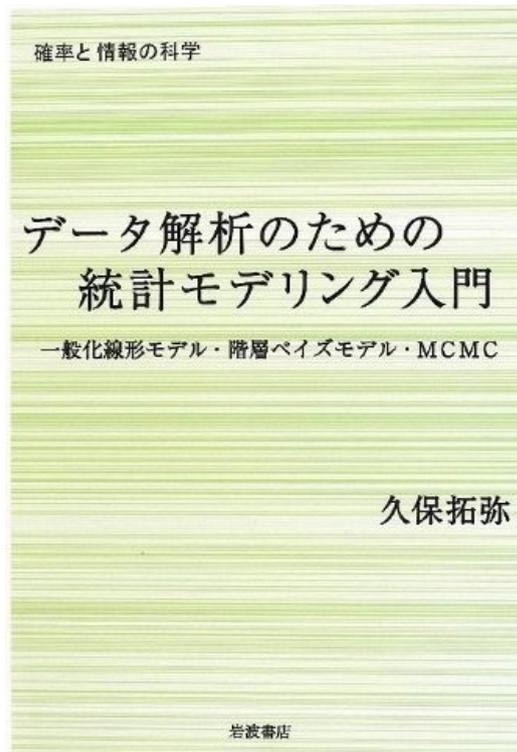


茨城大学集中講義・統計モデリング入門 (a)

# 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



茨城大学理工学研究科

- ・ 生物系特別講義 IV
- ・ 先端科学トピックス

2014 年 10/1-10/2

# 統計モデリング授業の web page

# <http://goo.gl/2QNwgl>

更新: 2014-09-28 22:03:54

## 生態学のデータ解析 - 茨城大学集中講義2014

まだ完成にはほどとおい状態……ですね

- 2014 年 10 月 1-2 日の茨城大学理学部での集中講義 ([統計学授業](#))
  - このページの短縮 URL: <http://goo.gl/2QNwgl>
- 参考文献: [統計モデリング入門](#)

(もくじ)

- [10 月 1 日 \(水\) の予定](#)
- [10 月 2 日 \(木\) の予定](#)
- [単位取得する人のための課題](#)

### 10 月 1 日 (水) の予定

場所: 理学部 ??? 室

- (a, b) 08:50-10:20 全体の流れ: 統計モデルと確率分布
  - 講義資料
    - [kuboIbaraki2014a.pdf](#)
    - [kuboIbaraki2014i2.pdf](#)

# 課題も!



# 自己紹介：久保拓弥

- 北大の環境科学学院という学部のない大学院
- 生態学に関するデータ解析とかやっています
  - 野外調査をしない生態学者
- データは誰か別の人がとってきてくれます

そもそも生態学って何？

- 生物の数の変化や分布や生活の様子を調べる
- いろいろな動植物が対象



# 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

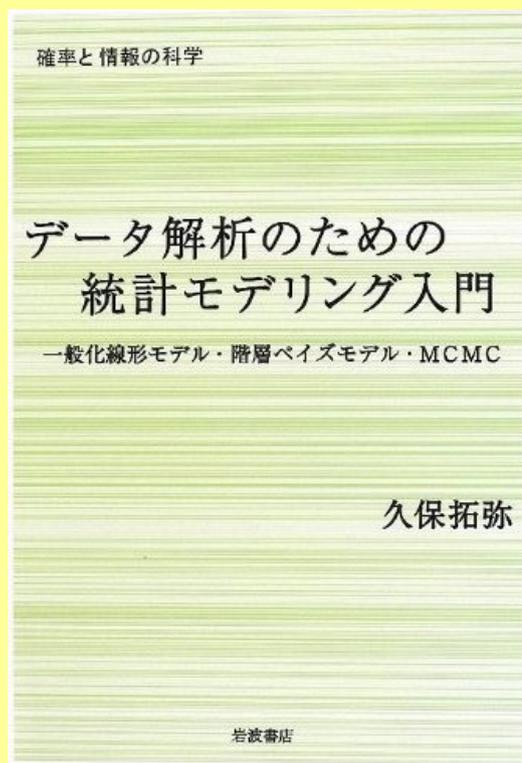
- ブラックボックス統計解析
- とにかく「ゆーい差」さえ出せばよいという発想になっている
- 大学・大学院でもあまりちゃんと教えられていない，教えられるヒトが少ない……とくに近年発達している統計モデリングについて

# 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!

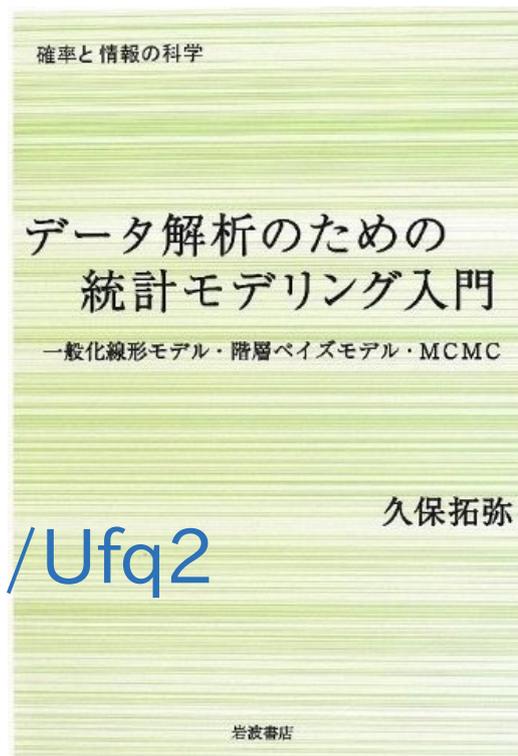
- データ解析で使われるの中でも比較的簡単な統計モデルを理解しよう
- 「ゆーい差」さえ出せばよいという発想をやめて、データと統計モデルの対応関係をよく見よう（作図重要）
- 統計ソフトウェア R を使い始めよう

