

統計モデリング入門 2019 (a)
 久保の多様性生物学の講義全体の流れの紹介
 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菜豆植物の葉：葉目の個体、この植物の体サイズ(個体の大きさ) x_i と肥料をやる無肥処理 J が種子数 y_i にどう影響しているのかを知りたい。

2019-07-22 統計モデリング入門 2019a 1/60

この統計モデリング授業の
 みなさんの質問に基づく成績評価

- 全 8 回の授業、それぞれについて
- 「よくわからなかった」箇所などなど…について、久保あてにメールを送ってください。質問の長さは自由
- 例「ポアソン分布と二項分布のちがいがわかりません」などなど…
- 「ここまでわかっているんですよ」といったアピールも可能
- よく考えられた質問には、よい点をあげたい…

2019-07-22 統計モデリング入門 2019a 2/60

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

- 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

2019-07-22 統計モデリング入門 2019a 3/60

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

この授業のポータル?

2019-07-22 統計モデリング入門 2019a 4/60

この授業の主題
 よい統計モデリング
 とは何だろう?

…あるいは…
 悪いモデルとは?

Ecologists' data-cooking
 nice data!



<http://www.wstephens.net/>

A caricatured example

Ecologists' data-cooking
 nice data!



<http://www.wstephens.net/>



data-chopping
<http://a1ispagnola.blogspot.com/>

Data-cooking in Ecology



<http://www.wstephens.net/>



data-chopping
<http://a1ispagnola.blogspot.com/>

data-flaming
<http://howdoifeelaboutthis.tumblr.com/>

Ecologists' data-cooking

Chainsaw Chop-Suey
data-chopping
data-flaming
your conclusion!

データ解析は統計モデルを作り!

- データ解析とは、統計モデルをデータにあてはめてみること!
- 統計モデルがおかしいと結論もおかしい
- データの性質・構造などをよくみて、よりよい統計モデルを作ろう (試行錯誤が重要)

まあ、授業しながら、おいしいと...

2019-07-22 統計モデリング入門 2019a 10/60

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

- (この授業限定の用語) 「ゆーい差」とは?
- 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
- ブラックボックス統計解析

2019-07-22 統計モデリング入門 2019a 11/60

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!

- 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

2019-07-22 統計モデリング入門 2019a 12/60

教科書とソフトウェア

2019-07-22

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!!

2019-07-22 統計モデリング入門 2019a 14/60

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!

2019-07-22 統計モデリング入門 2019a 15/60

統計モデルとは何か?

What? statistical modeling?

2019-07-22

「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
どれくらい良くあてはまっているかを
定量的に評価できる

2019-07-22 統計モデリング入門 2019a 17/60

「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

推定計算方法 MCMC
最尤推定法
最小二乗法

2019-07-22 統計モデリング入門 2019a 18/60

GLM and extended GLMs!

a better statistica model for better data analysis!

The Evolution of Linear Models

Hierarchical Bayesian Model (HBM) Parameter Estimation MCMC

Generalized Linear Mixed Model (GLMM) MLE

Generalized Linear Model (GLM) MSE

Linear Model

2019-07-22 統計モデリング入門 2019a 19/60

たとえばこんなデータがあったしましょう

An example

number of seeds
種子数

個体 i
種子数 y_i
体サイズ x_i

施肥処理 f_i
C: 肥料なし
T: 施肥処理

plant body size
体サイズ

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。

2019-07-22 統計モデリング入門 2019a 20/60

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!

正規分布

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

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。縦軸は y とともに変化する平均値。グレイで

2019-07-22 統計モデリング入門 2019a 21/60

全体の流れ (1/3)

第 1 回 観測されたパターンを説明する統計モデル
Introduction

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

第 3 回 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 (2/3)

第 4 回 モデル選択と検定
Model Selection and Statistical Test

第 5 回 一般化線形モデル: ロジスティック回帰
GLM: Logistic Regression

第 6 回 階層ベイズモデル 1
Hierarchical Bayesian Models (HBM) 1

全体の流れ (3/3)

第 7 回 繰り返し測定の階層ベイズモデル
Bayesian models for repeated measurements

第 8 回 時間変化データのベイズ統計モデル
Bayesian models for time series data

The End!

