

統計モデリング入門 2017 (a)

An overview for this Statistical Modeling class

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


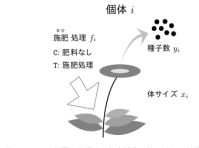





図 3.1 この例題に登場する架空植物の例。個体の個体、この植物の体サイズ(個体の大きさ) x_i と肥料をやる施肥処理 j が種子数 y_i にどう影響しているのかを知りたい。

2017-11-14統計モデリング入門 2017a (豊長研)1/58

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!**

2017-11-14統計モデリング入門 2017a (豊長研)2/58

この統計モデリング授業 改訂

- web 上の「課題」に回答してください
 - 回答もメールで送信してください
 - 課題 <https://goo.gl/z9yCJY>
- 成績評価は「課題」の回答
 - 出欠関係なし

subject: kubostat
to: kubo@ees.hokudai.ac.jp

2017-11-14統計モデリング入門 2017a (豊長研)3/58

Performance Rating 改訂

- E-mail assignment <https://goo.gl/z9yCJY>
 - **That's ALL!**
- Attendance? NOT care.

subject: kubostat
to: kubo@ees.hokudai.ac.jp

2017-11-14統計モデリング入門 2017a (豊長研)4/58

What for Statistical Modeling?

なぜデータ解析の方法を勉強しなければならないのか?


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" !

2017-11-14統計モデリング入門 2017a (豊長研)6/58

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

2017-11-14統計モデリング入門 2017a (豊長研)7/58

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

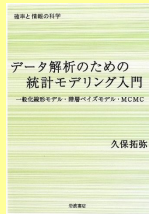
統計学の勉強には良い統計ソフトウェアが必要!

- 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

追記メモ: RStudio の紹介!

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

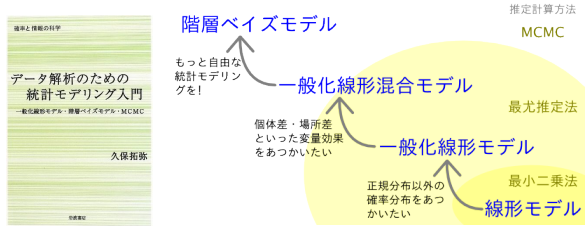
- 観察によってデータ化された現象を説明するために作られる
- 確率分布が基本的な部品であり、これはデータにみられるばらつきを表現する手段である
- データとモデルを対応づける手づきぎが準備されていて、モデルがデータにどれくらい良くあてはまっているかを定量的に評価できる



「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!

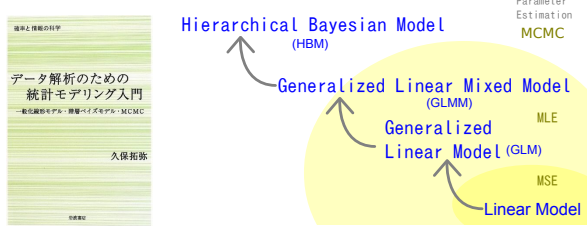
線形モデルの発展



GLM and extended GLMs!

a better statistica model for better data analysis!

The Evolution of Linear Models



たとえばこんなデータがあったらしよう

An example

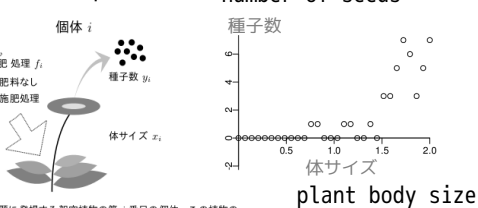
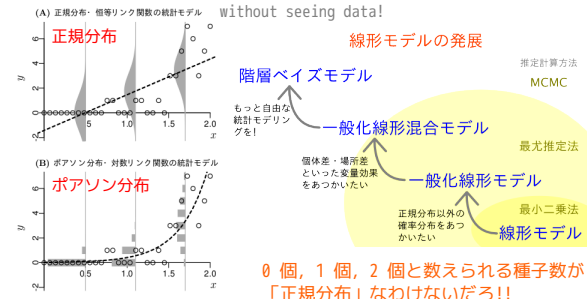


図 3.1 この例題に登場する架空植物の第 i 番目の個体、この植物の体サイズ (個体の大きさ) x_i と肥料をやる施肥処理 f_i が種子数 y_i にどう影響しているのかを知りたい。

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!



