

# 統計的方法のメモ

前川眞一

130607,140411,0417,24,0530

## 1 はじめに

## 2 標本分布

### 2.1 はじめに

母集団から大きさ  $n$  の標本を抽出する。

観測されるデータは、母集団からサンプルされたものであると考える。

$n$  が大きくなればなるほど、そのデータは母集団の状態をより正確に反映していると考えられる。

### 2.2 標本統計量

標本平均

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad (1)$$

標本分散

$$S_X^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \quad (2)$$

( $n-1$  で割る流儀もある。) 分散の平方根を標準偏差と呼ぶ。

標本共分散

$$c_{XY} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \quad (3)$$

標本相関係数

$$r_{XY} = \frac{c_{XY}}{S_X S_Y} \quad (4)$$

$$= \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}} \quad (5)$$

### 2.3 標本平均の分布

母集団において

$$X_i \sim N(\mu, \sigma^2), \quad i = 1, 2, \dots, n, \quad iid \quad (6)$$

である場合、独立な  $n$  個の観測値の標本平均

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad (7)$$

の分布は

$$\bar{X} \sim N(\mu, \sigma^2/n) \quad (8)$$

となる。

これは、母集団から、大きさ  $n$  の標本を採り、その平均を計算したものを  $\bar{X}_k$  とする、という作業を  $m$  回繰り返した時の  $\bar{X}_k, k = 1, 2, \dots, m$  の分布である。

これは、正規分布の性質から

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \quad (9)$$

とも書ける。

母分散が未知の場合は、上記の  $Z$  の中の  $\sigma$  を標本標準偏差

$$s^2 = \sum_{i=1}^n \frac{1}{n-1} (X_i - \bar{X})^2 \quad (10)$$

で置き換えたもの  $T$  が自由度  $n-1$  の  $t$  分布をすること

$$T = \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t(n-1) \quad (11)$$

が知られている。

したがって、母数の値  $(\mu, \sigma)$  がわかっていれば、母平均がある区間に入る確率を計算できる。

正規分布の場合、 $\mu, \sigma$  がわかっていれば、標準正規分布表から

$$\beta = \Pr(Z \leq z_\beta) \quad (12)$$

となる点  $z_\beta$  が計算できる。 $t$  分布の場合も同様に

$$\beta = \Pr(T \leq t_{\beta, n-1}) \quad (13)$$

となる点  $t_{\beta, n-1}$  が計算できる。

(両側検定の場合は  $\beta = \alpha/2$  とする。)

## 2.4 標本比率の分布

標本比率の分布は

$$p \sim N\left(\pi, \frac{\pi(1-\pi)}{n}\right) \quad (14)$$

で近似することが出来る。したがって、

$$Z = \frac{p - \pi}{\sqrt{\frac{\pi(1-\pi)}{n}}} \sim N(0, 1) \quad (15)$$

を利用する。

### 3 区間推定 (信頼区間)

母集団から  $n$  個のサンプルを取り、標本平均  $\bar{X}$  と標本標準偏差  $s$  を観測した時の母平均  $\mu$  の  $100 \times (1 - \alpha)\%$  信頼区間は以下の通りである。

$$m_L = \bar{X} - y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \quad (16)$$

$$m_U = \bar{X} + y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \quad (17)$$

ただし、 $y_{(1-\alpha/2)}$  はそれぞれの分布において  $100 \times (1 - \alpha/2)\%$  パーセンタイル、すなわち、 $100 \times (\alpha/2)\%$  の上側確率を与える点である。すなわち、正規分布の場合

$$y_{(1-\alpha/2)} = z_{(1-\alpha/2)} \quad (18)$$

t 分布の場合

$$y_{(1-\alpha/2)} = t_{(1-\alpha/2), n-1} \quad (19)$$

である。

母集団標準偏差  $\sigma$  がわかっている時、もしくは  $n$  が大きい時は正規分布を用いても良い。

#### 3.1 導出

$$1 - \alpha = \Pr \left( \mu - y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \leq \bar{X} \leq \mu + y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \right) \quad (20)$$

