

縦断データの分析

菅原ますみ（監訳）

2012 年（朝倉書店）

1 時間による変化を検討する際の枠組み

1.2 変化に関する 2 つの質問の違い

1. 個人内変化：個人が時間の経過とともにどう変化するか（線形か、非線形か、一貫性の有無、など記述的なもの）
2. 変化の個人差：変化の個人差を予測するものは何か（関係性：予測変数との関係）

1.3 変化に関する研究の 3 つの重要な特徴

変化を扱う際に注意すべきこと

1. データは複数回（multiple waves）収集されたものであること（3 波以上）
2. 実質的に意味のある時間軸（単位）を用いていること
3. 結果変数は時間とともに組織的に変化するものであること

1.3.1 複数回のデータ収集

横断研究の欠点 しばしば横断研究から時間的変化に関する一般的結論を導き出しがちだが、横断研究は年齢とコホートの効果が交絡しており、選択バイアスがかかりやすい傾向がある。

2 回のデータ収集では不十分な理由

1. 個人の発達の軌跡を示すことができない
2. 真の変化と測定誤差を区別することができない

データ収集の回数を増やすごとに、より柔軟なモデルをより制約の少ない中で当てはめることができる。

1.3.2 時間の適切な測定基準

測定が等間隔であるか否かとは別に、

- 時間構造化されている：全員が同じスケジュールで測定を行なった場合
- 時間構造化されていない：データ収集のスケジュールが個人ごとに異なる

という。

等間隔で収集されたデータは釣り合い性と対象性を備えているので魅力的ではあるが、それにこだわるよりも測定間隔や測定回数をうまく調節することの方が重要。

2 時間についての縦断データの探索

2.1 縦断データセットを作る

2.1.1 個人レベルデータセット（多変量フォーマット）

各行が各ヒトに相当（ヒト × 測定回ごとの値、属性変数など）

経験的成長記録が目で見えて簡単に確認できるが、

- 時点間の相関行列を算出しても、それは各時点間での順位の安定性を表しているに過ぎない
- 時間に関する情報が変数名に含まれているので、分析に使えない
- 測定時点の回数あるいは間隔が個人間で異なると役に立たない
- 予測変数が時間とともに変化するようなものであった場合、手に負えない¹

2.1.2 個人-時間データセット（単変量フォーマット）

各行が各測定時点のデータに相当 ← こちらの形式にすること

以下の4つのタイプの変数が含まれる。

1. 個人を識別する変数
2. 時間を識別する変数
3. 結果変数
4. 予測変数

2.2 個人の時間による変化の記述的分析

2.2.1 経験的成長プロット

個別に結果変数と時点をプロットする。ランダムに何人かを選択するか、重要な変数で層化して抽出するとよい。

2.2.2 個人の経験的成長記録の要約に曲線を使う

ノンパラメトリック・アプローチ

- スプライン平滑化
- Loess 平滑化
- カーネル平滑化
- 移動平均

パラメトリック・アプローチ 個人ごとのデータに最小2乗法による回帰モデルを当てはめる。（関数形を決める）

¹時不変 *time invariant* な変数なら良いが、時変の予測変数 *time-varying predictors* が含まれていれば、それごとに列を加えなくてはいけない。

2.3 変化の個人差を探る

2.3.1 全員分の平滑化曲線を検討する

1. 各測定時点での平均を用いて平滑化曲線を描く²
2. 各個人の平滑化曲線を重ねて描く

2.4 最小二乗法によって推定された変化率の精度と信頼性を改善する

精度 OLS 推定による変化率の精度は、

1. 残差分散（予測値と観測値とのズレ）
2. 観測回数とその間隔

に依存する。つまり、個人 i の OLS 推定における変化率の標本分散は、

$$\frac{\sigma_{\epsilon i}^2}{\sum_{j=1}^T (t_{ij} - \bar{t}_i)^2} = \frac{\sigma_{\epsilon i}^2}{CSST_i} \quad (2.1)$$

である。測定時点がより多様化すれば、変化の測定精度はより高くなる。

1. 測定時点が平均より離れたところまで拡大する
2. 測定時点を増やす

信頼性 精度は個人レベルで、信頼性はグループレベルで意味がある。

仮に完全に釣り合い型データであって、個人の残差は独立に共通の分散 σ_{ϵ}^2 の分布から得られたものであるとすると、OLS 推定による変化率の信頼性は、

$$\frac{\sigma_{trueslope}^2}{\sigma_{trueslope}^2 + \frac{\sigma_{\epsilon}^2}{CSST}} \quad (2.2)$$

と定義される。この量は、

- 変化率の精度（の逆数）
- 真の変化率の分散

の双方に左右される。つまり、個人の変化率の推定精度と変化の異質性とは交絡している。

メモ

変化率の信頼性を算出することが（グループレベルでも）どのような意義を生むのか、よくわからない。

テストであれば、個人差を識別すること自体が目的であるので、そもそも真の得点に個人差が含まれないことには意味がない。つまり、個人差がないことで信頼性が低くなること自体には、テスト得点の個人差は誤差由来でしかないといえる。

一方、変化の個人差については、それがあることが大前提なのではなく、あるとすればどれくらいあるのか現状を把握することが目的である。結果的に変化率に個人差がほとんどなければ、それはそれで現実を反映した知見であるのだから、そのことによって信頼性が低くても問題とは言えないと思う。

メモ

² 「平均の軌跡」であって、「軌跡の平均」ではない。曲線の平均と平均の曲線が一致するのは、パラメータに対して線形である場合（直線、2次曲線、3次曲線など）のみ。

3 変化についてのマルチレベルモデルの紹介

3.4 変化についてのマルチレベルモデルをデータに当てはめる

3.4.1 最尤推定法の利点

1. 漸近的に不偏
2. 漸近的に正規分布に従う
3. 漸近的に有効（他の方法で求められたものよりも標準誤差が小さい）
4. ML 推定量の任意の関数もまた ML 推定量である