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 図解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

全体の流れ Course Overview

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
京大豊長研の講義 <https://goo.gl/z9yCJY>
2017-11-14
ファイル更新時刻: 2017-11-11 16:02

単純化した例題

number of seeds per plant (individual)
こんなデータ (実数) があつたとしましょう

まあ、なんだかどういふような植物を調査していらっしゃいます

個体 i 種子数 y_i
全 50 個体 $i \in \{1, 2, 3, \dots, 50\}$ 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 8 2 4 7 2 2 6 2 4 5 4 6 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
となく図を細く!
```

Simplified examples to learn statistical modeling

2017-11-14 統計モデリング入門 2017a (豊長研) 19/58

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。観測データと平均 λ とその確率 prob の関係が示されている。図 5 の高さを同じにしたもの。右の prob は観測データの $\text{type} = "n"$ によって「高さ情報による表示」。 $\text{ity} = 2$ によって「折れ線は破線」とも表示している。

図 5 観測データと確率分布の対比をなごめる。ヒストグラムは図 4 と同じ。それに重ねられている高さと同じにしたもの。右の prob は観測データの $\text{type} = "n"$ によって「高さ情報による表示」。 $\text{ity} = 2$ によって「折れ線は破線」とも表示している。

2017-11-14 統計モデリング 2017 (b) 1 / 42

さいりゅう 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

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

$\lambda = 3.56$ がより

図 7 平均 λ (lambda) を変え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 $\log L(\lambda)$)。平均 λ が 3.56 のとき、対数尤度は最大になる。

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

2017-11-14 統計モデリング入門 2017a (豊長研) 21/58

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
京大豊長研の講義 <https://goo.gl/z9yCJY>
2017-11-14
ファイル更新時刻: 2017-11-11 16:02

2017-11-14 統計モデリング入門 2017 (c) 2017-11-14 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定植物の個体 i 。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 j が種子数 y_i にとり影響しているのかを知りたい。

the "normal distribution is NOT "normal"

(A) 正規分布・相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。また別の実定植物に対して GLM をあてはめたい。破線は y とともに変化する平均値。グレイで

2017-11-14 統計モデリング入門 2017a (豊長研) 23/58

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定植物の個体 i 。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 j が種子数 y_i にとり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, データモデル
  family = "poisson" link = "log",
  data = d, リンク関数の指定 (省略可)
) data.frame の指定
```

図 3.1 この例題に登場する実定植物の個体 i 。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 j が種子数 y_i にとり影響しているのかを知りたい。

図 17 平均種子数 λ の予測。図 17 に λ の予測値 (実測) を上げたもの。

2017-11-14 統計モデリング入門 2017a (豊長研) 24/58

統計モデリング入門 2017 (d)
 model selection and statistical test
 モデル選択と検定

久保拓弥 kubo@ees.hokudai.ac.jp
 京大豊長研の講義 <https://goo.gl/z9yCJY>
 2017-11-14

ファイル更新時刻: 2017-11-11 16:02

kubostat2017d (<https://goo.gl/z9yCJY>) 統計モデリング入門 2017 (d) 2017-11-14 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?

(A) パラメーター数 $k=1$

(B) パラメーター数 $k=7$

What is the "best?" parameter number k ?

2017-11-14 統計モデリング入門 2017a (豊長研) 26/58

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

統計モデリングの検定 そして、その中身も
 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 ←こっちは!

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

2017-11-14 統計モデリング入門 2017a (豊長研) 27/58

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 ($\beta_1 = 2.06$, $\beta_2 = 0.06$) が元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。

2017-11-14 統計モデリング入門 2017a (豊長研) 28/58

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
 京大豊長研の講義 <https://goo.gl/z9yCJY>
 2017-11-14

ファイル更新時刻: 2017-11-11 16:02

kubostat2017e (<https://goo.gl/z9yCJY>) 統計モデリング入門 2017 (e) 2017-11-14 1 / 47

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
- ゆーいがでたら万歳・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

2017-11-14 統計モデリング入門 2017a (豊長研) 30/58

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

NO! 何でも割算!