$$= \Pr \left( -y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \leq \bar{X} - \mu \leq y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \right) \quad (21)$$

$$= \Pr \left( -\bar{X} - y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \leq -\mu \leq -\bar{X} + y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \right) \quad (22)$$

$$= \Pr \left( \bar{X} + y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \geq \mu \geq \bar{X} - y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \right) \quad (23)$$

$$= \Pr \left( \bar{X} - y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} + y_{(1-\alpha/2)} \frac{s}{\sqrt{n}} \right) \quad (24)$$

#### 3.2 標本比率

$$m_L = p - z_{(1-\alpha/2)} \sqrt{\frac{p(1-p)}{n}} \quad (25)$$

$$m_U = p + z_{(1-\alpha/2)} \sqrt{\frac{p(1-p)}{n}} \quad (26)$$

## 4 統計的仮説検定の一般論

### 4.1 手続き

- 母集団の統計モデルを作る。  
正規分布、2項分布
- 検定仮説を立てる。  
 $H_0$  帰無仮説（何らかの制約が付いているもの）  
 $H_1$  制約を取ったものでふつうはこれを採択したい場合が多い。
- $H_0$  の下での検定統計量の分布を求める。  
標本平均、標本比率、平均平方和 その他  
正規分布、t 分布、F 分布 etc
- 棄却域を決める。  
 $H_0$  の下での検定統計量の分布で起こりそうもないこと（確率が小さいこと）を定義する。
- データを取る。  $X_i, i = 1, 2, \dots, n$
- 観測された検定統計量が棄却域に入るかどうかを調べる。  
棄却域に入ったということは  $H_0$  の下で起こりそうもないことが起こったということなので  $H_0$  を棄却する。  
もしくは  $p$ -値を計算する。 $p$ -値とは、 $H_0$  が真の時に、検定統計量が、実際にデータから計算された検定統計量の値よりも極端な値をとる確率。

### 4.2 検定統計量の選び方

Table 1: 仮説検定に伴う誤り

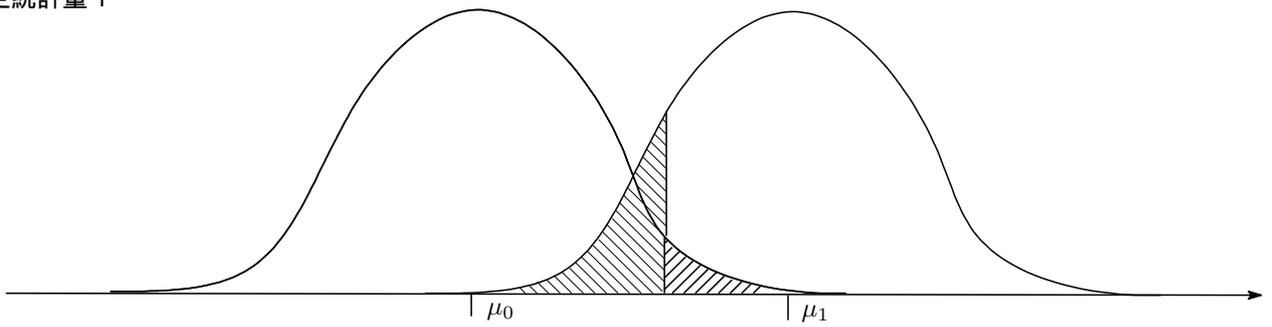
		真の状態	
		$H_0$	$H_1$
採 択 さ れ た 仮 説	$H_0$	正解	第2種の誤り Type II Error $\beta$
	$H_1$	第1種の誤り Type I Error $\alpha$	正解