また、

- 釣り合い型デザイン
- 測定時点が計画的で欠測値がない
- レベル 2 の各式で同じ予測変数が使われている

ならば、制限付き ML は小標本に対しても上記の性質は漸近的にではなく正確に当てはまる。

3.5 推定された固定効果の検討

最尤法によって推定を行なった場合、固定効果に関する帰無仮説 $H_0 : \gamma = 0$ の一母数検定は、漸近的な標準誤差を用いた検定統計量

$$z = \frac{\hat{\gamma}}{ase(\hat{\gamma})} \quad (3.1)$$

によって行う³。

ソフトウェアによって、同じ検定量が 準 t 統計量、 t 統計量、 t 比などと表記が異なる⁴。

3.6 推定された分散成分の検討

最尤法によって推定を行なった場合、固定効果に関する帰無仮説 $H_0 : \sigma^2 = 0$ の一母数検定は、漸近的な標準誤差を用いた検定統計量

$$z = \frac{\hat{\sigma}^2}{ase(\hat{\sigma}^2)} \quad (3.2)$$

もしくはこれを 2 乗した χ^2 統計量によって行われる。

ただし、この方法は

- 正規性からの逸脱に敏感
- 標本サイズやバランスの不釣り合い（個人ごとに測定回数が異なる）

の影響を受けやすく、使用されるべきではないと指摘する研究者もいる⁵。

³ただし、この検定の性質については漸近的にしか明らかにされていない。

⁴制限付き ML が正確検定になる条件のもとではこれは正確に t 統計量としての性質を持つ。

⁵Miller (1986); Raudenbush & Bryk (2002); Longford (1999)

4 変化についてのマルチレベルモデルのデータ分析

加工されていない変数に非線形モデルを当てはめるより、変換した変数に線形モデルを当てはめた方が、分析はしばしば明確なものになる。

メモ

データをやたら変換してからモデルに当てはめるのではなく、モデルによってデータの性質を表現すべきであるとする久保 (2012) とは対照的な意見だ。

メモ

4.2 合成的な定式化

$$Y_{ij} = \pi_{0i} + \pi_{1i} TIME_{ij} + e_{ij} \quad (4.1)$$

と

$$\pi_{0i} = \gamma_{00} + \gamma_{01} COA_i + \zeta_{0i} \quad (4.2)$$

$$\pi_{1i} = \gamma_{10} + \gamma_{11} COA_i + \zeta_{1i} \quad (4.3)$$

を合成すると、

$$Y_{ij} = [\gamma_{00} + \gamma_{10} TIME_{ij} + \gamma_{01} COA_i + \gamma_{11} (COA_i \times TIME_{ij})] + [\zeta_{0i} + \zeta_{1i} TIME_{ij} + e_{ij}] \quad (4.4)$$

となる。

この合成残差 $[\zeta_{0i} + \zeta_{1i} TIME_{ij} + e_{ij}]$ を見ると、時不偏な残差 ζ_{0i} と ζ_{1i} があることで、測定機会ごとの残差は個人内で自己相関していて等分散が仮定できないということが分かる。

4.3 推定法 (再考)

4.3.1 一般化最小二乗推定 GLS

OLS よりも残差に複雑な仮定をおくことが可能だが、真の誤差共分散行列の中身について知る必要がある。そこで、

1. OLS によってあてはめたモデルから残差を算出し、誤差共分散行列を推定
2. 推定された誤差共分散行列を真のものと見なしてモデルを再当てはめ (固定効果と SE を算出)

という 2 段階のアプローチを繰り返すことを考える (反復一般化最小二乗法 *IGLS*)。

4.3.2 完全最尤推定 *FML* と制限付き最尤推定 *RML*

FML 標本データを観測する尤度を最大化する母数を求める。

分散成分の FML 推定値は、固定効果の FML 推定値を含んだ形で与えられる。

→ 固定効果の値を既知として扱うことになる。

→ 分散成分の推定において、本来考えるべきであった固定効果の自由度を無視してしまう。

→ 分散成分の過小推定 ($N-1$ で割るべきところを N で割っているようなもの。)

モデル全体の当てはまりを検証できるので、適合度検定は固定効果も変量効果も両方対象にすることができる。

RML 標本データではなく、標本残差を観測する尤度を最大化する分散共分散を求める。つまり、

固定効果の推定 → 残差の算出 → この残差は、どのような分散成分のもとで最も得られやすいのか？
を考えるということ。

ただし、RML が FML よりも一方的に優れていることは証明されていない。

残差部分の当てはまりのみを議論しているので、適合度検定は変量効果しか対象にすることはできない。

4.4 2つの無条件マルチレベルモデル

4.4.1 無条件平均モデル

単に結果変数の変動を記述するもの。どのレベルにも予測変数は含まれない。

$$Y_{ij} = \pi_{0i} + e_{ij} \quad (4.6a)$$

$$\pi_{0i} = \gamma_{00} + \zeta_{0i} \quad (4.6b)$$

引き続き分析を行う価値があるだけの十分な変動がそのレベルにおいて見られているか確認する。

級内相関係数

1. 個人間変動が全体に占める割合を評価する
2. 合成無条件平均モデルの残差の自己相関の大きさを要約する

$$\rho = \frac{\sigma_0^2}{\sigma_0^2 + \sigma_e^2} \quad (4.7)$$

無条件平均モデルでは、残差の自己相関係数は級内相関係数に等しい。

$$\text{Cor}(Y_{ij}, Y_{ij'}) = \frac{\text{Cov}(Y_{ij}, Y_{ij'})}{\text{SD}(Y_{ij})\text{SD}(Y_{ij'})} = \frac{\sigma_0^2}{\sigma_0^2 + \sigma_e^2} = \rho$$

4.4.2 無条件成長モデル

レベル1モデルに予測変数 *TIME* を導入する。

$$\begin{aligned} Y_{ij} &= \pi_{0i} + \pi_{1i} \text{TIME}_{ij} + e_{ij} \\ \pi_{0i} &= \gamma_{00} + \zeta_{0i} \\ \pi_{1i} &= \gamma_{10} + \zeta_{1i} \end{aligned} \quad (4.9a)$$

すると、合成された残差は $\zeta_{0i} + \zeta_{1i} \text{TIME}_{ij} + e_{ij}$ であるから、異分散性と自己相関が導き出される⁶。

4.4.3 分散説明率の定量化

擬 R^2 統計量 *pseudo- R^2 statistics* 2つの方法がある。

1. 各個人の各時点での結果変数の予測値を計算し、観測値と予測値の相関を2乗する
2. 予測変数を追加することで $\hat{\sigma}_e^2$ がどれくらい減少したか、追加前の値との比を取る（残差分散減少率）

4.5 モデル構築のための実践的データ分析

- 主たる予測変数と、統制したい予測変数を区別すること
- 予測変数を投入したあとの残差分散を偏分散 *partial variance* や 条件付き分散 *conditional variance* と呼ぶ

⁶時変な変数 *TIME* によって分散が左右されること、また時不変な残差成分 ζ_{0i}, ζ_{1i} が存在することによる。

4.5.4 予測変数の中心化

- 時不変の変数を中心化する場合、標本平均を引くことが一般的だが、現実的に意味のある定数で中心化することが重要。
- ダミー変数であっても、平均（平均的な対象者）に意味があれば中心化することもある。
- レベル2の予測変数を全て中心化しておけば（レベル1の予測変数が *TIME* のみであれば）切片は無条件成長モデルに一致する。→ レベル2の予測変数は、2値変数でも全て中心化してしまうことが多い。

