

# 線形回帰モデルの推定

樋口さぶろお <https://hig3.net>

龍谷大学 先端理工学部 数理・情報科学課程

確率統計 I L12(2021-06-30 Wed)

最終更新: Time-stamp: "2021-06-30 Wed 18:23 JST hig"

## 今日の目標

- 回帰分析を確率変数の言葉で説明できる
- 隠されたパラメタを最尤推定できる
- 回帰分析で回帰係数を推定できる
- 多次元の確率変数の推定ができる

岩薩林 確率・統計 §9



## L11-Q1

Quiz 解答:カイ二乗分布の確率と  $\chi_k^2(\alpha)$ 

- ①  $z(0.025) = 1.960$ .
- ② 標準正規分布の確率密度関数は偶関数なので,  
 $z(1 - 0.025) = -z(0.025) = -1.960$ .
- ③  $\chi_1^2(0.05) = 3.841$ . 別解.  $0.05 = P(W > w_0) = P(Z^2 > w_0) =$   
 $P(Z < -\sqrt{w_0} \text{ or } Z > +\sqrt{w_0}) = 2 \times P(Z > \sqrt{w_0})$ . よって,  
 $\sqrt{w_0} = 1.960$ .
- ④  $\chi_1^2(1 - 0.05) = 0.00393$ .

## L11-Q2

## Quiz 解答:母分散の区間推定

標本サイズは  $n = 9$ , 自由度は  $9 - 1$ , 母分散  $\sigma^2$  の信頼係数 0.95 の信頼区間は,

$$\frac{n-1}{\chi_{n-1}^2(\alpha/2)} \times S^2 < \sigma^2 < \frac{n-1}{\chi_{n-1}^2(1-\alpha/2)} \times S^2$$

$$\frac{8}{17.53} \times 72 < \sigma^2 < \frac{8}{2.180} \times 72$$

$$32.85 < \sigma^2 < 264.2$$

## L11-Q3

Quiz 解答:t 分布の確率と  $t_k(\alpha)$ 

- ①  $z(0.025) = 1.960$ .
- ② 標準正規分布の確率密度関数は偶関数なので,  
 $z(1 - 0.025) = -z(0.025) = -1.960$ .

$$\textcircled{3} \quad t_{40}(0.025) = 2.021.$$

$$\textcircled{4} \quad t \text{ 分布の確率密度関数は偶関数なので,} \\ t_{40}(1 - 0.025) = -t_{40}(0.025) = -2.021.$$

## L11-Q4

## Quiz 解答:母平均値の区間推定 (母分散未知)

- ① 重さの標本平均値は  $\bar{X} = 50\text{g}$ . 不偏標本分散は  $S^2 = \frac{1}{4-1} \cdot 14\text{g}^2$ . t 分布表から 自由度  $k = n - 1 = 3$  の  $t_3(0.05/2)$  を参照して, 信頼係数 0.95 の信頼区間は

$$50 - 3.182 \times \sqrt{\frac{14}{3}/4} < \mu < 50 + 3.182 \times \sqrt{\frac{14}{3}/4}.$$

- ② 同様に, t 分布表から 自由度  $k = n - 1 = 3$  の  $t_3(0.01/2)$  を参照して,

$$50 - 5.841 \times \sqrt{\frac{14}{3}/4} < \mu < 50 + 5.841 \times \sqrt{\frac{14}{3}/4}.$$

## ここまで来たよ

11 母分散, 母平均値の区間推定

12 線形回帰モデルの推定

- 最尤推定
- 線形モデルとしての回帰分析
- 多次元の確率変数の母期待値の推定
- 標本共分散, 標本相関係数

## 線形モデル (統計モデルのある一族)

あるドーナツ製造機の作るドーナツの重さ  $Y$  は次のモデルに従う.

$$Y = \beta_0 + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

$Y, \epsilon$ : 連続型確率定数,  $\sigma > 0, \beta_0$ : 定数=パラメタ=母数

$\epsilon$ : 誤差, ノイズ. 小文字だけど確率変数.

(隠された) パラメタ母数  $\mu, \sigma^2$  を,  $n$  個のドーナツの重さのデータから推定したい.

$$E[Y] = \beta_0 + E[\epsilon] = \beta_0,$$

$$V[Y] = V[\epsilon] = \sigma^2.$$

正規分布と限定した以外は, ここしばらくやってた, 母平均値, 母分散の推定の言い換えに過ぎない.

