

統計モデリング入門 2013 (1)

観測されたパターンを説明する統計モデル

久保拓弥 (北海道大・環境科学)

kubo@ees.hokudai.ac.jp

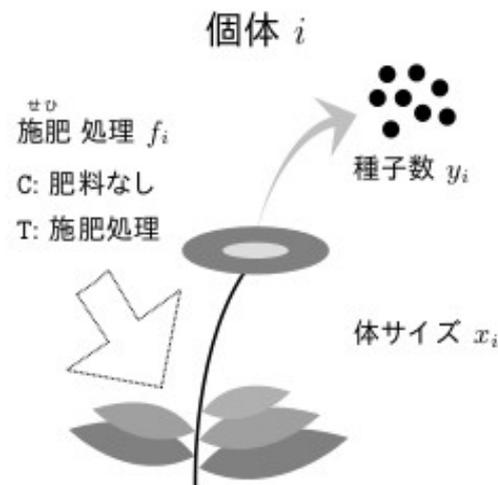


図 3.1 この例題に登場する架空植物の第 i 番目の個体. この植物の体サイズ (個体の大きさ) x_i と肥料をやる施肥処理 f_i が種子数 y_i にどう影響しているのかを知りたい.

統計モデリング授業の web page

<http://goo.gl/82dgC>

更新: 2013-06-27 15:24:25

生態学のデータ解析 - 統計学授業 2013

- [統計学の授業](#) やります (2013 年度前期後半, 2013 年 7 月)
 - 「植物生態学特論 I」の一部, 2013 年の 7/1 (月) から開始
 - 誰でも参加できます
 - 事前の申しこみなどは何も必要ありません (聴講のみで, 単位を必要としない場合)
 - 全部ではなく部分的に参加してもらってもかまいません
 - 教科書「[統計モデリング入門](#)」
 - 参考: [2008 年の講義の一と](#)
 - 月曜日・水曜日の 3 講目 (13:00-14:30)
 - 教室: [地球環境科学研究所 A 棟 A809](#) (エレベーターですすぐ前の部屋)
 - 短縮 URL: <http://goo.gl/82dgC>

【おもな内容】

- [第 1 回: 7/01 \(月\) 観測されたパターンを説明する統計モデル](#)
- [第 2 回: 7/03 \(水\) 確率分布と最尤推定](#)
- [第 3 回: 7/08 \(月\) 一般化線形モデル: ポアソン回帰](#)
- [第 4 回: 7/10 \(水\) モデル選択と検定](#)
- [第 5 回: 7/17 \(水\) 一般化線形モデル: ロジスティック回帰](#)

この統計モデリング授業の Mailing List (ML) **kubostat**

- 授業登録している人たちは自動的に ML に登録します
 - I will subscribe all registrants of the class to the **kubostat** ML.
- ML を使って各回の「課題」を出します
 - 回答もメールで送信してください
- 成績評価は「課題」の回答
 - 出欠関係なし（欠席の連絡いりません）
- 単位とらない人も ML 登録してください
 - 講義資料のダウンロード案内などあります

The main language of this class is
Japanese ... Sorry

- I didn't know foreign students are in the class.
- Why in Japanese? ... because even in Japanese, statistics is difficult for Japanese students to understand.
- I compensate for language disadvantages in foreign students when I give grades.
- Questions in English is welcomed.

統計モデルは データ解析の道具

なぜデータ解析の方法を
勉強しなければ
ならないのか？

科学のデータ解釈は統計的手法に依存

「データ→結論」のつなぎめ

- データ解析がおかしいと結論もおかしい
- データ解析を悪用して結論をねつぞうできる
- 論文を読むときにデータ解析の部分がわからないと「どうしてこのデータからこの結論が導かれたのか、妥当といえるのか」などがわからない→論文を批判的に読めない

データ解析はあまり重視されてなかった
内容がわからなくてもソフトウェアにまるなげ

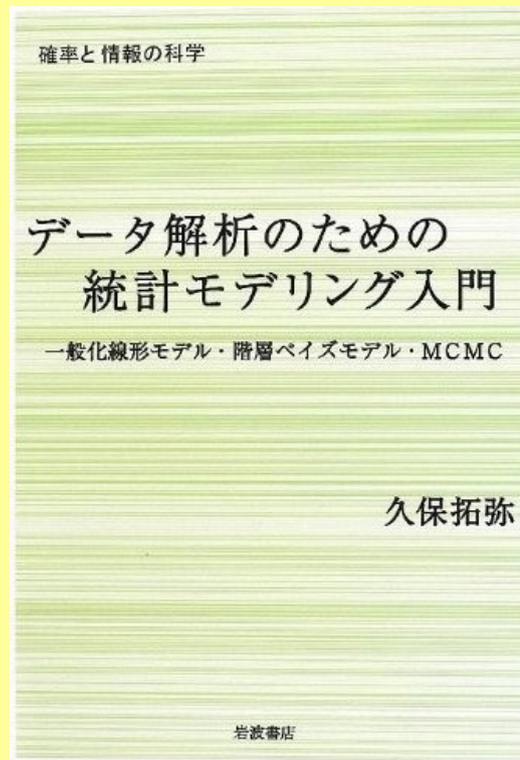
- ブラックボックス統計解析
- とにかく「ゆーい差」さえ出せばよいという
発想になっている
- 大学・大学院でもあまりちゃんと教えられて
いない, 教えられるヒトが少ない……とくに
近年発達している統計モデリングについて

この授業のねらい

できるだけ内容を理解して統計ソフトウェアを使おう！

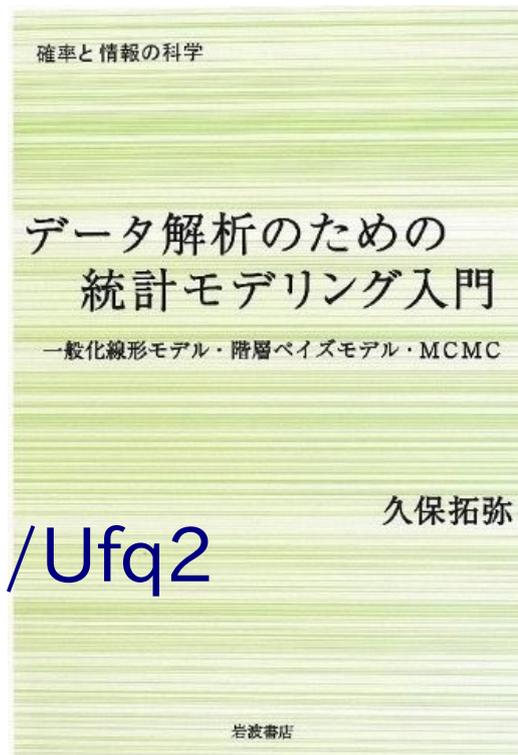
- データ解析で使われるの中でも比較的簡単な統計モデルを理解しよう
- 「ゆーい差」さえ出せばよいという発想をやめて、データと統計モデルの対応関係をよく見よう（作図重要）
- 統計ソフトウェア R を使い始めよう

