

統計モデリング入門 2015 (a)
 An Introduction to Statistical Modeling
 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

2015-07-06 統計モデリング入門 2015 (1) 1/46

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

- Why in Japanese? ... because even in Japanese, **statistics is difficult** for Japanese students to understand.
- I will **compensate for language disadvantages** in foreign students when I give grades.
- **Questions in English are always welcomed!**

2015-07-06 統計モデリング入門 2015 (1) 2/46

統計モデリング授業の web page
<http://goo.gl/76c4i>

植物生態学特論 I (Advanced Course of Plant Ecology I)
 生態学の統計モデリング 7 月 6 日から
 Statistical Modeling for Ecology, commence on July 6
 13:00 - 14:30, Monday and Wednesday
 担当: 久保拓弥 (KUBO Takuya) kubo@ees.hokudai.ac.jp

- 授業メーリングリスト (Course Mailing List)
 - <http://goo.gl/f0vCn8> で早めに登録してください。登録のユーザインターフェイスは日本語に必要できます。授業の資料ダウンロードの連絡などします。
 - subscribe at <http://goo.gl/f0vCn8> immediately, or you can not download course handouts.
- 授業 Web Site: <http://goo.gl/76c4i>

2015-07-06 統計モデリング入門 2015 (1) 3/46

統計モデリング授業 Mailing List
<http://goo.gl/f0vCn8>

植物生態学特論 I (Advanced Course of Plant Ecology I)
 生態学の統計モデリング 7 月 6 日から
 Statistical Modeling for Ecology, commence on July 6
 13:00 - 14:30, Monday and Wednesday
 担当: 久保拓弥 (KUBO Takuya) kubo@ees.hokudai.ac.jp

- 授業メーリングリスト (Course Mailing List)
 - <http://goo.gl/f0vCn8> で早めに登録してください。登録のユーザインターフェイスは日本語に必要できます。授業の資料ダウンロードの連絡などします。
 - subscribe at <http://goo.gl/f0vCn8> immediately, or you can not download course handouts.
- 授業 Web Site: <http://goo.gl/76c4i>

2015-07-06 統計モデリング入門 2015 (1) 4/46

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

- 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

2015-07-06 統計モデリング入門 2015 (1) 5/46

Performance Rating

- E-mail assignment (via Mailing List)
 - **That's ALL!**
- Attendance? NOT care.

2015-07-06 統計モデリング入門 2015 (1) 6/46

What for Statistical Modeling?
 なぜデータ解析の方法を勉強しなければならぬのか?

All you depend on statistics whenever you conclude something based on your data

- データ解析がおかしいと **結論もおかしい**
- Crazy data analysys → Crazy results
- 統計解析わからんと批判的に読めない
- A lack of statistical knowledge → no critical reading of papers

2015-07-06 統計モデリング入門 2015 (1) 8/46


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

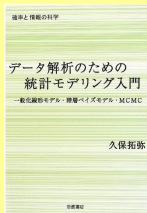
- ブラックボックス統計解析
- No “Blackbox” statistics!
- とにかく「ゆーい差」さえ出せばよいという発想になっている
- Don't blindly believe “Significance” !

この授業のねらい (aim)

できるだけ内容を理解して統計ソフトウェアを使おう!

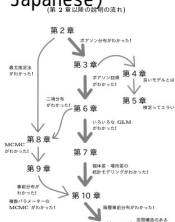
- Understand how to fit statistical models to your data データにあてはめられる統計モデルを作ろう
- Use the statistical software  to show your data structure

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

my text book (in Japanese)
 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



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

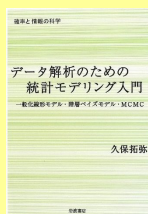
統計学の勉強には良い統計ソフトウェアが必要!

- 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



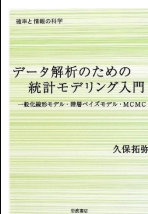
統計モデルとは何か?



「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づき準備されていて、モデルがデータにどれく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

階層ベイズモデル (MCMC)
 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つつきたい
 一般化線形モデル (最小二乗法)
 正規分布以外の確率分布をつつきたい
 線形モデル

2015-07-06 統計モデリング入門 2015 (1) 17/46

「統計モデリング入門」の主張

right probability distribution for right statistical modeling

The Evolution of Linear Models

Hierarchical Bayesian Model (HBM) (Parameter Estimation MCMC)
 Generalized Linear Mixed Model (GLMM) (最尤推定法)
 Generalized Linear Model (GLM) (MLE)
 Linear Model (MSE)

2015-07-06 統計モデリング入門 2015 (1) 18/46

たとえばこんなデータがあったしましょう

An example (次の時間の例題)

種子数 (number of seeds)
 個体サイズ (plant body size)

図 3.1 この初期に登場する架空植物の第1番目の個体。この植物の体サイズ(個体の大きさ)と、肥料をやる施肥処理が、種子数にどう影響しているのかを知りたい。

2015-07-06 統計モデリング入門 2015 (1) 19/46

一般化線形モデル - ばらつきをよく見る

Don't use the normal distribution without seeing data!