だけど, 多数のパラメタを含む一般的なモデルにも使える考え方をする.

## 尤度 likelihood

$\epsilon = Y - \beta_0 \sim N(0, \sigma^2)$  より, ドーナツの重さ  $y$  を得る確率密度は,

$$f(y|\beta_0, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\beta_0)^2}{2\sigma^2}}$$

サイズ  $n$  の標本が  $y_1, \dots, y_n$  である確率密度は, 独立同分布なので積で,

$$\begin{aligned} f(y_1, y_2, \dots, y_n | \beta_0, \sigma) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \beta_0)^2}{2\sigma^2}} \\ &= (2\pi\sigma^2)^{-n/2} e^{-\sum_{i=1}^n \frac{(y_i - \beta_0)^2}{2\sigma^2}} \end{aligned}$$

### $n$ 次元正規分布

多変量解析及び演習

この  $f$  を,  $y_i$  の確率密度関数と思わず, が測定済データ  $y_1, \dots, y_n$  が定数,  $\beta_0, \sigma$  が変数と思ったとき, **尤度** (ゆうど) 関数  $L(\beta_0, \sigma)$  という。

$$L(\beta_0, \sigma) = f(y_1, y_2, \dots, y_n | \beta_0, \sigma)$$

# 最尤推定

岩薩林 確率・統計なし

## 最尤推定

$\beta_0, \sigma$  の推定値として  $L(\beta_0, \sigma)$  が最大になる値を選ぶ

2変数関数の最大値  $\rightsquigarrow$  偏微分

微積分 II

$$0 = \frac{\partial L}{\partial \beta_0}(\beta_0, \sigma) = \frac{\partial L}{\partial \sigma}(\beta_0, \sigma)$$

ここでは最初の等式だけ解く。合成微分。

$$0 = (\text{定数})^{-n/2} \times \sum_{i=1}^n \frac{1}{\sigma^2} (y_i - \beta_0) \times e^{\text{同じ}}$$

$$0 = \sum_{i=1}^n (y_i - \beta_0)$$

$$\beta_0 = \frac{1}{n} \sum_{i=1}^n y_i$$

ここでは最初の式だけ。  
合成微分。

$\beta_0$  (母平均値) の最尤推定  
(MLE=maximum likelihood  
estimation) 値  $\hat{\beta}_0$  は標本平均値  
じゃん



## ここまで来たよ

11 母分散, 母平均値の区間推定

12 線形回帰モデルの推定

- 最尤推定
- 線形モデルとしての回帰分析
- 多次元の確率変数の母期待値の推定
- 標本共分散, 標本相関係数

## (確率変数でない) 変数 $x$ に依存する確率変数 $Y$

このドーナツ製造機で作るドーナツの重さ  $Y$  は、温度  $x$  によるらしい。  
次の線形回帰モデルを仮定する。

$$Y = \beta_0 + \beta_1 \cdot x + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

$Y, \epsilon$ : 連続型確率定数,  $\beta_0, \beta_1$ : 回帰係数,  $\sigma > 0$ : 定数 = パラメタ = 母数

$Y$ : 目的変数 (従属変数) ここでは確率変数

$x$ : 説明変数 (独立変数) ここでは確率変数でない

ノイズ・誤差  $\epsilon = Y - \beta_0 + \beta_1 \cdot x \sim N(0, \sigma^2)$ .

$\epsilon = Y - \beta_0 - \beta_1 x \sim N(0, \sigma^2)$  より、ドーナツの重さ  $y$  を得る確率密度は、

$$f(y|x, \beta_0, \beta_1, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\beta_0-\beta_1 \cdot x)^2}{2\sigma^2}}$$

(正確に) 指定した温度  $x_i$  ( $i = 1, \dots, n$ ) で製造したときの重さが  $y_i$  ( $i = 1, \dots, n$ ) である確率密度は, 独立分布なので積.

$$\begin{aligned} f(y_1, y_2, \dots, y_n | x_1, \dots, x_n, \beta_0, \beta_1, \sigma) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \beta_0 - \beta_1 \cdot x_i)^2}{2\sigma^2}} \\ &= (2\pi\sigma^2)^{-n/2} e^{-\sum_{i=1}^n \frac{(y_i - \beta_0 - \beta_1 \cdot x_i)^2}{2\sigma^2}} \\ &= L(\beta_0, \beta_1, \sigma) \end{aligned}$$