教科書とソフトウェア

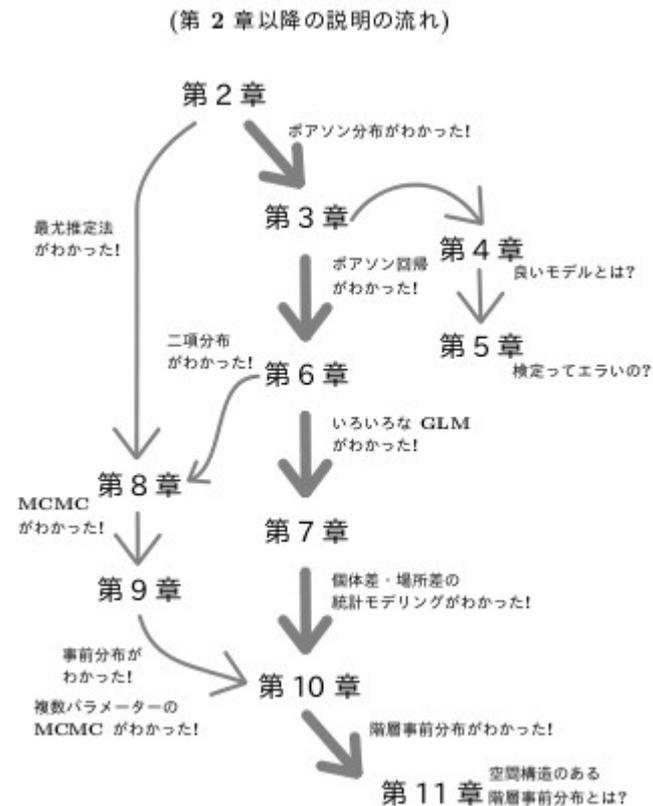


この授業は「統計モデリング入門」 にそった内容を説明します

著者：久保拓弥
出版社：岩波書店
2012-05-18 刊行
価格 3990 円

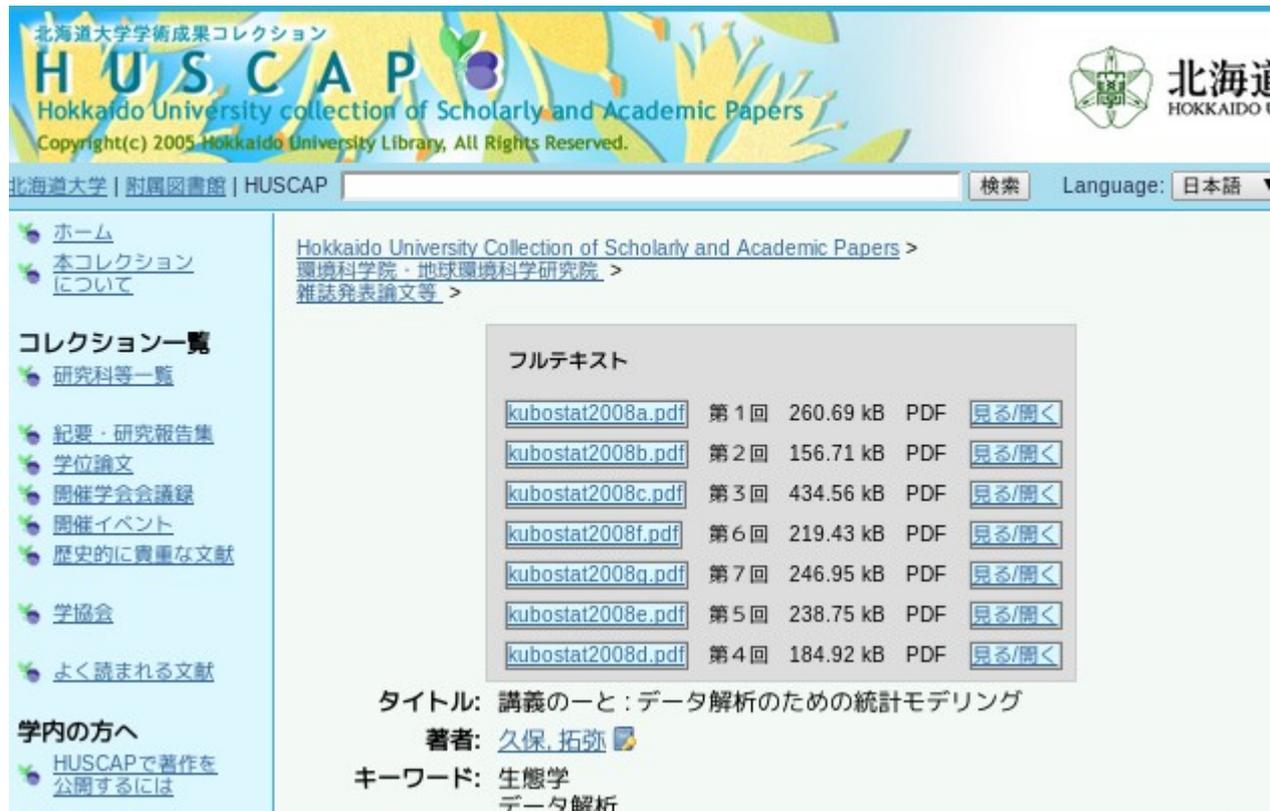


<http://goo.gl/Ufq2>



割引販売 3000 円!!

「統計モデリング入門」のもとになった「講義のーと」もあります



北海道大学学術成果コレクション
HUSCAP
Hokkaido University collection of Scholarly and Academic Papers
Copyright(c) 2005 Hokkaido University Library, All Rights Reserved.

北海道大学 | 附属図書館 | HUSCAP

検索 Language: 日本語

Hokkaido University Collection of Scholarly and Academic Papers >
環境科学院・地球環境科学研究所 >
雑誌発表論文等 >

コレクション一覧
研究科等一覧
紀要・研究報告集
学位論文
開催学会会議録
開催イベント
歴史的に貴重な文献
学協会
よく読まれる文献
学内の方へ
HUSCAPで著作を公開するには

フルテキスト

kubostat2008a.pdf	第1回	260.69 kB	PDF	見る/開く
kubostat2008b.pdf	第2回	156.71 kB	PDF	見る/開く
kubostat2008c.pdf	第3回	434.56 kB	PDF	見る/開く
kubostat2008f.pdf	第6回	219.43 kB	PDF	見る/開く
kubostat2008g.pdf	第7回	246.95 kB	PDF	見る/開く
kubostat2008e.pdf	第5回	238.75 kB	PDF	見る/開く
kubostat2008d.pdf	第4回	184.92 kB	PDF	見る/開く

タイトル: 講義のーと : データ解析のための統計モデリング
著者: 久保, 拓弥
キーワード: 生態学
データ解析

授業 web page に「講義のーと」へのリンクがあります! <http://goo.gl/82dgC>

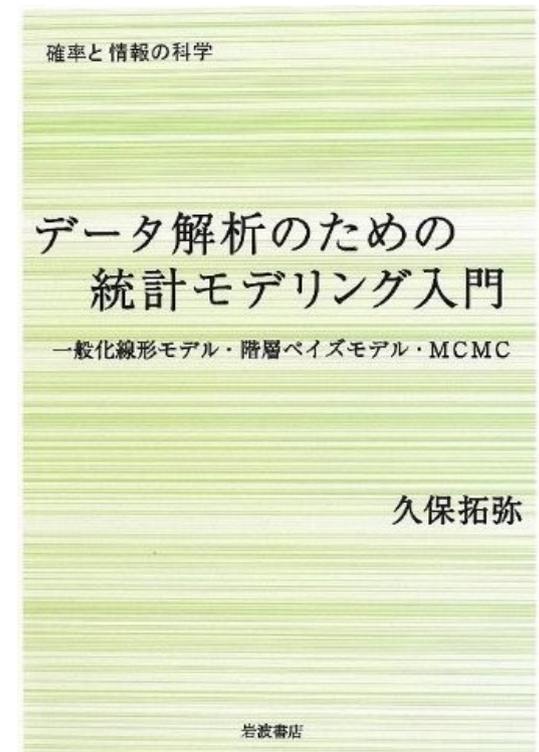
統計ソフトウェア R



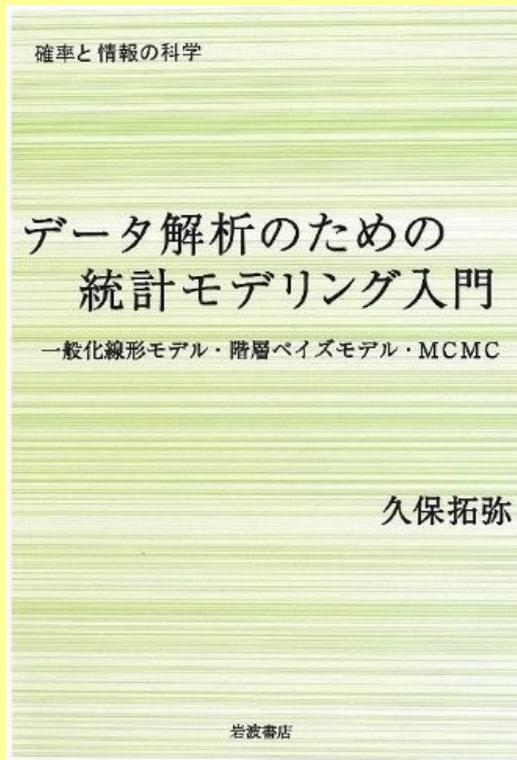
統計学の勉強には良い統計ソフトウェアが必要!

- 無料で入手できる
- 内容が完全に公開されている
- 多くの研究者が使っている
- 作図機能が強力

この教科書でも R を
使って問題を解決する
方法を説明しています



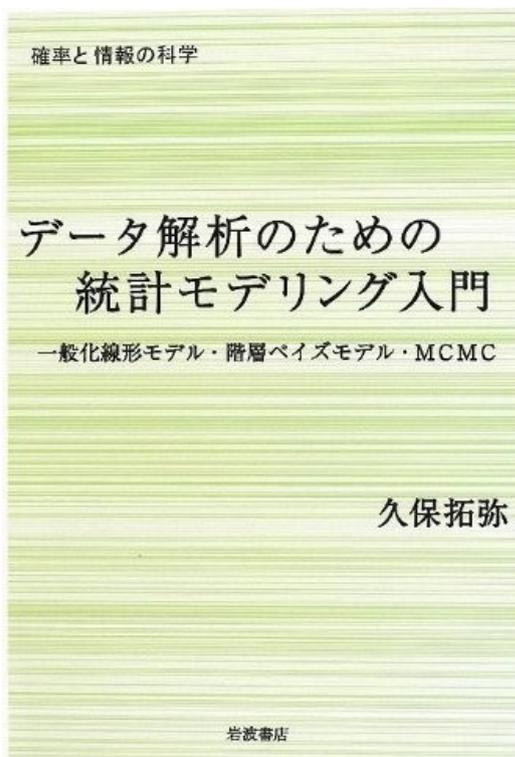
統計モデルとは何か？



「統計モデル」とは何か？

どんな統計解析においても
統計モデルが使用されている

- 観察によってデータ化された現象を説明するために作られる
- 確率分布が基本的な部品であり、これはデータにみられるばらつきを表現する手段である
- データとモデルを対応づける手つづきが準備されていて、モデルがデータにどれぐらい良くあてはまっているかを定量的に評価できる



「統計モデリング入門」の主張

「何でも正規分布」じゃないだろ!