メモ

標本平均や標本比率などを用いて中心化する場合、その値の信頼性や標本変動によって切片の有意性は変化すると思う。

メモ

4.6 乖離度統計量 *deviance statistics* を用いたモデルの比較

乖離度統計量

$$D = -2LL_{\text{current model}} \quad (4.10)$$

- 双方のモデルが同じデータを用いて推定されており⁷
- 片方のモデルがもう片方にネスとしている

場合、乖離度統計量を比較することができる。

情報量規準 ネストしていないモデルを比較する場合、AIC や BIC を用いる。

- AIC: モデルに含まれるパラメータの数⁸に基づいてペナルティを課す
- BIC: パラメータの数と標本サイズ⁹に基づいてペナルティを課す

4.7 複合仮説のワルド検定

帰無仮説を $H_0: C\gamma' = \mathbf{0}$ のように一般線形仮説の形で表し、パラメータの重み付き線形結合の2乗をその分散の推定値と比較する。なお、分散成分の複合帰無仮説の検定に利用することは推奨されない。

メモ 一般線形仮説のワルド検定については要確認。メモ

4.8 モデルの仮定の許容度の評価

得られた結果の妥当性は、モデルを当てはめる際の仮定をどれくらい許容できるかに依存する。

- 関数形: 結果変数と予測変数の散布図を各レベルで描いてみる
- 正規性:
 - 素残差 $\hat{e}_{ij}, \zeta_{0i}, \zeta_{1i}$ について、各値とそれに対応する正規スコアとの散布図（正規確率プロット）を描画する（→ 線形性から逸脱していないか？）
 - 標準化残差をプロットする（正規分布に従っているなら、約95%が中心 $\pm 2SD$ 以内に収まるはず）
- 等分散性: 予測変数と素残差のプロット

⁷欠測のあるデータの利用には特に注意。

⁸FML では全パラメータ、RML では分散成分のみ。

⁹人数なのか観測回数なのか定かではない。

4.9 経験ベイズ推定値

以下の2つの方法がある。

1. OLS 推定値と母平均の推定値の重み付き平均
2. 個人の予測変数から平均的な軌跡を得て、そこにレベル2の残差を加える

$$\begin{aligned}\tilde{\pi}_{0i} &= \hat{\pi}_{0i} + \hat{\zeta}_{0i} \\ \tilde{\pi}_{1i} &= \hat{\pi}_{1i} + \hat{\zeta}_{1i}\end{aligned}\tag{4.21}$$

属性を共有する他者からの情報を利用することで (borrowing strength)、個人の推定値が正確なものになる。経験ベイズ推定値では不偏性は犠牲になるが、OLS よりも精度は良い。ただし、推定値の質はモデル適合の質に大きく依存する。

5 時間的な変数 *TIME* をより柔軟に扱う

5.1 間隔が一定ではない測定時点

例えば、年齢が不均一なコホートを一定期間追いかける加速コホートデザインによって、より少ない測定時点数でより長い時間間隔を経た変化をモデリングできる。

時間構造化されていないデータを時間構造化されているように扱うことは、分析に誤差を持ち込んでいるようなもの。(測定間隔が違うのに、*TIME* を測定時期そのものではなく、1回目、2回目などとしてまとめてしまう。)

5.2 測定時点の数が異なる場合

測定時点の数が異なる非釣り合い型データであっても、マルチレベルモデルは問題なく当てはめられる。

メモ

結果変数が賃金のように正に歪んでいる場合、対数変換したものを予測の対象とすることがある。

$$\log Y_{ij} = \gamma_0 + \gamma_1 TIME_{ij} + e_{ij}$$

このとき、 $100(\exp(\gamma_1) - 1)$ をパーセント変化率として解釈する。つまり、*TIME* が1増加すれば、結果変数には $100(\exp(\gamma_1) - 1)$ 倍の増加が見込めるということ。

もし $\gamma_1 = 0$ であれば $\exp(\gamma_1) = 1$ であるから、*TIME* が変化しても結果変数は変化しないことになる (パーセント変化率は0)。

メモ

5.2.2 分析上の問題点

データが過度に非釣り合いであると、分散成分の推定に影響を与える可能性がある。特に、時点の少ない人が多すぎると、残差のばらつきを正常に評価できない。

- 境界制約を超えた推定値、非収束などの問題が生じる。
- 一部の残差をゼロに固定するなど、確率的な部分を単純化することで解決を図る。

5.2.3 欠測の様々なタイプを区別する

非釣り合いは計画的に生じたものではないことが多く、なぜ非釣り合いになったのか考えることが求められる。

以下の3つを無視可能な無回答と言い、いずれの条件下でもマルチレベルモデルを当てはめた結果は正しく一般化可能 (Laird, 1988)¹⁰。

完全にランダムな欠測: **missing completely at random (MCAR)** 観測される確率が

- 特定の時期
- 予測変数の値
- 結果変数の値

の3つと独立である場合。

共変量依存型の脱落: **covariate dependent dropout (CDD)** 観測される確率が

- 結果変数の値

と独立であるが、

- 特定の時期
- 予測変数の値

と関連している。

ランダムな欠測: **missing at random (MAR)** MCAR と CDD を示すには、「欠測確率は結果変数の同時期の値と関連しない」ことを示さなくてはならない。

MAR では、

- 欠測がいかなる観測データに依存していても構わない
- 欠測はいかなる観測されていない値に依存してはいけない

ことが必要。実際は、この仮定を満たすことすら難しい (Greenland & Finkle, 1995)。

無視できない欠測 欠測が上記のいずれでもない場合、

- 選択モデル (完全データのためのモデルと欠測メカニズムのモデルの2つを立てる)
- パターン混合モデル (欠測のパターンごとにマルチレベルモデルを当てはめる)

によって補正を行う必要がある。

5.3 時変の予測変数

5.3.1 事変の予測変数の主効果を含める

分散成分

- 時不変な予測変数を加えた場合: レベル2の分散成分は減少しても、レベル1の分散成分はあまり減少しない
- 時変な予測変数を加えた場合: レベル1も2も分散成分は変化する

※ただし、レベル2の分散成分の減少にはあまり意味がないことが多い。

時変の予測変数を追加すると、レベル1のパラメータ π_i は意味が変わる。

→ 追加前に比べて、レベル2の分散が増加することもある。

¹⁰一般化推定方程式 GEE では、MCAR が成り立っていることが必要。

5.3.2 時変の予測変数の効果が時間とともに変化することを許容する

時変予測変数間の交互作用を入れれば良い。

$$Y_{ij} = [\gamma_{00} + \gamma_{10} TIME_{ij} + \gamma_{20} UNEMP_{ij} + \gamma_{30} \underbrace{UNEMP_{ij} \times TIME_{ij}}_{\text{interaction}}] + [\zeta_{0i} + \zeta_{1i} TIME_{ij} + e_{ij}] \quad (5.7)$$