## 最尤推定

測定済のデータ  $x_i, y_i$   $i = 1, \dots, n$  を定数と思ったときの, 3 変数関数

$L(\beta_0, \beta_1, \sigma) = f(y_1, y_2, \dots, y_n | x_1, \dots, x_n, \beta_0, \beta_1, \sigma)$  の最大値は?  $\rightsquigarrow$  偏微分 微積分 II

$$0 = \frac{\partial L}{\partial \beta_0}(\beta_0, \beta_1, \sigma) = \frac{\partial L}{\partial \beta_1}(\beta_0, \beta_1, \sigma) = \frac{\partial L}{\partial \sigma}(\beta_0, \beta_1, \sigma)$$

ここでは最初の 2 つの等式だけ解く.

$$0 = \frac{\partial L}{\partial \beta_0} = (\text{定数}) \sum_i \frac{1}{\sigma^2} (y_i - \beta_0 - \beta_1 x_i) \times e^{\text{同じ}}$$

$$0 = \frac{\partial L}{\partial \beta_1} = (\text{定数}) \sum_i \frac{1}{\sigma^2} x_i (y_i - \beta_0 - \beta_1 x_i) \times e^{\text{同じ}}$$

しよせん,  $\beta_0, \beta_1$  の連立 1 次方程式

正規方程式 岩薩林 確率・統計 §9.2

$$n\beta_0 + \left(\sum_i x_i\right)\beta_1 = \sum_i y_i$$

$$\left(\sum_i x_i\right)\beta_0 + \left(\sum_i x_i^2\right)\beta_1 = \sum_i x_i y_i$$

加減法  $\rightsquigarrow$

$$\hat{\beta}_1 = \frac{s_{xy}}{s_{xx}},$$

$$\hat{\beta}_0 = \bar{y} - \frac{s_{xy}}{s_{xx}}\bar{x}$$

$$y = \frac{s_{xy}}{s_{xx}}x + \bar{y} - \frac{s_{xy}}{s_{xx}}\bar{x}$$

$$y - \bar{y} = \frac{s_{xy}}{s_{xx}}(x - \bar{x})$$

ここで,

$$\bar{x} = \frac{1}{n} \sum_i x_i \quad \text{平均値っぽい形}$$

$$\bar{y} = \frac{1}{n} \sum_i y_i \quad \text{平均値っぽい形}$$

$$s_{xy} = \frac{1}{n} \sum_i x_i y_i - \bar{x} \cdot \bar{y} \quad \text{岩薩林 確率・統計 定理 1.5} = \frac{1}{n} \sum_i (x_i - \bar{x})(y_i - \bar{y}) \quad \text{共分散っぽい形}$$

$$s_{xx} = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \quad \text{岩薩林 確率・統計 定理 1.2} = \frac{1}{n} \sum_i (x_i - \bar{x})^2 \quad \text{分散っぽい形}$$

$$Y = \hat{\beta}_0 + \hat{\beta}_1 \cdot x + \epsilon$$

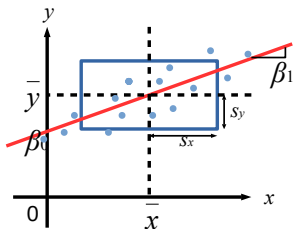
## 回帰直線

岩薩林 確率・統計 §9.1

推定結果  $\beta_0, \beta_1$  を係数とする  $xy$  平面の直線

$$y - \bar{y} = \frac{s_{xy}}{s_{xx}} (x - \bar{x})$$

$\beta_0, \beta_1$  回帰係数



## 回帰係数, 予測値の信頼区間

$Y_i$  は確率変数.

こうやって得られる推定結果を  $\hat{\beta}_0, \hat{\beta}_1$  と書くと, これらも確率変数.  
 $\hat{\beta}_0, \hat{\beta}_1$  にも信頼区間などが考えられる. さらに, (自分が新たに設定した温度  $x$  に対する) 予測値  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$  にも信頼区間などが考えられる.