線形モデルの発展

推定計算方法

MCMC

階層ベイズモデル

もっと自由な
統計モデリン
グを!

一般化線形混合モデル

最尤推定法

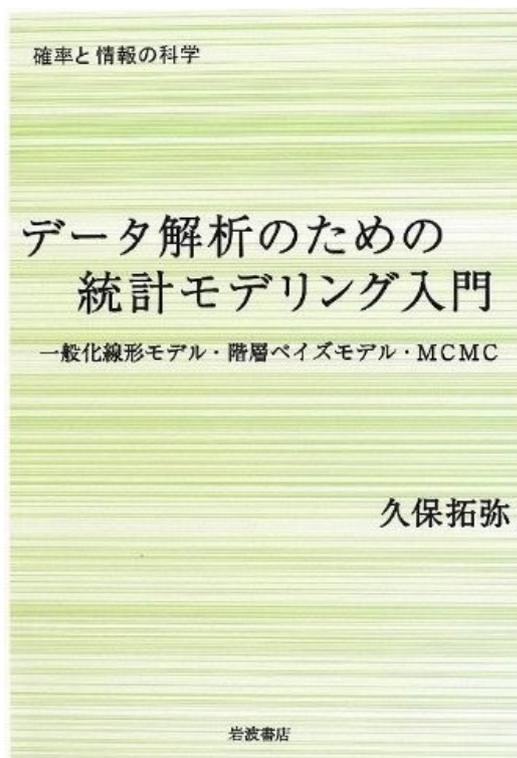
個体差・場所差
といった変量効果
をあつかいたい

一般化線形モデル

正規分布以外の
確率分布をあつ
かいたい

最小二乗法

線形モデル



たとえばこんなデータがあったしましょう

(次の時間の例題)

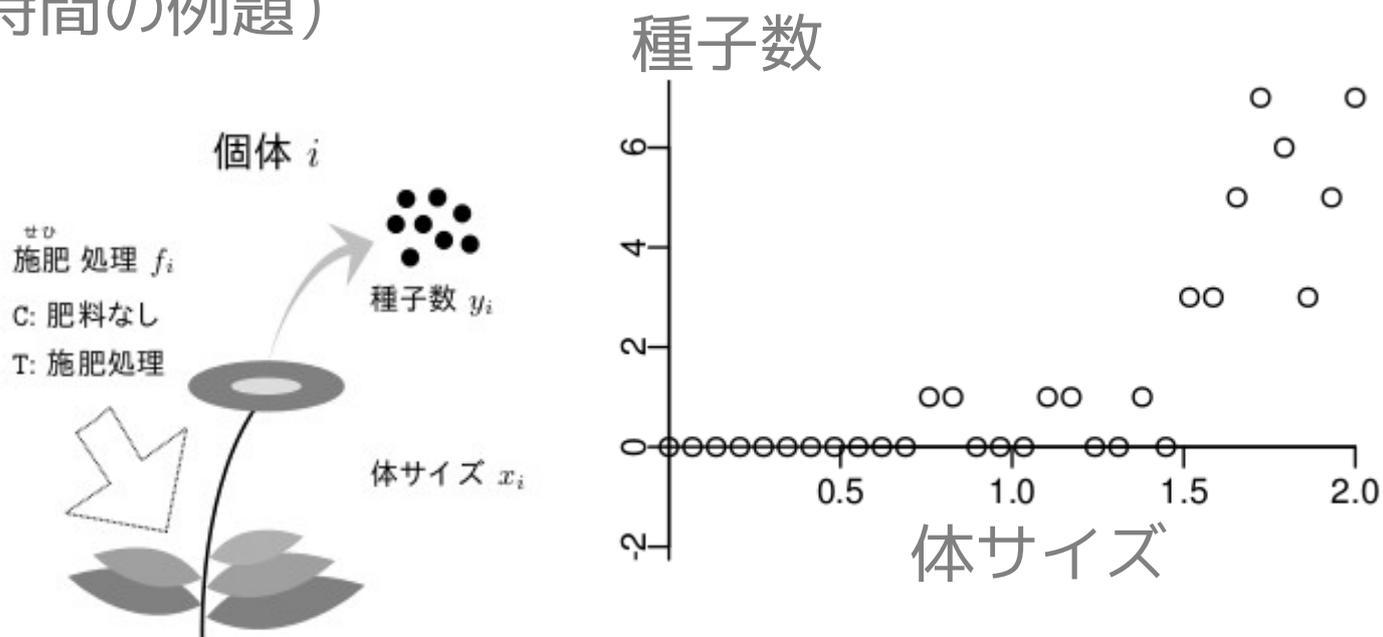
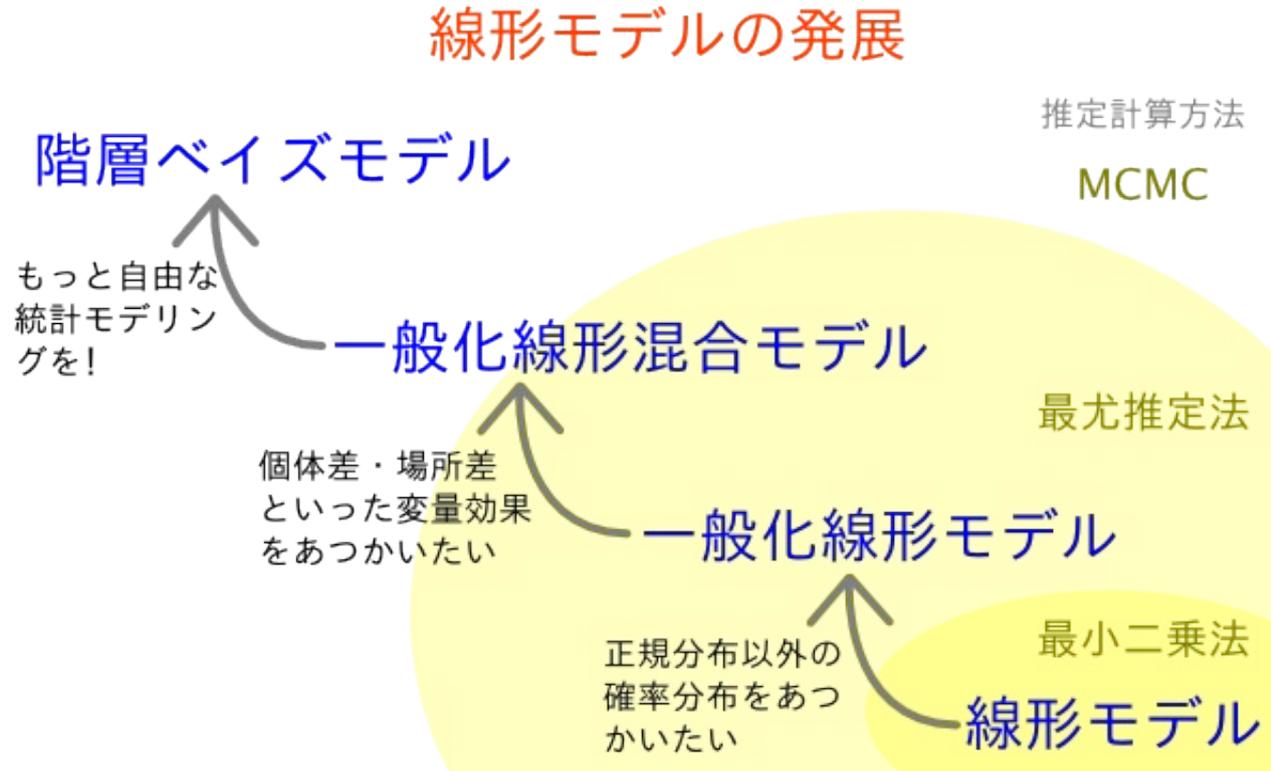
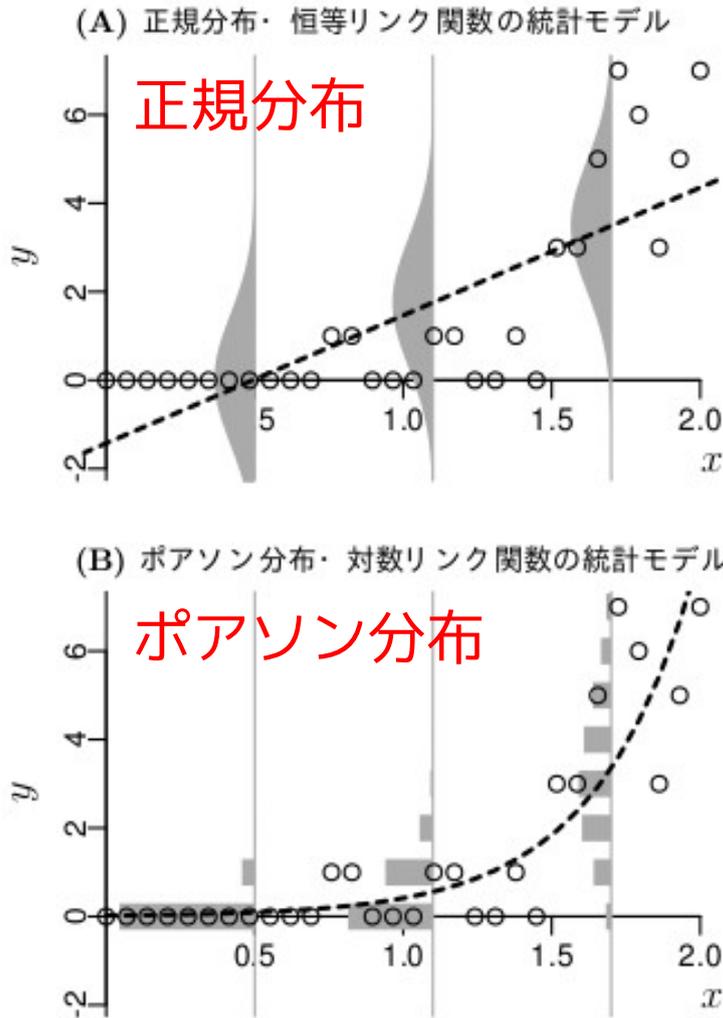


図 3.1 この例題に登場する架空植物の第 i 番目の個体. この植物の体サイズ(個体の大きさ) x_i と肥料をやる施肥処理 f_i が種子数 y_i にどう影響しているのかを知りたい.