# 教科書とソフトウェア

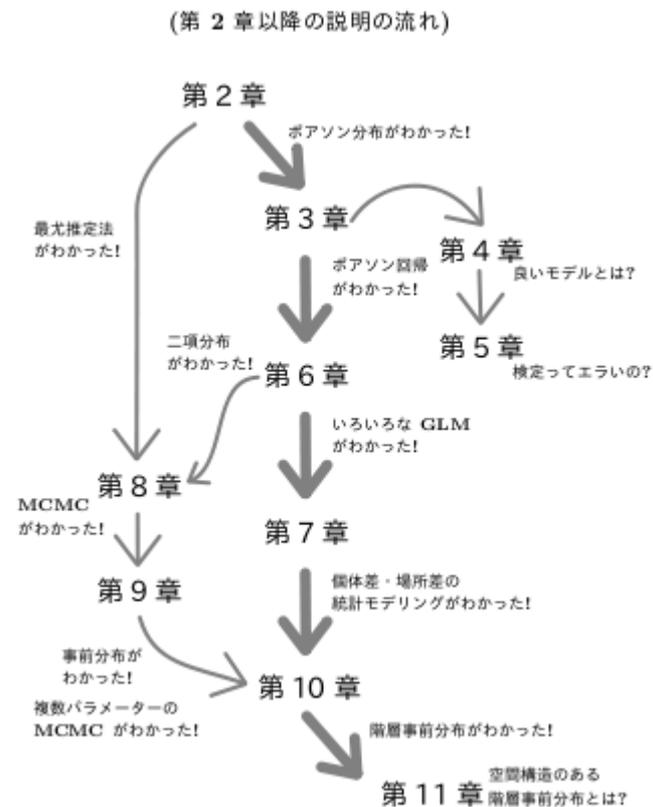


# この授業は「統計モデリング入門」 にそった内容を説明します

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

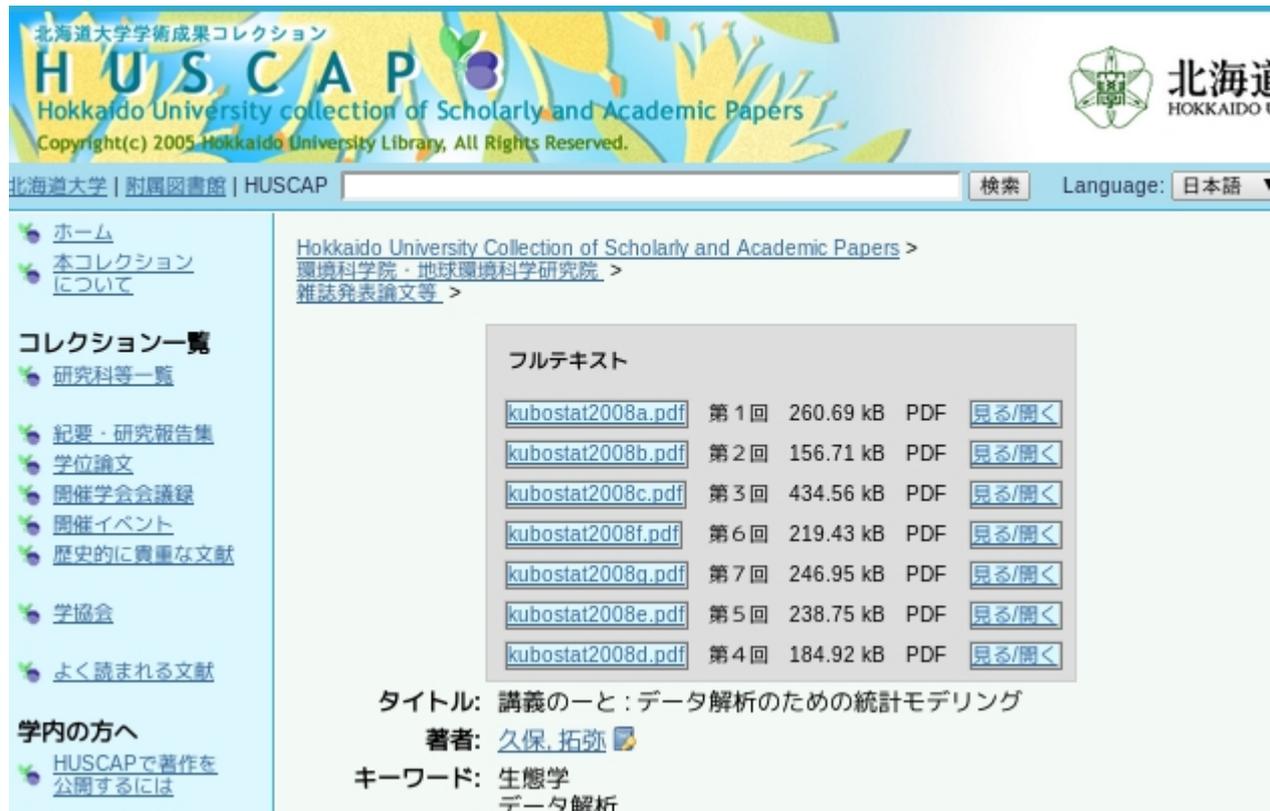


<http://goo.gl/Ufq2>



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 >

フルテキスト

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

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>

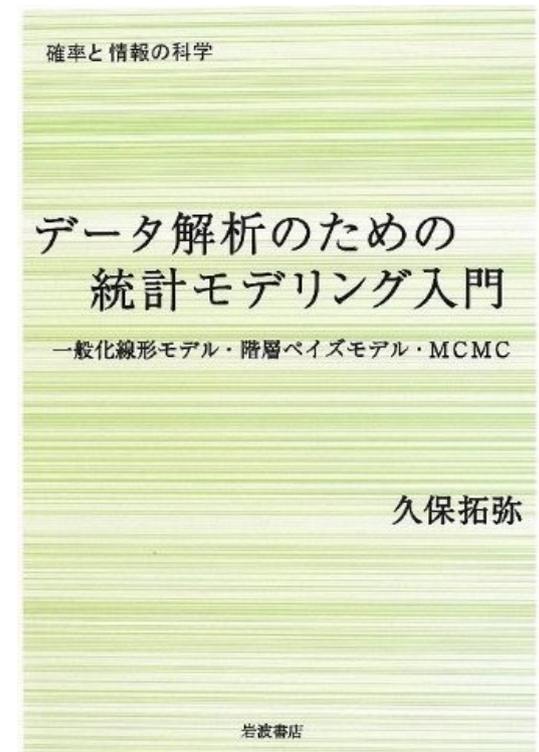
# 統計ソフトウェア R



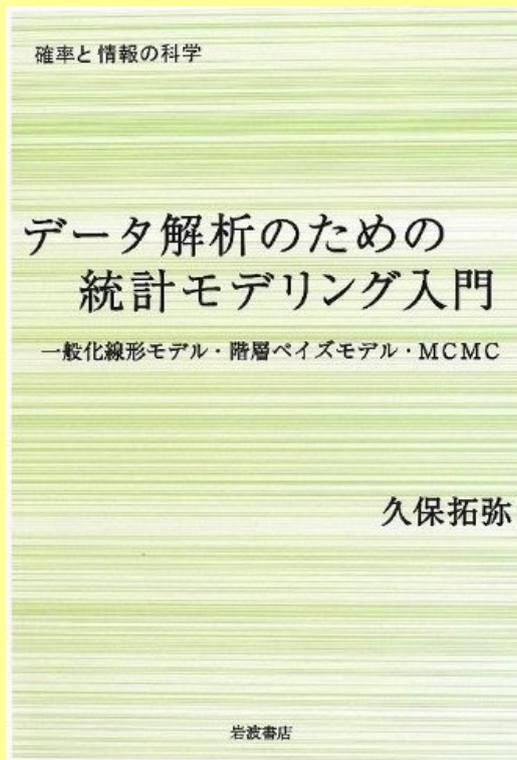
統計学の勉強には良い統計ソフトウェアが必要!