## 決定係数

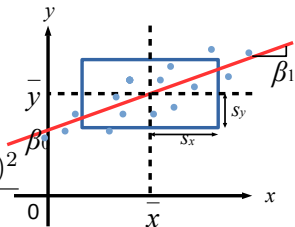
岩薩林 確率・統計 §9.2

残差 直線 (予測値) からの上下方向のずれ,

$$e_i = y_i - \hat{y}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i).$$

決定係数

$$R^2 = 1 - \frac{\sum_i e_i^2}{\sum_i (y_i - \bar{y})^2} = 1 - \frac{\sum_i (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2}{\sum_i (y_i - \bar{y})^2}$$



$0 \leq R^2 \leq 1$  で, 1 に近いほど, あてはまりがよい.

実は  $R^2$  は, 2次元データ  $(x_i, y_i)$  の相関係数

岩薩林 確率・統計 p.22

データ分析



岩薩林 確率・統計 例題 9.2, 9.3, §9 問題 3,4,5, §9 練習問題 1

## L12-Q1

### Quiz(回帰係数と回帰直線)

$x$  を説明変数,  $y$  を目的変数とする線形モデル  $Y = \beta_0 + \beta_1 \cdot x + \epsilon$ ,  $\epsilon \sim N(0, \sigma^2)$  において,  $x$  を  $x = x_1, \dots, x_n$  と変えて  $Y$  を測定したところ,  $y_1, \dots, y_n$  となった. それをまとめたのが下である.

$(x, y)$  のデータの個数  $n$       16

$$\bar{x} = \frac{1}{n} \sum_i x_i \quad 9$$

$$\bar{y} = \frac{1}{n} \sum_i y_i \quad -4$$

$$\frac{1}{n} \sum_i x_i^2 - \bar{x}^2 \quad 49$$

$$\frac{1}{n} \sum_i y_i^2 - \bar{y}^2 \quad 36$$

$$\frac{1}{n} \sum_i x_i y_i - \bar{x} \cdot \bar{y} \quad -25$$

このとき, 回帰直線の式を,  $x, y$  で書こう. 整理しなくてよい.

## ここまで来たよ

11 母分散, 母平均値の区間推定

12 線形回帰モデルの推定

- 最尤推定
- 線形モデルとしての回帰分析
- 多次元の確率変数の母期待値の推定
- 標本共分散, 標本相関係数

# (復習) 多次元の確率変数

確率統計 I(2021)L04

## 確率分布 (母分布, 母集団)

### 離散型

確率関  
数  
 $p(x, y)$

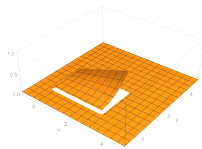
$y \backslash x$	7	8	9	計
0	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{6}$	
1	0	0	$\frac{1}{3}$	
計				

### 標本 (サイズ $n$ )

#### 離散型

データ番号 $i$	$x$	$y$
1	8	0
2	9	1
$\vdots$	$\vdots$	$\vdots$
$n$	8	0

### 連続型



確率密  
度関数

$f(x, y) =$

### 連続型

データ番号 $i$	$x$	$y$
1	1.2	0.95
2	2.5	1.24
$\vdots$	$\vdots$	$\vdots$
$n$	0.9	0.04

## 多次元の母期待値の推定は1次元と同じのりで

同時分布が与えられたときの母期待値

岩薩林 確率・統計 (3.16)p.60

確率統計 I(2021)L04

$$\text{連続型 } E[g(X, Y)] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} p(x, y) \cdot g(x, y) dx dy$$

1次元 確率統計 I(2021)L10  $\rightsquigarrow$  多次元

標本期待値

$$g(X, Y) \text{ の標本期待値 } \overline{g(X, Y)} = \frac{1}{n} [g(X_1, Y_1) + \cdots + g(X_n, Y_n)]$$

が,  $E[g(X, Y)]$  の 'よい' 推定量になっている.

$g(X, Y)$  が  $X$  のみによる場合  $E[g(X)]$ , 例えば  $E[X]$  だったら, 周辺分布  $f_X(x) = \int f(x, y) dy$  で考えられれば, 確率統計 I(2021)L11 そのままじゃん.