検定力とは  $H_1$  が正しい時にそれを正しく採択する確率のこと、すなわち  $1 - \beta$ 。

以下の図では、右上から左下への斜線部分が第1種の誤り Type I Error の確率  $\alpha$  を表し、左上から右下への斜線部分が第2種の誤り Type II Error の確率  $\beta$  を表している。したがって、検定力は右側の分布において、白い部分からこの斜線部分を引いたものとなる。

この図では Type I Error の大きさは同じだが、下の統計量の方が検定力が高いことがわかる。

検定統計量 1



検定統計量 2

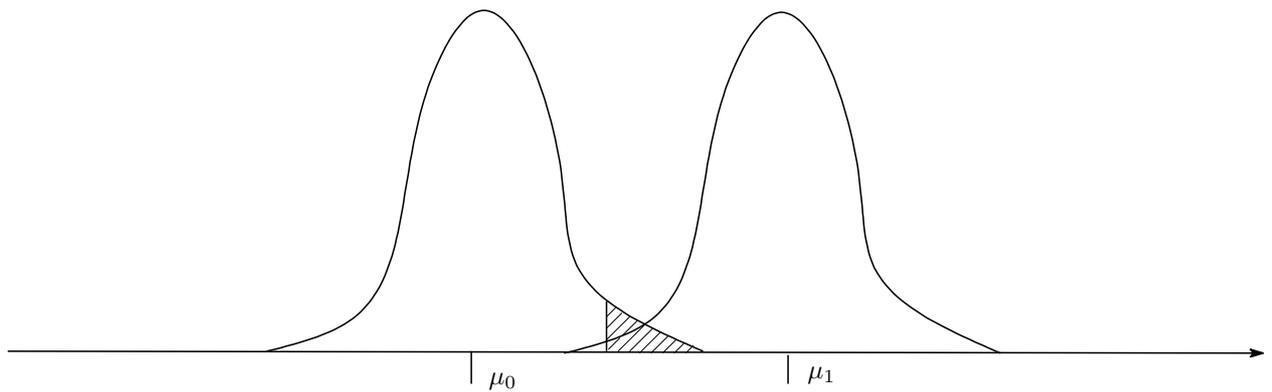


Figure 1: Power Comparison

## 5 仮説検定

### 5.1 一つの母集団の母数に関する仮説検定

- 母集団の統計モデルを作る。  
母集団においてデータは正規分布をする。

$$X_i \sim N(\mu, \sigma^2), \quad i = 1, 2, \dots, n, \quad iid \quad (27)$$

- 検定仮説を立てる。

両側検定

$H_0: \mu = \mu_0$  母平均  $\mu$  は  $\mu_0$  に等しい。

$H_1: \mu \neq \mu_0$  母平均  $\mu$  は  $\mu_0$  に等しくない。

片側検定 1

$H_0: \mu = \mu_0$  母平均  $\mu$  は  $\mu_0$  に等しい。

$H_1: \mu \geq \mu_0$  母平均  $\mu$  は  $\mu_0$  より大きい。

片側検定 2

$H_0: \mu = \mu_0$  母平均  $\mu$  は  $\mu_0$  に等しい。

$H_1: \mu \leq \mu_0$  母平均  $\mu$  は  $\mu_0$  より小さい。

- $H_0$  の下での検定統計量の分布を求める。

$$\bar{X} \sim N(\mu_0, \sigma^2/n) \quad \text{or} \quad Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim N(0, 1) \quad (28)$$

もしくは  $\sigma$  を  $s$  で置き換えて

$$T = \frac{\bar{X} - \mu_0}{s/\sqrt{n}} \sim t(n-1) \quad (29)$$

- 棄却域を決める。

$H_0$  の下での検定統計量の分布で起こりそうもないこと（確率が小さいこと）を定義する。

両側検定

$$\text{棄却域} = (T < t_{(\alpha/2), n-1}) \quad \text{or} \quad (t_{(1-\alpha/2), n-1} < T) \quad (t_{(\alpha/2)} = -t_{(1-\alpha/2), n-1})$$