- 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？



# たとえばこんなデータがあったしましょう

(次回の講義の例題)

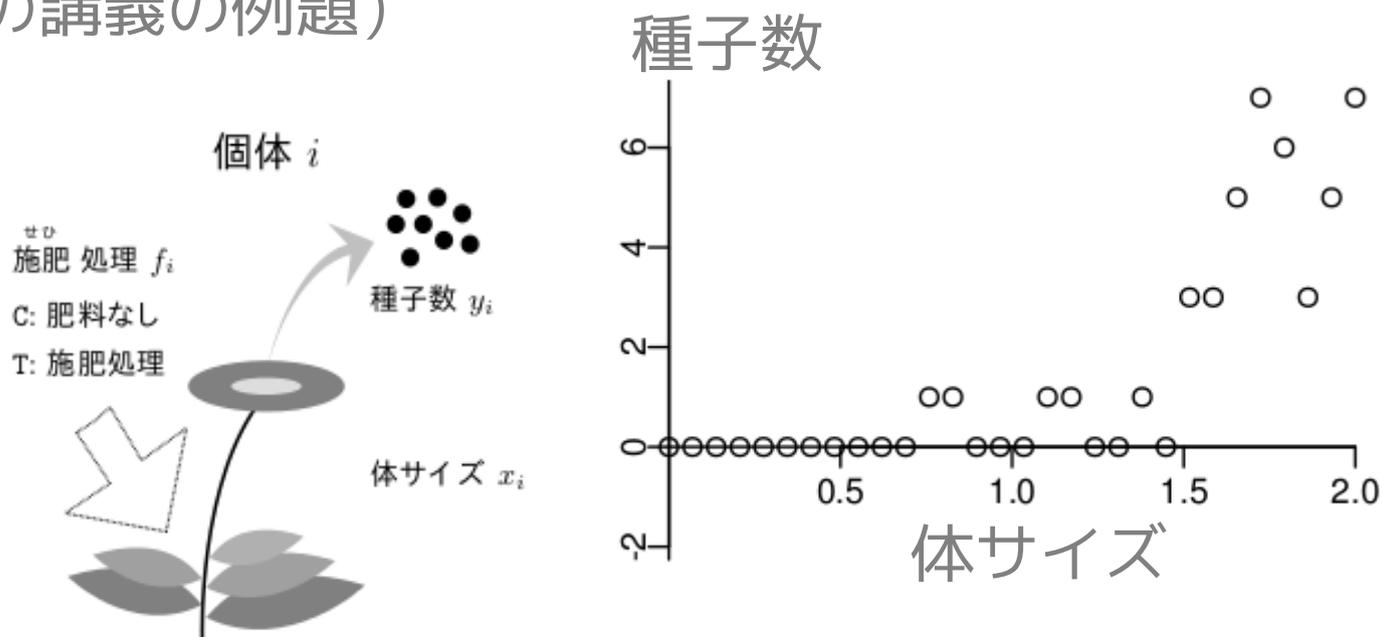
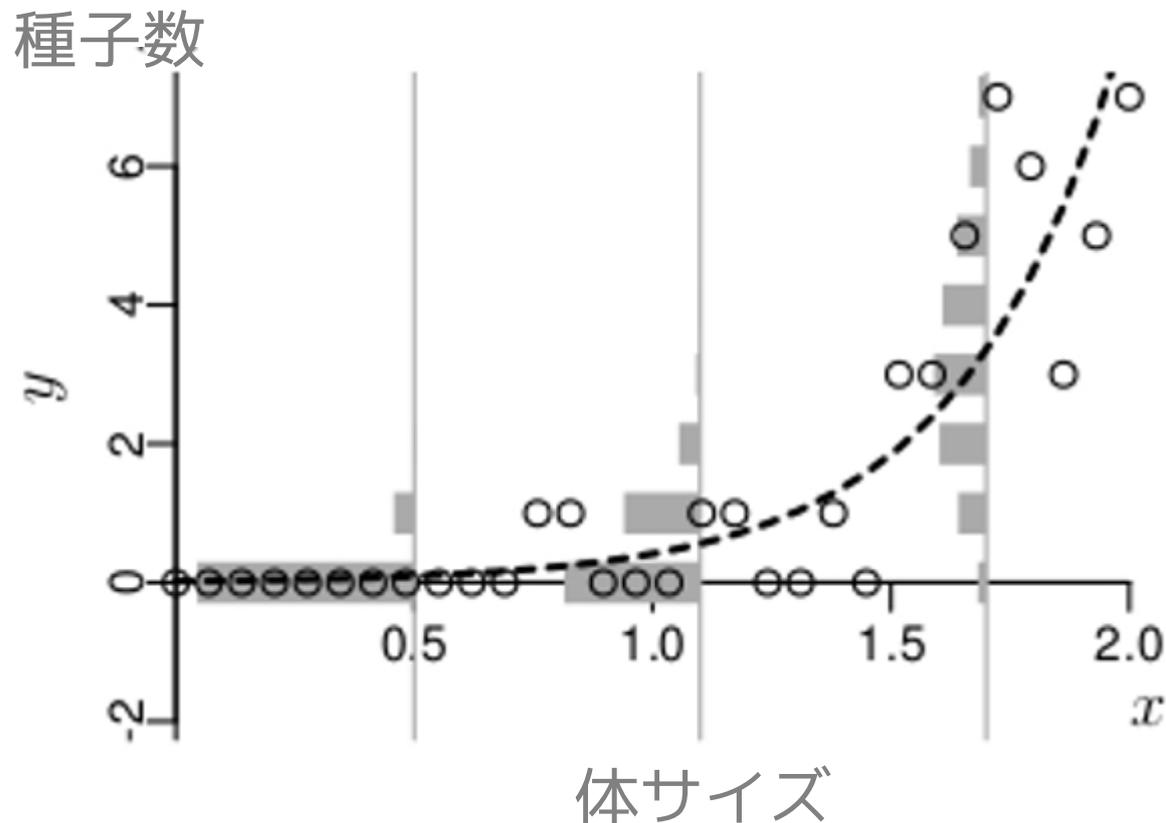
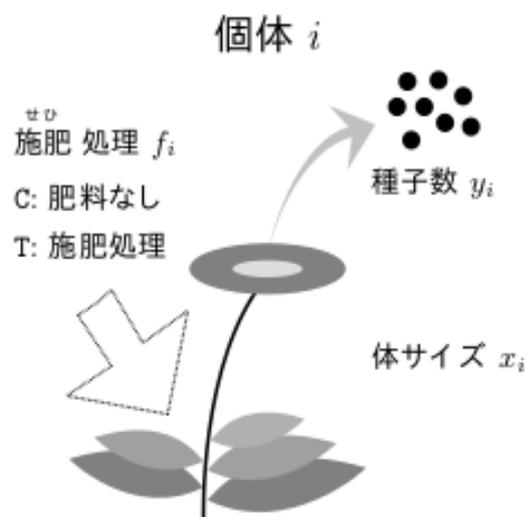


図 3.1 この例題に登場する架空植物の第  $i$  番目の個体. この植物の体サイズ(個体の大きさ)  $x_i$  と肥料をやる施肥処理  $f_i$  が種子数  $y_i$  にどう影響しているのかを知りたい.

こういう「定量的な説明」があったらいいな……

なんとなく「わかった」ような気分?



数式で書かれた  
「統計モデル」を  
準備する  
それをデータに  
あてはめる



種子数の平均値はサイズ  $x$  とともに増大する  
平均値が増大するとばらつきが変化する

……などなど……

# 「統計モデル」のしくみを理解しよう!

もうすこし「わかった」ような気分?