Overview Statistical Modeling 2019 (b)

Probability distributions and maximum likelihood estimation

さまざまな確率分布と最尤推定

単純化した例題

こんなデータ (架空) があってとしましょう

まあ、なんだからこういうヘンな植物を調査しているときは

個体 i 種子数 y_i

全 50 個体 $y = (1, 2, 3, \dots, 50)$ この $\{y_i\}$ が観測データ!

このデータ $\{y_i\}$ がすでに R という統計ソフトウェアに格納されていた、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 4 6 1 3 2 3
```

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))
```

Histogram of data

データ解析における重要な事項とにかく図を張れ!

Simplified examples to learn statistical modeling

2019-07-22
統計モデリング入門 2019a
26/60

カウントデータとは、たとえば植物の種子が1個、2個、3個、...と数えられるデータのことをさす。

人間の身長などはカウントデータではない

2019-07-22 統計モデ...

さいゆう 最尤推定という考えかたを説明します

対数尤度を最大化する λ をさがす

対数尤度 $\log L(\lambda) = \sum (y_i \log \lambda - \lambda - \sum \log A)$

2019-07-22 統計モデリング入門 2019a 28/60

Overview Statistical Modeling 2019 (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

個体 i

施肥処理 f_i
○: 肥料なし
□: 施肥処理

種子数 y_i

体サイズ x_i

(A) 正規分布・相関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

the "normal distribution is NOT "normal"

2019-07-22 統計モデリング入門 2019a 30/60

Free の統計ソフトウェア R で統計モデリング

```
結果を格納するオブジェクト
関数名      確率分布の指定
fit <- glm(y ~ x, family = poisson(link = "log"), data = d)
) data.frame の指定      リンク関数の指定 (省略可)
```

2019-07-22 統計モデリング入門 2019a 31/60

Overview Statistical Modeling 2019 (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てもヘン?

What is the "best?" parameter number k ?

2019-07-22 統計モデリング入門 2019a 33/60

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

統計学は簡単 そして、その評価も簡単
But their procedures are similar
しかしモデル選択と検定の手順は途中で同じ

統計モデルの検定 ←こっちだ!
AICによるモデル選択

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

2019-07-22 統計モデリング入門 2019a 34/60

統計学って「検定」のこと?

「検定」って何なの?

fallacy of statistical significance?

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 ($\beta_0 = 2.06$, β_1 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。

2019-07-22 統計モデリング入門 2019a 35/60

Overview

Statistical Modeling 2019 (e)

Logistic regression, a generalized linear model

ロジスティック回帰

measurement / measurement?... 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

1. データどどん割算・割算
2. 何でもいからひたすらセンをひく
3. ゆーいいでたら万歳・万歳
4. 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

2019-07-22 統計モデリング入門 2019a 37/60

Important: statistical education for graduate students in ecology

Photo taken by Dr. GM ISM public lecture "MCMC and Bayesian modeling" in 2009, 2010

統計モデリング教育が重要

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: 不自然にデータをごねくりまわさない
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

2019-07-22 統計モデリング入門 2019a 39/60

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な?

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) 観測データの一例 ($y=3$) (B) 商業されるモデル

a statistical model for fractions using binomial distributions

二項分布: N 回のうち y 回, となる確率

2019-07-22 統計モデリング

8/5

Overview Statistical Modeling 2019 (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……

GLMs do NOT work?!

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

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 σ は hyper parameter

二項分布 生存確率 $q[i]$ 植物の個体差 $r[i]$

全体の平均 a 事前分布 個体差のばらつき σ

無情報事前分布 無情報事前分布 (超事前分布) σ 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

2019-07-22 統計モデリング入門 2019a 43/60

なぜ階層ベイズモデルまで勉強するの?

たとえば生態学とかでは…

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

生態学にかぎらず、実験などにおける反復でも、必要

2019-07-22 統計モデリング入門 2019a 44/60

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!

