

統計モデリング入門 2018 (a)  
 生物多様性学特論  
 An overview: 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



図 3.1 この問題に登場する架空植物の高さ(茎の長さ)、この植物の体サイズ(茎の長さ)  $y_1$  と肥料をやる量(施肥量)  $x_1$  が種子数  $y_2$  にどう影響しているのかを知りたい。

2018-06-18 統計モデリング入門 2018a 1/56

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!**

2018-06-18 統計モデリング入門 2018a 2/56

Performance Rating

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

2018-06-18 統計モデリング入門 2018a 3/56

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
  - **Send your assignment via the class ML**
- 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

2018-06-18 統計モデリング入門 2018a 4/56

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

mailing list  
<http://goo.gl/f0vCn8>

2018-06-18 統計モデリング入門 2018a 5/56

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 analysis → 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

2018-06-18 統計モデリング入門 2018a 7/56


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" !**

2018-06-18 統計モデリング入門 2018a 8/56

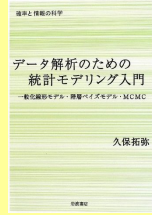

この授業のねらい (aim)

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

2018-06-18 統計モデリング入門 2018a 9/56

教科書とソフトウェア


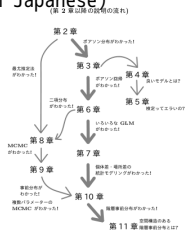



2018-06-18 統計モデリング入門 2018a 10/56

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!!


2018-06-18 統計モデリング入門 2018a 11/56

Statistical software for this course

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

追記メモ: RStudio の紹介!

2018-06-18 統計モデリング入門 2018a 12/56

統計モデルとは何か?

What? statistical modeling?




2018-06-18 統計モデリング入門 2018a 13/56

「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

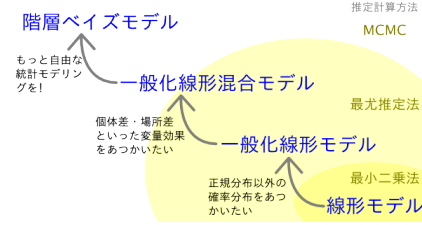


2018-06-18 統計モデリング入門 2018a 14/56

「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  
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  
最小二乗法

最尤推定法

推定計算方法 MCMC

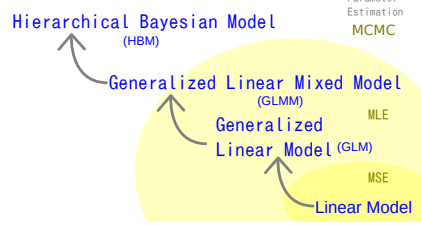


2018-06-18 統計モデリング入門 2018a 15/56

GLM and extended GLMs!

a better statistica model for better data analysis!

The Evolution of Linear Models



Hierarchical Bayesian Model (HBM)

Generalized Linear Mixed Model (GLMM)


Generalized Linear Model (GLM)

Linear Model

Parameter Estimation MCMC

MLE

MSE



2018-06-18 統計モデリング入門 2018a 16/56



さいゆう  
最尤推定という考えかたを説明します

対数尤度  $\log L(\lambda) = \sum (n_i \log \lambda - \lambda - \sum x_i \log x_i)$

How to fit the distribution to the observation?  
Maximum likelihood estimation!

2018-06-18 統計モデリング入門 2018a 25/56

6/25  
Overview  
Statistical Modeling 2018 (c)  
Poisson regression  
and generalized linear model  
ポアソン回帰と GLM

ここで登場する ---  
「何でも正規分布」ではダメ! という発想

the "normal" distribution is NOT "normal"

2018-06-18 統計モデリング入門 2018a 27/56

Free の統計ソフトウェア R で統計モデリング

2018-06-18 統計モデリング入門 2018a 28/56

6/25  
Overview  
Statistical Modeling 2018 (d)  
Model Selection  
and Statistical Test  
モデル選択と統計学的検定

statistical model selection  
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ともヘン?  
What is the "best?" parameter number k?  
2018-06-18 統計モデリング入門 2018a 30/56

model selection for better predictions  
A. より良い予測をする統計モデルを探すこと

統計学の方法が... 異なるが、その目的は同じ  
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

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 モデル選択標準 AICの評価

2018-06-18 統計モデリング入門 2018a 31/56

統計学って「検定」のこと?  
「検定」って何なの?  
fallacy of statistical significance?