種子数の平均値はサイズ  $x$  とともに増大する

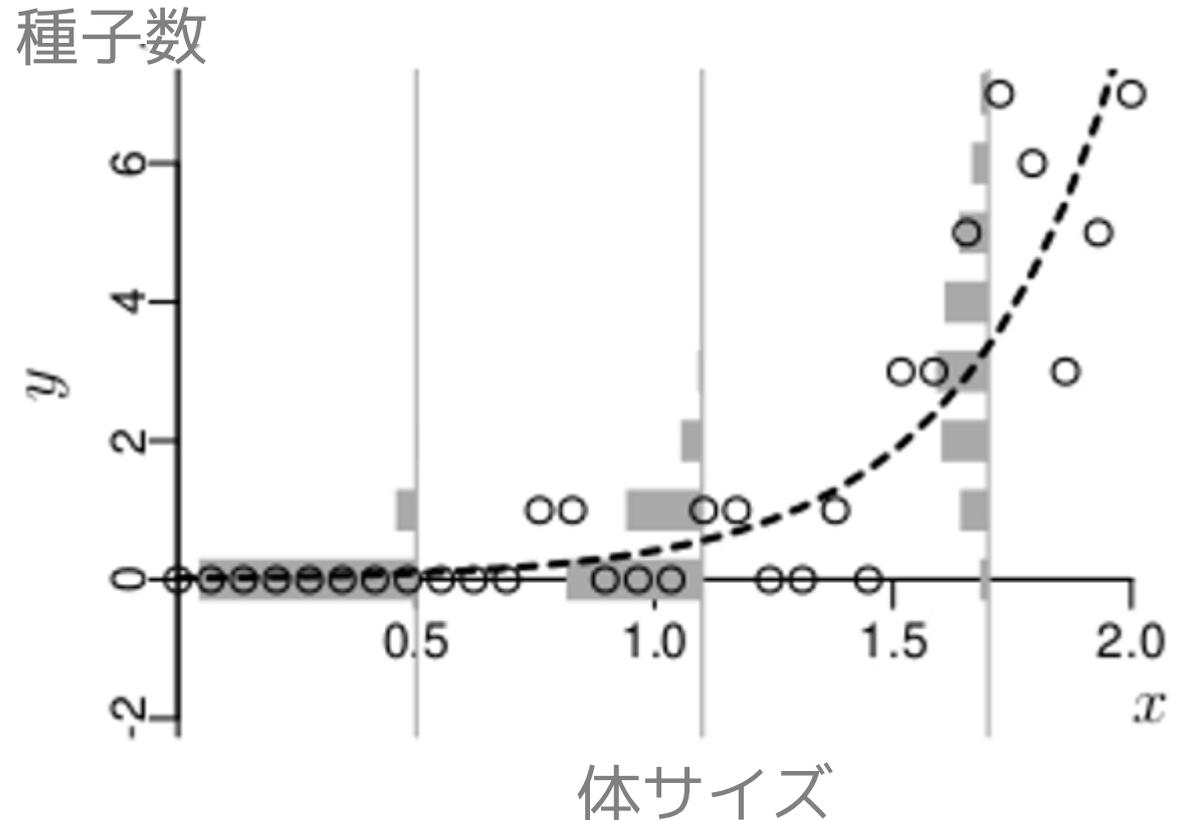
➡ **どのように変化するのか?**  
数式で書くとどうなる?

平均値が増大するとばらつきが変化する

➡ **どのようにばらつくのか?**  
確率分布?

