

統計モデリング入門 2017 (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


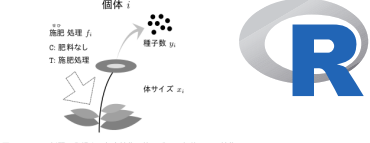



図 3.1 この問題に登場する架空植物の例。緑色の楕円、この植物の体サイズ(個体の大きさ) x_i と肥料をやる施肥処理 f_i が種子数 y_i にどう影響しているのかわかりたい。

2017-06-05 統計モデリング入門 2017a 1/59

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!**

2017-06-05 統計モデリング入門 2017a 2/59

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 6月5日から
Statistical Modeling for Ecology, commence on June 5
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> に登録してください。授業の資料ダウンロードの連絡などします。単位を取得する院生は必須。登録のユーザーインターフェイスは日本語に必要です。
 - subscribe at <http://goo.gl/f0vCn8> immediately, or you can not download course handouts.
- 授業 Web Site: <http://goo.gl/76c4i>

2017-06-05 統計モデリング入門 2017a 3/59

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 6月5日から
Statistical Modeling for Ecology, commence on June 5
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> に登録してください。授業の資料ダウンロードの連絡などします。単位を取得する院生は必須。登録のユーザーインターフェイスは日本語に必要です。
 - subscribe at <http://goo.gl/f0vCn8> immediately, or you can not download course handouts.
- 授業 Web Site: <http://goo.gl/76c4i>

2017-06-05 統計モデリング入門 2017a 4/59

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 ML**
- 成績評価は「課題」の回答
 - 出欠関係なし (欠席の連絡りません)
- 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

2017-06-05 統計モデリング入門 2017a 5/59

Performance Rating

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

2017-06-05 統計モデリング入門 2017a 6/59

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

2017-06-05 統計モデリング入門 2017a 8/59


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

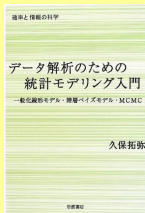
- ブラックボックス統計解析
- No “Blackbox” statistics!
- とにかく「ゆーい差」さえ出せばよいという発想になっている
- Don't blindly believe “Significance” !

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

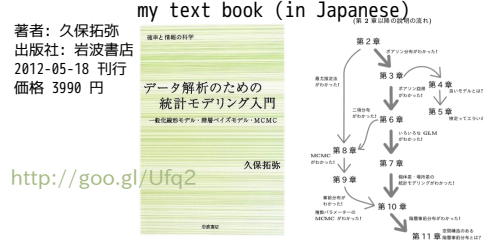
教科書とソフトウェア



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!!

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!

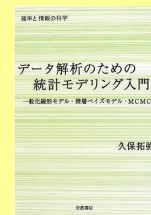
- 無料で入手できる
- 内容が完全に公開されている
- 多くの研究者が使っている
- 作図機能が強力

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



統計モデルとは何か?

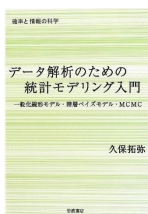
What? statistical modeling?



「統計モデル」とは何か?

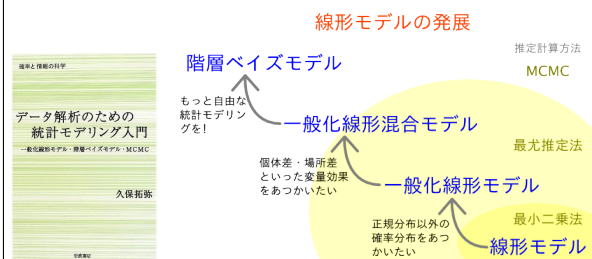
どんな統計解析においても統計モデルが使用されている

- 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ぎが準備されていて、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!



「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) ← MLE

Generalized Linear Model (GLM) ← MSE

Linear Model

データ解析のための統計モデリング入門
一般化線形モデル・階層ベイズモデル・MCMC

久保拓弥

2017-06-05 統計モデリング入門 2017a 17/59

たとえばこんなデータがあったらしよう

An example

number of seeds

種子数

個体 i

種子数 y_i

体サイズ x_i

体サイズ

plant body size

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を知りたい。

2017-06-05 統計モデリング入門 2017a 18/59

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

一般化線形混合モデル

一般化線形モデル

線形モデル

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!!