一般化線形モデル - ばらつきをよく見る

Generalized Linear Model, GLM



0 個, 1 個, 2 個と数えられる種子数が「正規分布」なわけないだろ!!

3.9 回帰モデルと確率分布の関係. また別の架空データに対して GLM をあてはめた例. 破線は x とともに変化する平均値. グレイで

全体の流れ

- 第 1 回: 7/01 (月) 観測されたパターンを説明する統計モデル
- 第 2 回: 7/03 (水) 確率分布と最尤推定
- 第 3 回: 7/08 (月) 一般化線形モデル: ポアソン回帰
- 第 4 回: 7/10 (水) モデル選択と検定
- 第 5 回: 7/17 (水) 一般化線形モデル: ロジスティック回帰
- 第 6 回: 7/22 (月) 一般化線形混合モデル
- 第 7 回: 7/24 (水) 階層ベイズモデル

統計モデリング入門 2013 (2)

確率分布と最尤推定

久保拓弥 kubo@ees.hokudai.ac.jp

北大環境科学院の講義 <http://goo.gl/82dgC>

2013-07-03

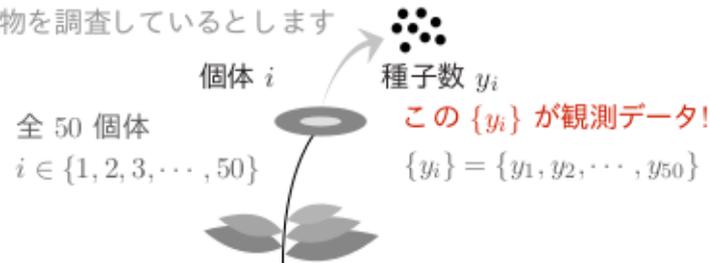
ファイル更新時刻: 2013-06-25 15:00

単純化した例題

例題: 種子数の統計モデリング まあ、かなり単純な例から始めましょう

こんなデータ (架空) があったとしましょう

まあ、なんだかこういうヘンな植物を調査しているとします



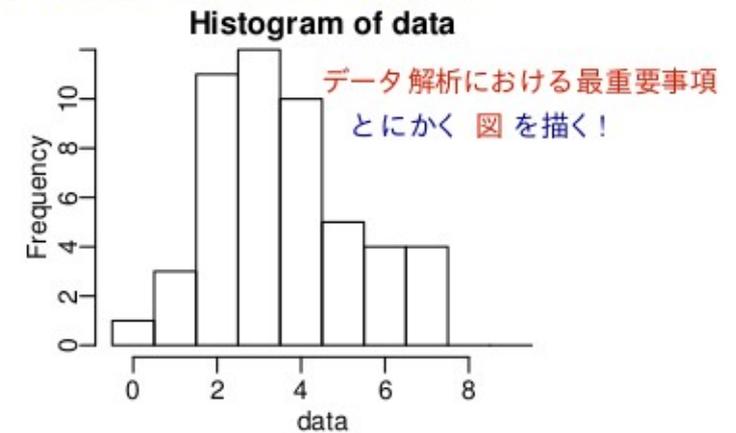
このデータ $\{y_i\}$ がすでに R に格納されていた、としましょう

```
> data
[1] 2 2 4 6 4 5 2 3 1 2 0 4 3 3 3 3 4 2 7 2 4 3 3 3 4
[26] 3 7 5 3 1 7 6 4 6 5 2 4 7 2 2 6 2 4 5 4 5 1 3 2 3
```

例題: 種子数の統計モデリング まあ、かなり単純な例から始めましょう

とりあえずヒストグラムを描いてみる

```
> hist(data, breaks = seq(-0.5, 9.5, 1))
```



カウントデータはポアソン分布を使って説明できないかを調べる

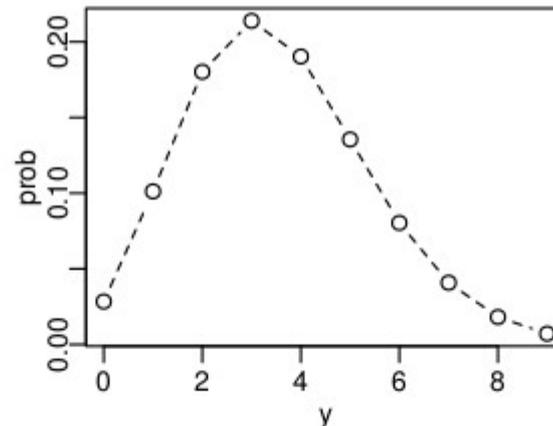


図 4 平均 $\lambda = 3.56$ のポアソン分布. 種子数 y とその確率 prob の関係が示されている. 図 4 の表を図にしたもの. R の `plot()` 関数の引数, `type = "b"` によって「丸と折れ線による図示」, `lty = 2` によって「折れ線は破線で」と指示している.

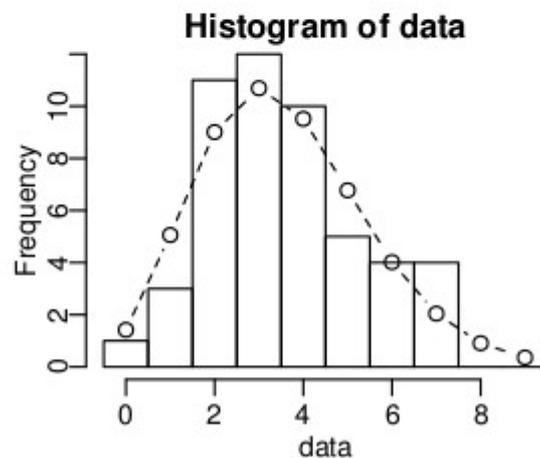
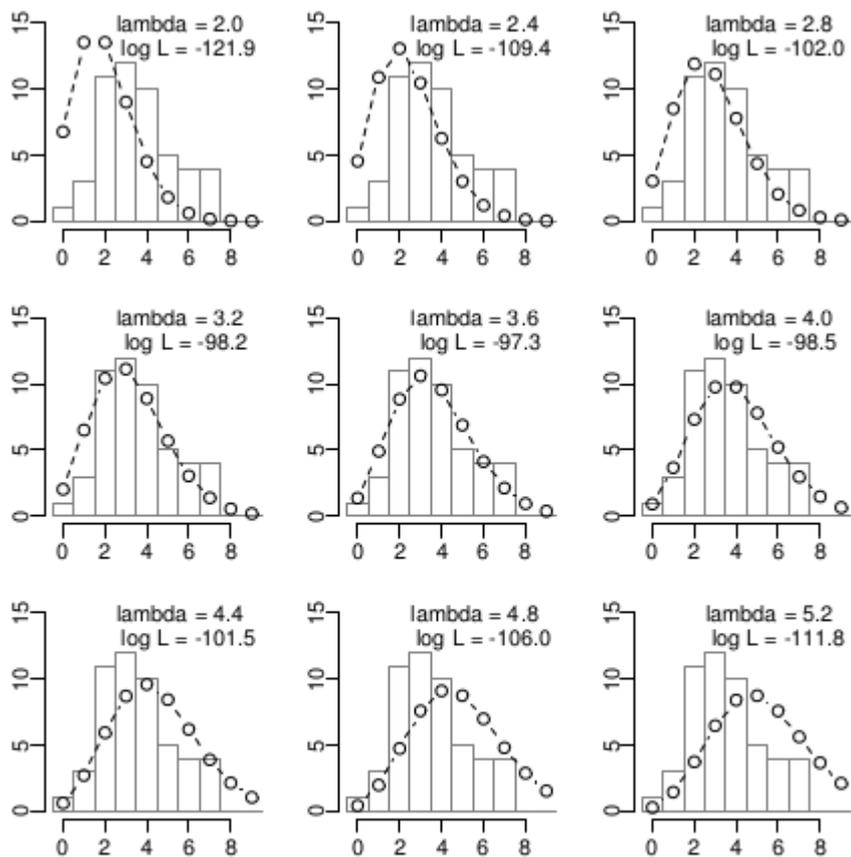


図 5 観測データと確率分布の対応をながめる. ヒストグラムは図 4 と同じ. それに重ねられている丸と破線は y 個の種子をもつ個体数の予測. 平均 3.56 の図 4 のポアソン分布の確率分布に全個体数 50 をかけて得られる.

