

データ解析のための統計モデリング (2008 年 10-11 月)

全 5 (+2) 回中の第 1 回 (2008-10-27)

# 生態学データ解析の統計モデリングとは?

久保拓弥 kubo@ees.hokudai.ac.jp

<http://hosho.ees.hokudai.ac.jp/~kubo/ce/EesLecture2008.html>

この講義のーとが「データ解析のための統計モデリング入門」として出版されました!

<http://hosho.ees.hokudai.ac.jp/~kubo/ce/IwanamiBook.html>

まちがいを修正し, 詳しい解説・新しい内容をたくさん追加したものです

## 今日のもくじ

- |                          |   |
|--------------------------|---|
| 1. 環境科学にとって統計学的なデータ解析とは? | 2 |
| 2. 統計モデルって何?             | 4 |
| 3. データ解析の勉強のやりかた         | 6 |
| 4. データ解析ソフトウェア R って何?    | 8 |
| 5. この全 5 回の授業では何をあつかうのか? | 8 |

(Appendix)

- |                           |    |
|---------------------------|----|
| A1. これまでの生態学における統計学の使われかた | 12 |
| A2. 私はどのようにデータ解析の専門家になったか | 14 |

この統計学講義のーと<sup>1</sup>は

- 私 (久保) が講義でしゃべることを忘れずのを防ぐ
- 皆さんがノートをとる手間をはぶく<sup>2</sup>
- 皆さんがあとで読みかえして復習する
- ネット上で誰でもダウンロードできるようにする<sup>3</sup>

ために準備しました。まあ, この講義のーともいろいろと不備なんかがあると思いますが, ネット上のこの講義ページ (URL は目次の上に, 入力がめんどくさければ「久保 講義のーと 2008」で検索してください) でまちがい修正だとか随時更新していければ, と考えています。また, 次回以降の講義ノートもこのへんにアップロードするので, それをあらかじめダウンロード

1. 教科書というよりもっといいかげんな講義のーとです。

2. 左右の余白は注釈・メモ・らくがきのためのものです。手を動かしてないと眠くなりますよ。

3. しかしこの教室にいないとよくわからない内容かもしれませんねえ。

& 印刷出力してもらおうと、当方としても助かります。これはこの講義専用のメイリングリストで連絡します。

この「生態学の統計モデリング」って講義でいったいぜんたい何をやるつもりなのか？ そうですねえ、スローガンの的には

- これからデータ解析をする大学院生たちのために
- 実際のデータ解析で使う統計モデルの考えかたを
- 基礎から最先端につながるような「地図」を示しつつ
- 理解できる統計学めざして

といったモノになってくれればなあ、と。こうやって皆さんの前でしゃべってる私のほうはそういうつもりなんだけど.....実際のところ、聴いてる皆さんのほうはなかなかすんなりとは理解できないかもしれません。ということで、わからないことがあったらどんどん質問してください。<sup>4</sup>

4. しかし院生の皆さんに質問を發してもらうのは、ときとしてかなりの難事だよなあ、という気分だったり。

## 1. 環境科学にとって統計学的なデータ解析とは？

この統計学モデリングの講義では、おもに生物の環境科学<sup>5</sup> というか生態学 (ecology) なんかの研究をしている大学院生に説明しているつもりでハナシをススメます。生態学固有の知識なんかはできるだけ持ちださないでハナシをしたい、と考えています。しかしながら、ときとして「個体群」とか生態学用語なんぞを使ってしまうかもしれないので、これまたわからなかったら臆せず質問してください。

5. これはいちおう「環境起学」専攻とやらの講義なんで.....

さて、ハナシの順番からいって、まずは「そもそも何で統計学だの統計モデリングだのが必要なのか？」といったところから始めてみましょうか。

生態学では野外・室内での観測や実験によって、自然現象から何か情報を取りだそうとするんですけど、このときに二段階の情報消失・情報圧縮が発生します (図1)。つまり.....

- 第一段階: 自然現象から観測・実験といった手段で情報をとりだして「データ」(数値・記号の集まり) に変えるときに多くの情報が失われる
- 第二段階: データを統計解析するときに多くの情報が失われる

まず第一段階では、生きている生物がもつ莫大な情報を人間が観測・測定してみることはほぼ不可能、それを実験室の簡単な条件下においても、ということ。たとえばそこらへんの地面に生えてる草について考えてみましょ