はてな？

$UNEMP$ と $TIME$ の交互作用を考える際に、センタリングしておかなくて良いのだろうか？

はてな？ここまで

5.3.3 時変の予測変数を再中心化する

中心化は大きな関心を集めているテーマ (Kreft et al., 1995; Hofmann & Gavin, 1998)。

全平均中心化 予測変数の全標本平均を引き算すること。

ただし、測定時点と測定回数が人によって異なる場合、この「平均」には意味がないかもしれない。
→ 実質的な意味のある定数で中心化することがよくある。

文脈内中心化 (群平均の中心化) 個人ごとに与えられた定数で中心化すること。

※ 内生性の解釈の問題が生じる。

5.3.4 重要な注意：逆方向因果の問題

逆方向因果 (内生性) X と Y に相関があるとき、どちらが原因で結果か断定できない問題。縦断データだけで逆方向因果の問題が解決されるわけではない。

予測変数の種類によってどれくらい問題が生じるかが異なる。

1. 確定した defined 変数

- 時変の予測変数の値があらかじめ決定している
- ほとんどの場合は時間の関数
- 逆方向因果は問題にならない

2. 付属する ancillary 変数

- 時変の予測変数の値が参加者とは無関係な外部の確率的な過程によって (不規則に) 決まる
- 被験者の生きている物理的もしくは社会的な環境の、変化する特性¹¹
- 逆方向因果は問題にならない

3. 文脈的な contextual 変数

- 確率的に値が決まる時変の予測変数だが、参加者と関係がある¹²
- 結果変数の値によって影響されることがあり、逆方向因果が問題になる

4. 内的な internal 変数

- 個人の時間とともに潜在的に変化する時変の予測変数
- 心理的、身体的、社会的な状態など¹³

¹¹ その地域の失業率、天気、ランダムな処遇への割当など

¹² 両親の離婚、通っている保育所の質、など

¹³ 気分、血圧、既婚/未婚など

- 逆方向因果が問題になる

解決の指針として、

1. 理論を指針して、最も厳しい批判を考える
2. 予測変数について、1 時点前の値を結果変数に対応させるようにする

など。

5.4 *TIME* の効果の再中心化

時間変数を T と表し、定数 c で中心化することを考える。

$$\begin{aligned} Y_{ij} &= \pi_{0i} + \pi_{1i}(T_{ij} - c) + e_{ij} \\ \pi_{0i} &= \gamma_{00} + \gamma_{01} TREAT_i + \zeta_{0i} \\ \pi_{1i} &= \gamma_{10} + \gamma_{11} TREAT_i + \zeta_{1i} \end{aligned} \quad (5.12)$$

なるモデルを考える。このとき、 $Var(\zeta_{0j}) = \sigma_0^2$ が c によって変化するのと合わせて、 ζ_{0i} と ζ_{1i} の共分散も（結果的に相関も）大きく変化する（Rogosa & Willett, 1985）。
切片と傾きの相関が強くなりすぎると、安定した推定値を得られにくくなる。
→ 切片パラメータの必要性をなくすような c を探すことも有効。

メモ

上記のモデルについて、中心化定数 c の値と ζ_{0i} と ζ_{1i} の相関の変化について図的に描いたものが図 5.1。

メモここまで

次のようなモデルを考えることで、初期状態と最終状態に対する問いに同時に応えることができる。

$$\begin{aligned} Y_{ij} &= \pi_{0i} \left(\frac{T_{max} - T_{ij}}{T_{max} - T_{min}} \right) + \pi_{1i} \left(\frac{T_{ij} - T_{min}}{T_{max} - T_{min}} \right) + e_{ij} \\ \pi_{0i} &= \gamma_{00} + \gamma_{01} TREAT_i + \zeta_{0i} \\ \pi_{1i} &= \gamma_{10} + \gamma_{11} TREAT_i + \zeta_{1i} \end{aligned} \quad (5.13)$$

すると、

- π_{0i} : 初期値
- π_{1i} : 最終値

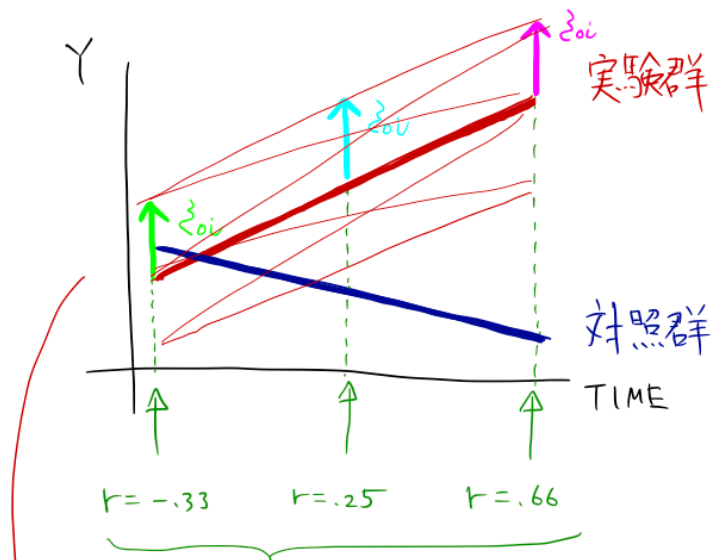
となる。

はてな？

TIME が人によってばらばらである場合、「初期値」とか「最終値」はあくまでサンプル全体での値であるから、「その人の初期値」「その人の最終値」というわけではない。外挿となる可能性が高いのではないか。