さいゆう

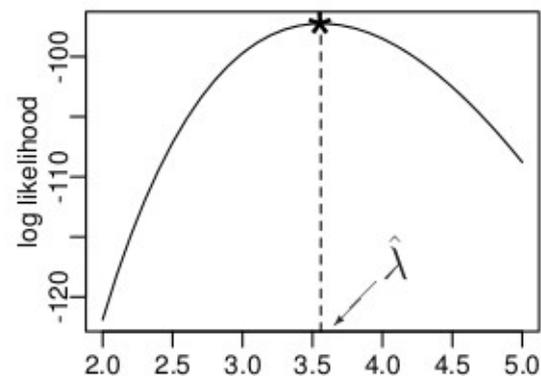
最尤推定という考えかたを説明します



ポアソン分布のパラメータの 最尤推定 もっとももらしい推定?

対数尤度を最大化する $\hat{\lambda}$ をさがす

$$\text{対数尤度 } \log L(\lambda) = \sum_i (y_i \log \lambda - \lambda - \sum_k y_i \log k)$$



kubostat2013a (<http://goo.gl/82dgC>)

統計モデリング入門 2013 (2)

2013-07-03

23 / 28

図 7 平均 λ (lambda) を変化させていったポアソン分布と、観測データへのあてはまりの良さ (対数尤度 $\log L$)。すべてのヒストグラムは図 2 と同じ。

統計モデリング入門 2013 (3)

一般化線形モデル: ポアソン回帰

久保拓弥 kubo@ees.hokudai.ac.jp

北大環境科学院の講義 <http://goo.gl/82dgC>

2013-07-08

ファイル更新時刻: 2013-06-25 15:24

ここで登場する ---

「何でも正規分布」ではダメ! という発想

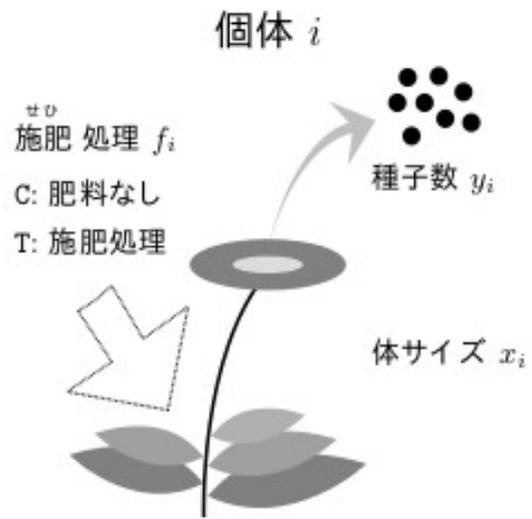


図 3.1 この例題に登場する架空植物の第 i 番目の個体。この植物の体サイズ（個体の大きさ） x_i と肥料をやる施肥処理 f_i が種子数 y_i にどう影響しているのかを知りたい。

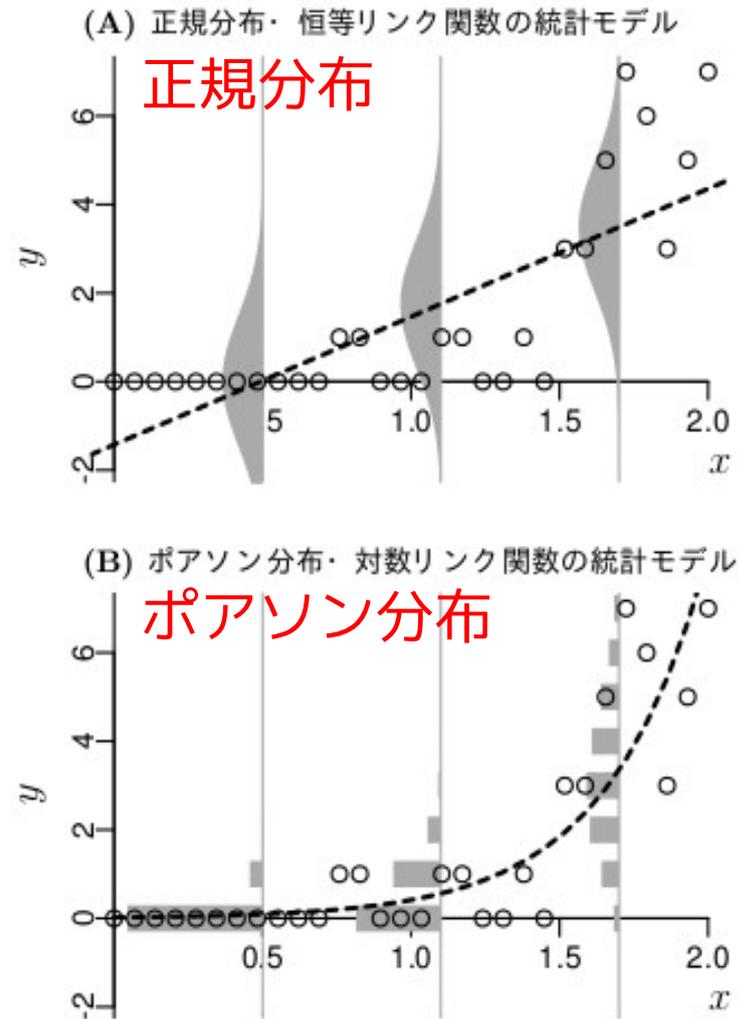


図 3.9 回帰モデルと確率分布の関係。また別の架空データに対して GLM をあてはめた例。破線は x とともに変化する平均値。グレイで

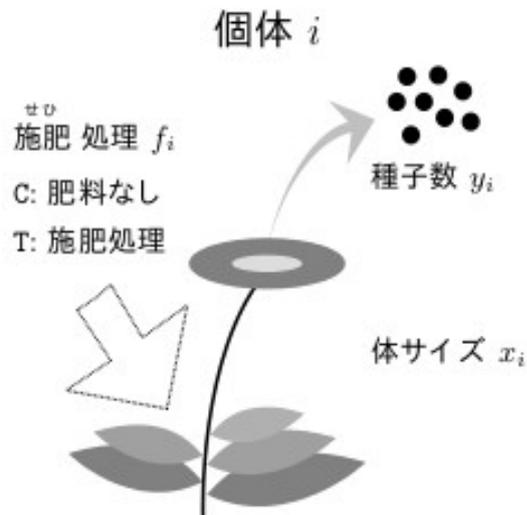


図 3.1 この例題に登場する架空植物の第 i 番目の個体
体サイズ(個体の大きさ) x_i と肥料をやる施肥処理、
にどう影響しているのかを知りたい。

結果を格納するオブジェクト

```
fit <- glm(
  y ~ x,
  family = poisson(link = "log"),
  data = d
)
```

関数名
確率分布の指定
モデル式
リンク関数の指定 (省略可)
) data.frame の指定

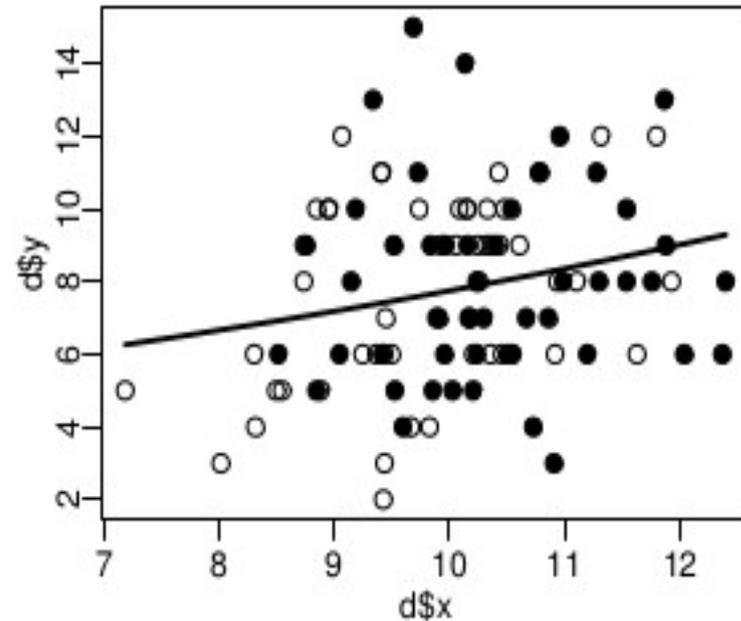


図 17 平均種子数 λ の予測. 図 12 に λ の予測値 (実線) を上かきしたもの。

7/10 (水)

統計モデリング入門 2013 (4)