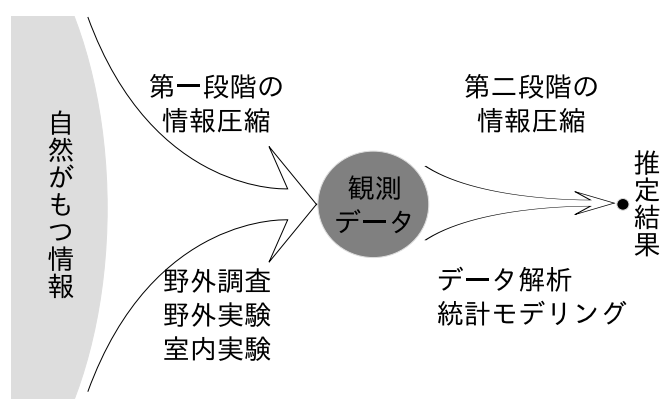


図 1: 生態学における二段階の情報消失・情報圧縮

う。高さを測定したり葉の枚数を数えたり光合成速度を計ったり，といったことは可能でしょうし，刈り取ってしまえば重量を測定したり，遺伝子を測定したりといったこともできるでしょう。

しかし何もかも調べることはできません。野外観測データをあつかったことがある研究者なら，われわれ人間はさまざまな観測記録機器を駆使しても自然がもつ情報のごく一部しか「データ化」できない，ということに同意してくれるだろうと思います。

何やらくたくたく説明してきましたが，この講義ではこの「第一段階の情報消失」については検討しません。そこまでいくと各分野での観測の方法論だの「研究によって明らかにしたいこと」だのにまで踏みこまねばならないでしょうし，そんな面倒は私も皆さんも回避したいところでしょう。

この講義では第二段階の情報消失，つまり統計解析によって情報が失われる，良くいえば「圧縮」される問題に集中したい，と考えています。<sup>6</sup> 自然現象に含まれていたはずの莫大な情報のごくごく一部が数値と記号の羅列すなわち「観測データ」とやらにいったん変換されてしまったならば，どんな学問分野で得られたデータであろうと，たいていの場合には統計学的手法を使ってそれを調べるほかないからです。これはじつに驚くべきことなのかもしれないですね。

で，豊潤な自然現象の「しぼりかす」みたいな観測データを何でまたわざわざ統計学的手法とかで情報消失・圧縮してやらねばならないのか，という問題なんですけど……ヒトこと言えば，われわれ人間はアタマが悪すぎてデータなる数値・記号の羅列をうまく把握できないからなんです。私なんかだとデータ量が 3 行 × 3 列ぐらいをこえると，もう何がなんだかよくわからなくなったりします。

ということで，「おそらく自然現象の一部を反映してくれているであろう<sup>7</sup> 観測データ，その観測データに見られる<sup>8</sup> さまざまな『パターン』をうま

6. ホントはこの統計解析の問題は第一段階にも影響する，つまりどう自然を観測・測定するか，にも密接に関わるのですが，ハナシが広がりすぎるのでやはり省略します。

7. これは単なる願望であることが多いんですけどね

8. あるいは見えてるのかどうかよくわからない

く説明できるような数理モデルがほしい」という状況になります。数理モデルってのは数式で書かれた比較的簡単なモデル，さらにそのモデルに「測定誤差」だの「個体差」だの「不確定性」だの「変異性」だのが混入してしまっているもの，つまり確率論的な数理モデルのことをこの講義では統計モデルとよぶことにしましょう。私が言いたいことをまとめてみると，

- 「自然科学という方法で自然現象を理解する」なる営為は結局のところ「情報を捨てることで人間でも理解できる<sup>9</sup> ようにする」こと
- とくにいったん「データ化」してしまえば，そこからの「情報捨て」は統計学的方法であつかうしかない—そしてここから先の情報処理は“客観的”，つまりあなたがやったことを私もまた同じように再現してみせることが簡単である，ということです
- とくに（実験による現象の再現が難しい）生態学など環境科学の研究者の世界では「自然を観察・観測して言えること」の少なくない部分は統計学的方法に依存していて，<sup>10</sup> 研究の信頼度を上げるための研究者の相互検証はこのデータ解析の再検証によることが多い

9. もしくは理解した「気分」になれる

後述するように計算機ハードウェア・ソフトウェアの進歩によって，1990年代から統計学的方法は大きく発展しました。生態学まわりでのデータ解析の方法も大きく変わってきました。いわば「りすとらされたデータ解析の世界」で，何から勉強していけばいいのかをこの講義で考えていきましょう。

11

10. もちろん現場で自然現象をみて「自分はここでこういうことが生じていると思った」という直接観察は重要です……しかしながら，その思索だけでは自然科学には含まれないですね。自然科学では数値化・記号化されたデータこそが（共有できるという意味で）「客観的なものだ」という共同幻想を維持していますので，こいつを「客観的」にモデリングしてパターンを抽出する必要があります。