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

2017-11-14 統計モデリング入門 2017a (豊長研) 31/58

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) 観測データの一例 ($y=0$) (B) 推定されるモデル

a statistical model for fractions using binomial distributions

二項分布: N 回のうち y 回, とする確率

2017-11-14 統計モデリング入門 2017a (豊長研) 32/58

統計モデリング入門 2017 (f)
Generalized Linear Mixed Model (GLMM)
一般化線形混合モデル

久保拓弥 kubo@ees.hokudai.ac.jp
京大豊長研の講義 <https://goo.gl/z9yCJY>
2017-11-14
ファイル更新時刻: 2017-11-11 16:02

kubostat2017f (<https://goo.gl/z9yCJY>) 統計モデリング入門 2017 (f) 2017-11-14 1 / 35

GLM ではうまく説明できないデータ!
GLMM は階層ベイズモデルの一種 事前分布をどう選ぶかが重要

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?!

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 $q[i]$ ← $r[i]$ ← 植物の個体差
事前分布 個体差のばらつき σ は σ と思ってください
無情報事前分布 全体平均 a ← 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

2017-11-14 統計モデリング入門 2017a (豊長研) 35/58

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

The Evolution of Linear Models
Linear Model → Generalized Linear Model (GLM) → Generalized Linear Mixed Model (GLMM) → Hierarchical Bayesian Model (HBM)

Parameter Estimation MCMC
MLE
MSE

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

2017-11-14 統計モデリング入門 2017a (豊長研) 36/58

統計モデリング入門 2017 (g)
階層ベイズモデル Hierarchical Bayesian Model

久保拓弥 kubo@ees.hokudai.ac.jp
京大豊長研の講義 <https://goo.gl/z9yCJY>
2017-11-14
ファイル更新時刻: 2017-11-11 16:02

kubostat2017g (<https://goo.gl/z9yCJY>) 統計モデリング入門 2017 (g) 2017-11-14 1 / 68

Metropolis Method のルールで q を動かす
もっともっと長くサンプリングしてみる

最尤推定法
メトロポリス法 (MCMC)

MCMC は何をサンプリングしている?
対数尤度 $\log L(q)$ 尤度 $L(q)$ に比例する確率分布
尤度に比例する確率分布からのランダムサンプル
最尤推定はパラメータの値の点推定
MCMC は「パラメータの事後分布」(←推定したいこと) はこういう分布ですと推定している