片側検定

$$\text{棄却域} = (t_{(1-\alpha), n-1} < T)$$

- データを取る。  $X_i, i = 1, 2, \dots, n$  から  $z$  や  $t$  を計算する。（小文字はデータから計算された  $Z$  や  $T$  の値を示す。）

- 観測された検定統計量が棄却域に入るかどうかを調べる。

棄却域に入ったということは  $H_0$  の下で起こりそうもないことが起こったということなので  $H_0$  を棄却する。

両側検定の場合、  $t < -t_{(1-\alpha/2), n-1}$  または  $t_{(1-\alpha/2), n-1} < t$  ならば  $H_0$  を棄却する。

p-値の計算

$H_0$  の下で  $T > t$  または  $T < t$  となる確率。この確率が小さい場合、  $H_0$  の下で起こりにくいことが起こったとして  $H_0$  を棄却する。

## 5.2 二つの母集団の母数（の差）に関する仮説検定

$\sigma_1 = \sigma_2 (= \sigma^2)$  が仮定出来る場合

母集団の統計モデルは、

$$X_{ij} \sim N(\mu_j, \sigma^2), \quad i = 1, 2, \dots, n_j, \quad j = 1, 2 \quad iid \quad (30)$$

両側検定

$H_0 : \mu_1 = \mu_2$     ふたつ母平均  $\mu_1, \mu_2$  は等しい。

$H_1 : \mu_1 \neq \mu_2$     ふたつ母平均等しくない。

片側検定 1

$H_0 : \mu_1 = \mu_2$     ふたつ母平均  $\mu_1, \mu_2$  は等しい。

$H_1 : \mu_1 \geq \mu_2$     母平均  $\mu_1$  は  $\mu_2$  より大きい。

片側検定 2

$H_0 : \mu_1 = \mu_2$     ふたつ母平均  $\mu_1, \mu_2$  は等しい。

$H_1 : \mu_1 \leq \mu_2$     母平均  $\mu_1$  は  $\mu_2$  より小さい。

検定統計量は

$$\bar{X}_1 \sim N(\mu_1, \sigma^2/n_1) \quad \text{and} \quad \bar{X}_2 \sim N(\mu_2, \sigma^2/n_2) \quad (31)$$

であるから

$$\bar{X}_1 - \bar{X}_2 \sim N(\mu_1 - \mu_2, \sigma^2/n_1 + \sigma^2/n_2) = N\left(\mu_1 - \mu_2, \sigma^2\left(\frac{1}{n_1} + \frac{1}{n_2}\right)\right) \quad (32)$$

したがって

$$Z = \frac{\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)}{\sqrt{\sigma^2\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \sim N(0,1) \quad (33)$$

もしくは  $\sigma$  を  $s$  で置き換えて

$$T = \frac{\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)}{\sqrt{\frac{s^2}{n_1} + \frac{s^2}{n_2}}} = \frac{\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)}{s\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t(n_1 + n_2 - 2) \quad (34)$$

ただし、ふたつのサンプルのデータをまとめて（プールして）推定した  $\sigma^2$  の推定値を

$$s^2 = \frac{1}{n_1 + n_2 - 2} \left( \sum_{i=1}^{n_1} (X_{i1} - \bar{X}_1)^2 + \sum_{i=1}^{n_2} (X_{i2} - \bar{X}_2)^2 \right) \quad (35)$$

とする。

棄却域の計算は以下の通り。ただし、 $t_{(\beta)}$  を自由度  $n_1 + n_2 - 2$  の  $t$  分布において下側確率が  $\beta$  となる点とする。なお、 $t_{(1-\beta)} = -t_{(\beta)}$  である。