11. このへんの導入部分では「ヤミにほうむられる修士論文」とデータ解析の関係，といったことも検討したのですが，時間ぎれで書けませんでした。これはまた別の機会に……

**現状** — 生態学研究まわりにおける

- 軽視されている（授業でも適切な方法を教えない）
- そもそも何やってるか**わかってない**ヒトが多い
- まちがっている方法に**固執**する（指摘すると逆ぎれ）

**理想** — この統計学授業のネらい

- 理念: スジのとあった**合理的**な統計解析をめざそう
- 手段: データの**性質・構造**によくあった手法を（データの有効利用）
- 目的: 自然現象うまく説明できる**モデリング**になってれば

図 2: これまでの状況，この講義でめざすところ

## 2. 統計モデルって何?

「観測データにもとづいて『何か』を主張する」ためには統計学的なデータ解析がたいへん重要な役割を果たすはず— にもかかわらず、生態学まわりでの統計学の使われかたを見ていると、ちょっと前までは何ともおそまつなものでした (図2 の左, あるいは Appendix でぐだぐだと書いていること) .

現状は必ずしもよろこばしいものではない, それがゆえに, この講義では「理念: スジのとおった」「手段: データをよくみて」「目的: 現象を説明する」統計モデルをめざしましょう (図2 の右), といったあたりを身につけてもらいたいデータ解析の姿勢としてかかげてみたいと思います. なんだかおかたいハナシのようですが, 観測データをこうやって解析していくのはじつに楽しいことという点も強調しておきます. <sup>12</sup>

さて, ここで重要になるのは統計モデルといった考えかたです (図3) . 統計モデルってのは何なのかを説明してみるとこんなかんじになるかと思ひます .

- 観測データを解析する, ってことは統計モデリングである
- 統計モデルは観測データのパターンをうまく説明できるようなモデル
- 統計モデルの基本的な部品は確率分布, 確率分布のカタチはパラメーターによって決まる
- 観測データをうまく説明できるようにパラメーターの値を決めることを「統計モデルのあてはめ」または「統計モデルによる推定」という

13

各項がのべようとしていることは, この講義の中でおいおい説明してみてもりです .

### データ解析は統計モデリング

- 統計モデルは観測データのパターンをうまく説明できるようなモデル
- 基本的部品: 確率分布 (とそのパラメーター)
- データにもとづくパラメーター推定, あてはまりの良さを定量的に評価できる

12. これに対して, スジのとおらない, 何をやってるのかワケのわからない, 現象を説明できてるのか判然とせぬデータ解析はまったくカナしいものです .

13. 観測データとの対応関係が明確ではない「数理モデル」をあつかっているヒトたちもいます. こういったモデルは, (1) 観測データからパラメーターを推定する手段, (2) モデルの予測とデータの乖離を定量化する手段, のあいだに整合性がありません. つまり観測データにあらわれるパターンを説明したいのであれば, (1) と (2) を最初から考えている統計モデルのほうがずっと有用なのです .

図 3: 統計モデルとは何か?

最初に一点だけ強調しておくなら、誰もがデータ解析をやるときに統計モデリングを回避できない、というあたりでしょうか。「もできる？ オレは自然現象そのものを調べてるから数理モデルなんかとは何の関係ないもんね」てなことを言ってるカンちがいぎみの研究者であっても、「まあこの論文にのせるデータ、そのまま table にしても reviewer に文句言われるから『いちおう検定にかけて』<sup>14</sup> 『ゆーい差』でも出しておくか」などと行動するのであれば.....このヒトは、確率分布をくみあわせた統計モデルを観測データにあてはめてみて、パラメータだの検定統計量だのを推定している、ということになるのです。つまり自分が何やってるのかわかっていない、そういうヒトたちが統計学的手法ユーザーのたいはんをしめています。

14. これはじつにひんぱんに聞かされる意味不明な表現.....というか「ゆーい差」にかかる枕詞?