(A) 正規分布・恒等リンク関数の統計モデル
 正規分布
 線形モデルの発展
 階層ベイズモデル (MCMC)
 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つつきたい
 一般化線形モデル (最小二乗法)
 正規分布以外の確率分布をつつきたい
 線形モデル

(B) ポアソン分布・対数リンク関数の統計モデル
 ポアソン分布

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をあてはめたとき、誤差はほとんどに変化する平均値、グレイで

2015-07-06 統計モデリング入門 2015 (1) 20/46

全体の流れ (1/2)

- 第 1 回: 7/06 (月) 観測されたパターンを説明する統計モデル
 Introduction
- 第 2 回: 7/08 (水) 確率分布と最尤推定
 Probability Distributions and Maximum Likelihood Estimation (MLE)
- 第 3 回: 7/13 (月) 一般化線形モデル: ポアソン回帰
 Generalized Linear Model (GLM): Poisson Regression

全体の流れ (2/2)

- 第 4 回: 7/15 (水) モデル選択と検定
 Model Selection and Statistical Test
- 第 5 回: 7/22 (水) 一般化線形モデル: ロジスティック回帰
 GLM: Logistic Regression
- 第 6 回: 7/27 (月) 一般化線形混合モデル
 Generalized Linear Mixed Model (GLMM)
- 第 7 回: 7/29 (水) 階層ベイズモデル
 Bayesian GLMM and Markov Chain Monte Carlo

7/8 (水)

統計モデリング入門 2015 (b)

probability distribution and maximum likelihood estimation
 確率分布と最尤推定

久保拓弥 kubo@ees.hokudai.ac.jp

北大環境科学の講義 <http://goo.gl/76c4i>

2015-07-08

ファイル更新時刻: 2015-07-02 16:53

kubostat2015b (<http://goo.gl/76c4i>) 統計モデリング入門 2015 (b) 2015-07-08 1 / 37

単純化した例題

```

> data
[1] 2 2 4 4 5 2 3 1 2 0 4 3 3 3 3 4 2 7 2 4 3 3 4
[26] 3 7 9 3 1 7 6 6 5 2 4 7 2 2 6 2 4 6 4 5 1 3 2 3
    
```

Histogram of data
 data

2015-07-06 統計モデリング入門 2015 (1) 24/46

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

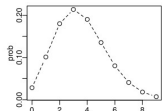


図 4 平均 $\lambda = 3.56$ のポアソン分布。種子数 y とその確率 $prob$ の関係が示されている。図 5 の表を用いたもの。表の `plot()` 関数の引数: `type = "p"` によって「点と折れ線による表示」。`log = TRUE` によって「折れ線は緑色で」も表示している。

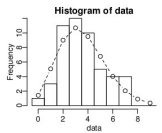
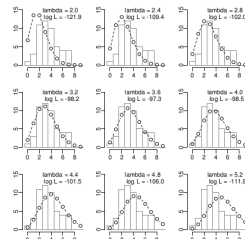


図 5 観測データと確率分布の対応をみるため、ヒストグラムと `plot()` を用い、それに重ねられている `lambda = 3.56` のポアソン分布の確率分布の表。平均 $\lambda = 3.56$ の確率分布のポアソン分布の確率分布に全観測データ y_i を重ねて表示される。

2015-07-06

25/46

さいゆう 最尤推定という考えかたを説明します



seek the maximum (Maximum estimate, λ)
対数尤度を最大化する λ をさがす

$$\text{対数尤度 } \log L(\lambda) = \sum (\ln \log \lambda - \lambda + \sum \log k)$$

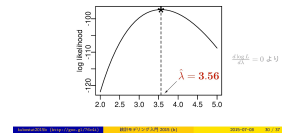


図 7 平均 λ (`lambda`) を変化させていったポアソン分布と、観測データへのあてはまりの良さ (対数尤度 `logL`) を示すヒストグラムと対数尤度関数。

2015-07-06

統計モデリング入門 2015 (1)

26/46

7/13 (月)

統計モデリング入門 2015 (c) Poisson regression, a generalized linear model (GLM) 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/78c41>

2015-07-13

ファイル更新時刻: 2015-07-02 16:24

kubostat2015c (<http://goo.gl/78c41>) 統計モデリング入門 2015 (c) 2015-07-13 1 / 46

ここで登場する --- 「何でも正規分布」ではダメ! という発想

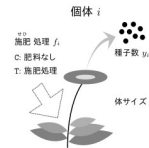
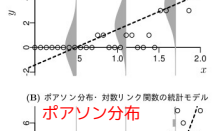


図 3.1 この例題に登場する実在植物の類: 番目の個体。この植物の体サイズ (個体の大きさ) x_i と肥料をやる施肥処理 f_i が種子数 y_i にどう影響しているのかを知りたい。