両側検定

$$\text{棄却域} = (T < t_{(\alpha/2)}) \text{ or } (t_{(1-\alpha/2)} < T)$$

片側検定 1

$$\text{棄却域} = (t_{(1-\alpha)} < T)$$

片側検定 2

$$\text{棄却域} = (T < t_{(\alpha)})$$

### 5.3 二つの母集団の母数（の差）に関する仮説検定

$\sigma_1 = \sigma_2$  が仮定出来ない場合

$$\bar{X}_1 \sim N(\mu_1, \sigma_1^2/n_1) \quad \text{and} \quad \bar{X}_2 \sim N(\mu_2, \sigma_2^2/n_2) \quad (36)$$

したがって、

$$\bar{X}_1 - \bar{X}_2 \sim N(\mu_1 - \mu_2, \sigma_1^2/n_1 + \sigma_2^2/n_2) \quad (37)$$

しかし、これの  $\sigma_1, \sigma_2$  を標本標準偏差  $s_1, s_2$  で置き換えたものは  $t$  分布をしない。したがって、検定統計量とその自由度を以下のように変更する。（Welch）

$$T = \frac{\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \sim t(\nu) \quad (38)$$

ただし、修正された自由度を

$$\nu = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{\left(\frac{s_1^2}{n_1}\right)^2}{n_1-1} + \frac{\left(\frac{s_2^2}{n_2}\right)^2}{n_2-1}} = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{s_1^4}{n_1^2(n_1-1)} + \frac{s_2^4}{n_2^2(n_2-1)}} \quad (39)$$

とする。

棄却域の計算は上記の通り。ただし、 $t_{(\beta)}$  を自由度  $\nu$  の  $t$  分布において下側確率が  $\beta$  となる点とする。

## 6 3 母集団以上の平均値の差の検定

t 検定は、一つの母集団の平均値が、特定の値であるという仮説や、二つの母集団の平均値が同じであるという仮説を検定する際に用いられる。

母集団が三つ以上になった場合、t 検定は使わない。

その理由は第 1 種の過誤の確率をコントロールできないからである。

いま、三つの母集団があり、その平均を  $\mu_1, \mu_2, \mu_3$  とする。数学的には  $\mu_1 = \mu_2, \mu_2 = \mu_3$  が  $\mu_1 = \mu_2 = \mu_3$  を意味するので、検定仮説として次の二つを考える。

$$\begin{aligned} H_{01} : \mu_1 &= \mu_2 \\ H_{11} : \mu_1 &\neq \mu_2 \end{aligned} \tag{40}$$

and

$$\begin{aligned} H_{02} : \mu_2 &= \mu_3 \\ H_{12} : \mu_2 &\neq \mu_3 \end{aligned} \tag{41}$$

それぞれの仮説を有意水準  $\alpha$  で検定すると、二つの検定で少なくとも一つの type I error を起こす確率は、これらの仮説が独立であるとしても、 $1 - (1 - \alpha)^2 \approx 2\alpha$  となる。独立でない場合はどうなるか解らない。

一般的に、母集団が  $p$  個あり、 $p - 1$  個の仮説群の検定を行った場合、それらが独立の場合には全体としての type I error の確率 (少なくとも一つの type I error を起こす確率) は  $1 - (1 - \alpha)^{p-1} \approx (p - 1)\alpha$  となる。

したがって、母集団が三つ以上ある場合には、別の方式を考えなければならない。

一般的に  $p$  個の母集団があり、母分散は共通であると仮定できる場合、

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_p \quad (\text{すべての母平均は等しい}) \tag{42}$$

$$H_1 : \text{not } H_0 \quad (\text{少なくとも一つの母平均は他とことなる}) \tag{43}$$

を適切な検定統計量を用いて有意水準  $\alpha$  で検定することになる。(通常は  $\sigma$  の値を未知として、F 統計量を用いる。)

このような手法を (一元配置の) 分散分析と呼ぶ。(後述)