この講義の目的としては、皆さんがいかなる状況にあっても「自分はいまどういう統計モデルをどうあつかっているのか、観測データとはどう対応しているのか」を明確に意識しつつデータ解析にとりくめるようになる、てなことがこの講義のゴールで実現できればよいのだろう、と考えています。

### 3. データ解析の勉強のやりかた

近ごろの北大の修論発表会や生態学会の大会の研究発表をみていると、大学院生たちのほうが統計学をよくわかっている、その指導教員たちにくらべて、といった気分させられます。そうなる理由としては、大学院生のほうが統計学をよく勉強していて、しかも図 4 に示しているように現在は統計学の勉強をすすめていくのに良い条件がそろっているためではないでしょうか。

教員世代より院生世代のほうが「わかっている」理由:  
新しいデータ解析を学ぶときの勉強環境がすぐれている  
(教員世代はもはや勉強しない..... ヒトが多い)

- 計算機のハードウェアがよい ..... 格段に速い!
- 統計解析のソフトウェアがよい ..... R!!
- インターネット上にいくらでも教材がある

皆さんは良い環境で勉強できる、ということです

2008-10-27

11 / 22

図 4: 今どきのめぐまれた勉強環境

簡単に説明してみますと.....そも

そもデータ解析は計算機のパワーが必要であり、この授業であつかうような手法は現代の「パソコン」ぐらいの計算能力がどうしても必要となります。

<sup>15</sup> また計算機のパワーを使いこなすためには良いソフトウェアが必要であり、(後述する) R のような優れた統計ソフトウェアが無料で簡単に入手できる、というのも昔なら信じられないことです。さらにインターネットを介して「世界中の統計学講義のーと」が読めるようになりました。<sup>16</sup> 数式ばかりの統計学教科書にひたすらとりくむ他なかつた教員世代とくらべると、現在はずいぶん勉強しやすい環境になったといえます。とくにデータを使っ

15. これは 20 年前であれば数億円の大型計算機に匹敵するものです。

16. これまでの久保「講義のーと」のたぐいも数百人(もっとたくさん?)のヒトにダウンロードしていただきました。

て試行錯誤，つまり手をうごかしてデータの図を作ったりしながら概念を理解していくことはとても有用です．ですから，皆さんもこの勉強環境をうまく活用してみてください．

この講義では「ふつーの統計学の講義」では扱わない「基本から最先端までを橋渡りする統計モデリングの流れ」を重視します．一方で，ふつーの講義で扱われることを，あまり説明しない傾向があるかもしれません（図 5）．そこでこの「欠けた部分」をおぎなうために，皆さんに統計学の本を買ってもらってそのあたりを自分である程度は勉強する必要があります．

これは「ヘン」な統計学授業です

- とにかく「統計モデリング」偏重!
- $t$  検定とか分散分析とかやりません
- 「正規分布」もほとんどでません

皆さん，他の統計学本とかも読んでね……

2008-10-27

14 / 22

図 5: 「ヘン」な講義によるこそ

さきほど述べたようにインターネット上でいろいろな統計学の講義資料などが配布されています．この講義の欠落部分をおぎなうためにはそういったものを読むことも有用だろうと思います．

また世の中には多数のすぐれた統計学本がでています．これらもぜひ活用してください．図書館や本屋でいろいろとばらばらと見てもらって，自分にとってわかりやすそうな入門的な教科書が一冊でもとにあればいざと勉強の役にたつだろうと思います．とくに図の多い統計学教科書は重要で，その意味では（もしわかりやすければ）「マンガでわかる……」のたぐいの統計学本も良いかと思います．

ここでは二冊ほど紹介してみます．まず一冊目は「統計のはなし」（粕谷英一，1998，文一総合出版），生態学研究者のあいだでは「ぴんく本」として知られるものです．これは読みやすく，また生態学まわりでの統計学的手法が概観できて，とてもおすすめの一冊です．巻末付録「君にもできるごまかし」もいろいろな意味でたいへん勉強になります．<sup>17</sup>

もう一冊は「統計（第二版）」（竹村彰通，2007，共立出版）です．統計学の専門家が書いた「統計的手法の使いかたを示しながら理論を説明」するもので「統計学を学ぶ際の動機づけ」を重視して「統計学的手法を用いるとどんなことができるか」を示しています．またここでは統計ソフトウェア R による解析例が示されています．

データ解析の手法の勉強するときに実践したら良いことを図 6 にまとめてみました．私がおすすめしたいのは実際にデータをあつかいながら，とくに

17. このぴんく本に影響されたヒトはたくさんいるので，皆さんの指導教員のアタマの中がどうなってるのか，それを予測するときにもたいへん有用です．

皆さんに実践してもらいたいこと

- 観測データの図をたくさん作ろう
- データ解析で試行錯誤，他人と議論しよう
- 以上を実現すべく R を使いこなそう

2008-10-27

15/ 22

### 図 6: データ解析の勉強するときに...

図を描くという作業をしながら勉強をすすめていくことです。そして他人との議論も重要です。それは人によって「データを使っての試行錯誤」体験がずいぶんと異なっているからで、そのあたりの自分・他人の「ずれ」から多くのことが学べます。<sup>18</sup> ひきこもっているのはデータ解析の勉強はできません。このように勉強をすすめるとき、そしてもちろん研究のデータ解析をすすめるときに有用なのが R と呼ばれるソフトウェアです。