2017-06-05 統計モデリング入門 2017a 19/59

全体の流れ (1/3)

第 1 回: 6/05 (月) 観測されたパターンを説明する統計モデル
Introduction

6/07 (水) (授業なし no class) ← 奇怪な学年暦

第 2 回: 6/12 (月) 確率分布と最尤推定
Probability Distributions and Maximum Likelihood Estimation (MLE)

第 3 回: 6/14 (水) 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: 6/19 (月) モデル選択と検定
Model Selection and Statistical Test

第 5 回: 6/21 (水) 一般化線形モデル: ロジスティック回帰
GLM: Logistic Regression

第 6 回: 6/26 (月) 一般化線形混合モデル
Generalized Linear Mixed Model (GLMM)

全体の流れ (3/3)

第 7 回: 6/28 (水) 階層ベイズモデル
Bayesian GLMM and Markov Chain Monte Carlo

第 8 回: 7/03 (月) 時間変化データの統計モデル
Time change data analysis: common mistakes

(tentative)

6/12 (月)

統計モデリング入門 2017 (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

北大環境学院の講義 <http://goo.gl/T6c4i>

2017-06-12

ファイル更新時刻: 2017-05-24 16:30

kubostat2017b (<http://goo.gl/T6c4i>) 統計モデリング入門 2017 (b) 2017-06-12 1 / 42

単純化した例題

number of seeds per plant (individual)

こんなデータ (架空) があつたしましょう

まあ、なんだかどういうへんな植物を調査しているたします

個体 i のこの $\{a_i\}$ が観測データ!
 $f \in \{1, 2, 3, \dots, 50\}$
 $\{a_i\} = \{0, 20, \dots, 80\}$