モデル選択と検定

久保拓弥 kubo@ees.hokudai.ac.jp

北大環境科学院の講義 <http://goo.gl/82dgC>

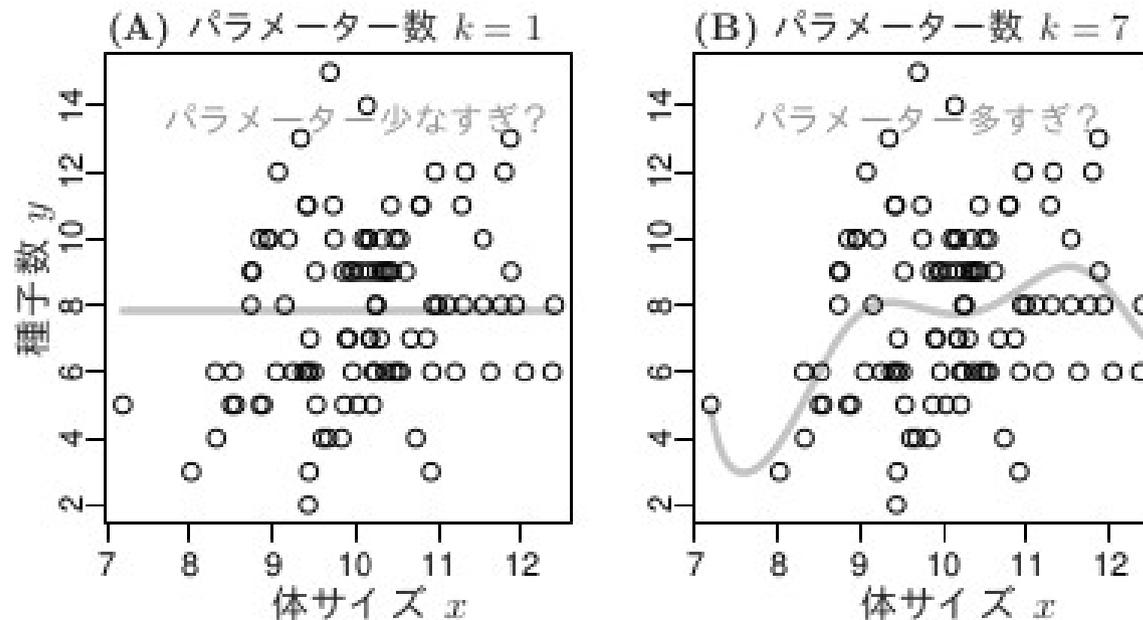
2013-07-10

ファイル更新時刻: 2013-06-25 15:44

Q. モデル選択とは何か？

データと確率分布の対応 どういう関係なのか図示してながめる

パラメーター数は多くても少なくてもヘン？



A. より良い予測をする統計モデルを探すこと

もくじ

モデル選択と検定の手順

統計モデルの検定

AICによるモデル選択

←こっちだ!

検定はモデル選択じゃない!

解析対象のデータを確定



データを説明できるような統計モデルを設計

(帰無仮説・対立仮説)

(単純モデル・複雑モデル)



ネストした統計モデルたちのパラメーターの最尤推定計算



帰無仮説棄却の危険率を評価 モデル選択規準 AIC の評価

統計学って「検定」のこと?

「検定」って何なの?

「検定」ってエラいの?

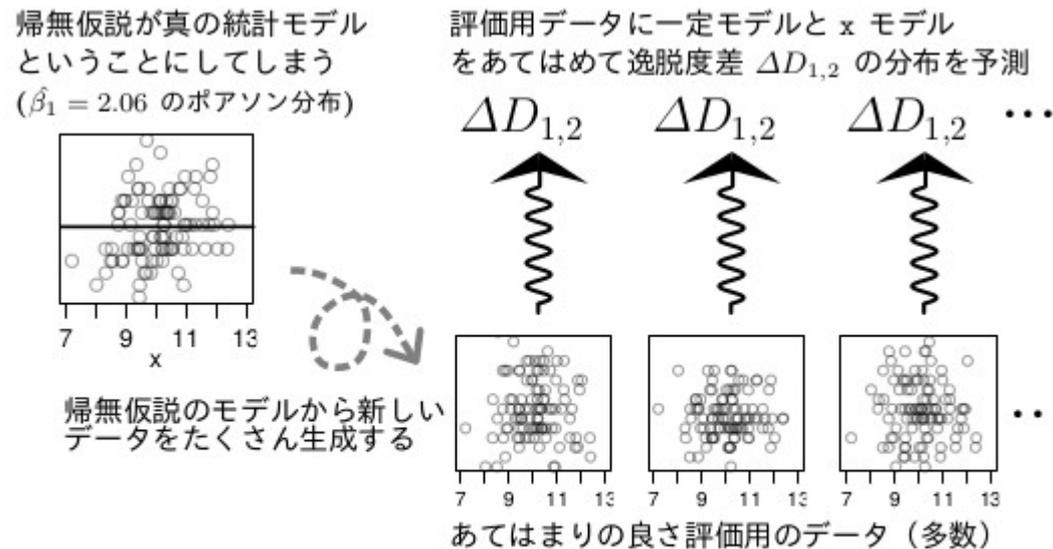


図 6 尤度比検定に必要な $\Delta D_{1,2}$ の分布の生成。まず帰無仮説である一定モデル ($\hat{\beta}_1 = 2.06$, p. 参照) が真の統計モデルだと仮定し、そこから得られるデータを使って逸脱度差 $\Delta D_{1,2}$ がどのような分布になるかを調べる。

7/17 (水)

統計モデリング入門 2013 (5)

一般化線形モデル: ロジスティック回帰

久保拓弥 kubo@ees.hokudai.ac.jp

北大環境科学院の講義 <http://goo.gl/82dgC>

2013-07-17

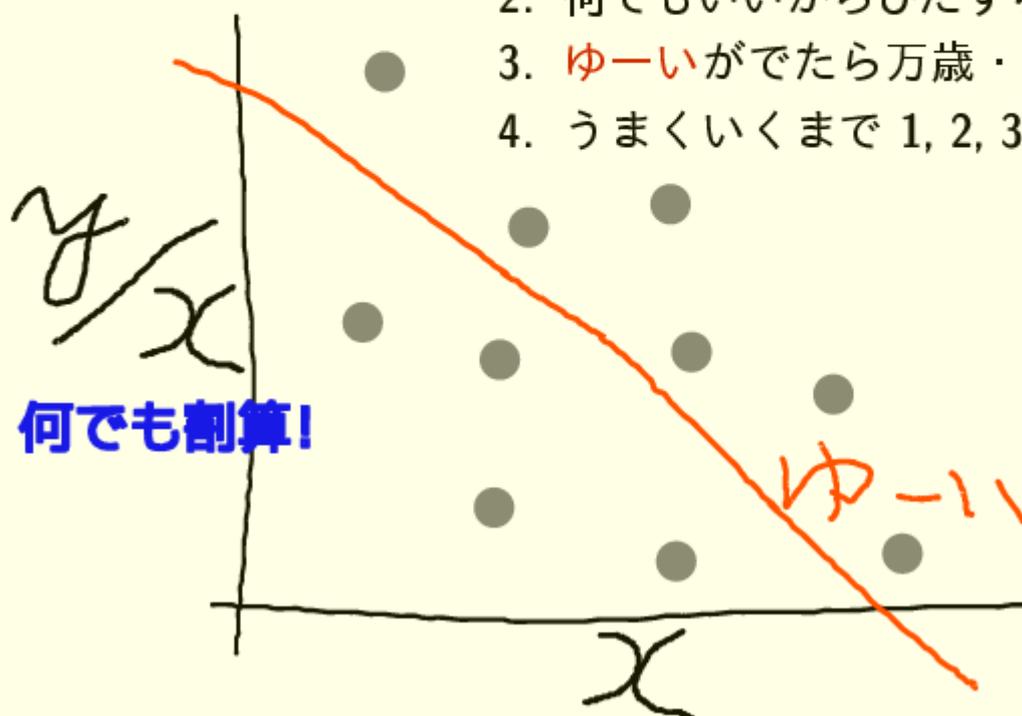
ファイル更新時刻: 2013-06-25 16:36

生物学のデータ解析は「割算」しまくり!!

この悪しき「割算」な統計学

世間でよくみかけるおススメできない作法の例

1. データどんどん割算・割算
2. 何でもいいからひたすらセンをひく
3. ゆーいがでたら万歳・万歳
4. うまくいくまで 1, 2, 3 ぐるぐる

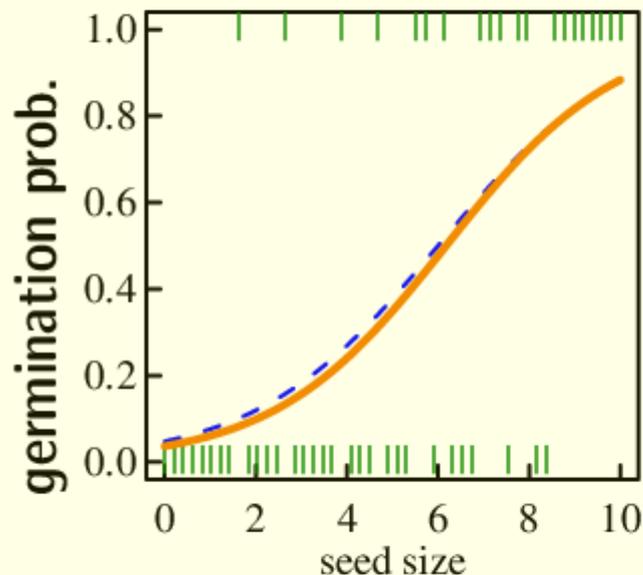


ちなみにこれは w と $0/w$ を比較してるんだから、反比例みたいな偽「負の相関」ができるのはあたりまえ

GLM のひとつ, ロジスティック回帰を使おう

データにあわせたより良い統計モデリングを!