JAGS を R の「しとうけ」して使う
MCMC sampling from posterior distributions
事後分布からのランダムサンプル

モデルの構造
データとパラメータの初期値
サンプリングの仕組み
Input Output

2017-11-14 統計モデリング入門 2017a (豊長研) 38/58

Applications of Hierarchical Bayesian Model
- nested data and "paired data" -

統計モデリング入門 2017 (h)
階層ベイズモデルの応用と“対応”のあるデータ解析

久保拓弥 kubo@ees.hokudai.ac.jp
2017-11-15

統計モデリング入門 2017

1. 複数ランダム効果の階層ベイズモデル
個体差 + 植木鉢差 など

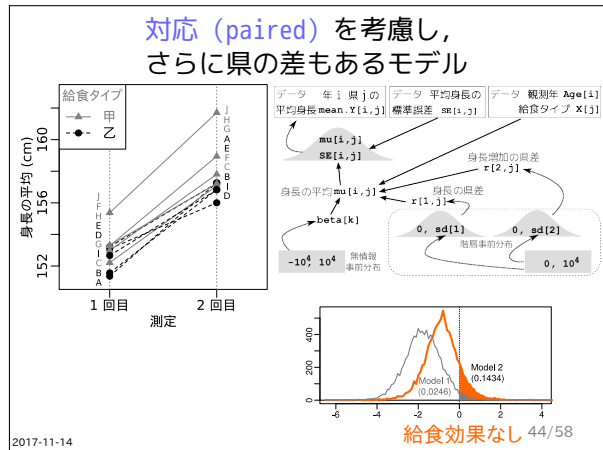
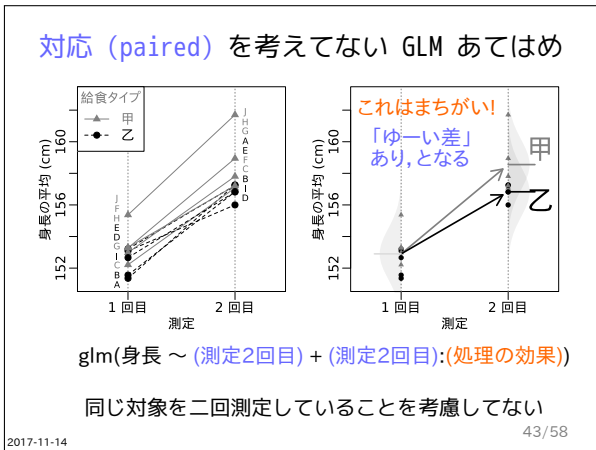
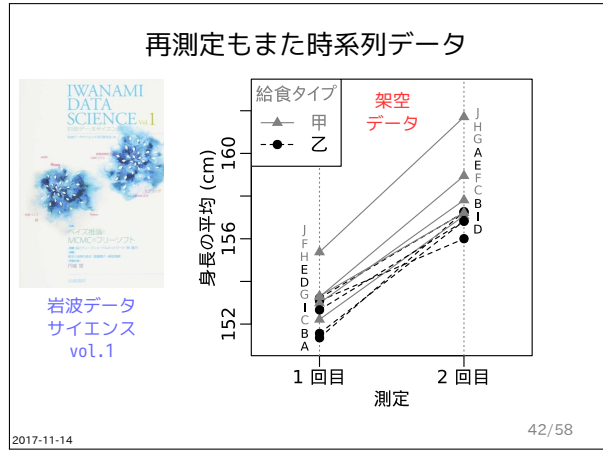
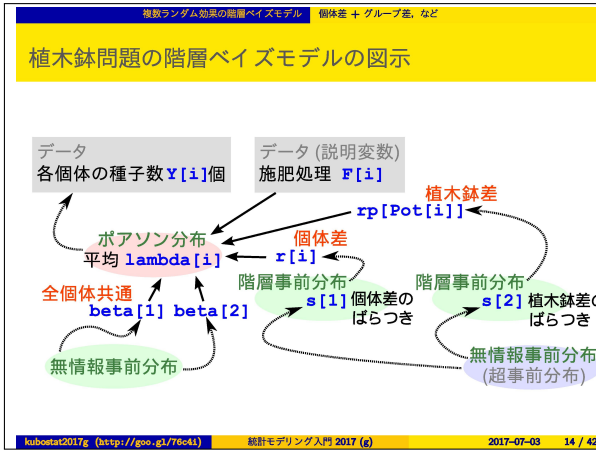
individual effects + pot effects ... etc

response variable 応答変数 種子数 $\{y_i\}$
explanatory variable 説明変数 施肥処理 $\{f_i\}$

seed number 種子数 $\{y_i\}$
fertilisation 施肥処理 $\{f_i\}$

seeds 種子数
C: 肥料なし
T: 施肥処理

個体 i



統計モデリング入門 2017 (i)