このデータ $\{a_i\}$ がすでに何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 4
[26] 3 7 5 3 1 7 6 4 6 5 2 4 7 2 2 6 2 4 6 4 5 1 3 2 3
```

start with data plotting, always

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最重要事項
とにかく 図を感!

Simplified examples to learn statistical modeling

2017-06-05 統計モデリング入門 2017a 24/59

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

Find some appropriate probability distributions to fit the observed distributions

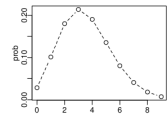


図4 平均 $\lambda = 3.56$ のポアソン分布、種子数 y とその確率 pmf の関係が示されている。図5の表を照らしよ。左の $plot()$ 関数の引数: $lambda = \lambda$ によって「丸と折れ線による図表」、 $log = TRUE$ によって「折れ線図表」が出力される。

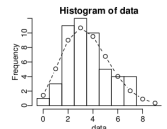
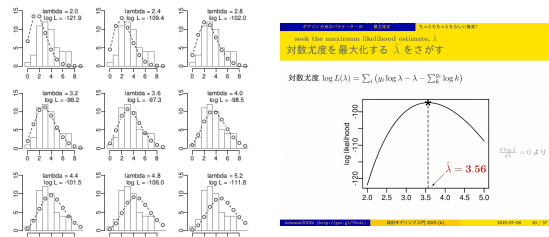


図5 観測データと確率分布の対応を定める。ヒストグラム $plot()$ を用いて、それに重ねた上で、左から順に y 軸の種子数 y と x 軸の頻度の関係、平均 $\lambda = \lambda$ の関係のポアソン分布の確率分布と観測データ y を重ねて見られる。

2017-06-05

統計モデリング入門 2017a

さいゆう 最尤推定という考えかたを説明します



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

2017-06-05

統計モデリング入門 2017a

26/59

6/14 (水)

統計モデリング入門 2017 (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/76c4i>

2017-06-14

ファイル更新時刻: 2017-05-24 16:30

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2017-06-14 1 / 47

ここで登場する --- 「何でも正規分布」ではダメ! という発想

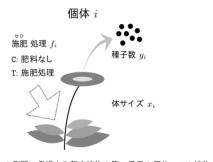


図 3.1 この例題に登場する実在植物の根; 番目の個体、この植物の体サイズ(個体の大きさ) x_i と肥料をやる施肥処理 f_i が種子数 y_i にどう影響しているのかを知りたい。

the "normal distribution is NOT "normal"

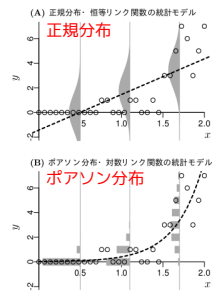


