

統計モデリング入門 2017 (b)

確率分布と最尤推定

久保拓弥 kubo@ees.hokudai.ac.jp

北大環境科学院の講義 <http://goo.gl/76c4i>

2019-07-22

ファイル更新時刻: 2019-07-20 20:55

今日のハナシ I

- ① 例題: 種子数の統計モデリング
An example: a distribution of seed number
- ② データと確率分布の対応
probability distribution, the core of statistical model
- ③ ポアソン分布のパラメーターの さいゆうすいてい 最尤推定
もっとももっともらしい推定?
- ④ 統計モデルの要点
乱数発生・推定・予測

本題にはいる前に 統計モデリング授業前半の 「テーマ」を 再確認しておきましょう

この授業であつかう統計モデルたち

0 個, 1 個, 2 個と数えられるデータ

カウントデータ ($y \in \{0, 1, 2, 3, \dots\}$ なデータ)

(たとえば卵数) (たとえば体重) 説明変数

- たとえば x は植物個体の大きさ, y はその個体の花数
- 体サイズが大きくなると花数が増えるように見えるが.....
- この現象を表現する統計モデルは?

正規分布を使った統計モデル ムリがある?

正規分布・恒等リンク関数の統計モデル

説明変数

- タテ軸のばらつきは「正規分布」なのか?
- y の値は 0 以上なのに
- 平均値がマイナス?

GLM という統計モデル

ポアソン分布を使った統計モデルなら良さそう?!

ポアソン分布・対数リンク関数の統計モデル

応答変数 y

説明変数 x

- タテ軸に対応する「ばらつき」 fair distribution
- 負の値にならない「平均値」 non-negative mean
- 正規分布を使ってるモデルよりましだね bye-bye, the normal distribution

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 7 / 42

GLM という統計モデル

データの性質をよくみる

Plot your data and observe it

確率分布という部品を選ぶ

Choose proper distributions

「正規分布」は万能ではない!

the normal distribution is NOT good at everything

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 8 / 42

GLM という統計モデル

今日の内容と「統計モデリング入門」との対応

<http://goo.gl/Ufq2>

今日はおもに「第2章 確率分布と統計モデルの最尤推定」の内容を説明します。

- 著者: 久保拓弥
- 出版社: 岩波書店
- 2012-05-18 刊行

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 9 / 42

例題: 種子数の統計モデリング An example: a distribution of seed number

2. 例題: 種子数の統計モデリング

An example: a distribution of seed number

R でデータをあつかいつつ

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 10 / 42

例題: 種子数の統計モデリング An example: a distribution of seed number

この授業では架空植物の架空データをあつかう

number of seeds taken from 50 imaginary plants

理由: よけいなことは考えなくてすむので

現実のデータはどれも授業で使うには難しすぎる……

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 11 / 42

例題: 種子数の統計モデリング An example: a distribution of seed number

こんなデータ (架空) があったとしましょう

まあ、なんだかこういうヘンな植物を調査しているとします

seed number of i 種子数 y_i

この $\{y_i\}$ が観測データ!
 $\{y_i\} = \{y_1, y_2, \dots, y_{50}\}$

全 50 個体 $i \in \{1, 2, 3, \dots, 50\}$

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 5 4 5 1 3 2 3
```

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 12 / 42

例題: 種子数の統計モデリング An example: a distribution of seed number

R: a free statistical software
これ使いましょう: 統計ソフトウェア R
<http://www.r-project.org/>

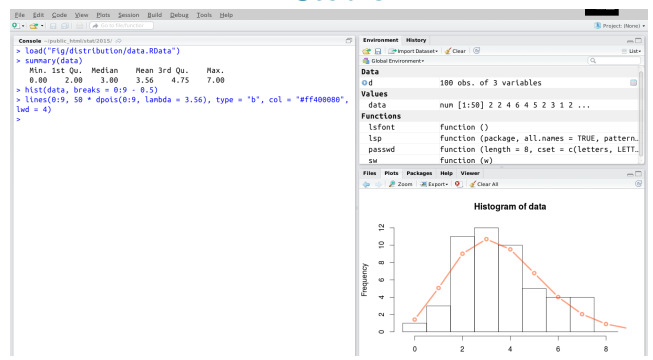


- いろいろな OS で使える **free software**
- 使いたい機能が充実している
- **作図**機能も強力
- S 言語による **プログラミング**可能
- **RStudio** <http://www.rstudio.com/>

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 13 / 42

例題: 種子数の統計モデリング An example: a distribution of seed number

RStudio



<http://www.rstudio.com/>

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 14 / 42

例題: 種子数の統計モデリング An example: a distribution of seed number

R でデータの様子をながめる



の table() 関数を使って種子数の頻度を調べる