統計モデルをデータにうまくあてはめる

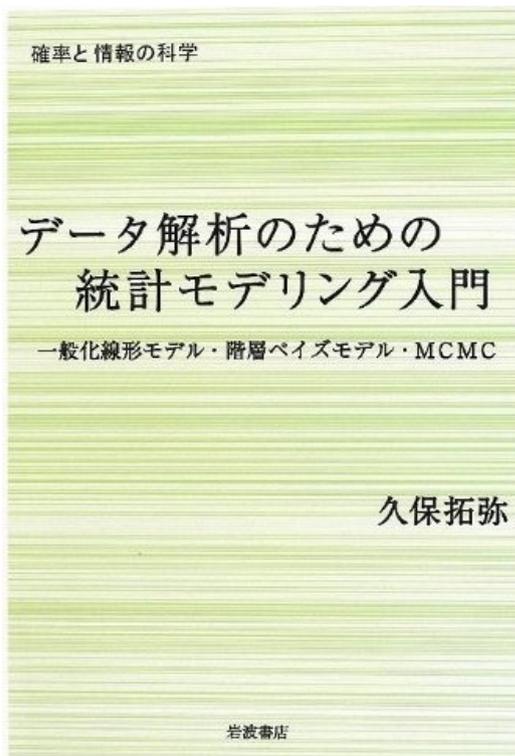
➡ **どのようにあてはめるのが妥当なのか? パラメーター推定法?**



# 「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

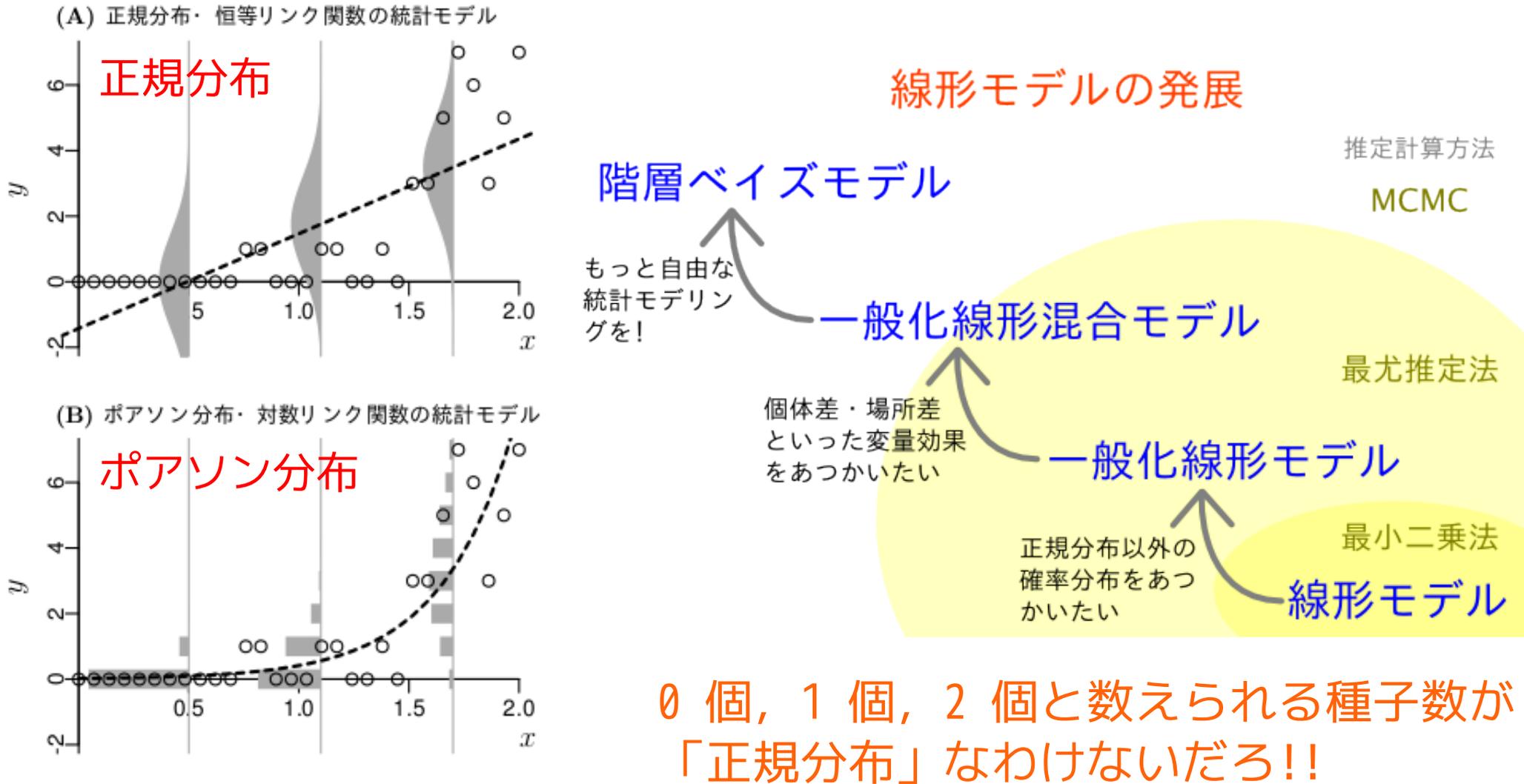
最小二乗法

線形モデル



# 一般化線形モデル - ばらつきをよく見る

Generalized Linear Model, GLM



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

10/1 (水)

第 1 講時 (a) 観測されたパターンを説明する統計モデル  
+ (b) 確率分布と最尤推定

第 2 講時: (c) 一般化線形モデル: ポアソン回帰

第 3 講時: (d) モデル選択と検定

第 4 講時: (e) 一般化線形モデル: ロジスティック回帰

第 5 講時: (f) 階層ベイズモデル

10/2 (木) 第 1,2 講時: 情報基盤センターで R 実習 (r1, b-f)

# 統計モデリング入門 2014 (2)

確率分布と最尤推定

久保拓弥 `kubo@ees.hokudai.ac.jp`

茨城大集中講義 <http://goo.gl/2QNwgl>

2014-10-01

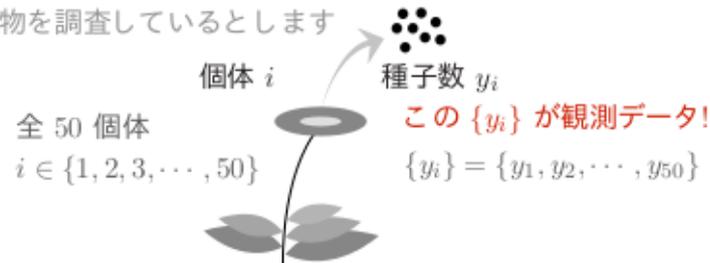
ファイル更新時刻: 2014-09-25 14:39

# 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

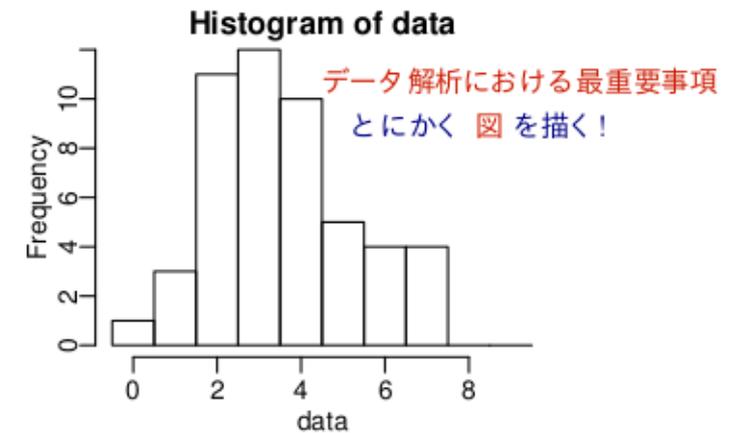
```
> data
[1] 2 2 4 6 4 5 2 3 1 2 0 4 3 3 3 3 4 2 7 2 4 3 3 3 4
[26] 3 7 5 3 1 7 6 4 6 5 2 4 7 2 2 6 2 4 5 4 5 1 3 2 3
```