18. 私はいつも大学院生たちからの質問、とくにふうがわりな質問に勉強させられます..... いや、ホントにそう思います。

## 4. データ解析ソフトウェア R って何?

この授業では統計ソフトウェア R の利用を前提にハナシをすすめていきます。<sup>19</sup> みなさんも自分の PC にさっそく R をインストールしてください。インストール方法は

<http://hosho.ees.hokudai.ac.jp/~kubo/ce/HowtoInstallR.html>

なんか参考になるだろうと思います。<sup>20</sup>

R は統計ソフトウェアは実際の解析だけでなく、統計学的手法を勉強するときにも必要なものです。<sup>21</sup>

いろいろあるソフトウェアの中で R を皆さんにおすすめします。というか、この講義の中では R を使うことを前提にハナシをすすめます。R の利点は

- free software である — つまり誰でも無料で入手でき、しかもそのしくみ (source code) が完全に公開されている<sup>22</sup>
- 強力無比な作図機能がある — 「データを図にする」ことはデータ解析においてとてもとても重要<sup>23</sup>
- プログラミング可能

19. R については、今回 (印刷出版前にこっそり) 配布する「R で改善する生態学のデータ解析」も参照してください。

20. このあたりの URL などとはあとでメイリングリストでお知らせします。

21. いろいろな方法を適用する「実験」ができますからね。

22. このことは学術研究においては特に重要です。

23. この講義の中でもそれを強調していきたいと思います。



- add-on package library が充実していて，機能拡張が容易

といったことです．データ解析の勉強のおともとして，あるいは自分の研究の強力な仲間としてたいへん頼りになるソフトウェアです．

## 5. この全 5 回の授業では何をあつかうのか?

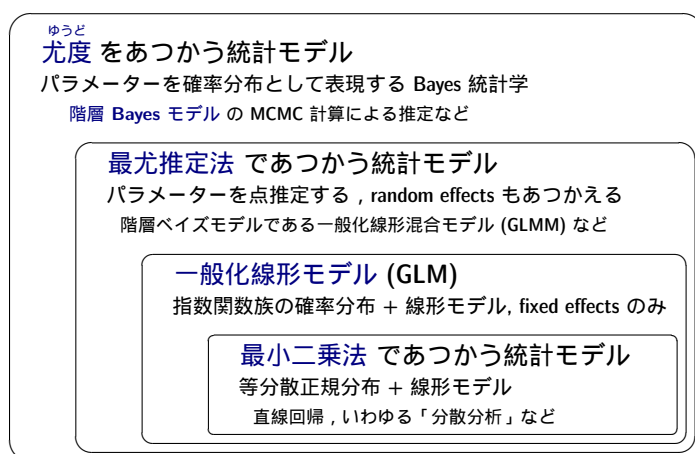
この講義では R を駆使して，「自分のアタマで理解して使う統計学」をめざしたい，と考えています．同時に，観測データを不自然にひねくりまわさない「素材のよさ」をいかしたデータ解析ができるようになってくれば，とも考えています．

そのあたりを実現するために，今回の全 5 回の統計学の講義は

- 確率分布で表現される統計モデルを重視する
- 統計モデルがデータにどれだけあてはまっているかを <sup>ゆうど</sup>尤度 で評価する
- 線形モデルを中心にあつかい，これを一般化線形モデル → 一般化線形混合モデルと拡張していく <sup>24</sup>

という基本方針で，図 7 に示しているような統計モデルの「地図」の内側から外側にむかってすすめていく予定です．この図をみれば何となくわか

24. といった予定だったのですが，後述しているように一般化線形混合モデルの部分は消えました．



2008-10-27

12 / 22

図 7: 統計モデルの「地図」

るように，上でちょっとマンガふうに描写した「指導教員たちの統計学の世

界」ってのは図の中のごく狭い領域「最小二乗法であつかう統計モデル」に限定されている，といったことを示しているつもりです．この講義では，そういう「最小二乗法であつかう」狭い領域は一般化線形モデル (generalized linear model; GLM) の一部にすぎない，というあたりも説明できれば，と考えています．

とりあえず全 7 回の流れに対応させてこのあたりの概要を説明してみましよう．

### 1. 10/27 (月) 生態学データ解析の統計モデリングとは?

今回ですね．統計モデリングの必要性，生態学における今までの統計学的手法の使われかたについて説明してみました．このあと，統計学勉強に必要な本とソフトウェアについて説明するつもりです．