```
> table(data)
 0  1  2  3  4  5  6  7
 1  3 11 12 10  5  4  4
```

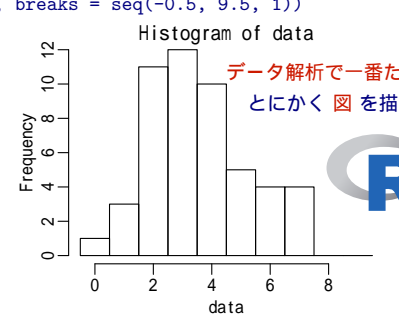
(種子数 5 は 5 個体, 種子数 6 は 4 個体)

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 15 / 42

例題: 種子数の統計モデリング An example: a distribution of seed number

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))
```



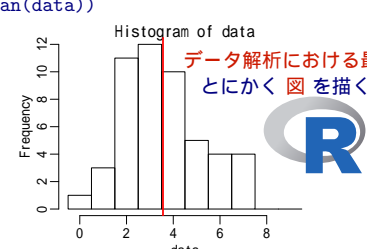
データ解析で一番たいせつなことに
とにかく **図** を描く!

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 16 / 42

例題: 種子数の統計モデリング An example: a distribution of seed number

How to evaluate mean value using R?

```
> mean(data)
[1] 3.56
> abline(v = mean(data))
```



データ解析における最重要事項
とにかく **図** を描く!

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 17 / 42

例題: 種子数の統計モデリング An example: a distribution of seed number

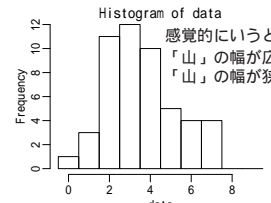
「ばらつき」の統計量

あるデータの **ばらつき** をあらわす標本統計量の例: **標本分散**

```
> var(data)
[1] 2.9861
```

標本標準偏差 とは標本分散の平方根 ($SD = \sqrt{\text{variance}}$)

```
> sd(data)
[1] 1.7280
> sqrt(var(data))
[1] 1.7280
```



感覚的にいうと
「山」の幅が広い: 分散が大きい
「山」の幅が狭い: 分散が小さい

kubostat2017b (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (b) 2019-07-22 18 / 42

データと確率分布の対応 probability distribution, the core of statistical model

3. データと確率分布の対応

probability distribution, the core of statistical model

確率分布は統計モデルの重要な部品

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 19 / 42

データと確率分布の対応 probability distribution, the core of statistical model

Empirical VS Theoretical Distributions

統計モデルの部品である **確率分布** には
 “データそのまま” な **経験分布** と
 数式で定義される **理論的な分布** がある

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 20 / 42

データと確率分布の対応 probability distribution, the core of statistical model

“データそのまま” な経験分布

```
> data.table <- table(factor(data, levels = 0:10))
> cbind(y = data.table, prob = data.table / 50)
```

y	prob
0	1 0.02
1	3 0.06
2	11 0.22
3	12 0.24
4	10 0.20
5	5 0.10
6	4 0.08
7	4 0.08
8	0 0.00
9	0 0.00
10	0 0.00

- 確率分布とは **発生する事象** と **発生する確率** の対応づけ
- 確率分布のひとつである **経験分布** とは
 “たまたま手もとにある” データから
 “発生確率” を決める確率分布

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 21 / 42

データと確率分布の対応 probability distribution, the core of statistical model

なるほど**経験分布**は“直感的”かもしれないが.....

- データが変わると確率分布が変わる?
- 種子数 $y = \{0, 1, 2, \dots\}$ となる確率が、個々におたがい無関係に決まる?
- パラメーターは $\{p_0, p_1, p_2, \dots, p_{99}, p_{100}, \dots\}$ 無限個ある?
 道具として使うには、ちょっと不便かもしれない.....

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 22 / 42

データと確率分布の対応 probability distribution, the core of statistical model

なにか理論的に導出された確率分布のほうが便利ではないか?