図 3.9 回帰モデルと確率分布の関係。また別の実在データに対して GLM をあてはめた例。縦軸は y とともに変化する平均値。グレイで

2017-06-05

統計モデリング入門 2017a

28/59

Free の統計ソフトウェア R で統計モデリング

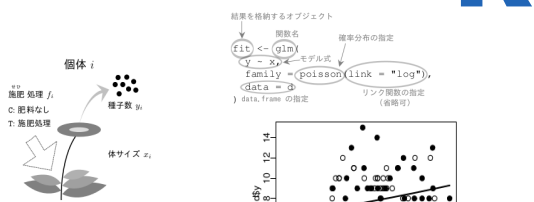


図 3.1 この例題に登場する実在植物の根; 番目の個体、この植物の体サイズ(個体の大きさ) x_i と肥料をやる施肥処理 f_i が種子数 y_i にどう影響しているのかを知りたい。

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

2017-06-05

統計モデリング入門 2017a

29/59

6/19 (月)

統計モデリング入門 2017 (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/76c4i>

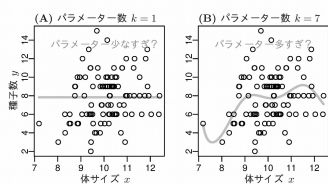
2017-06-19

ファイル更新時刻: 2017-05-24 16:32

kubostat2017d (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (d) 2017-06-19 1 / 44

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 ?

kubostat2017d (<http://goo.gl/76c4i>) 統計モデリング入門 2017a (d) 2017-06-19 4 / 37

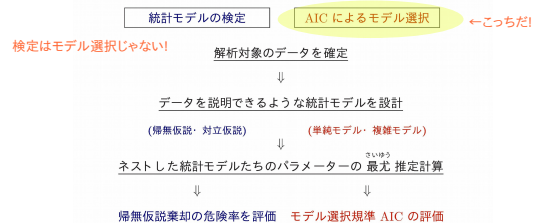
2017-06-05

統計モデリング入門 2017a

31/59

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 けれど、その詳細は異なる。しかしモデル選択と検定の手順は途中で同じ



kubostat2017d (<http://goo.gl/76c4i>) 統計モデリング入門 2017a (d) 2017-06-19 26 / 37

2017-06-05

統計モデリング入門 2017a

32/59

統計学って「検定」のこと?

「検定」って何なの?

「検定」ってエライの?

線形回帰が真の統計モデル
ということにしてしまう
($\beta_1 = 2.0$ のガウス分布)

観測データに一定モデルと ϵ モデル
をあてはめて選別差 $\Delta D_{1,2}$ の分布を予測

観測データのモデルから新しい
データをたくさん生成する

あてはまりの良い評価用のデータ (多数)

図 6 尤度最大化に必要な $\Delta D_{1,2}$ の分布の生成。まず観測データである一定モデル ($\beta_1 = 2.0$)、 ϵ (標準) 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選別差 $\Delta D_{1,2}$ がどのような分布になるかを調べる。

2017-06-05 統計モデリング入門 2017a 33/59

6/21 (水)

統計モデリング入門 2017 (e)

GLM logistic regression
