n は，いずれかの潜在クラス c に確率 s_c で帰属する，と考えるモデル

- 「主体 n がどのクラスにどれくらい帰属しているか」は観測データからはわからない=帰属確率 s_c が隠れ変数

- 主体の潜在クラスへの帰属確率 s_c (隠れ変数)とクラスごとのパラメータ β_{kc} (顕在変数)を同時に推定する。 k : 説明変数 x の番号

- E, Train Kenneth (1st 2003, 2nd 2009) " Discrete choice methods with simulation" の"Chapter 14. EM Algorithms" を参考に実装.

ロジットモデルへの適用：定式化

- 潜在クラス c ごとの選択確率

$$P'_{ni}(\beta_{kc}) = \frac{\exp(\beta_{kc}x_{kc})}{\sum_{i \in I} \exp(\beta_{kc}x_{kc})}$$

(主体の数) N 行 \times (クラス数) C 列

- 条件なしの選択確率

$$P_{ni}(s) = \sum_{c \in C} s_c P'_{ni}(\beta_c)$$

(主体の数) N 行 \times 1列

- モデルの尤度関数

$$LL(s) = \sum_n \log P_{ni}(s)$$

- 確率 s_c でクラス c に所属する主体 n が
選択肢 i を選ぶ場合の条件付き確率 (s は固定)

$$h(\beta_{kc}; i_n, s_c) = \frac{s_c P'_{ni}(\beta_c)}{P_{ni}(s)}$$

(主体の数) N 行 \times (クラス数) C 列

ロジットモデルへの適用：定式化

- この問題におけるQ関数は

$$Q(\theta^0|\theta) = \sum_n \sum_c h(\beta_{kc}; i_n, s_c) \log(s_c P'_{ni}(\beta_c))$$

- 潜在クラスごとの選択確率には隠れ変数 s_c

が含まれていないので以下の二式をそれぞれ扱ってやれば良い
隠れ変数を含む/含まない式で分離

$$Q(s_c|\theta) = \sum_n \sum_c h(\beta_{kc}; i_n, s_c) \log s_c$$

$$Q(\beta_{kc}|\theta) = \sum_c h(\beta_{kc}; i_n, s_c) \log P'_{ni}(\beta_c)$$

ロジットモデルへの適用

1. 初期値パラメータ $\theta^0 = \langle \beta_c^0, s^0 \rangle$ を設定

この時、パラメータ β_c^0 は $C = 1$ の時の結果(理由は後述)、 $s^0 = \frac{1}{C}$ とする。

2. E-step : 負担率の計算

$$h^{old}(\beta_{kc}^{old}; i_n, s^{old}) = \frac{s_c^{old} P'_{ni}(\beta_c^{old})}{P_{ni}(s^{old})} \quad \text{old : 一つ前のstep}$$

3. M-step: Q関数を最大化するパラメータを推定

$Q(s|\theta)$ の場合

$$s_c^{new} = \frac{\sum_n h^{old}(\beta_{kc}^{old}; i_n, s_{c'}^{old})}{\sum_n \sum_{c'} h^{old}(\beta_{kc}^{old}; i_n, s_{c'}^{old})}$$

前のstepでQ関数を
最大化するパラメータ
 β を代入して更新する

$Q(\beta_{kc}|\theta)$ の場合

$$\beta_{kc}^{new} = \operatorname{argmax} \sum_c h^{old}(\beta_{kc}^{old}; i_n, s_c^{old}) \log P'_{ni}(\beta_c)$$

対数尤度最大化に対して
最適化問題を解くやり方で
 β を推定する

4. 対数尤度 $LL = \sum_n \log P_{ni}(s^{new})$ を計算

5. 1-4のステップを繰り返し、収束するまでパラメータを更新

実装例：推定コード

```
# EMアルゴリズム
##結果を格納するマトリクスの設定 nrow=計算回数
### (省略) ###

##初期尤度の算出
p0<-matrix(1/table(Data[,12]))
pc0<-sapply(1:hh,function(x){sum(sc*p0[x,])}) # 条件なしの選択確率
L0<-sum(log(pc0)) # 初期尤度

##アルゴリズムの本編
for(i in 1:num){
  # クラスごとの確率Pncを計算する (p12.段階1)
  pnc<-sapply(1:c,function(y){
  # 帰属確率scを重みにして、条件なしの選択確率Pを計算する (p12.段階2)
  p<-sapply(1:hh,function(x){sum(sc*pnc[x,])})
  # クラスごとのパラメータβcで計算された選択確率にクラス帰属確率scを重みづけした事後確率hを計算する(E-step:期待値の計算)
  h<-t(sapply(1:hh,function(x){sc*pnc[x,]/p[x]})
  # 次のステップのクラス帰属確率sc_newを計算しscを更新する (M-step:Q関数を最大化するパラメータsc_newを算出)
  sc<-apply(h,2,sum)/sum(apply(h,2,sum))
  # 次のステップのパラメータβを最尤推定でoptim関数を使って計算する (M-step:Q関数を最大化するパラメータb_newを算出)
  ## パラメータを更新する
  tval<-matrix(0,nrow=1,ncol=length(b0))
  for(j in 1:c){
  # 更新したパラメータで対数尤度を計算する(p.14 step4)
  L_pnc<-sapply(1:c,function(y){
  L_p<-sapply(1:hh,function(x){sum(sc*L_pnc[x,])})
  LL<-sum(log(L_p)) # 対数尤度の計算

#結果を格納
### (省略) ###
if(abs(resLL[i+1]-resLL[i])<0.001){ #収束条件
  break
}
```

メリット／デメリット

- 比較的扱いやすい／理解しやすいアルゴリズム
- 計算に時間がかかる
 - E-stepとM-stepを繰り返すのでそれはそう
 - 1回のstepごとに、クラス数分のモデルの顕在パラメータ β を最適化関数で推定しているため、クラス数が増えた分計算時間がかかる
今回は2クラスが最適なクラス数だったのでなんとかなった
- 初期値によっては局所最適解に陥る可能性がある
 - 潜在クラスを仮定しない $C=1$ の場合のモデル構造で推定したパラメータを使用

さいごに



さいごに

- 重要な点として…

- 観測データの特徴を吟味し，理解すること
- 観測データの不確実性や潜在データを考慮して，モデリングすること
- モデル構造にとって最適な推定方法や推定アルゴリズムを吟味すること
 - 今日紹介したEMアルゴリズムはその一例です

- 史料も観測された実データの種類です

- 史料批判とは確からしい叙述への第一歩
- どの範囲で価値があるのかを念頭に置きつつ，発見をやっていきましょう

楽しく研究をやっていきましょう！