- 少数のパラメーターで分布の“カタチ”が決まる
- “なめらかに” 確率が変化する
- いろいろと数理的な道具が準備されている (パラメーター推定方法など)

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 23 / 42

データと確率分布の対応 probability distribution, the core of statistical model

確率分布 (ポアソン分布) を数式で決めてしまう

種子数が y である確率は以下のように決まる、と考えている

$$p(y | \lambda) = \frac{\lambda^y \exp(-\lambda)}{y!}$$

- $y!$ は y の階乗で、たとえば $4!$ は $1 \times 2 \times 3 \times 4$ をあらわしています.
- $\exp(-\lambda) = e^{-\lambda}$ のこと ($e = 2.718\dots$)
- ここではなぜポアソン分布の確率計算が上ようになるのかは説明しません— まあ、こういうもんだと考えて先に進みましょう

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 24 / 42

データと確率分布の対応 probability distribution, the core of statistical model

数式で決められたポアソン分布?

とりあえず R で作図してみる

```
> y <- 0:9 # これは種子数 (確率変数)
> prob <- dpois(y, lambda = 3.56) # ポアソン分布の確率の計算
> plot(y, prob, type = "b", lty = 2) # cbind で「表」作り
> cbind(y, prob)
```

y	prob
1	0.02843882
2	0.10124222
3	0.18021114
4	0.21385056
5	0.19032700
6	0.13551282
7	0.08040427
8	0.04089132
9	0.01819664
10	0.00719778

平均 (λ) が 3.56 である
Poisson distribution

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 25 / 42

データと確率分布の対応 probability distribution, the core of statistical model

データとポアソン分布を重ね合わせる

```
> hist(data, seq(-0.5, 8.5, 0.5)) # まずヒストグラムを描き
> lines(y, prob, type = "b", lty = 2) # その「上」に折れ線を描く
```

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 26 / 42

データと確率分布の対応 probability distribution, the core of statistical model

パラメーター λ はポアソン分布の平均

```
> # cbind で「表」作り
> cbind(y, prob)
```

y	prob
1	0.02843882
2	0.10124222
3	0.18021114
4	0.21385056
5	0.19032700
6	0.13551282
7	0.08040427
8	0.04089132
9	0.01819664
10	0.00719778

- 平均 λ はポアソン分布の唯一のパラメーター
- 確率分布の平均は λ である ($\lambda \geq 0$)
- 分散と平均は等しい: $\lambda = \text{平均} = \text{分散}$
- $y \in \{0, 1, 2, \dots, \infty\}$ の値をとり、すべての y について和をとると 1 になる

$$\sum_{y=0}^{\infty} p(y | \lambda) = 1$$

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 27 / 42

データと確率分布の対応 probability distribution, the core of statistical model

どういった場合にポアソン分布を使う?

統計モデルの部品としてポアソン分布が選んだ理由:

- データに含まれている値 y_i が $\{0, 1, 2, \dots\}$ といった非負の整数である (カウントデータである)
- y_i に下限 (ゼロ) はあるみただけで上限はよくわからない
- この観測データでは平均と分散がだいたい等しい
 - このだいたい等しいがあやしいのだけど、まあ気にしないことにしよう

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 28 / 42

データと確率分布の対応 probability distribution, the core of statistical model

ポアソン分布の λ を変えてみる

$$p(y | \lambda) = \frac{\lambda^y \exp(-\lambda)}{y!}$$

λ は平均をあらわすパラメーター

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 29 / 42

ポアソン分布のパラメーターの 最尤推定 もっとももらしい推定?

さいゆうすいいてい

4. ポアソン分布のパラメーターの最尤推定

もっとももらしい推定?

「あてはめる」ことは推定すること

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 30 / 42

ポアソン分布のパラメータの 最尤推定 もっとももっともらしい推定?

ゆうど
尤度 (likelihood) とは何か?

- 最尤推定法では、**尤度** という**あてはまりの良さ**をあらわす統計量に着目
- 尤度は**データが得られる確率**をかけあわせたもの
- この例題の場合、パラメータ λ を変えると尤度が変わる
- もっとも「あてはまり」が良くなる λ を見つけたい
- たとえば、いまデータが 3 個体ぶん、たとえば、 $\{y_1, y_2, y_3\} = \{2, 2, 4\}$ 、これだけだった場合、尤度はだいたい $0.180 \times 0.180 \times 0.19 = 0.006156$ といった値になる

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 31 / 42

ポアソン分布のパラメータの 最尤推定 もっとももっともらしい推定?