Overview Statistical Modeling 2019 (g)

Modeling time change data (short term)

短い時系列データの統計モデル

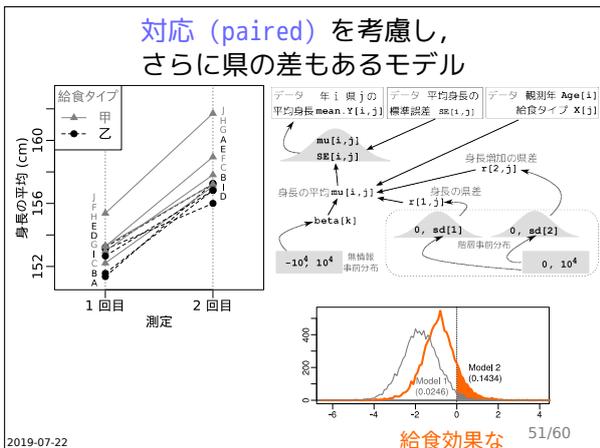
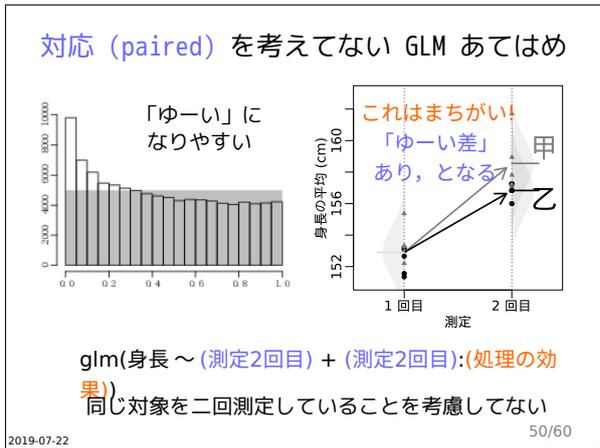
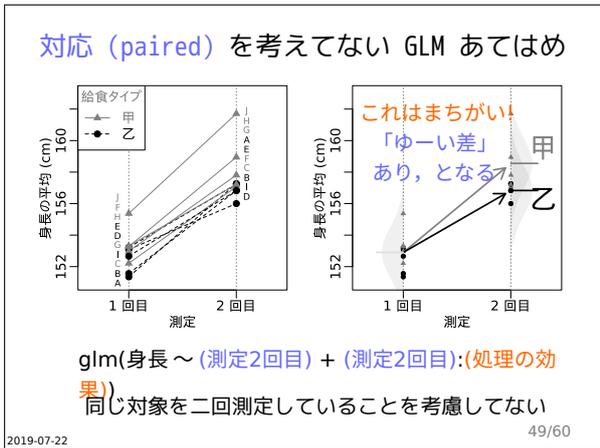
A Time series model for single step data 短い時系列データ