おススメできないデータ解析を回避するための注意点



- むやみに 区画わけしない!
- 何でも 割り算するな!
- たくさん 図を描く
- 「観測データを説明する 確率分布は何か?」を考える



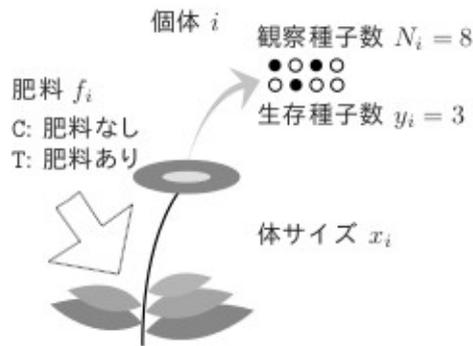
コツ: 不自然にデータをこねくりまわさない
データの性質・構造にあったモデリングを!

GLM のひとつ, ロジスティック回帰を使おう

データと確率分布の対応 どういう関係なのか表示してながめる

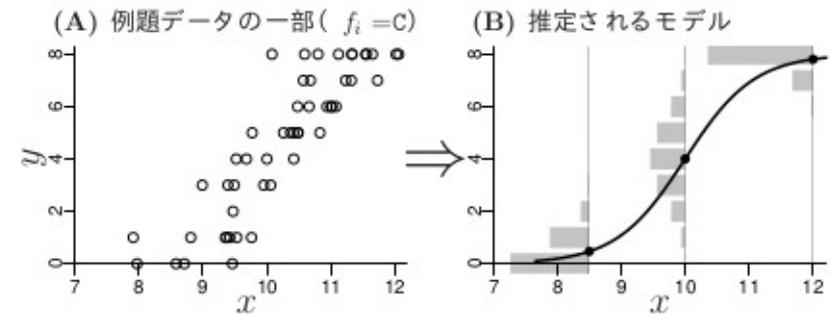
またいつもの例題? …… ちょっとちがう

8 個の種子のうち y 個が **発芽可能** だった! …… というデータ



データと確率分布の対応 どういう関係なのか表示してながめる

ロジスティック回帰とは何なのか?



kubostat2013a (<http://goo.gl/82dgC>)

統計モデリング入門 2013 (5)

2013-07-17 4 / 16

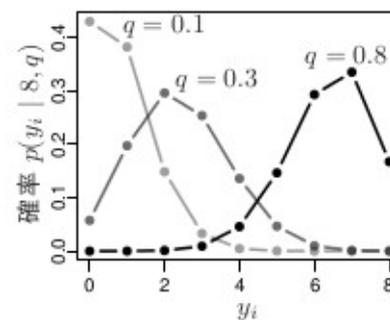
kubostat2013a (<http://goo.gl/82dgC>)

統計モデリング入門 2013 (5)

2013-07-17 9 / 16

データと確率分布の対応 どういう関係なのか表示してながめる

二項分布: N 回のうち y 回, となる確率



7/22 (月)

統計モデリング入門 2013 (6)

一般化線形混合モデル

久保拓弥 kubo@ees.hokudai.ac.jp

北大環境科学院の講義 <http://goo.gl/82dgC>

2013-07-22

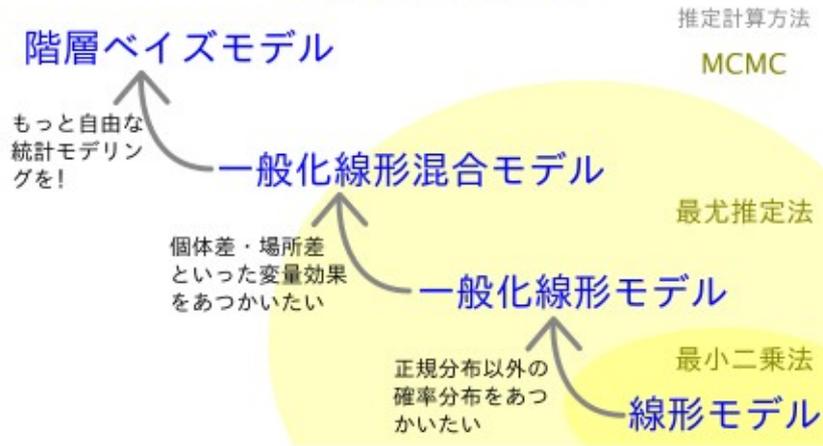
ファイル更新時刻: 2013-06-27 12:17

GLM ではうまく対処できない問題

GLM では説明できない種子データ 「ばらつき」が大きすぎる!

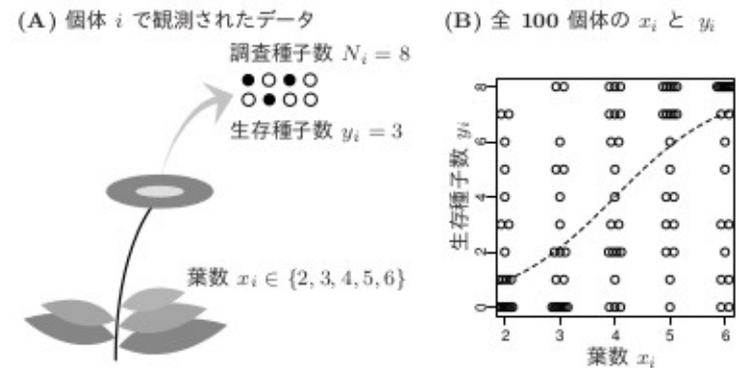
この授業であつかう統計モデルたち

線形モデルの発展



GLM では説明できない種子データ 「ばらつき」が大きすぎる!

今日の例題: 種子の生存確率, ただし……



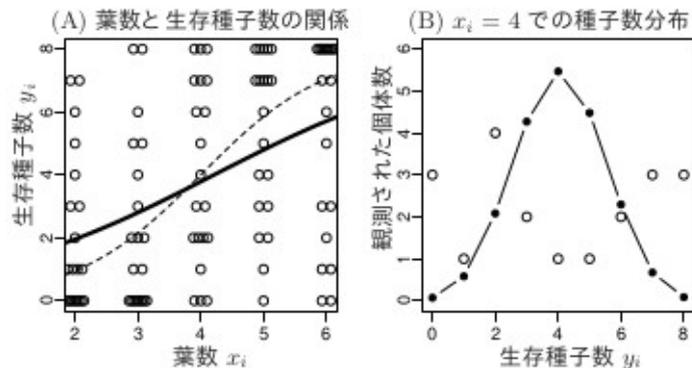
統計モデル勉強のプラン: 線形モデルを発展させる

kubostat2013a (<http://goo.gl/82dgc>) 統計モデリング入門 2013 (6) 2013-07-22 5 / 21

GLM では説明できない種子データ 「ばらつき」が大きすぎる!

GLM では説明できないばらつき!

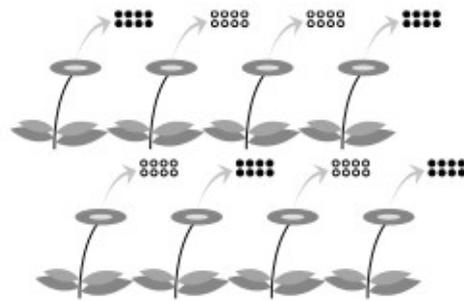
kubostat2013a (<http://goo.gl/82dgc>) 統計モデリング入門 2013 (6) 2013-07-22 6 / 21



今まで「個体差」無視した統計モデル (GLM) を使っていた!

過分散と個体差 観測されていない個体差をもたらす過分散

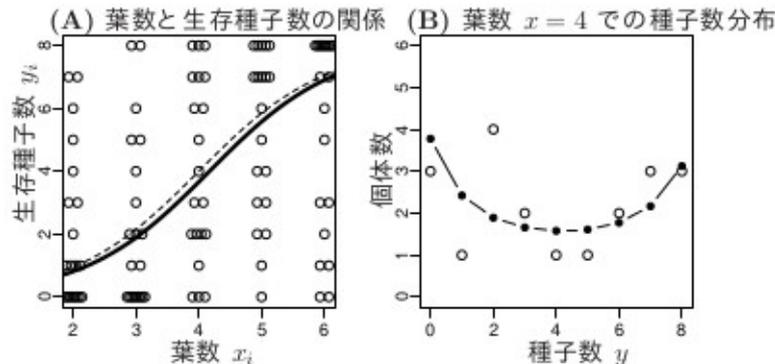
過分散 (overdispersion) とは何か?