尤度 $L(\lambda)$ はパラメータ λ の関数

この例題の尤度

$$L(\lambda) = (y_1 \text{ が } 2 \text{ である確率}) \times (y_2 \text{ が } 2 \text{ である確率}) \times \dots \times (y_{50} \text{ が } 3 \text{ である確率})$$

$$= p(y_1 | \lambda) \times p(y_2 | \lambda) \times p(y_3 | \lambda) \times \dots \times p(y_{50} | \lambda)$$

$$= \prod_i p(y_i | \lambda) = \prod_i \frac{\lambda^{y_i} \exp(-\lambda)}{y_i!}$$

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 32 / 42

ポアソン分布のパラメータの 最尤推定 もっとももっともらしい推定?

尤度はしんどのいで対数尤度を使う

尤度は確率 (あるいは確率密度) の積であり、あつかいがふべん (大量のかけ算!)

そこで、パラメータの最尤推定では、**対数尤度関数** (log likelihood function) を使う

$$\log L(\lambda) = \sum_i \left(y_i \log \lambda - \lambda - \sum_k \log k \right)$$

対数尤度 $\log L(\lambda)$ の最大化は尤度 $L(\lambda)$ の最大化になるから
 まずは、平均をあらわすパラメータ λ を変化させていったときに、ポアソン分布のカタチと対数尤度がどのように変化するかを調べてみましょう

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 33 / 42

ポアソン分布のパラメータの 最尤推定 もっとももっともらしい推定?

λ を変えるとあてはまりの良さが変わる

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 34 / 42

ポアソン分布のパラメータの 最尤推定 もっとももっともらしい推定?

対数尤度を最大化する $\hat{\lambda}$ をさがす

$$\text{対数尤度 } \log L(\lambda) = \sum_i (y_i \log \lambda - \lambda - \sum_k y_k \log k)$$

- 最尤推定量 (ML estimator): $\sum_i y_i / 50$ 標本平均値!
- 最尤推定値 (ML estimate): $\hat{\lambda} = 3.56$ ぐらい

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 35 / 42

ポアソン分布のパラメータの 最尤推定 もっとももっともらしい推定?

最尤推定を使っても真の** λ は見つからない**

真の λ が 3.5 と設定して架空データを生成

データは有限なので**真の** λ はわからない
 標本サイズが 50 の場合、「平均値の推定」すらなかなかうまくできない

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 36 / 42

統計モデルの要点 乱数発生・推定・予測

5. 統計モデルの要点

乱数発生・推定・予測

統計モデルとデータの対応づけ

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 37 / 42

統計モデルの要点 乱数発生・推定・予測

確率分布: 乱数発生 と 推定

(人間には見えない) 真の統計モデル $\lambda = 3.5$ のポアソン分布

観測データから推定された $\lambda = 3.56$ のポアソン分布

パラメータ推定

観測されたデータ

データをサンプル

確率分布から乱数を発生

データ?...ここでは確率・統計モデルが生成していると仮定

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 38 / 42

統計モデルの要点 乱数発生・推定・予測

推定されたモデルを使った 予測

(人間には見えない) 真の統計モデル $\lambda = 3.5$ のポアソン分布

観測データから推定された $\lambda = 3.56$ のポアソン分布

予測: 新しいデータにあてはまるのか? (予測) の良さを調べている

新しいデータをサンプル

同じ調査方法で得られた新データ

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 39 / 42

統計モデルの要点 乱数発生・推定・予測

この講義で登場する確率分布

- ポアソン分布: $y \in \{0, 1, 2, 3, \dots\}$ となるデータ
- 二項分布: $y \in \{0, 1, 2, \dots, N\}$ となるデータ
- 正規分布: $-\infty < y < \infty$ の連続値をとるデータ
- その他あれこれ — 一様分布とかも登場

そんなに多くの確率分布は登場しません

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 40 / 42

統計モデルの要点 乱数発生・推定・予測

いろいろな確率分布があるけれど.....

- この講義では多種多様な確率分布を[あつかいません](#)
- しかし [確率分布を混ぜあわせる](#) ことによって、自分で確率分布を作り出すことができます
- ハナシの後半に登場する GLMM や階層ベイズモデル

線形モデルの発展

階層ベイズモデル (HBM)

もっと自由な統計モデリングを!

一般化線形混合モデル (GLMM)

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 (GLM)

正規分布以外の確率分布をあつかいたい

線形モデル (LM)

最尤推定法

最小二乗法

推定計算方法 MCMC

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 41 / 42

統計モデルの要点 乱数発生・推定・予測

次回予告

The next topic

YES!

一般化線形モデルのひとつ: ポアソン回帰

Poisson Regression, a Generalized Linear Model (GLM)

kubostat2017b (http://goo.gl/76c4i) 統計モデリング入門 2017 (b) 2019-07-22 42 / 42