はてな？ここまで

このモデルは、上記の T_{ij} を c でセンタリングしたモデルとまったく同じ乖離度を持つ。



中心化するcの値により、切片とTIMEの傾きの
レベル2の共分散 (相関) は大きく変化する

6つの個人直線の平均とを比較すると、

cが大きいと、切片が平均から高い直線は、

傾きが急なものが多い。→ 相関(強)

cが中くらいだと、平均な直線と同じ傾きのものが

切片 + 傾きを通る。→ 相関(中)

cが小さいと、平均的な直線よりゆるやかなものが

切片 + 傾きを通る。→ 負の相関

Figure 5.1: 中心化定数とレベル2の残差相関の変化

6 非連続あるいは非線形の変化のモデリング

初期値と変化率という問いから、増加、減少、転換点、移行、漸近線を考えることへと変わる。

- *TIME* の変換
- 関数形式の仮定

6.2 個人の非線形の変化を変換によってモデリングする

- 変換された尺度でもパラメータの解釈は明瞭
- 多くの変数の測度はそもそも後づけなので、オリジナルで分析しようが変換して分析しようが問題は生じない

べき乗のはしご (Mosteller & Tukey, 1977) データの経験的プロットのでっぱりをみながら、

- *Y* の次数を上げるか/下げるか
- *TIME* の次数を下げるか/上げるか

探索的に決める。

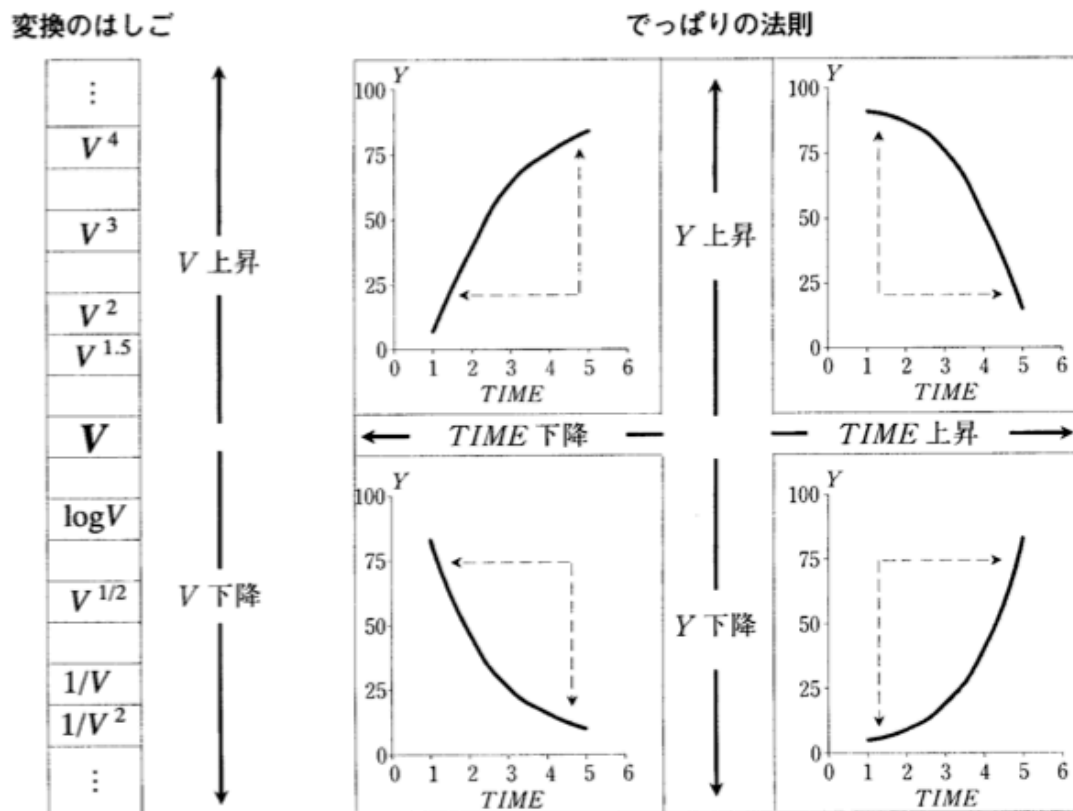


図 6.5 変換のはしごとでっばりの法則

個人の成長の軌跡の賢明な変換方法選択のためのガイドライン。

- 平均値のプロットから判断するのではない
- 変換することで意味がわかりにくくなるのであれば、別の測度を変換する

6.3 多項式成長モデル

$TIME$ の 2 次、3 次の項を含める。例えば、

$$Y_{ij} = \pi_{0i} + \pi_{1i} TIME_{ij} + \pi_{2i} TIME_{ij}^2 + e_{ij}$$

この場合、

- π_{1i} は変化率ではなく、 $TIME = 0$ の瞬間的な変化率
- π_{2i} は曲率

となる。時間による微分 $= 0$ を解くことで、曲線の頂点は $-\pi_{1i}/2\pi_{2i}$ であることがわかる。
多項式では、各パラメータの解釈が難しくなる。

6.4 真に非線形な軌跡

動的的一致性 (Keats, 1983)

1. 「平均の曲線」と「曲線の平均」が一致する
2. 個人の軌跡と平均の軌跡が同じ関数になる¹⁴

パラメータに対して線形な個人の成長モデルは、動的的一致性を満たす。

6.4.3 真に非線形な軌跡を調べる

双曲線型 以下の曲線を直角双曲線と呼び、生物学的成長や農作物の成長をモデリングするのに用いられる。

$$Y_{ij} = \alpha_i - \frac{1}{\pi_{1i} TIME_{ij}} + e_{ij}$$

ここで、

- α_i は上方漸近線
- π_{1i} は漸近線に近づく割合（小さいほど速く近づく）

を表す。ただし、 $TIME \rightarrow 0$ のとき、 $Y \rightarrow -\infty$ となることもあり、あまり用いられない。

逆多項式型 次の式で表される曲線を、**2 次の逆多項式**という（直角双曲線の拡張）。

$$Y_{ij} = \alpha_i - \frac{1}{(\pi_{1i} TIME_{ij} + \pi_{2i} TIME_{ij}^2)} + e_{ij}$$

- π_{2i} は上方漸近線に近づくのを抑える働きをしている（大きく負であれば、一度近づいた漸近線から下降してゆく）

指数型 生物学的、農学的、物理的成長のモデリングに最も広く使われている。

¹⁴異なる 2 次曲線を持つグループの平均は 2 次曲線になる。一方、ロジスティック曲線の集合の平均は平滑化されたステップ関数になる。

単純な指数成長モデル 無限の栄養がある状態での細菌の増殖など、爆発的軌跡として知られている。結果変数が対数変換されたものと見るとわかりやすい。

$$Y_{ij} = \pi_{0i} e^{\pi_{1i} TIME_{ij}} + e_{ij}$$

ここで、

- π_{0i} は切片に相当（値が大きいほどより高い位置から直線が始まる）
- π_{1i} は傾きに相当（値が大きいほどより速く無限大に近づいてゆく）

である。

負の指数成長モデル 上方漸近線がある点で直角双曲線や2次の逆多項式に形状は似ているが、 $TIME = 0$ で $-\infty$ にならない。農作物の収穫量、感染症の患者数など一定を超えると横ばいになるようなもの。

$$Y_{ij} = \alpha_i - (\alpha_i - \pi_{0i}) e^{-\pi_{1i} TIME_{ij}} + e_{ij}$$

ここで、

- α は上方漸近線
- π_{0i} は切片に相当（値が大きいほどより高い位置から直線が始まる）
- π_{1i} は傾きに相当（値が大きいほどより速く漸近線に近づいてゆく）

