

## 統計モデリング入門 2017 (c)

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>

2019-07-22

ファイル更新時刻: 2019-07-20 08:16

もくじ

### 今日のハナシ I

- ① ポアソン回帰の統計モデル  
応答変数  $y$  と説明変数  $x$
- ② ポアソン回帰の例題: 架空植物の種子数データ  
植物個体の属性, あるいは実験処理が種子数に影響?
- ③ GLM の詳細を指定する  
確率分布・線形予測子・リンク関数
- ④ R で GLM のパラメーターを推定  
あてはまりの良さは 対数尤度関数で評価
- ⑤ 処理をした・しなかった 効果も統計モデルに入れる  
GLM の因子型説明変数

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 1 / 45

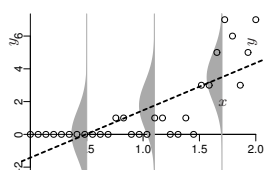
kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 2 / 45

もくじ

### 今日のハナシ II

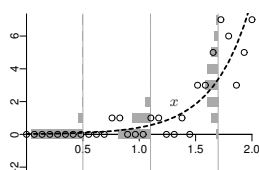
Normal distribution  
and identity link function

正規分布・恒等リンク関数の  
統計モデル



Poisson distribution  
and log link function

ポアソン分布・log リンク関数の  
統計モデル



kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 3 / 45

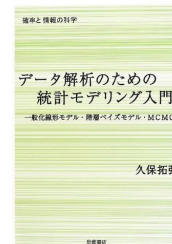
もくじ

### 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3章 一般化線形モデル (GLM)」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



kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 4 / 45

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 4 / 45

もくじ

### 一般化線形モデルって何だろう?

#### 一般化線形モデル (GLM)

- **ポアソン回帰** (Poisson regression)
- ロジスティック回帰 (logistic regression)
- 直線回帰 (linear regression)
- .....

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 5 / 45

ポアソン回帰の統計モデル

応答変数  $y$  と説明変数  $x$

#### 1. ポアソン回帰の統計モデル

応答変数  $y$  と説明変数  $x$

一般化線形モデルにとりくんでみる

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 5 / 45

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 6 / 45

ポアソン回帰の統計モデル 応答変数  $y$  と説明変数  $x$

### この授業であつかう統計モデルたち

The development of linear models

Hierarchical Bayesian Model (parameter estimation MCMC)

Generalized Linear Mixed Model (GLMM) (Be more flexible)

Generalized Linear Model (GLM) (MLE)

Linear model (MSE)

Incorporating random effects such as individuality

Always normal distribution? That's non-sense!

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 7 / 45

ポアソン回帰の統計モデル 応答変数  $y$  と説明変数  $x$

### 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)

応答変数 (たとえば卵数)

説明変数 (たとえば体重)

- 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?

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 8 / 45

ポアソン回帰の統計モデル 応答変数  $y$  と説明変数  $x$

### 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

応答変数

説明変数

**NO!**

とにかくセンヒキゃいいんでしょ  
傾き「ゆーい」ならいいんでしょ  
...という安易な発想のデータ解析

- 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?

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 9 / 45

ポアソン回帰の統計モデル 応答変数  $y$  と説明変数  $x$

### 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

説明変数

**YES!**

- 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

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 10 / 45

ポアソン回帰の例題: 架空植物の種子数データ 植物個体の属性, あるいは実験処理が種子数に影響?

## 2. ポアソン回帰の例題: 架空植物の種子数データ

植物個体の属性, あるいは実験処理が種子数に影響?

Modeling number of seeds of plants using GLM

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 11 / 45

ポアソン回帰の例題: 架空植物の種子数データ 植物個体の属性, あるいは実験処理が種子数に影響?

### 個体サイズと実験処理の効果を調べる例題

- 応答変数: 種子数  $\{y_i\}$
- 説明変数:
  - 体サイズ  $\{x_i\}$
  - 施肥処理  $\{f_i\}$

標本数

- 無処理 ( $f_i = C$ ): 50 sample ( $i \in \{1, 2, \dots, 50\}$ )
- 施肥処理 ( $f_i = T$ ): 50 sample ( $i \in \{51, 52, \dots, 100\}$ )

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 12 / 45



ポアソン回帰の例題: 架空植物の種子数データ 植物個体の属性, あるいは実験処理が種子数に影響?