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

再測定もまた時系列データ

岩波データサイエンス vol.1

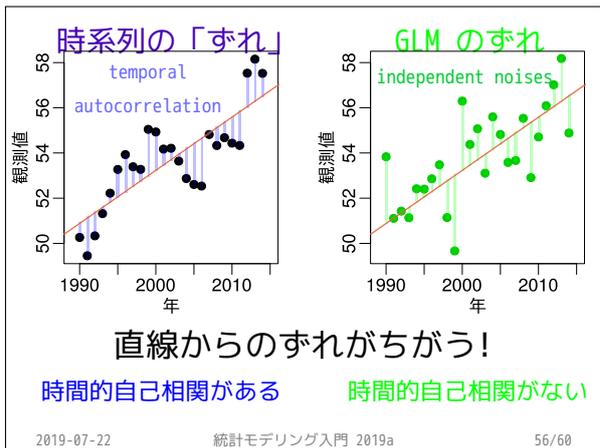
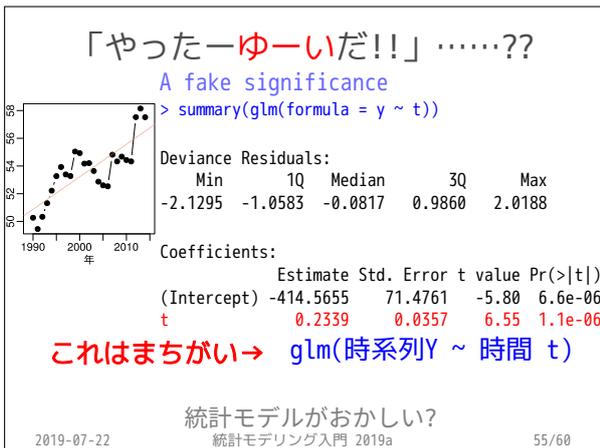
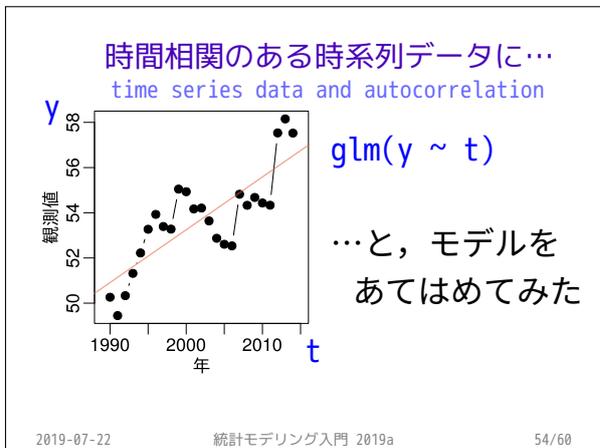
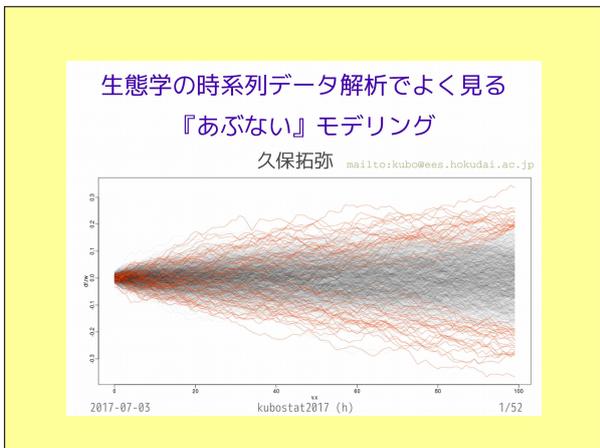
2019-07-22 48/60



Overview
Statistical Modeling 2019 (h)

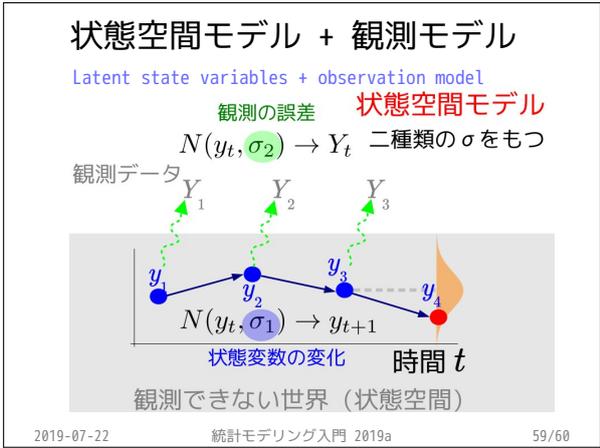
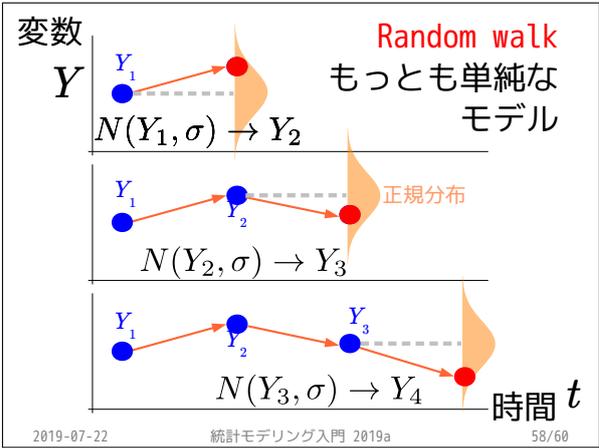
Modeling time series data (long term)

長い時系列データの統計モデル



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



この時間はここまで

any questions?