## 7 回帰分析

### 7.1 はじめに

単一の連続的な変数  $Y$  (従属変数、基準変数、被説明変数) を一つ、または、複数の変数群  $X$  (独立変数、説明変数、予測変数) で説明する方法。

一般的な回帰曲線の定義は、説明変数  $X$  を与えた時の基準変数  $Y$  の条件付き期待値  $\mathcal{E}(Y|X)$  の描く曲線(曲面)のこと。この曲線が、直線(平面)であると仮定したものが、以下に示す回帰分析である。

$$y_i = \alpha + \beta x_i + \varepsilon_i \quad \text{単回帰分析} \quad (44)$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \varepsilon_i \quad \text{重回帰分析(説明変数ふたつ)} \quad (45)$$

$$y_i = \beta_0 + \sum_{j=1}^q \beta_j x_{ij} + \varepsilon_i \quad \text{重回帰分析(一般形)} \quad (46)$$

ただし、 $i = 1, 2, \dots, N$ , として

$$\varepsilon_i \sim N(0, \sigma^2) \quad (47)$$

また、誤差を除いた部分、すなわち、 $X$  を与えた時の  $Y$  の期待値の部分

$$\hat{y}_i = \beta_0 + \sum_{j=1}^q \beta_j x_{ij} \quad (48)$$

を予測値と呼ぶ。

回帰係数  $\beta_j$  は、他の説明変数の値を固定した時に、第  $j$  番目の変数  $x_{ij}$  を一単位だけ動かした時の  $y$  の変化量を表す。

### 7.2 回帰係数の推定

単回帰分析の場合、回帰係数の推定値は

$$\hat{\beta} = \frac{\text{Cov}(Y, X)}{\text{var}(X)} = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} \quad (49)$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X} \quad (50)$$

となり、一般的には

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (51)$$

となる。

### 7.3 分散説明率・決定係数・ $R^2$

$$R^2 = \frac{\text{var}(\hat{Y})}{\text{var}(Y)} \quad (52)$$

を分散説明率と呼ぶ。この値が大きい方がよい。ちなみに、単回帰の場合、 $R^2$  は  $Y$  と  $X$  の相関係数の2乗となっている。重回帰の場合は  $\text{Corr}(Y, \hat{y})$  の2乗である。

どのような説明変数を用いても、説明変数を増やすことにより、 $R^2$  は増加する。

## 7.4 回帰分析の行列表記

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (53)$$

$$= \hat{\mathbf{y}} + \boldsymbol{\varepsilon}, \quad \hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\beta} \quad (54)$$

$\mathbf{y} \rightarrow N \times 1$	基準変数
$\mathbf{X} \rightarrow N \times (q+1)$	説明変数 + 1
$\boldsymbol{\beta} \rightarrow (q+1) \times 1$	回帰係数 (切片を含む)
$\hat{\mathbf{y}} \rightarrow \mathbf{X}\boldsymbol{\beta} = N \times 1$	モデルの予測値

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_{21} \\ \vdots \\ \hat{y}_N \end{pmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1q} \\ 1 & x_{21} & x_{22} & \dots & x_{2q} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{Nq} \end{bmatrix} \times \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_q \end{pmatrix} \quad (55)$$

$$y_i = \beta_0 + \sum_{j=1}^q x_{ij}\beta_j + \varepsilon_{ij} \quad (56)$$

ここで

$$x_{i0} = 1 \quad (57)$$

とすると

$$y_i = \sum_{j=0}^q x_{ij}\beta_j + \varepsilon_{ij} \quad (58)$$

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (59)$$

$$\hat{y}_i = \sum_{j=0}^q x_{ij}\beta_j \quad (60)$$

残差の2乗和 RSS を回帰係数  $\boldsymbol{\beta}$  の要素で偏微分すると

$$\frac{\partial \text{RSS}}{\partial \beta_k} = 2 \sum_{i=1}^N (y_i - \hat{y}_i) \times \frac{\partial (y_i - \hat{y}_i)}{\partial \beta_k} \quad (61)$$

$$= 2 \sum_{i=1}^N (y_i - \hat{y}_i) \times \frac{\partial \sum_{j=0}^q x_{ij}\beta_j}{\partial \beta_k} \quad (62)$$