あぶない時系列データ解析

久保拓弥 kubo@ees.hokudai.ac.jp
 京大豊原研の講義 <https://goo.gl/z9yCJY>
 2017-11-15

ファイル更新時刻: 2017-11-11 17:56

統計モデリング入門 2017 (i) 2017-11-15 1 / 1

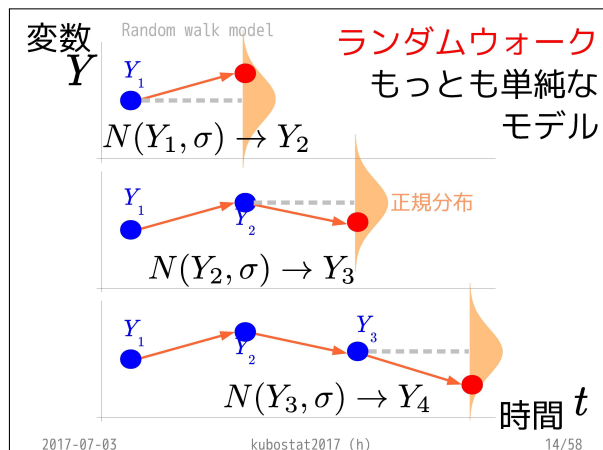
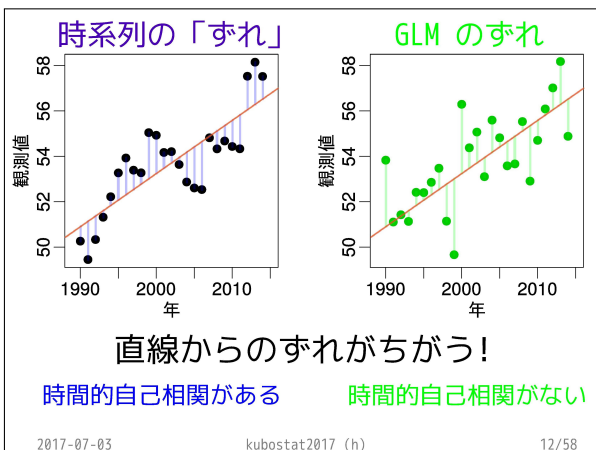
(危1) 時系列データを GLM で

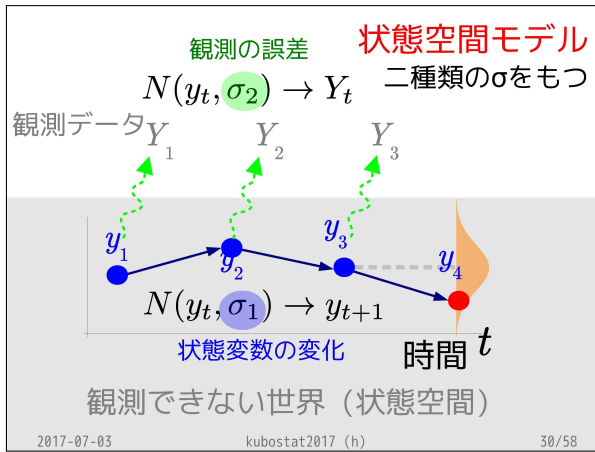
Do NOT apply GLM to time-series data!

「ゆーいな傾き」をねつぞうする原因

- 傾きの検定やめて
- AIC モデル選択
- しても同様になる

検定とかモデル選択とかそういう問題ではない
 統計モデルがおかしい?





統計モデリング入門 2017 (j)
 分割の統計モデル Categorical Data Analysis

久保拓弥 kubo@ees.hokudai.ac.jp
 北大環境科学院の講義 <http://geo.g1/76c41>
 2017-11-15
 ファイル更新時刻: 2017-11-08 17:11

kubostat2017 (<http://geo.g1/76c41>) 統計モデリング入門 2017 (j) 2017-11-15 1 / 63

xtabs: 分割表をあつかう R のクラス

```

y x Spc
1 286 0 A
2 85 0 B
4 378 1 A
5 148 1 B

> (ct2 <- xtabs(y ~ x + Spc, data = d2))
Spc
x    A  B
0 286 85
1 378 148
    
```

2017-11-15 9 / 63

xtabs: 分割表の図示

```

Spc
x    A  B
0 286 85
1 378 148

> plot(ct2, col = c("orange", "blue"))
    
```

2017-11-15 11 / 63

xtabs: 分割表の図示

```

Spc
x    A  B  C  D  E  F  G  H  I
0 62 21 14 11 10 10 2 0 2
1 48 34 22 17 16 7 2 1 1

> plot(ct9, col = c("ごちゃごちゃと指定"))
    
```

2017-11-15 37 / 63

重量分割モデル (階層ベイズモデル): そのプロセス

2017-11-15 47 / 63

例題 1
 「アロメトリーな回帰」はヤメて
 重量分割モデルを作ってみよう

2017-11-15 44 / 63

このモデルで複雑な重量分配を表現できる

- オレンジ色の線は中央値 (median)
- 黄色の領域は個体差による予測のばらつき (95% CI)

2017-11-15 60 / 63

A long trailer, done

any questions?

11/15 (明日) 午後は
「統計モデリングの個別相談」
を実施します (先着順)

<https://goo.gl/z9yCJY> の下のほう