kubostat2013a (<http://goo.gl/82dgC>) 統計モデリング入門 2013 (2) 2013-07-03 5 / 28

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))
```



kubostat2013a (<http://goo.gl/82dgC>) 統計モデリング入門 2013 (2) 2013-07-03 7 / 28

# 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

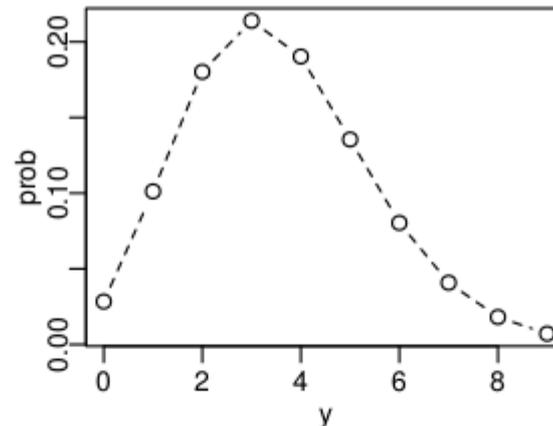


図 4 平均  $\lambda = 3.56$  のポアソン分布. 種子数  $y$  とその確率  $\text{prob}$  の関係が示されている. 図 3 の表を図にしたもの. R の `plot()` 関数の引数, `type = "b"` によって「丸と折れ線による図示」, `lty = 2` によって「折れ線は破線で」と指示している.

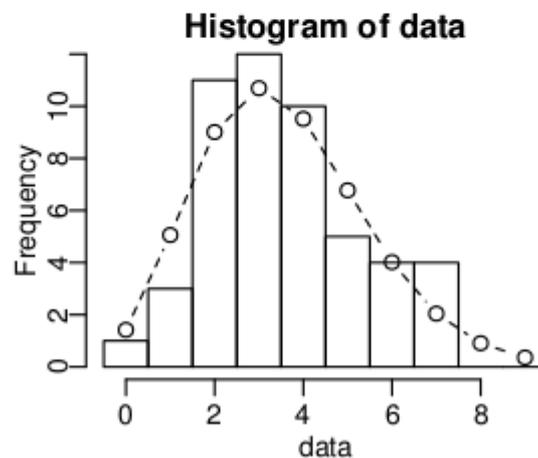
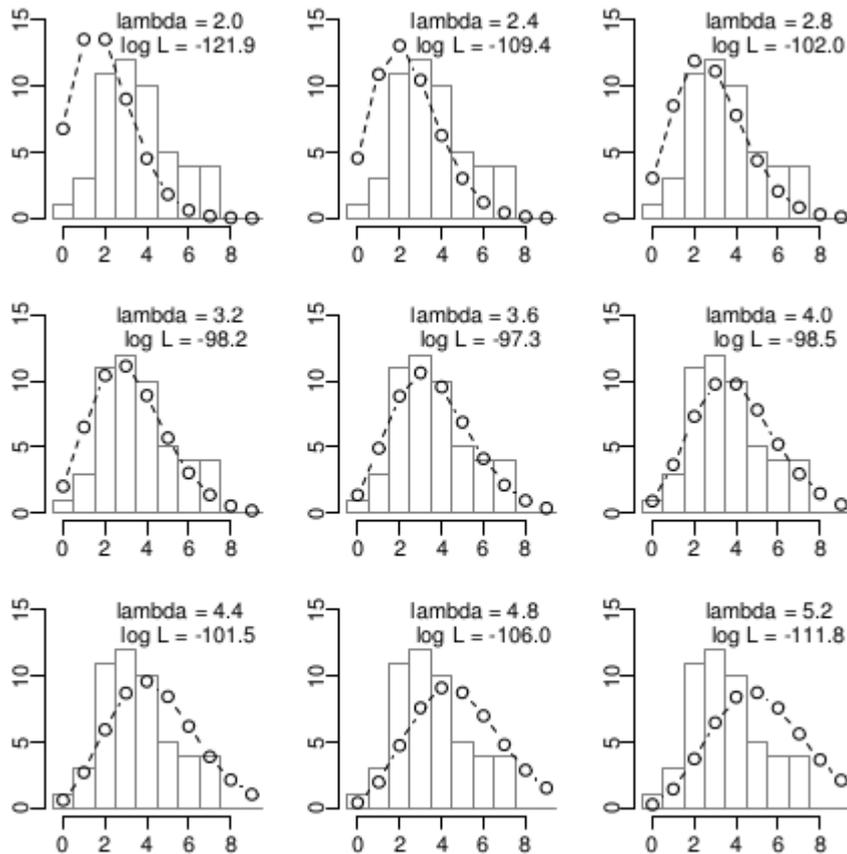


図 5 観測データと確率分布の対応をながめる. ヒストグラムは図 3 と同じ. それに重ねられている丸と破線は  $y$  個の種子をもつ個体数の予測. 平均 3.56 の図 3 のポアソン分布の確率分布に全個体数 50 をかけて得られる.

さいゆう

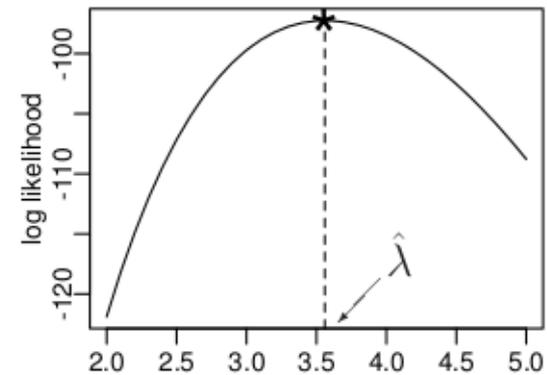
# 最尤推定という考えかたを説明します



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  $\lambda$  (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 2014 (3)

一般化線形モデル: ポアソン回帰

久保拓弥 `kubo@ees.hokudai.ac.jp`

茨城大集中講義 <http://goo.gl/2QNwgl>

2014-10-01

ファイル更新時刻: 2014-09-29 12:44

# ここで登場する ---

## 「何でも正規分布」ではダメ! という発想

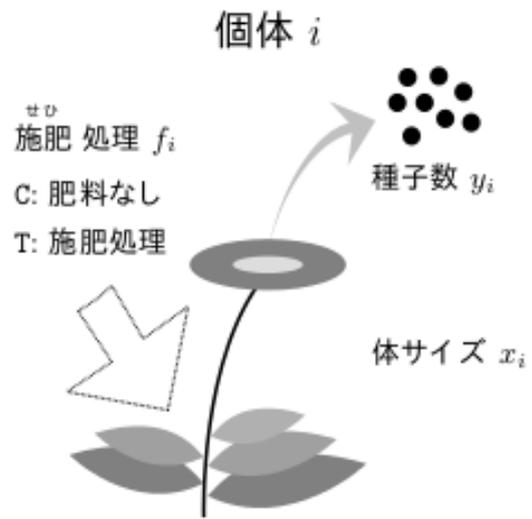