### 2. 10/30 (木) さまざまな確率分布と最尤推定

今回は統計モデルの「基本部品」というべき確率分布について説明します．世の中にはたくさんの確率分布があって，ホントはいろいろ紹介したいのですが，この講義で使う 3 つの確率分布，すなわち「ポアソン分布 (Poisson distribution)」「二項分布 (binomial distribution)」「正規分布 (normal distribution / Gaussian distribution)」についてだけ説明する予定です．またデータへの「あてはまり」をあらわす尤度<sup>ゆうど</sup>について説明し，これまでの統計学のパラメーター推定の基本的な考えかた，最尤推定法<sup>さいゆう</sup>についても説明します．<sup>25</sup>

25. 11/3 (月) は祝日なので講義はありません．

### 3. 11/06 (木) 一般化線形モデル (GLM) 1

ポアソン分布で説明されるであろう現象<sup>26</sup>の統計モデリングについて説明します．ポアソン分布モデルのパラメーターを最尤推定することを，ポアソン回帰といたりします．一般化線形モデル (GLM) というのは (線形) ポアソン回帰とかこの次の回の (線形) ロジスティック回帰など線形モデルのたぐいをいっしょくたにして呼ぶときの総称です．

26. 観測データが「0 個, 1 個, 2 個, ...」というふうに整数として数えられるカウントデータであるとき，など．

### 4. 11/10 (月) 一般化線形モデル (GLM) 2

よく使われるもうひとつの GLM であるロジスティック回帰について説明します．これは「生物が刺激に应答した・しなかった」とか「観察された  $N$  個体のうち  $y$  匹が死亡した」といった「あり・なし」データを解析する統計モデルにもとづいていて，この統計モデルは二項分布を基本部品とする統計モデルです．

## 5. 11/13 (木) 検定とモデル選択

ある観測データを説明できる統計モデルはたいてい複数あります。このときに複数の統計モデルをどう関係づけるのか、伝統的な Neyman-Pearson わくぐみの「検定」の正体を説明します。さらにまた別のモデル評価法、情報量規準をつかったモデル選択についても説明します。

## 6. (消された講義その 1) 一般化線形混合モデル (GLMM)

データ解析で GLM を使いこなせるようになってくると、GLM ではうまく説明できない現象あるいは観測データ内のパターンがあることに気づきます。これらは「個体差」「場所差」など人間が測定してない / 測定できない要因によって生じるばらつきだったりします。そこで「個体差」「場所差」を random effects として GLM に組みこんで一般化線形混合モデル (generalized linear mixed model; GLMM) に拡張します。

## 7. (消された講義その 2) 階層ベイズモデル

データ解析に GLMM が使えるようになると、今度は観測データのあちこちに random effects になりうる要因があることに気づきます。たとえば野外観測データの場合「場所差」があってその中に「個体差」がある、といった状況です。このように統計モデルが複雑になってくると、最尤推定法によるパラメータ推定がだんだん「苦しく」なります。そこで、GLMM を階層ベイズモデルの一種だとみなして、Markov Chain Monte Carlo (MCMC) 法を使ってパラメータの事後分布をもとめる方法を紹介します。

### 説明したい統計モデリングのお作法

- 観測データをどんな確率分布で表現できるか考えよう
- 「割算値」の統計モデリングはやめよう
- むやみに「グループ平均」とるのはやめよう

つまり観測データの「もち味をいかした」  
「ひねくりまわさない」統計モデリング

図 8: この授業で強調する統計モデリングの「お作法」

最後にもう一度この授業のネライみたいなものをまとめてみましょう。私がおススメしてみたいデータ解析の姿勢(図 2 の右)としては、「理念: スジ

のとあった」「手段: データをよくみて」「目的: 現象を説明する」統計モデリングをめざしましょうといったあたりです。データ解析の勉強(図 6)としておススメしたいのはデータを使った試行錯誤と他人との議論です。そして統計モデリングお作法としては、図 8 にまとめているように、「確率分布の重視」「なんでも割算するな」「なんでもグループ化するな」といったところ です。次回からはこれら三つの注意点を実現できるようにハナシをすすめていきます。

とりあえず、今日はここまでにしませう。

## A1. これまでの生態学における統計学の使われかた

本文で「もっと観測データと対応のつく数理モデルにしる」とか言ってたヒトたち自身がやっているデータ解析は、けっこう悲惨なモノとのべました。実際のところどういうものなのか、「何をやってるか」というより「どのようにやってるのか」あたりの行動観察の結果を示してみましよう。以下の流れは、かなり多くの生態学研究者で見られるものです。<sup>27</sup>