$$= 2 \sum_{i=1}^N (y_i - \hat{y}_i) \times x_{ik} \quad (63)$$

$$= 2 \sum_{i=1}^N (y_i - \sum_{j=0}^q x_{ij}\beta_j) \times x_{ik} \quad (64)$$

$$= 2 \sum_{i=1}^N y_i x_{ik} - 2 \sum_{i=1}^N \left( \sum_{j=0}^q x_{ij} \beta_j \right) x_{ik} \quad (65)$$

$$= 2 \sum_{i=1}^N y_i x_{ik} - 2 \sum_{i=1}^N \left( \sum_{j=0}^q x_{ij} x_{ik} \beta_j \right) \quad (66)$$

$$= 2 \sum_{i=1}^N y_i x_{ik} - 2 \sum_{j=0}^q \left( \sum_{i=1}^N x_{ij} x_{ik} \right) \beta_j \quad (67)$$

$\sum_{i=1}^N y_i x_{ik}$  は  $X'y$  の第  $k$  要素、 $\sum_{i=1}^N x_{ij} x_{ik}$  は  $X'X$  の第  $jk$  要素  
したがって、方程式  $\frac{\partial \text{RSS}}{\partial \beta_k} = 0, k = 0, 1, 2, \dots, q$  は、まとめて

$$X'y - X'X\beta = 0 \quad (68)$$

と書ける。従って、解は

$$\hat{\beta} = (X'X)^{-1} X'y \quad (69)$$

となる。

## 7.5 仮説検定

モデル全体の検定、すなわち、 $H_0: \beta_0 = 0, \beta_1 = 0, \beta_2 = 0, \dots, \beta_q = 0$  の検定には F 分布を使う。  
個々の回帰係数が 0 出ないという仮説  $H_0: \beta_j = 0$  の検定には t 分布を使う。

$R^2$  が小さくてもモデルの帰無仮説が棄却されることがある。(統計的には回帰モデルを作る意味はあるが、あまり役に立つモデルではない。)

統計的仮定：誤差項は正規分布 (平均は 0、分散は共通)

$$y_i = \beta_0 + \sum_{j=1}^q x_{ij} \beta_j + \varepsilon_{ij} = \sum_{j=0}^q x_{ij} \beta_j, \quad \varepsilon_{ij} \sim N(0, \sigma^2) \quad (70)$$

行列で書くと

$$y = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I) \quad (71)$$

このモデルに関して、以下の仮説検定を考える。

$$H_0: \beta_{(1)} = \beta_{(2)} = \dots = \beta_{(r)} = 0 \quad (72)$$

$$H_1: \quad \text{上記の制約がない} \quad (73)$$

すなわち、ある特定の  $r$  個の回帰係数が 0 であるという帰無仮説である。このとき、帰無仮説の下で

$$F = \frac{(\text{RSS}_0 - \text{RSS}_1)/r}{\text{RSS}_1/(N - q - 1)} \quad (74)$$

が第 1 自由度  $r$ 、第 2 自由度  $N - q - 1$  の F 分布をすることが分かっている。ただし、 $\text{RSS}_k$  は、仮説  $H_k$  の元でパラメタを推定した際の残差の 2 乗和

$$\text{RSS}_k = \sum_{i=1}^N (y_i - \widehat{y}_{(k)_i})^2, \quad k = 0 \text{ or } 1 \quad (75)$$

である。(  $\text{RSS}_0$  は  $H_0$  で 0 と指定された回帰係数に対応する説明変数を  $X$  行列から除いて回帰係数を推定した時の残差の 2 乗和。その時の予測値が  $\widehat{y}_{(0)_i}$  。また、 $\text{RSS}_1$  は全ての説明変数を用いて計算した残差の 2 乗和で、その時の予測値が  $\widehat{y}_{(1)_i}$  ) 。