(A) 正規分布: 恒等リンク関数の統計モデル

正規分布



(B) ポアソン分布: 対数リンク関数の統計モデル

ポアソン分布

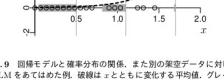


図 3.9 回帰モデルと確率分布の関係。また別の実在データに対して GLM をあてはめた例。縦軸 y とともに変化する平均値、グレイで

2015-07-06

統計モデリング入門 2015 (1)

28/46

Free の統計ソフトウェア R で統計モデリング

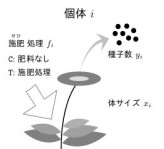


図 3.1 この例題に登場する実在植物の類: 番目の個体。この植物の体サイズ (個体の大きさ) x_i と肥料をやる施肥処理 f_i が種子数 y_i にどう影響しているのかを知りたい。

```
結果を格納するオブジェクト
fit <- glm(
  y ~ x, # モデル式
  family = poisson(link = "log"), # リンク関数の指定 (省略可)
  data = train) # データフレームの指定
```

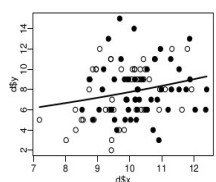


図 17 平均種子数の予測。図 13 に x の予測値 (実験) を上がしたものを。

2015-07-06

統計モデリング入門 2015 (1)

29/46

7/15 (水)

統計モデリング入門 2015 (d)

モデル選択と検定

久保拓弥 kubo@ees.hokudai.ac.jp

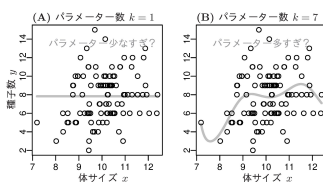
北大環境科学院の講義 <http://goo.gl/78c41>

2015-07-15

ファイル更新時刻: 2015-07-02 16:24

Q. モデル選択とは何か?

パラメーター数は多くても少なくてもヘン?



What is the "best?" parameter number k ?

kubostat2015d (<http://goo.gl/78c41>) 統計モデリング入門 2015 (d) 2015-07-15 4 / 37

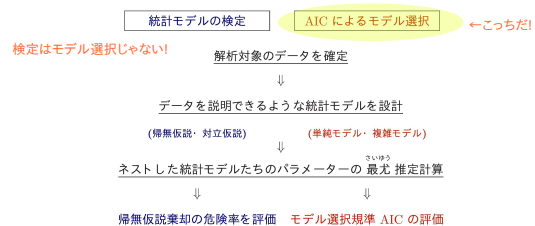
2015-07-06

統計モデリング入門 2015 (1)

31/46

A. より良い予測をする統計モデルを探すこと

But their procedures are similar
しかしモデル選択と検定の手順は途中まで同じ



2015-07-06

統計モデリング入門 2015 (1)

32/46

7/29 (水)

統計モデリング入門 2015 (e)

階層ベイズモデル

久保拓弥 kubo@ees.hokudai.ac.jp

北大環境科学院の講義 <http://goo.gl/7f6c41>

2015-07-29

ファイル更新時刻: 2015-07-02 16:24

kubostat2015e (<http://goo.gl/7f6c41>) 統計モデリング入門 2015 (e) 2015-07-29 1 / 87

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……

さっきの例題と同じようなデータなのに?
(「統計モデリング入門」第 10 章の最初の例題)

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

データ: 8 個中の $Y[i]$ 個の種子が生存 σ は hyper parameter

二項分布: 生存確率 $q[i]$ ← $r[i]$ ← 植物の個体差

事前分布: σ ← σ は σ ばらつき

無情報事前分布: a ← σ は σ と思ってください

無情報事前分布 (超事前分布)

矢印は手順ではなく、依存関係をあらわしている

2015-07-06 統計モデリング入門 2015 (1) 43/46

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

そういう難しい状況では……

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

2015-07-06 統計モデリング入門 2015 (1) 44/46

全体の流れ (1/2)

第 1 回: 7/06 (月) 観測されたパターンを説明する統計モデル
Introduction

第 2 回: 7/08 (水) 確率分布と最尤推定
Probability Distributions and Maximum Likelihood Estimation (MLE)

第 3 回: 7/13 (月) 一般化線形モデル: ポアソン回帰
Generalized Linear Model (GLM): Poisson Regression

全体の流れ (2/2)

第 4 回: 7/15 (水) モデル選択と検定
Model Selection and Statistical Test

第 5 回: 7/22 (水) 一般化線形モデル: ロジスティック回帰
GLM: Logistic Regression

第 6 回: 7/27 (月) 一般化線形混合モデル
Generalized Linear Mixed Model (GLMM)

第 7 回: 7/29 (水) 階層ベイズモデル
Bayesian GLMM and Markov Chain Monte Carlo