### 施肥処理 f した・しないの箱ひげ図 (box-whisker plot)

```
> plot(d$f, d$y) # note that d$f is factor type!
```

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 19 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### 3. GLM の詳細を指定する

確率分布・線形予測子・リンク関数

ポアソン回帰では log link 関数を使うのが便利

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 20 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### 一般化線形モデルを作る

一般化線形モデル (GLM)

- 確率分布は?
- 線形予測子は?
- リンク関数は?

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 21 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### GLM のひとつである直線回帰モデルを指定する

直線回帰のモデル

- 確率分布: 正規分布
- 線形予測子: e.g.,  $\beta_1 + \beta_2 x_i$   
直線の式: (切片) + (傾き)  $\times x_i$
- リンク関数: 恒等リンク関数

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 22 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### 線形モデルの予測子, predictor of linear model

- 応答変数 (response variable)
- 説明変数 (explanatory variable)
- 係数 (coefficient)
- 線形予測子 (linear predictor):  

$$(\text{応答変数}) = \text{定数 (切片, intercept)} \\ + (\text{係数 1}) \times (\text{説明変数 1}) \\ + (\text{係数 2}) \times (\text{説明変数 2}) \\ + (\text{係数 3}) \times (\text{説明変数 3}) \\ + \dots$$

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 23 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### GLM のひとつであるポアソン回帰モデルを指定する

ポアソン回帰のモデル

- 確率分布: ポアソン分布
- 線形予測子: e.g.,  $\beta_1 + \beta_2 x_i$
- リンク関数: 対数リンク関数

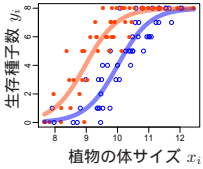
kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 24 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### GLM のひとつである logistic 回帰モデルを指定する

ロジスティック回帰のモデル

- 確率分布: 二項分布
- 線形予測子: e.g.,  $\beta_1 + \beta_2 x_i$
- リンク関数: logit リンク関数



kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 25 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### R で一般化線形モデル (GLM) の推定を.....

	確率分布	乱数発生	GLM あてはめ
(離散)	ベルヌーイ分布	rbinom()	glm(family = binomial)
	二項分布	rbinom()	glm(family = binomial)
	ポアソン分布	rpois()	glm(family = poisson)
	負の二項分布	rnbinom()	glm.nb() in library(MASS)
(連続)	ガンマ分布	rgamma()	glm(family = gamma)
	正規分布	rnorm()	glm(family = gaussian)

- glm() で使える確率分布は上記以外にもある
- GLM は直線回帰・重回帰・分散分析・ポアソン回帰・ロジスティック回帰その他の「よせあつめ」と考えてもよいかも

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 26 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### さて、種子数の例題にもどって...

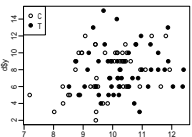
種子数  $y_i$  は平均  $\lambda_i$  のポアソン分布にしたがうとしましょう

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

個体  $i$  の平均  $\lambda_i$  を以下のようにおいてみたらどうだろう.....?

$$\lambda_i = \exp(\beta_1 + \beta_2 x_i)$$

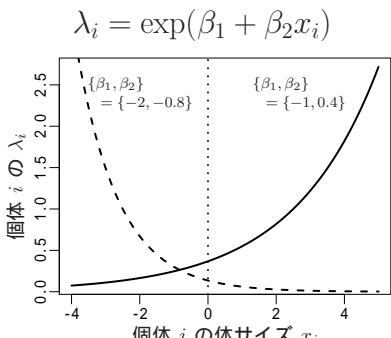
- $\beta_1$  と  $\beta_2$  は係数 (パラメーター)
- $x_i$  は個体  $i$  の体サイズ,  $f_i$  はとりあえず無視



kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 27 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### 指数関数ってなんだっけ?