kubostat2013a (<http://goo.gl/82dgC>) 統計モデリング入門 2013 (6) 2013-07-22 9 / 21

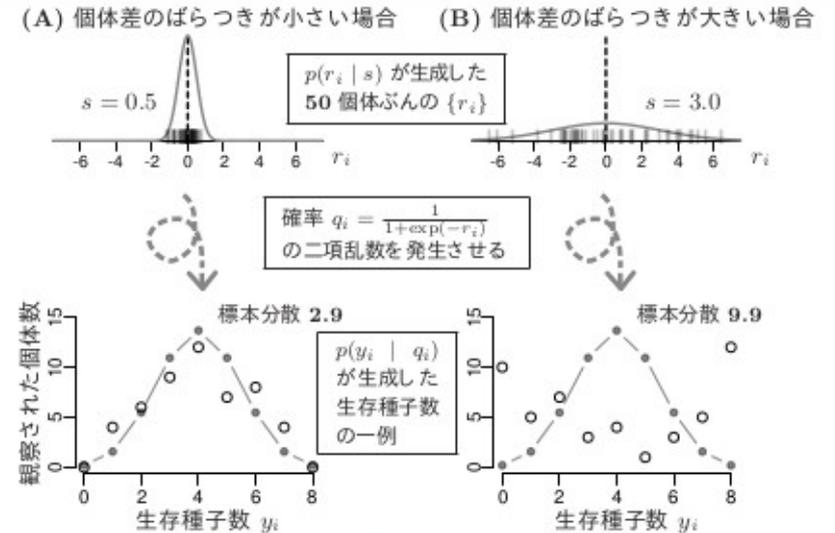
一般化線形混合モデルの最尤推定 「積分する」とは分布を混ぜること

推定された GLMM を使った予測



個体差のばらつきをあらわす確率分布 平均的な個体や「異端」な個体のばらつき

個体差 r_i の分布と過分散の関係



kubostat2013a (<http://goo.gl/82dgC>) 統計モデリング入門 2013 (6) 2013-07-22 14 / 21

一般化線形混合モデル
(Generalized Linear Model,
GLMM) を使って問題解決

7/24 (水)

2013-07-24

全部で 7 回中の 7 回目

統計モデリング入門 2013 (7)

階層ベイズモデル

久保拓弥 kubo@ees.hokudai.ac.jp

<http://goo.gl/0yB2k>

2013-07-24

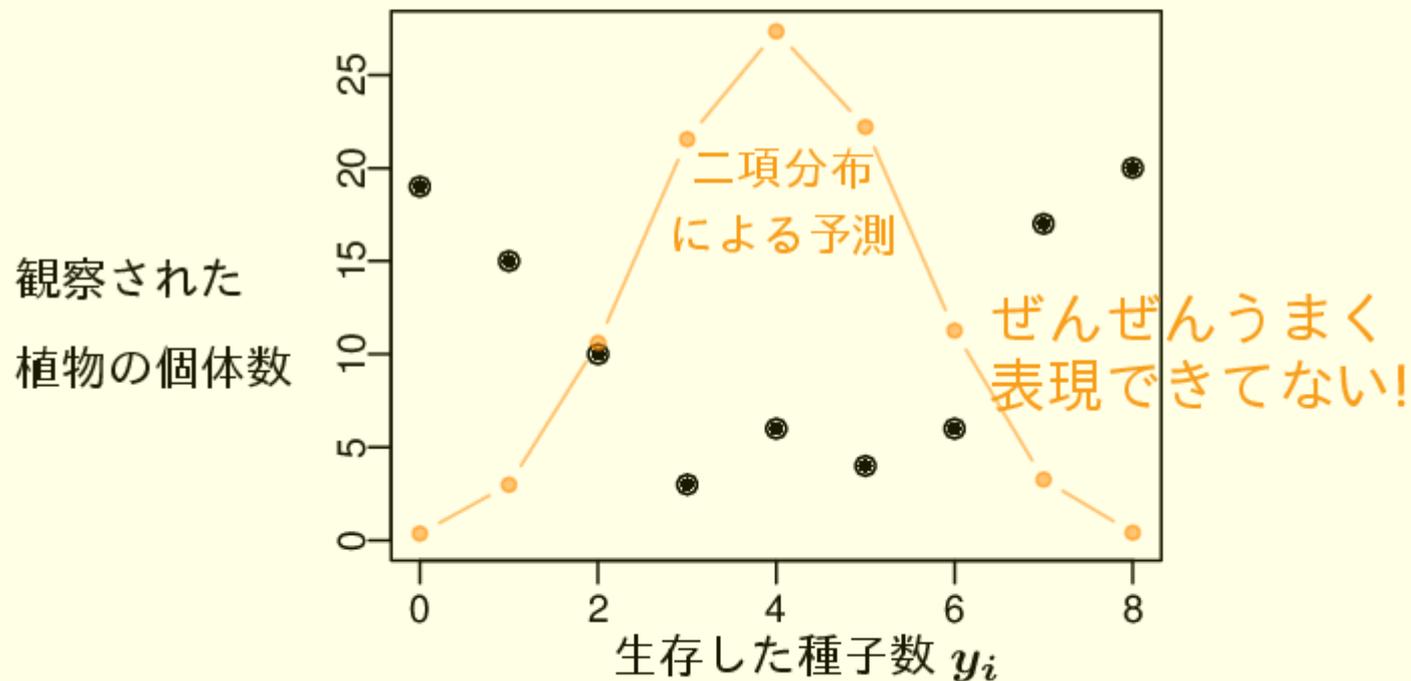
(2013-06-27 13:21 修正版)

1

GLM ではうまく説明できないデータ!?

また別の観測データ：二項分布だめだめ?!

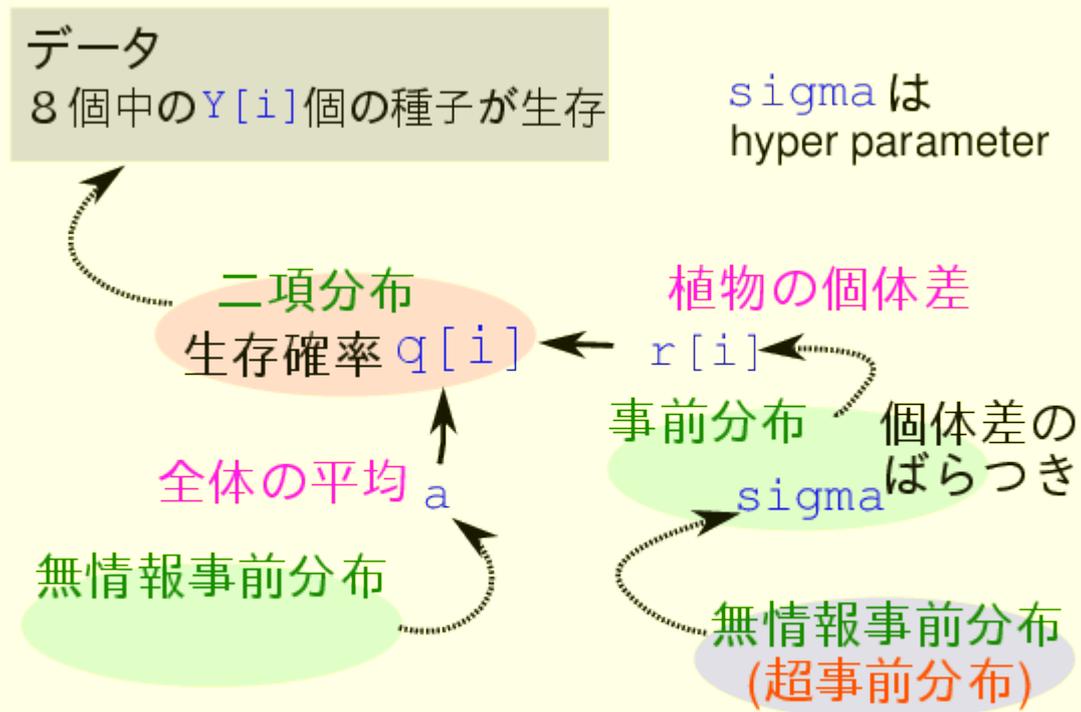
100 個体の植物の合計 800 種子中 **403 個** の生存が見られたので，平均生存確率は 0.50 と推定されたが……



第 6 回と同じような例題を，こんどはベイズモデルを使ってモデリングします

GLM を階層ベイズモデル化して対処

なぜ「階層」ベイズモデルと呼ばれるのか？



超事前分布 → 事前分布という階層があるから

なぜ階層ベイズモデルまで勉強するの？

• 生態学や漁業のデータ

解析は難しいから！

- ✓ 個体差・エリア差・空間相関・時間相関・種差などめんどろなことをあつかわないといけない

線形モデルの発展

階層ベイズモデル

もっと自由な統計モデリングを！

一般化線形混合モデル

個体差・場所差といった変量効果をあつかいたい

一般化線形モデル

正規分布以外の確率分布をあつかいたい

線形モデル

推定計算方法
MCMC

最尤推定法

最小二乗法

そういう難しい状況では……

- ベイズモデル化
- そのパラメーターの事後分布を MCMC 法を使って推定するのが無難

今日のハナシはここまで

時間があれば R インストール実演

- 第 1 回: 7/01 (月) 観測されたパターンを説明する統計モデル
- 第 2 回: 7/03 (水) 確率分布と最尤推定
- 第 3 回: 7/08 (月) 一般化線形モデル: ポアソン回帰
- 第 4 回: 7/10 (水) モデル選択と検定
- 第 5 回: 7/17 (水) 一般化線形モデル: ロジスティック回帰
- 第 6 回: 7/22 (月) 一般化線形混合モデル
- 第 7 回: 7/24 (水) 階層ベイズモデル

次回以降も楽しく勉強をすすめてみましょう!