2018-06-18 統計モデリング入門 2018a 32/56

7/2

## Overview

### Statistical Modeling 2018 (e)

Logistic regression,  
a generalized linear model

ロジスティック回帰

measurement / mesurement?... sounds bad!  
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

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

2012-11-02 k4 (2012-10-26 17:07 修正版) 14 / 44

2018-06-18 統計モデリング入門 2018a 34/56

Use logistic regressions!  
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!

2012-11-02 k4 (2012-10-26 17:07 修正版) 43 / 44

2018-06-18 統計モデリング入門 2018a 35/56

GLM のひとつ, ロジスティック回帰を使おう

またいつもの例題? ..... ちょっとちがう

ロジスティック回帰とは何なのか?

8個の種子のうち y 個が発芽可能だった! ..... というデータ

(A) 観測データの一部 (y=0) (B) 推定される確率

二項分布: N 回のうち y 回, となる確率

a statistical model  
for fractions  
using binomial distributions

2018-06-18 統計モデリング

7/2

## Overview

### Statistical Modeling 2018 (f)

Hierarchical Bayesian model  
and MCMC sampling

階層ベイズモデルと MCMC

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.....

GLM does NOT work?!

ぜんぜんうまく表現できてない!

さっきの例題と同じようなデータなの? (「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, ここんどはベイズモデルを使ってモデリングします

A solution: Hierarchical Bayesian GLM  
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

sigma は hyper parameter

二項分布: 生存確率 q[i]

植物の個体差: r[i]

事前分布: 個体差の sigma は a ばらつき

全体平均: a

無情報事前分布: 無情報事前分布 (超事前分布)

sigma は a と思ってください

矢印は手順ではなく, 依存関係をあらわしている

2012-11-02 k4 (2012-10-26 17:07 修正版) 43 / 44

2018-06-18 統計モデリング入門 2018a 39/56

なぜ階層ベイズモデルまで勉強するのか?

生態学!

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

The Evolution of Linear Models

- Linear Model
- Generalized Linear Model (GLM)
- Generalized Linear Mixed Model (GLMM)
- Hierarchical Bayesian Model (HBM)
- Parameter Estimation: MCMC

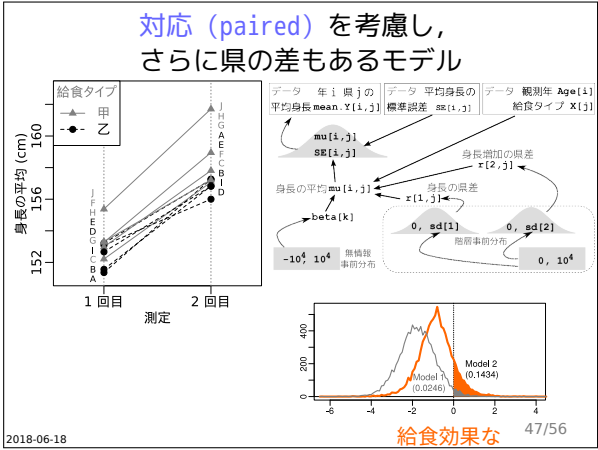
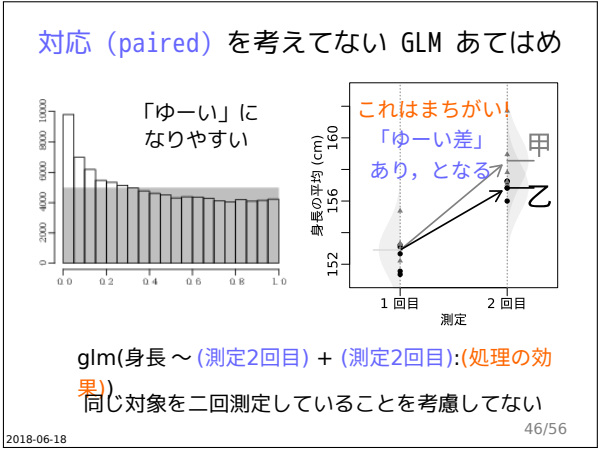
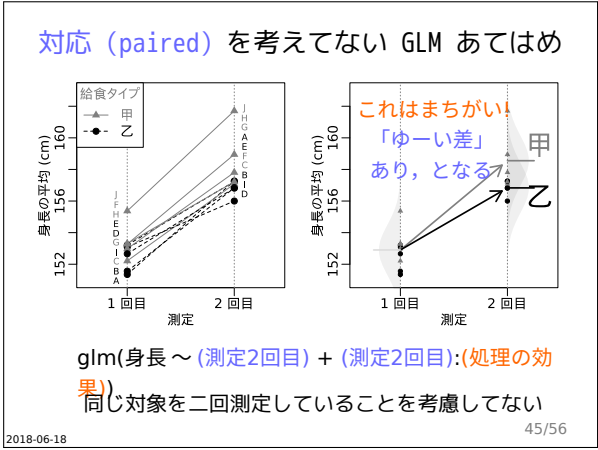
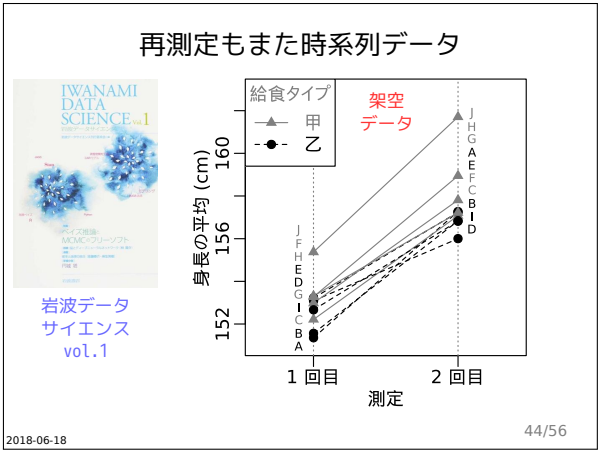
What for hierarchical Bayesian modeling? --- to detect interesting effects embedded in noisy & dirty data in the field of Ecology!