$$\lambda_i = \exp(\beta_1 + \beta_2 x_i)$$


kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 28 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### GLM のリンク関数と線形予測子 ← (直線の式)

個体  $i$  の平均  $\lambda_i$

$$\lambda_i = \exp(\beta_1 + \beta_2 x_i)$$

$$\Downarrow$$

$$\log(\lambda_i) = \beta_1 + \beta_2 x_i$$

**log(平均) = 線形予測子**

log リンク関数とよばれる理由は、上のようになっているから

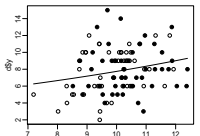
kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 29 / 45

GLM の詳細を指定する 確率分布・線形予測子・リンク関数

### この例題のための統計モデル

#### ポアソン回帰のモデル

- 確率分布: ポアソン分布
- 線形予測子:  $\beta_1 + \beta_2 x_i$
- リンク関数: 対数リンク関数



kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 30 / 45

で GLM のパラメーターを推定 あてはまりの良さは 対数尤度関数で評価

### 4. Rで GLM のパラメーターを推定

あてはまりの良さは 対数尤度関数で評価

推定計算はコンピューターにおまかせ

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 31 / 45

で GLM のパラメーターを推定 あてはまりの良さは 対数尤度関数で評価

### glm() 関数の指定

```
> d
  y   x f
1  6 8.31 C
2  6 9.44 C
3  6 9.50 C
... (中略) ...
99 7 10.86 T
100 9 9.97 T
```

これだけ!

```
> fit <- glm(y ~ x, data = d, family = poisson)
```

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 32 / 45

で GLM のパラメーターを推定 あてはまりの良さは 対数尤度関数で評価

### glm() 関数の指定の意味

```
fit <- glm(
  y ~ x,
  family = poisson(link = "log"),
  data = d
)
```

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

- モデル式 (線形予測子  $z$ ): どの説明変数を使うか?
- link 関数:  $z$  と応答変数 ( $y$ ) 平均値 の関係は?
- family: どの確率分布を使うか?

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 33 / 45

で GLM のパラメーターを推定 あてはまりの良さは 対数尤度関数で評価

### glm() 関数の出力

```
> fit <- glm(y ~ x, data = d, family = poisson)

all: glm(formula = y ~ x, family = poisson, data = d)

Coefficients:
(Intercept)          x
      1.2917      0.0757

Degrees of Freedom: 99 Total (i.e. Null); 98 Residual
Null Deviance: 89.5
Residual Deviance: 85      AIC: 475
```

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 34 / 45

で GLM のパラメーターを推定 あてはまりの良さは 対数尤度関数で評価

### glm() 関数のくわしい出力

```
> summary(fit)
Call:
glm(formula = y ~ x, family = poisson, data = d)

Deviance Residuals:
   Min       1Q   Median       3Q      Max
-2.368  -0.735  -0.177   0.699   2.376

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)  1.2917     0.3637    3.55  0.00038
x              0.0757     0.0356    2.13  0.03358

..... (以下, 省略) .....
```

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 35 / 45

で GLM のパラメーターを推定 あてはまりの良さは 対数尤度関数で評価

### 推定値と標準誤差のイメージ (かなりいいかげんな説明)

- 確率  $p$  は ゼロからの距離 をあらわしている
- $p$  がゼロに近いほど 推定値  $\hat{\beta}$  はゼロから離れている
- $p$  が 0.5 に近いほど 推定値  $\hat{\beta}$  はゼロに近い

(注: 頻度主義的な信頼区間の正しい解釈はもっとめんどくさい)

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 36 / 45

GLM のパラメーターを指定 あるいはまりの良さは 対数尤度関数で評価

### モデルの予測

```
> fit <- glm(y ~ x, data = d, family = poisson)
...
Coefficients:
(Intercept)          x
      1.2917         0.0757

> plot(d$x, d$y, pch = c(21, 19)[d$f]) # data
> xp <- seq(min(d$x), max(d$x), length = 100)
> lines(xp, exp(1.2917 + 0.0757 * xp))
```