図 3.1 この例題に登場する架空植物の第  $i$  番目の個体。この植物の体サイズ（個体の大きさ） $x_i$  と肥料をやる施肥処理  $f_i$  が種子数  $y_i$  にどう影響しているのかを知りたい。

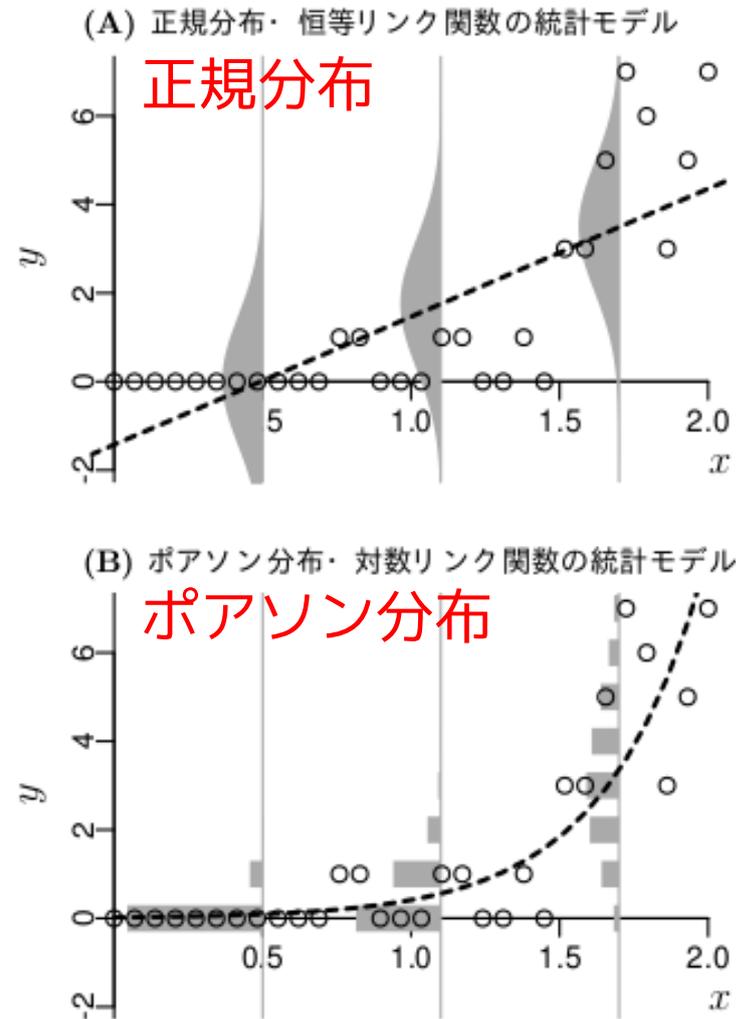


図 3.9 回帰モデルと確率分布の関係。また別の架空データに対して GLM をあてはめた例。破線は  $x$  とともに変化する平均値。グレイで

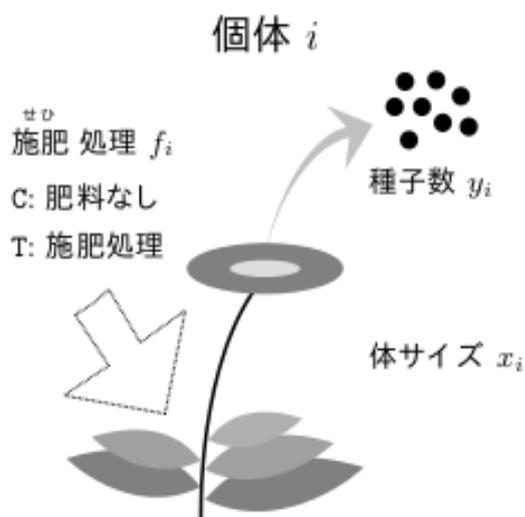


図 3.1 この例題に登場する架空植物の第  $i$  番目の個体  
体サイズ(個体の大きさ)  $x_i$  と肥料をやる施肥処理、  
にどう影響しているのかを知りたい。

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

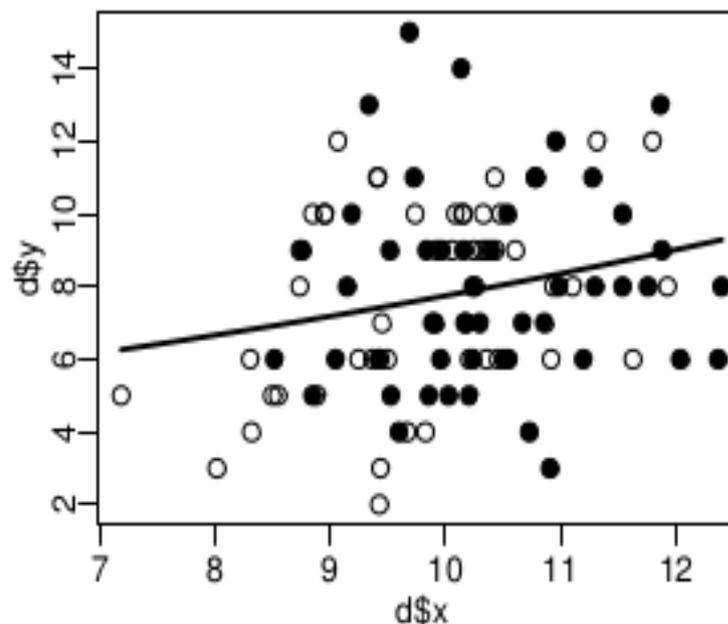


図 17 平均種子数  $\lambda$  の予測. 図 12 に  $\lambda$  の予測値 (実線) を上げきしたものの。

# 統計モデリング入門 2014 (4)

モデル選択と検定

久保拓弥 [kubo@ees.hokudai.ac.jp](mailto:kubo@ees.hokudai.ac.jp)

茨城大集中講義 <http://goo.gl/2QNwgl>

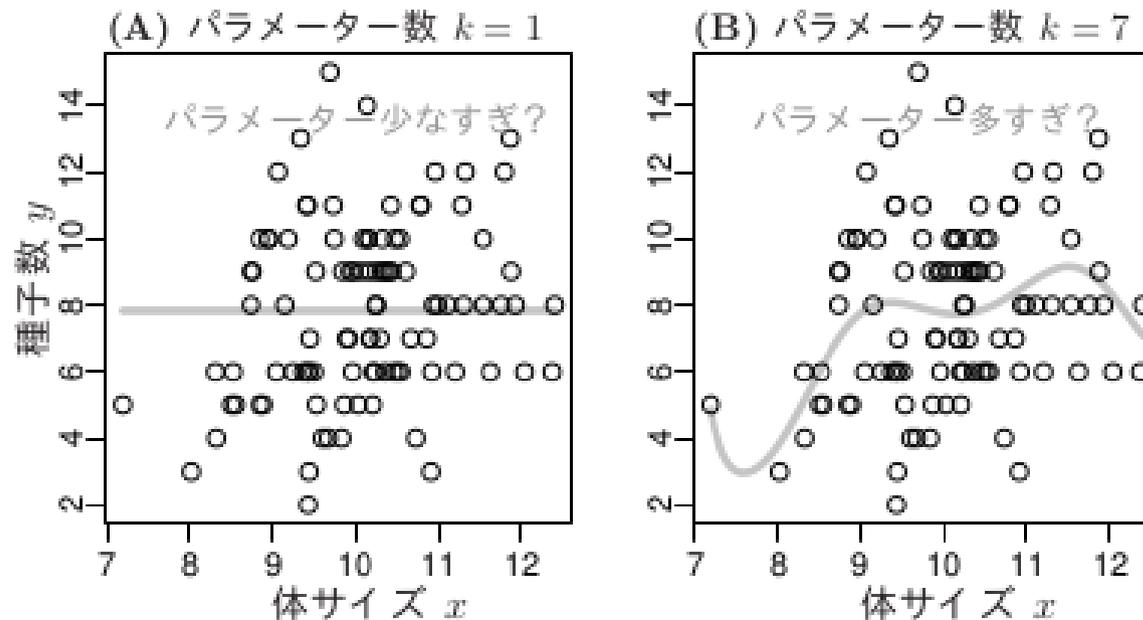
2014-10-01

ファイル更新時刻: 2014-09-29 12:44

# 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?

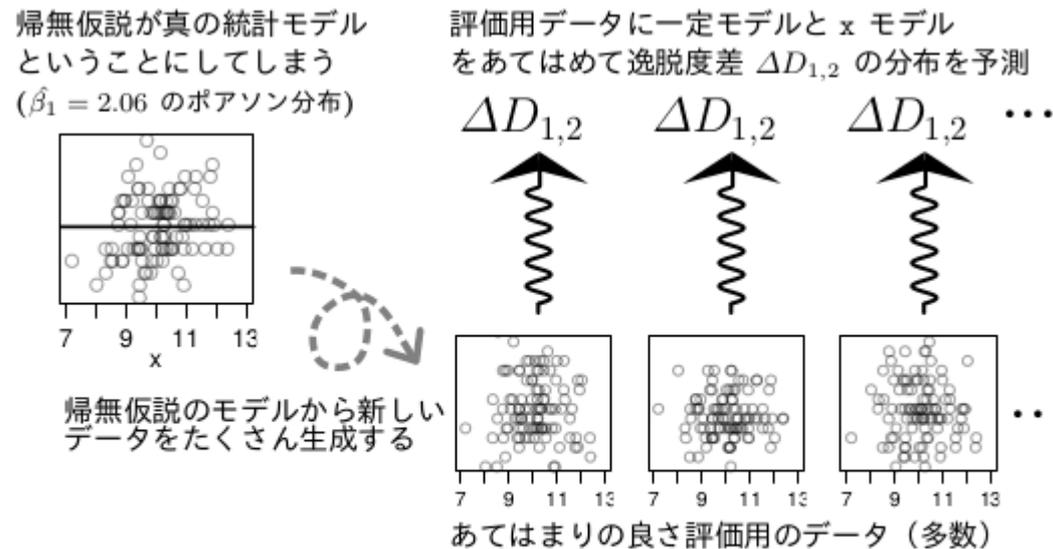


図 6 尤度比検定に必要な  $\Delta D_{1,2}$  の分布の生成。まず帰無仮説である一定モデル ( $\hat{\beta}_1 = 2.06$ , p. 参照) が真の統計モデルだと仮定し、そこから得られるデータを使って逸脱度差  $\Delta D_{1,2}$  がどのような分布になるかを調べる。

# 統計モデリング入門 2014 (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`

茨城大集中講義 <http://goo.gl/2QNwgl>

2014-10-01