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

北大環境科学の講義 <http://goo.gl/776c4i>

2017-06-21

ファイル更新時刻: 2017-05-24 16:32

2017-06-05 統計モデリング入門 2017 (e) 1 / 43

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. ゆーいいでたら万歳・万歳
4. うまくいくまで 1, 2, 3 ぐるぐる

ゆーいい

何でも割算!

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

2017-06-05 統計モデリング入門 2017a 35/59

GLM のひとつ, ロジスティック回帰を使おう

データにあわせて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

NO!

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

2017-06-05 統計モデリング入門 2017a 36/59

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? ちよつとちがう

8 種の種子のうち y 個が 発芽可能 だった! というデータ

ロジスティック回帰とは何なのか?

二項分布: N 回のうち y 回, とする確率

2017-06-05 37/59

6/26 (月)

統計モデリング入門 2017 (f)

階層ベイズモデル Hierarchical Bayesian Model

久保拓弥 kubo@ees.hokudai.ac.jp

北大環境科学の講義 <http://goo.gl/776c4i>

2017-06-26

ファイル更新時刻: 2017-05-24 16:32

2017-06-05 統計モデリング入門 2017 (f) 1 / 66

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.....

観察された植物の個体数

生存した種子数

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か?
(「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

2017-06-05 統計モデリング入門 2017a 40/59

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

二項分布
生存確率 q[i]

事前分布
全体の平均 a

無情報事前分布
(超事前分布)

sigma は hyper parameter

植物の個体差
r[i]

個体差のばらつき
sigma

sigma は a と思ってください

無情報事前分布
(超事前分布)

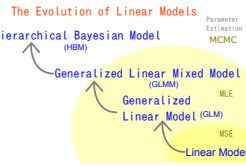
矢印は手順ではなく, 依存関係をあらわしている

2017-06-05 統計モデリング入門 2017a 40/59

なぜ階層ベイズモデルまで勉強するの？

• **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