ここでは観測データと予測の関係  
を見ているだけ、なのだが

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 37 / 45

処理をした・しなかった 効果も統計モデルに入れる GLM の因子型説明変数

### 5. 処理をした・しなかった 効果も統計モデルに入れる

GLM の因子型説明変数

数量型 + 因子型 という組み合わせで

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 38 / 45

処理をした・しなかった 効果も統計モデルに入れる GLM の因子型説明変数

### 肥料の効果 $f_i$ もいれましょう

種子数  $y_i$  は平均  $\lambda_i$  のポアソン分布にしたがうと  
しましょう

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

個体  $i$  の平均  $\lambda_i$  を次のようにする

$$\lambda_i = \exp(\beta_1 + \beta_2 x_i + \beta_3 d_i)$$

- $\beta_3$  は施肥処理の効果の係数
- $f_i$  のダミー変数

$$d_i = \begin{cases} 0 & (f_i = C \text{ の場合}) \\ 1 & (f_i = T \text{ の場合}) \end{cases}$$

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 39 / 45

処理をした・しなかった 効果も統計モデルに入れる GLM の因子型説明変数

### glm(y ~ x + f, ...) の出力

```
> summary(glm(y ~ x + f, data = d, family = poisson))
... (略) ...

Coefficients:
      Estimate Std. Error z value Pr(>|z|)
(Intercept)  1.2631     0.3696   3.42  0.00063
x              0.0801     0.0370   2.16  0.03062
fT            -0.0320     0.0744  -0.43  0.66703

..... (以下, 省略) .....
```

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 40 / 45

処理をした・しなかった 効果も統計モデルに入れる GLM の因子型説明変数

### x + f モデルの予測

```
> plot(d$x, d$y, pch = c(21, 19)[d$f]) # data
> xp <- seq(min(d$x), max(d$x), length = 100)
> lines(xp, exp(1.2631 + 0.0801 * xp), col = "blue", lwd = 3) # C
> lines(xp, exp(1.2631 + 0.0801 * xp - 0.032), col = "red", lwd = 3) # T
```

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 41 / 45

処理をした・しなかった 効果も統計モデルに入れる GLM の因子型説明変数

### 複数の説明変数をいれた場合の統計モデル

- $f_i = C: \lambda_i = \exp(1.26 + 0.0801x_i)$
- $f_i = T: \lambda_i = \exp(1.26 + 0.0801x_i - 0.032)$   
 $= \exp(1.26 + 0.0801x_i) \times \exp(-0.032)$

施肥効果である  $\exp(-0.032)$  は  
かけ算できくことに注意!

kubostat2017c (http://goo.gl/76c4i) 統計モデリング入門 2017 (c) 2019-07-22 42 / 45

処理をした・しなかった 効果も統計モデルに入れる GLM の因子型説明変数

### リンク関数が違うとモデルの解釈が異なる

(A) 対数リンク関数  
 $\lambda = \exp(\beta_1 + \beta_2 x + \dots)$

相乗的

平均種子数  $\lambda_i$

体サイズ  $x_i$

(B) 恒等リンク関数  
 $\lambda = \beta_1 + \beta_2 x + \dots$

相加的

平均種子数  $\lambda_i$

体サイズ  $x_i$

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 43 / 45

処理をした・しなかった 効果も統計モデルに入れる GLM の因子型説明変数

### GLM: 適切な確率分布 とリンク関数を選ぶ

正規分布・恒等リンク関数の統計モデル

ポアソン分布・log リンク関数の統計モデル

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 44 / 45

処理をした・しなかった 効果も統計モデルに入れる GLM の因子型説明変数

### この授業であつかう統計モデルたち

The development of linear models

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 45 / 45

処理をした・しなかった 効果も統計モデルに入れる GLM の因子型説明変数

### 次回予告

The next topic

(A)  $k = 1$

Too simple?

(B)  $k = 7$

Too complex?

### モデル選択と統計学的検定

Model selection and statistical test

kubostat2017c (<http://goo.gl/76c4i>) 統計モデリング入門 2017 (c) 2019-07-22 46 / 45