## L12-Q2

## Quiz(多次元の確率変数の母期待値母共分散の推定)

以下は、2次元の確率変数  $X, Y$  のサイズ  $n = 5$  の標本である

$X$	$Y$
1	5
3	15
4	14
5	11
7	20

- ① 母平均値  $E[X]$  を推定しよう.
- ② 母期待値  $E[X^2Y]$  を推定しよう.
- ③ 母分散  $V[X]$  を推定しよう.
- ④ 母共分散  $\text{Cov}[X, Y]$  を推定しよう.
- ⑤ 母相関係数  $\rho[X, Y]$  を推定しよう.

## ここまで来たよ

11 母分散, 母平均値の区間推定

12 線形回帰モデルの推定

- 最尤推定
- 線形モデルとしての回帰分析
- 多次元の確率変数の母期待値の推定
- 標本共分散, 標本相関係数

## 母分散の推定値

確率統計 I(2021)L04

不偏標本分散 確率統計 I(2021)L10

$$\begin{aligned} \text{不偏標本分散 } S_x^2 = S_{xx} &= \frac{1}{n-1} [(X_1 - \bar{X})^2 + \cdots + (X_n - \bar{X})^2] \\ &= \frac{n}{n-1} \left[ \frac{1}{n} \sum_i X_i^2 - (\bar{X})^2 \right] \end{aligned}$$

が, 母分散  $V[X]$  の 'よい' 推定値になっている.

$V[Y]$  の推定値は,  $S_y^2 = S_{yy}$  と書く.

## 母共分散の推定

母共分散 covariance 岩薩林 確率・統計 (3.19)p.61,(3.20)p.62

$X, Y$  が確率変数で,  $\mu_X = E[X], \mu_Y = E[Y]$  とおいたとき,

$$\text{母共分散 } \text{Cov}[X, Y] \stackrel{\text{定義}}{=} E[(X - \mu_X)(Y - \mu_Y)] \quad (3.19)$$

$$= \text{岩薩林 確率・統計 (3.20)} \cdots = E[XY] - E[X] \times E[Y].$$

(C1,(3.20))

不偏標本共分散 岩薩林 確率・統計 なし

不偏標本共分散

$$S_{xy} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \frac{n}{n-1} \left[ \frac{1}{n} \sum_{i=1}^n X_i Y_i - \bar{x} \cdot \bar{y} \right]$$

が母共分散  $\text{Cov}[X, Y]$  の 'よい' 推定値になっている.



## 母相関係数の推定

母相関係数 correlation

岩薩林 確率・統計 (3.22)p.64

$$\text{母相関係数 } \rho[X, Y] \stackrel{\text{定義}}{=} \frac{\text{Cov}[X, Y]}{\sqrt{V[X]}\sqrt{V[Y]}} \quad (3.22)$$

標本相関係数

岩薩林 確率・統計 §1.3(1.19)p.22

標本相関係数

$$r_{xy} = \frac{S_{xy}}{\sqrt{S_x^2}\sqrt{S_y^2}}$$

が母相関係数  $\rho[X, Y]$  に対応する.

( $n$  や  $n - 1$  は分子分母でキャンセルしているから, 分子分母一斉に変えてもよい)

## フィッシャーの $z$ -変換

### フィッシャーの $z$ -変換

$(X, Y)$  を 2次元正規分布にしたがう確率変数とするとき,  
サイズ  $n$  の標本の標本相関係数を確率変数  $R$  と考える.  
確率変数

$$Z = \frac{1}{2} \log \frac{1+R}{1-R}$$

は, 近似的に  $N(\text{Cov}[X, Y], \frac{1}{n-3})$  にしたがう.

逆変換は  $R = \tanh Z$ .

これで, 相関係数の区間推定や, 検定 (典型的には, 帰無仮説  $r = 0$ ) ができると.

## L12-Q2

## Quiz(多次元の確率変数の母期待値母共分散の推定)

以下は, 2次元の確率変数  $X, Y$  のサイズ  $n = 5$  の標本である

$X$	$Y$
1	5
3	15
4	14
5	11
7	20

- ① 母平均値  $E[X]$  を推定しよう.
- ② 母期待値  $E[X^2Y]$  を推定しよう.
- ③ 母分散  $V[X]$  を推定しよう.
- ④ 母共分散  $\text{Cov}[X, Y]$  を推定しよう.
- ⑤ 母相関係数  $\rho[X, Y]$  を推定しよう.