1. 自分と同じようなデータをとってる論文をひたすらさがす
2. 見つけたら「どういう『検定』を使っているか」をチェックする（これを検定 X としましよ）
3. 「統計ソフトウェアで検定 X の結果として という output が得られた → その場合は と断定してよい」というパターンをひたすらまる暗記する<sup>28</sup>
4. 検定 X ができる統計ソフトウェアを買ってくる
5. 統計ソフトウェアにデータを流しこむ
6. 検定 X つかった output が得られたら 3. で暗記したパターンを適用して論文をかき、投稿する
7. reiviewer が「検定に文句をつけたら」、reviewer が指定する検定を新しく検定 (新) X と指定しなおして、ひたすら 3.- 7. をくりかえす<sup>29</sup>

何とも「悪しき行動ばたーん」と言いたいところですが、<sup>30</sup> ここまでくるとむしろ何かこっけいなもののように思えてしまいます。もちろんやってる当人はごく真摯なもので、むしろ上のような行動パターンこそが「データ解析の必勝パターン」とでも信じているのでしょう。

上の一連の行動では何やら「検定」 & 「ゆーい差」なるもの<sup>31</sup> ばかりが重視されているみたいですが、この部分を極限まで押しすすめてしまうヒトたち<sup>32</sup> なんかもいたりしまして、そういうヒトたちのが示している奇妙な行動をちょっと列挙してみましよう。これらもまたことごとく私が実際に目撃・観察した事例の一部です。

- 「ゆーい差」が出るまで「検定」手法をひたすらとりかえる
- データ中の観測値どうしを割算したりして新しい「指標」をでっちあげる…… 「ゆーい差」がでるまで新発明・珍発明をくりかえす
- 都合の悪いデータをとりのぞく
- 「ゆーい差」が出るまでデータサンプリングをひたすら続ける

27. そしてこの教室にいる大学院生の多くも、修論とか書くときにはこういう行動とるんじゃないか、という悲観もあります。

28. そして、これらはだいたいにおいてまちがった「チャート式」になってたりするんですよ。

29. reviewer が指定する検定 (新) X をやるためだけに、十数万円もする統計ソフトウェアを新しく買う……私が目撃した事例のひとつです。

30. あるいは「検定 X の正体は何なのか、まったく理解してなくても論文が書けてしまう」ひどく巧妙なる方法論、と言えるのかもしれない。

31. これらはあとの回の講義でとりあげます

32. 「ゆーい差」決戦主義者…… まあ、これにもいろいろなタイプがあるのですが。

このヒトたちがここまでして「ゆーい差」にこだわるのは

- 「ゆーい差」さえ「出せば」、あとは何でもかんでも好き勝手に断定してみせてよい、という思いこみがある— 逆に「ゆーい」じゃないと不安で不安でしょうがない
- 統計学のこととかよく知らないけど、「ゆーい差」ならわかる<sup>33</sup>

33. 実際には「ゆーい差」についてもたいしてわかってないことが多いです。

といったキモチがあるようです。

また、大学院生にとって多くの場合、自分の研究の方向性を左右したり、データをどうあつかうかを決めるのは指導教員ですよね……ということで、その「ありがちな統計学的世界」を知っておくのは有用だろうと思います。

以下に示されてる図式はほとんどまちがいのなので、その意味するところがわからなければマネしないほうが安全だろうと思います。

- 「検定」はたいていの場合、正規分布を仮定すればよい
  - 値 X vs 値 Y みたいなデータはことごとく直線回帰すればよい
  - 「処理区」「無処理区」というデータは ANOVA すればよい
  - しかしこれらはどう違うのか理解できない
  - それはともかく、こういう解析やったときに統計ソフトウェアが吐き出す  $R^2$  値は「説明力」で、それが大きけりゃなんでもいい
- 「等分散じゃない」とか文句をつけられたら、データを対数変換だの arcsin 変換して回帰や ANOVA すればよい
  - あるいはめんどうになったら割算値とかぼんぼん作って「のんぱらめとりっく検定」やればよい
- 「検定を何度もやってる多重比較だ」と文句つけられれば、何でもかんでも多重検定法による  $p$  値補正をやればよい
- データがちょっと複雑な table であれば、何でもかんでも  $\chi^2$  検定やればよい
- もっとめんどうなデータは何でもかんでも多変量解析やればよい

実際のところ、こういうヒトたちのたいへんなところは、上にあげたことの当否というより、

- 上のようにデータ解析して論文が accept されてきた、つまり自分自身の「成功体験」だと信じている

- とにかく上のような「世界」を固守していればよい
- それゆえに、もはや統計学の勉強なんかする必要がない

つまり「アタマがかたい」というべき状態にあることかもしれません。

大学院生にとって重要なことのひとつは「大学院という牢獄」から脱出することなので、指導教員には逆らわないでその言いなりになって研究をまとめて卒業する、というのかなり現実的な方策だろうとは思いますが。そういう大学院生たちには、残念ながらこの講義はそれほど有用ではないかもしれません……というか統計学の勉強なんかする必要ありませんよね。