である。

ロジスティックモデル 下方漸近線と上方漸近線があり、生物学でよく使われる。

$$Y_{ij} = \alpha_{1i} + \frac{\alpha_{2i} - \alpha_{1i}}{1 + \pi_{0i} e^{-\pi_{1i} TIME_{ij}}} + e_{ij}$$

ここで、

- α_{1i} は下方漸近線
- α_{2i} は上方漸近線
- π_{0i} は切片に相当（値が大きいほどより高い位置から直線が始まる）
- π_{1i} は傾きに相当（値が大きいほどより速く漸近線に近づいてゆく）

である。

はてな？

例では、 α_{1i}, α_{2i} を課題の条件から定め、 π_{0i} と π_{1i} にレベル2の予測変数を含んだ線形モデル（誤差は2変量正規分布する）を仮定している。

この、レベル2のモデル自体は、非線形にできないのだろうか？あるいは、そのようなことが必要になる文脈はないのだろうか？

はてな？ここまで

6.4.4 実質科学的理論から個人の成長の数学的表現へ

自触媒の原理を応用した学習理論 (Robertson, 1909) 学習が進む割合は、

1. すでに生じた学習の量
2. まだ生じていない学習の量

に比例する。

すなわち、 Y を 時点 t で学習した量、 α を学習可能な量の上限、 k を比例定数とすると、

$$\frac{dY}{dt} = kY(\alpha - Y)$$

と表せる。

ここから、 $\pi_0 = e^{-c\alpha}$ 、 $\pi_1 = k\alpha$ とおくと、この微分方程式は、

$$Y = \frac{\alpha}{1 + \pi_0 e^{-\pi_1 \text{TIME}}}$$

と解け、ロジスティック関数となる。

これをデータに当てはめるため、添え字 i, j をつけ、さらに誤差項を加えることで、

$$Y_{ij} = \frac{\alpha_i}{1 + \pi_{0i} e^{-\pi_{1i} \text{TIME}_{ij}}} + e_{ij}$$

とモデル化すれば良い。

はてな？

誤差項は、単に足すだけで良いのか？

はてな？ここまで

反復回数と記憶量 (Robertson, 1908) 仮に人間が保持できる記憶の量は常に一定であるとする、時間 t における記憶量 Y の変化は、

$$\frac{dY}{dt} = kY$$

と表せる。

メモ

すでに記憶した量によって、今記憶できる速さは左右されるということ。(最初はたくさん記憶できても、記憶した量がたくさんになるとそれ以上記憶できなくなってくる。)

メモここまで

これをデータに当てはめるため、微分方程式を解いた上で添え字 i, j をつけ、さらに誤差項を加えることで、

$$Y_{ij} = \pi_{0i} e^{-\pi_{1i} \text{TIME}_{ij}} + e_{ij}$$

とモデル化すれば良い。丸暗記のモデルとしてよく用いられる。

習慣強度の発達 (Hull, 1943; 1952) 学習された反応がどう想起されるのかをもとに、負の指数型の軌跡

$$Y_{ij} = \alpha_i (1 - e^{-\pi_{1i} \text{TIME}_{ij}}) + e_{ij}$$

を仮定した。

学習関数に対する理論的な方程式 (Thurstone, 1917; 1930) 獲得 Y が練習 (それまでの時間) に依存するというモデルを作り、ネズミの迷路課題学習に適用した。

$$Y_{ij} = \pi_{0i} + \frac{(\alpha_i - \pi_{0i})TIME_{ij}}{\pi_{1i} + TIME_{ij}} + e_{ij}$$

7 マルチレベルモデルの誤差共分散構造を検討する

7.1 変化についての「標準的な」定式化

今、以下のモデルを考える。

$$\begin{aligned} Y_{ij} &= \pi_{0i} + \pi_{1i} TIME_{ij} + e_{ij} \\ e_{ij} &\stackrel{iid}{\sim} N(0, \sigma_e^2) \\ \pi_{0i} &= \gamma_{00} + \gamma_{01}(COG_i - \overline{COG}) + \zeta_{0i} \\ \pi_{1i} &= \gamma_{10} + \gamma_{11}(COG_i - \overline{COG}) + \zeta_{1i} \\ \begin{bmatrix} \zeta_{0i} \\ \zeta_{1i} \end{bmatrix} &\stackrel{iid}{\sim} N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_0^2 & \sigma_{01} \\ \sigma_{10} & \sigma_1^2 \end{bmatrix}\right) \end{aligned} \quad (7.1)$$

7.2 合成モデルにおける誤差分散共分散行列

このとき、これを合成モデルで表すと、その確率部分は

$$r_{ij} = [e_{ij} + \zeta_{0i} + \zeta_{1i} TIME_{ij}] \quad (7.5)$$

となる。

すると、 r_{ij} は、個人間では独立であっても、個人内では時点間で相関しており、かつ非等質である。

つまり、各個人単位の分散共分散行列を Σ_r と表すと、 \mathbf{r} の分布は

$$\mathbf{r} \sim N\left(\mathbf{0}, \begin{bmatrix} \Sigma_r & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \Sigma_r & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \Sigma_r \end{bmatrix}\right) \quad (7.9)$$

という、ブロック対角行列になる。

7.2.1 合成残差の分散

分散共分散行列 Σ_r の対角部分、つまり残差分散は、

$$\sigma_{r_j}^2 = \text{Var}(e_{ij} + \zeta_{0i} + \zeta_{1i} t_j) \quad (7.11)$$

$$= \left(\sigma_e^2 + \frac{\sigma_0^2 \sigma_1^2 - \sigma_{01}^2}{\sigma_1^2}\right) + \sigma_1^2 \left(t_j + \frac{\sigma_{01}}{\sigma_1^2}\right)^2 \quad (7.12)$$

となり、時間 t に対して 2 次の関係性を持つことになる。つまり、1 つの最小値から放射線状に増加すると仮定していることになる。

7.2.2 合成残差の共分散

分散共分散行列 Σ_r の非対角部分、つまり残差共分散は、

$$\sigma_{r_j r_{j'}} = \sigma_0^2 + \sigma_{01}(t_j + t_{j'}) + \sigma_1^2 t_j t_{j'} \quad (7.13)$$

となり、強い時間依存性を持つことがわかる。

もし $\sigma_0^2 = \sigma_1^2 = \sigma_{01} = 0$ なら、 Σ_r は対角行列となり、OLS を当てはめることができる。