ファイル更新時刻: 2014-09-29 12:44

# 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



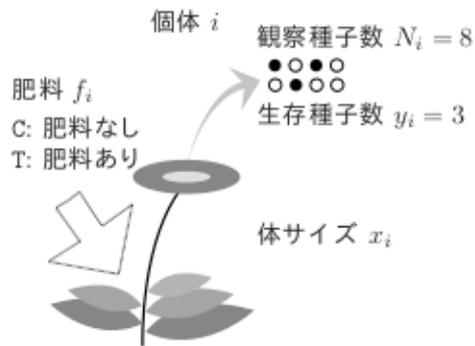
ちなみにこれは  $w$  と  $0/w$  を比較してるんだから、反比例みたいな偽「負の相関」ができるのはあたりまえ

# 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

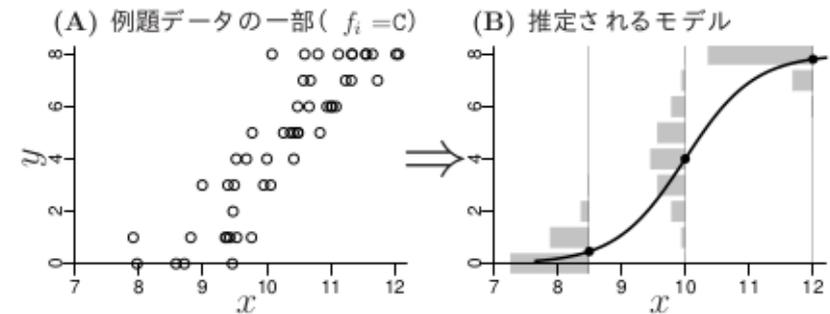
またいつもの例題? …… ちょっとちがう

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図示してながめる

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

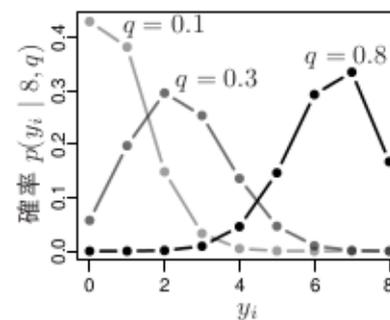
kubostat2013a (<http://goo.gl/82dgC>)

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図示してながめる

二項分布:  $N$  回のうち  $y$  回, となる確率



# 統計モデリング入門 2014 (6)

階層ベイズモデル (先端科学トピックス)

久保拓弥 [kubo@ees.hokudai.ac.jp](mailto:kubo@ees.hokudai.ac.jp)

茨城大集中講義 <http://goo.gl/2QNwg1>

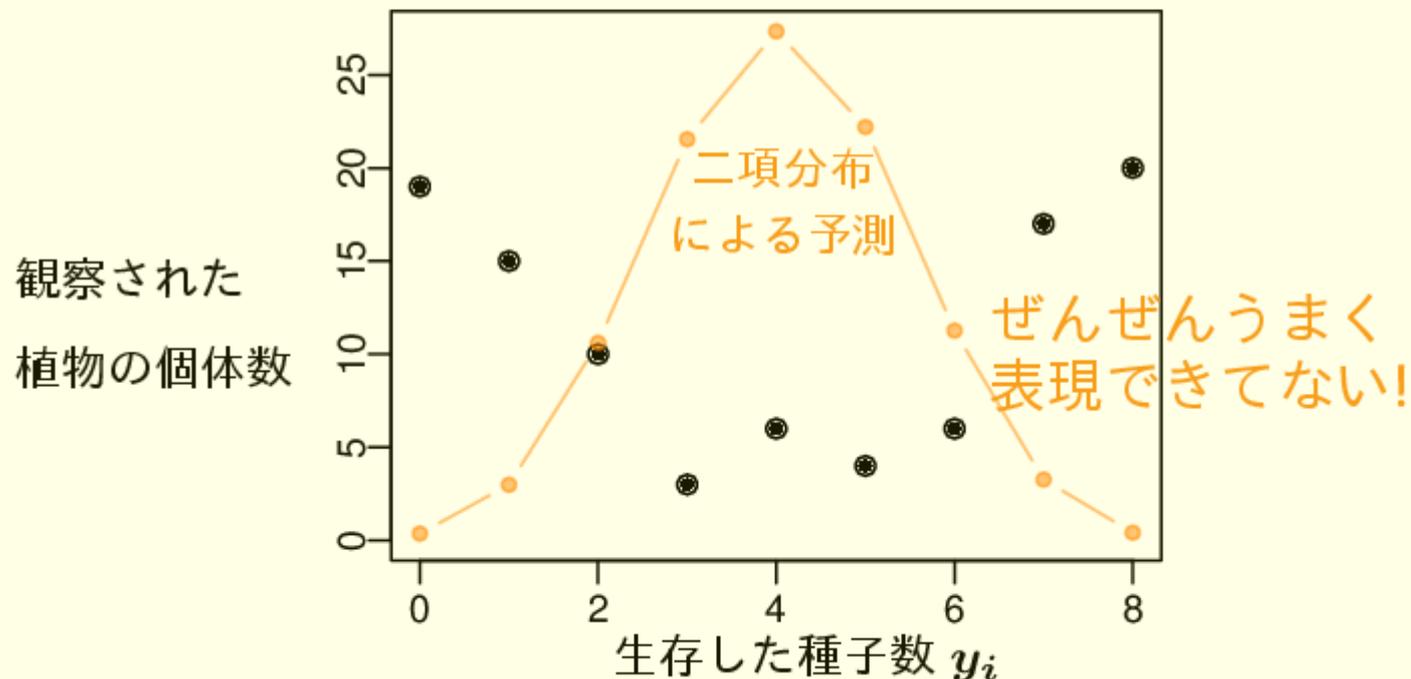
2014-10-01

ファイル更新時刻: 2014-09-29 12:44

# 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……



要点： **現実のデータ** を解析するためには，さらに **工夫が必要!**

# 現実のデータの複雑さを表現できる 統計モデルをつくる!

個体差・地域差・生物種差・  
空間相関・時間相関など  
めんどろなことをあつかわない  
といけない

GLM にそういう要因を組みこむ

データに複雑なモデルをあてはめる  
工夫をする (パラメーター推定法の  
改善)

## 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