## A2. 私はどのようにデータ解析の専門家になったか

この節ではどう勉強したかというより、まあ、私の研究者としての背景、みたいところを説明してみまじょうかね。

私はデータ解析つまり統計モデリングをもっぱらとしていて、観測データにみられるいろいろなパターンをうまく説明する、という部分が現在の専門といえるでしょう。

しかしながら、もともとそういう統計学なんかの教育を受けたわけではなくて、大学院生として所属していたのは「数理生物学」なるものを専門とする研究室でした。この数理生物学ってのは現実ばなれというか、観測データとの定量的な対応とかぜんぜん気にしてない分野です。<sup>34</sup> もちろん数理モデルを使えばいろいろなパターンを生成できますから、「やあ、何だかこれは自然現象に見られるパターンと『似てる』ぞ」とか学会発表してみたりしたんですけど、実際にホンモノの生物・生物集団をあいてにしているヒトたちからは「ふーん、へー」とあまり感心されない。親切なヒトなんかはわざわざ「もっと観測データと対応のつく数理モデルにしる」とか言ってくれたりしました。

そりゃ、たしかに当方が「似てる、似てる」などと騒いでみせたところで、それはぜんぜん定量的なハナシでも何でもありません。そして、観測データなるものを定量的にあつかおうとすると……これはやはり統計モデリングするしかない、ということに大学院生だった私でもいやおうなく気づかされたわけです。で、独学で統計学の勉強なんぞやったりすると、数理生物な研究室のヒトたちからは「データ解析なんかやって何がわかるんだ!」とか言われたりするんですね。<sup>35</sup>

ともかく、統計学を勉強して野外なんかの観測データをあつかった論文をすみずみまできちんと理解して読める<sup>36</sup> ようになってくると、ですねえ……

34. なんというか「身内」だけに分かるいろいろな方程式のたぐいをヒネくりまわすことばかりやっている分野、と戯画化できるのかもしれませんが。

35. 実話です。

36. ありていに言ってデータ解析の手法が理解できていない人は自然現象をあつかった論文をごく不十分にしか読めていません。

「もっと観測データと対応のつく数理モデルにしる」とか言っていたヒトたち自身がやっているデータ解析ってのは、じつはかなり悲惨なモノだといったことが理解できるようになりました。このヒトたちはなるほどデータを得るのに多大な苦労したんだろうけど、それを解析の段階で「めちゃくちゃ」にしている、と。<sup>37</sup>

このあたり、つまり生態学研究者たちの「統計学の使いかた」といった行動観察についてはまたあとで述べますが……個人的には大学院生ぐらいのころは「データ解析なんかやって何がわかるんだ!」と決めつけられる理不尽とか「もっと観測データと対応のつく数理モデルにしる」と言いつつだめだめなデータ解析ばかり横行してる現状、に対する何というか反抗的なココロがデータ解析・統計モデリングの勉強をすすめていく動機の一部でした。

しかし統計モデリングに関する理解が深まり、野外調査や実験をもっぱらとする研究者たちとの共同研究ができるようになって、その観測データからいろいろなパターンが抽出できたり、それをうまく統計モデルで簡単に説明できるようになってくると……まあ、上のようなヒトたちは取るにたらない、わざわざ相手するほどでもないヒトたちだとわかってきました。

そんなことよりも、データをとっているヒトたちとこういう統計モデリングの考えかたを共有していっしょに改善していく、ということのほうがよほど重要だろう、と。とくに大学院生なんかはまだアタマが柔軟なヒトが多いですから、<sup>38</sup> いろいろ議論したりやりとりしてるうちにその指導教員なんぞはあつというまに抜きさるほど統計モデリングに対する理解を深めたりしてくれる。これは何度みてもホントにおもしろいことです……ということで、私としてはこの講義で何人かの大学院生やひとつデータ解析を勉強してやろうと思ってるヒトたちが、統計モデリングに関して今までよりすっきりとした理解ができるようになってくれれば、私にとっても皆さんにとってもおもしろいことだろう……といったコトなんかこの講義をやる動機になっています。

37. じつは他ならぬ私自身なんかも、昔の自分の論文とかみると「これはいくらなんでもヒドすぎる」と思うことが多々あるのですが……

38. 実際のところ、私のところにデータ解析質問メールなんかをくれるのは、たいてい大学院生だったりします。