もし $\sigma_{01} = 0$ なら、 Σ_r は複合対称性を持つことになる。

7.2.3 合成残差の自己相関

合成残差の自己相関は、

$$\rho_{r_j r_{j'}} = \frac{\sigma_{r_j r_{j'}}}{\sqrt{\sigma_{r_j}^2 \sigma_{r_{j'}}^2}}$$

7.3 誤差共分散構造の別の仮定の仕方

7.3.1 非構造的誤差共分散行列

乖離度統計量は、あらゆる誤差共分散構造の中で常に最小になる。

$$\Sigma_r = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} \\ \sigma_{21} & \sigma_2^2 & \sigma_{23} \\ \sigma_{31} & \sigma_{32} & \sigma_3^2 \end{bmatrix}$$

7.3.2 複合対称的誤差共分散行列

標準的モデルの特殊ケース。

$$\Sigma_r = \begin{bmatrix} \sigma^2 + \sigma_1^2 & \sigma_1^2 & \sigma_1^2 \\ \sigma_1^2 & \sigma^2 + \sigma_1^2 & \sigma_1^2 \\ \sigma_1^2 & \sigma_1^2 & \sigma^2 + \sigma_1^2 \end{bmatrix}$$

7.3.3 異分散複合対称的誤差共分散行列

自己相関 ρ は共通。各時点における分散が非等質となるため、共分散も非等質となる。

$$\Sigma_r = \begin{bmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \rho & \sigma_1 \sigma_3 \rho \\ \sigma_2 \sigma_1 \rho & \sigma_2^2 & \sigma_2 \sigma_3 \rho \\ \sigma_3 \sigma_1 \rho & \sigma_3 \sigma_2 \rho & \sigma_3^2 \end{bmatrix}$$

7.3.4 (1 次の) 自己回帰的誤差共分散行列

自由度を相当節約するが、かなり厳しい制約ではある。

$$\Sigma_r = \begin{bmatrix} \sigma^2 & \sigma^2 \rho & \sigma^2 \rho^2 \\ \sigma^2 \rho & \sigma^2 & \sigma^2 \rho \\ \sigma^2 \rho^2 & \sigma^2 \rho & \sigma^2 \end{bmatrix}$$

7.3.5 異分散自己回帰的誤差共分散行列

$$\Sigma_r = \begin{bmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \rho & \sigma_1 \sigma_3 \rho^2 \\ \sigma_2 \sigma_1 \rho & \sigma_2^2 & \sigma_2 \sigma_3 \rho \\ \sigma_3 \sigma_1 \rho^2 & \sigma_3 \sigma_2 \rho & \sigma_3^2 \end{bmatrix}$$

7.3.6 トープリッツ誤差共分散行列

自己回帰的であるが、対角成分に対する各帯は同一比に制約されない。

$$\Sigma_r = \begin{bmatrix} \sigma^2 & \sigma_1 & \sigma_2 \\ \sigma_1 & \sigma^2 & \sigma_1 \\ \sigma_2 & \sigma_1 & \sigma^2 \end{bmatrix}$$

7.3.7 「正しい」誤差共分散構造を選ぶことは本当に重要なのか？

- 誤差構造に何を選ぶかに関わらず、固定効果のパラメータ推定値を本質的に変化させることは滅多にない
- 誤差共分散構造についての仮説を洗練させることは、固定効果の推定精度に影響を与える

8 共分散構造分析を用いて変化のモデリングを行う

8.1 一般的な共分散構造モデル

8.1.1 X 測定モデル

外生的構成概念に関する観測変数 X_1 の i 番目の観測値について、

$$X_{1i} = \tau_{x_1} + \lambda_{11}^x \xi_{1i} + \delta_{1i} \quad (8.1)$$

を考える。

- τ_{x_1} : X_{1i} の母集団平均、同じ構成概念に対する異なる指標の観測スコアが異なる平均を取ることができる
- λ_{11}^x : 尺度因子、ある構成概念が異なる尺度で測定されることが可能
- ξ_{1i} : (外生的な) 構成概念の値
- δ_{1i} : 指標 X のうち、仮定された構成概念に依存しない部分

外生的構成概念すべてについて測定モデルをまとめると、

$$X = \tau_x + \Lambda_x \xi + \delta \quad (8.4)$$

と表記できる。ここで、 δ の母分散共分散行列を Θ_δ とおく。ここには、あらゆる構造を仮定することが可能。また、外生的構成概念 ξ については、平均 κ 、母共分散行列 Φ を仮定する。

8.1.2 Y 測定モデル

内生的構成概念に関する測定モデルを、

$$Y = \tau_y + \Lambda_y \eta + \epsilon \quad (8.11)$$

と表す。また、 ϵ の母共分散行列は Θ_ϵ とする。

外生的構成概念の場合と異なり、内生的構成概念の平均と分散は、構造モデルによって表現される。

8.1.3 構造モデル

外生的構成概念と内生的構成概念の関係性を、

$$\underbrace{\eta}_{\text{endogenous factors}} = \underbrace{\alpha}_{\text{intercept}} + \underbrace{\Gamma \xi}_{\text{loadings} \times \text{exogenous factors}} + \underbrace{B \eta}_{\text{loadings} \times \text{endogenous factors}} + \underbrace{\zeta}_{\text{residuals}} \quad (8.15)$$

と表す。ただし、 ζ の母共分散行列は Ψ とする。

8.2 潜在成長モデリングの基礎

- Meredith & Tisak (1984; 1990), Tisak & Meredith (1990) による開発
- McArdle らによる拡張、心理・社会学への応用
- Muthén らによる、時間構造化されていないデータや欠測値の扱いに関する開発

潜在成長モデルを当てはめるには、多変量フォーマットに則った個人データセットの形で用意すること。つまり、時間構造化されたデータが最も扱いやすく、望ましい。

8.2.1 レベル 1 モデルの Y 測定モデルへの移植

今、レベル 1 モデルとして

$$Y_{ij} = \pi_{0i} + \pi_{1i}t_j + e_{ij} \quad (8.19)$$

を考える。

これは、仮に 3 時点あるとすると

$$\begin{bmatrix} Y_{i1} \\ Y_{i2} \\ Y_{i3} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 & t_1 \\ 1 & t_2 \\ 1 & t_3 \end{bmatrix} \begin{bmatrix} \pi_{0i} \\ \pi_{1i} \end{bmatrix} + \begin{bmatrix} e_{i1} \\ e_{i2} \\ e_{i3} \end{bmatrix} \quad (8.21)$$

$$\mathbf{Y} = \boldsymbol{\tau}_y + \boldsymbol{\Lambda}_y \boldsymbol{\eta} + \boldsymbol{\epsilon}$$