2018-06-18 統計モデリング入門 2018a 40/56

第 7, 8 回は  
 「時間変化」するデータの  
 統計モデリング  
 (階層ベイズモデルの応用)  
 Modeling of time-series data as  
 an application of hierarchical  
 Bayesian modeling!

7/9  
 Overview  
 Statistical Modeling 2018 (g)  
 Modeling time change data  
 (short term)  
 短い時系列データの統計モデル

A Time series model  
 for single step data  
 短い時系列データ  
 時系列の長短に関係なく  
 「対応のある」データ点が  
 どうか本質的な問題



7/9  
 Overview  
 Statistical Modeling 2018 (h)  
 Modeling time series data  
 (long term)  
 長い時系列データの統計モデル



7/27 (水)

生態学の時系列データ解析でよく見る  
『あぶない』モデリング

久保拓弥 <mailto:kubo@ees.hokudai.ac.jp>

2017-07-03 kubostat2017 (h) 1/52

時間相関のある時系列データに…  
time series data and autocorrelation

$Y$

$glm(y \sim t)$

…と、モデルを  
あてはめてみた

2018-06-18 統計モデリング入門 2018a 50/56

「やったーゆーいだ!!」……??  
A fake significance

`> summary(glm(formula = y ~ t))`

Deviance	Residuals:
	Min
	1Q
	Median
	3Q
	Max

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-414.5655	71.4761	-5.80	6.6e-06
t	0.2339	0.0357	6.55	1.1e-06

これはまちがい →  $glm(\text{時系列} Y \sim \text{時間 } t)$

統計モデルがおかしい?

2018-06-18 統計モデリング入門 2018a 51/56

時系列の「ずれ」  
temporal autocorrelation

GLM のずれ  
independent noises

直線からのずれがちがう!

時間的自己相関がある      時間的自己相関がない

2018-06-18 統計モデリング入門 2018a 52/56

統計モデルづくりの要点

時系列データの解析は  
階層ベイズモデル化した  
状態空間モデルを使うのが便利

Latent state model is a better model to know the characteristics of time-series data

変数  $Y$

Random walk  
もっとも単純な  
モデル

$N(Y_1, \sigma) \rightarrow Y_2$

$N(Y_2, \sigma) \rightarrow Y_3$

正規分布

$N(Y_3, \sigma) \rightarrow Y_4$

時間  $t$

2018-06-18 統計モデリング入門 2018a 54/56

状態空間モデル + 観測モデル  
Latent state variables + observation model

観測の誤差      状態空間モデル

$N(y_t, \sigma_2) \rightarrow Y_t$       二種類の  $\sigma$  をもつ

観測データ  $Y_1, Y_2, Y_3$

$N(y_t, \sigma_1) \rightarrow y_{t+1}$

状態変数の変化      時間  $t$

観測できない世界 (状態空間)

2018-06-18 統計モデリング入門 2018a 55/56

今日はここまで

any questions?