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

2017-06-05

統計モデリング入門 2017a

41/59

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)

6/ (水)

統計モデリング入門 2017 (g)

階層ベイズモデルと時間変化モデル

久保拓弥 kubo@ees.hokudai.ac.jp

北大環境科学院の講義 <http://goo.gl/76c41>

2017-07-03

ファイル更新時刻: 2017-05-24 16:32

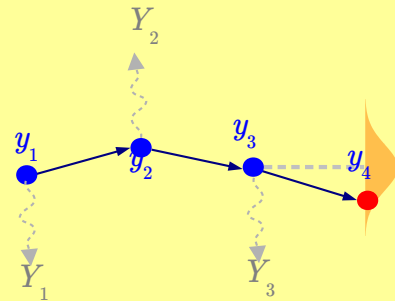
kubostat2017g (http://goo.gl/76c41)

統計モデリング入門 2017 (g)

2017-07-03 1 / 42

時間変化のデータ解析

時系列データ解析



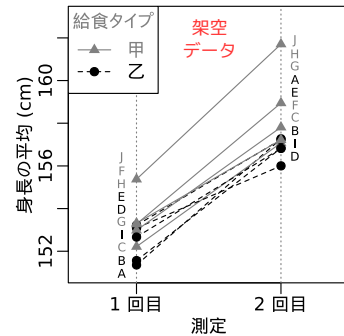
短い時系列データ

時系列の長短に関係なく
「対応のある」データ点か
どうかが本質的な問題

再測定もまた時系列データ



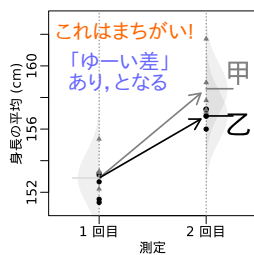
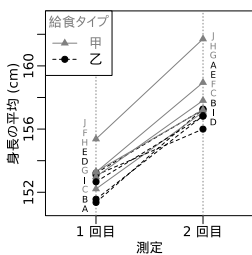
岩波データサイエンス vol.1



2017-06-05

46/59

対応 (paired) を考えてない GLM あてはめ



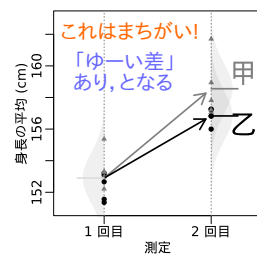
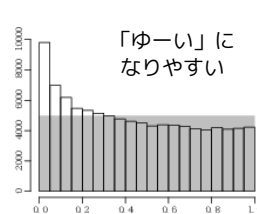
$glm(\text{身長} \sim (\text{測定2回目}) + (\text{測定2回目}):(\text{処理の効果}))$

同じ対象を二回測定していることを考慮してない

2017-06-05

47/59

対応 (paired) を考えてない GLM あてはめ

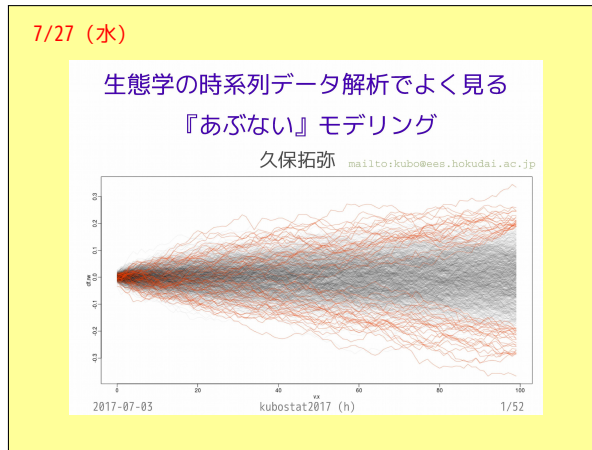
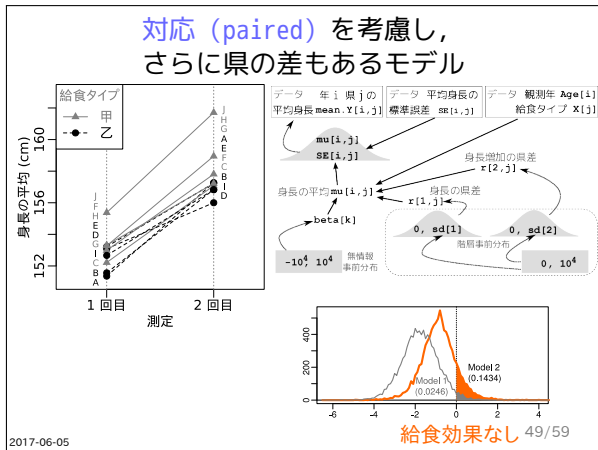


$glm(\text{身長} \sim (\text{測定2回目}) + (\text{測定2回目}):(\text{処理の効果}))$

同じ対象を二回測定していることを考慮してない

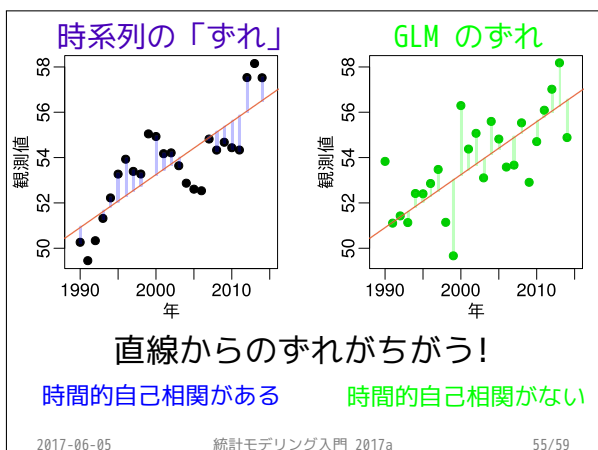
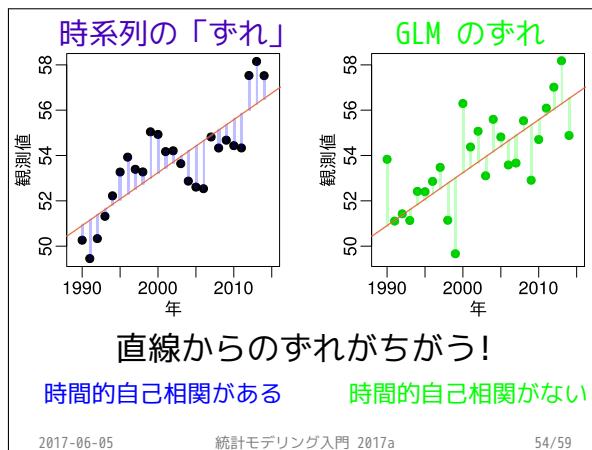
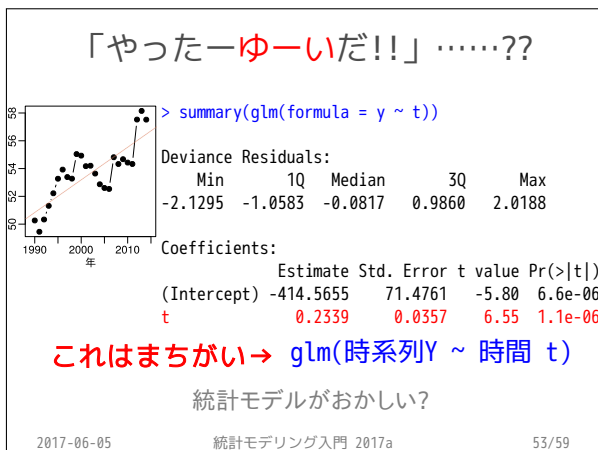
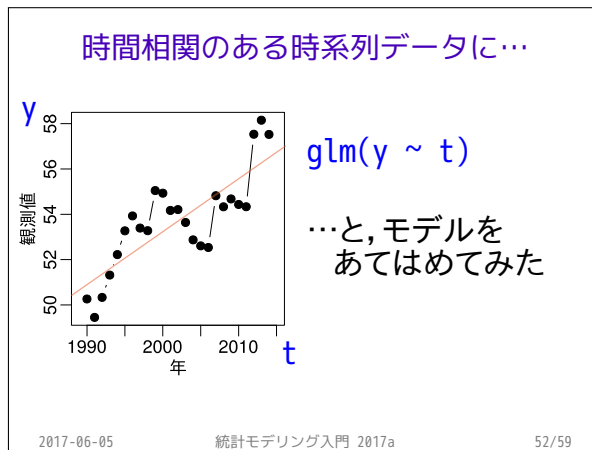
2017-06-05

48/59



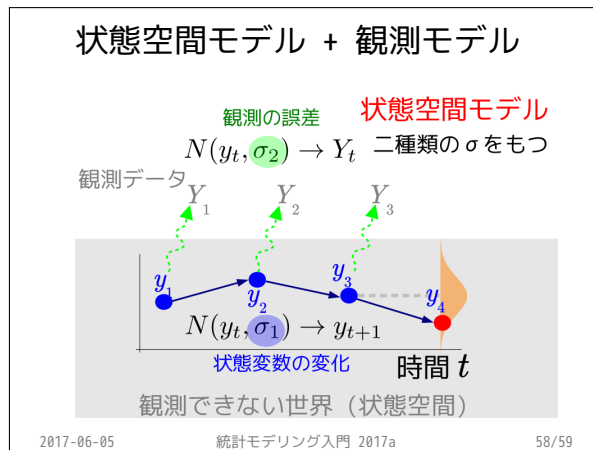
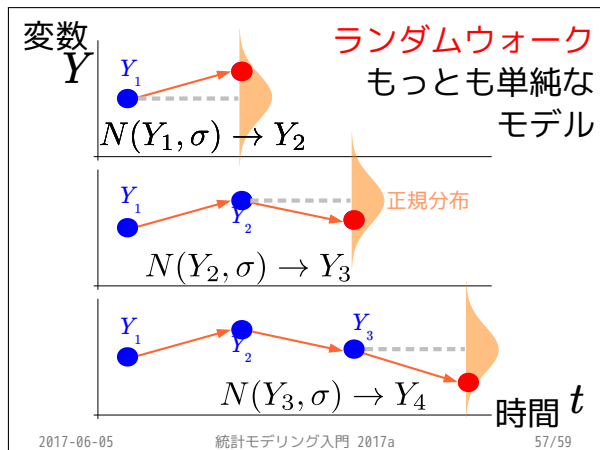
長い時系列データ

データ上で「時間相関」が見える
「時間相関」のモデリングが必要



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



今日はここまで

any questions?