と Y 測定モデルの形で書ける。ただし、負荷 $\boldsymbol{\Lambda}_y$ はデータとして固定されていることに注意。この固定された負荷に合うように、 $\boldsymbol{\eta}$ が調整されることになる。

8.2.2 レベル 2 モデルの構造モデルへの移植

無条件成長モデル 今、レベル 2 モデルを説明変数のない無条件成長モデルとすると、

$$\begin{bmatrix} \pi_{0i} \\ \pi_{1i} \end{bmatrix} = \begin{bmatrix} \mu_{\pi_{0i}} \\ \mu_{\pi_{1i}} \end{bmatrix} + \begin{bmatrix} \zeta_{0i} \\ \zeta_{1i} \end{bmatrix} \quad (8.26)$$

は

$$\begin{bmatrix} \pi_{0i} \\ \pi_{1i} \end{bmatrix} = \begin{bmatrix} \mu_{\pi_{0i}} \\ \mu_{\pi_{1i}} \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \pi_{0i} \\ \pi_{1i} \end{bmatrix} + \begin{bmatrix} \zeta_{0i} \\ \zeta_{1i} \end{bmatrix} \quad (8.27)$$

$$\boldsymbol{\eta} = \boldsymbol{\alpha} + \boldsymbol{\Gamma} \boldsymbol{\xi} + \boldsymbol{B} \boldsymbol{\eta} + \boldsymbol{\zeta}$$

と表すことができる。

X 測定モデルへの時不変な予測変数の追加（レベル 2 の説明変数） 今、レベル 2 を予測変数 X を含んだモデルにすることを考える。すなわち、

$$\begin{aligned} \pi_{0i} &= \mu_{\pi_{0i}} + \gamma_{\pi_0} X_i + \zeta_{0i} \\ \pi_{1i} &= \mu_{\pi_{1i}} + \gamma_{\pi_1} X_i + \zeta_{1i} \end{aligned}$$

である。これは、

$$\begin{bmatrix} \pi_{0i} \\ \pi_{1i} \end{bmatrix} = \begin{bmatrix} \mu_{\pi_{0i}} \\ \mu_{\pi_{1i}} \end{bmatrix} + \begin{bmatrix} \gamma_{\pi_0} \\ \gamma_{\pi_1} \end{bmatrix} X_i + \begin{bmatrix} \zeta_{0i} \\ \zeta_{1i} \end{bmatrix}$$

とも表わせるから、式 (8.3.2) に合わせると、

$$\begin{bmatrix} \pi_{0i} \\ \pi_{1i} \end{bmatrix} = \begin{bmatrix} \mu_{\pi_{0i}} \\ \mu_{\pi_{1i}} \end{bmatrix} + \begin{bmatrix} \gamma_{\pi_0} \\ \gamma_{\pi_1} \end{bmatrix} [X_i] + \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \pi_{0i} \\ \pi_{1i} \end{bmatrix} + \begin{bmatrix} \zeta_{0i} \\ \zeta_{1i} \end{bmatrix}$$

$$\boldsymbol{\eta} = \boldsymbol{\alpha} + \boldsymbol{\Gamma} \boldsymbol{\xi} + \boldsymbol{B} \boldsymbol{\eta} + \boldsymbol{\zeta}$$

と表すことができる。

すると、X 測定モデルとして

$$\underbrace{\mathbf{X}}_{X_i} = \underbrace{\boldsymbol{\tau}_x}_0 + \underbrace{\boldsymbol{\Lambda}_x}_1 \underbrace{\boldsymbol{\xi}}_{X_i} + \underbrace{\boldsymbol{\delta}}_0 \quad (8.31)$$

ただし $\text{Mean}(\zeta) = \kappa, \text{Cov}(\zeta) = \Phi$ 、もしくは、

$$\underbrace{X_i}_{X_i} = \underbrace{\tau_x}_{\mu_x} + \underbrace{\Lambda_x}_1 \underbrace{\xi}_{X_i - \mu_x} + \underbrace{\delta}_0$$

を考えているとすれば良いことになる。後者の方法では、 $\text{Mean}(X_i) = \tau_x (= \mu_x)$ と指定しており、 X_i が右辺において中心化されていることになる。

8.3 変数横断的な変化の分析

8.3.1 X, Y 測定モデルの両方で個々人の変化をモデリングする

内生的構成概念も、それを予測する外生的構成概念も個人内で変化している場合、測定モデルは

$$\begin{aligned} X &= \tau_x + \Lambda_x \xi + \delta \\ Y &= \tau_y + \Lambda_y \eta + \epsilon \end{aligned}$$

と表すとして、 δ の共分散行列 Θ_δ と ϵ の共分散行列 Θ_ϵ に加え、 δ と ϵ の共分散行列 $\Phi_{\delta\epsilon} (= \text{diag}\{\sigma_{\delta_i\epsilon_i}\})$ を指定する必要がある。

メモ

この共分散行列は、外生的構成概念における測定モデルの誤差と内生的構成概念における測定モデルの誤差間の相関関係を表したものの。仮にすべてゼロとおいたとしても、分析自体は行えるはず。

メモここまで

8.3.2 構造モデルで変化の軌跡間の関係性をモデリングする

今、外生的構成変数と内生的構成概念が1つずつあるモデルを考える。外生的構成変数の切片と傾きを $[\pi'_{0i}, \pi'_{1i}]^T$ 、内生的構成概念の切片と傾きを $[\pi_{0i}, \pi_{1i}]^T$ とおく。

すると、切片から切片、切片から傾き、傾きから切片、傾きから傾き、へのパス係数を Γ に含んだモデル

$$\begin{aligned} \begin{bmatrix} \pi_{0i} \\ \pi_{1i} \end{bmatrix} &= \begin{bmatrix} \mu_{\pi_{0i}} \\ \mu_{\pi_{1i}} \end{bmatrix} + \begin{bmatrix} \gamma_{\pi_0\pi'_0} & \gamma_{\pi_0\pi'_1} \\ \gamma_{\pi_1\pi'_0} & \gamma_{\pi_1\pi'_1} \end{bmatrix} \begin{bmatrix} \pi'_{0i} \\ \pi'_{1i} \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \pi_{0i} \\ \pi_{1i} \end{bmatrix} + \begin{bmatrix} \zeta_{0i} \\ \zeta_{1i} \end{bmatrix} \\ \eta &= \alpha + \Gamma \xi + B \eta + \zeta \end{aligned}$$

を考えれば良い。

8.4 潜在成長モデリングの拡張

不規則な間隔で収集されたデータ 行列 λ_y と λ_x の負荷を適切な値に指定すれば良い。

媒介効果の検討 行列 B を導入することで、内生的構成概念を互いに予測し合うことが可能になる。