これは  $H_0$  が切片を除くすべての回帰係数が 0 という  $r=q$  個の制約の場合

$$F = \frac{\sum_{i=1}^N (\hat{y}_i - \bar{y})^2 / q}{\sum_{i=1}^N (y_i - \hat{y}_i)^2 / (N - q - 1)} \quad (76)$$

となる。ただし  $\hat{y}_i$  は  $H_1$  の下での予測値、 $\bar{y}$  は  $H_0$  の下での予測値である。

また、一つの回帰係数  $\beta_{(1)}$  のみが 0 という  $r=1$  個の制約の場合も同様の F 検定 (第 1 自由度 1、第 2 自由度  $N - q - 1$ ) が使えるが、通常は、それと同値の t 検定 (自由度  $N - q - 1$ ) を用いることが多い。

$RSS_1 = \sum_{i=1}^N (y_i - \hat{y}_i)^2$  を誤差平方和 (error SS)、 $RSS_0 - RSS_1$  を回帰平方和 (model SS)、 $\sum_{i=1}^N (y_i - \bar{y})^2$  を全平方和 (total SS) と呼ぶことがある。また、

$$R^2 = \frac{\sum_{i=1}^N (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^N (y_i - \bar{y})^2} = \frac{\text{var}(\hat{\mathbf{y}})}{\text{var}(\mathbf{y})} = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2} \quad (77)$$

を決定係数と呼ぶ。

## 7.6 回帰分析のバリエーション

### 1. 単回帰分析

基準変数 1 個、説明変数 1 個

$$\mathbf{y} = \beta_0 \mathbf{1} + \beta \mathbf{x} + \boldsymbol{\varepsilon} \quad (78)$$

### 2. 重回帰分析

基準変数 1 個、説明変数  $q$  個

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (79)$$

### 3. 多変量回帰分析

基準変数  $p$  個、説明変数  $q$  個

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}' + \boldsymbol{\varepsilon}, \quad \boldsymbol{\beta} = (p \times q) \quad (80)$$

$$\widehat{\boldsymbol{\beta}}' = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \quad (81)$$

### 4. 多項式回帰分析

$$\hat{y}_i = \sum_{m=0}^M \beta_m x_i^m \quad (82)$$

$$\mathbf{Y} = \mathbf{G}^{(M)}(\mathbf{x})\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \mathbf{G}^{(M)}(\mathbf{x}) = (n \times M + 1) \quad (83)$$

$$\mathbf{G}(\mathbf{x}) = [\mathbf{g}_m(\mathbf{x})], \quad m = 0, 1, \dots, M \quad (84)$$

$$\mathbf{g}_m(\mathbf{x}) = [x_i^m] \quad (85)$$

## 5. 非線形回帰としてのスプライン回帰分析 (多項式 + 切斷冪関数)

$$\hat{y}_i = \sum_{m=0}^M \beta_m x_i^m + \sum_{\ell=1}^K \beta_\ell (x_i - c_\ell)_+^M \quad (86)$$

$$(x_i - c_\ell)_+^M = \begin{cases} 0 & (x_i < c_\ell) \\ (x_i - c_\ell)^M & (x_i \geq c_\ell) \end{cases} \quad (87)$$

$$\mathbf{Y} = \mathbf{G}^{(K,M)}(\mathbf{x})\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad \mathbf{G}^M(\mathbf{x}) = (n \times M + K + 1) \quad (88)$$

ただし、 $K$  は節点の数、 $c_\ell$  は節点

## 6. ニューラルネットワーク

$$y_i = \beta_0 + \sum_{j=1}^Q \psi(a_j(x_i - b_j))\beta_j + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (89)$$

ただし、

$$\psi(x) = \frac{1}{(1 + \exp(-x))} \quad (90)$$

は係数  $a_j, b_j$  を持つロジスティク関数である。従って

$$\mathbf{y} = \mathbf{G}(\mathbf{x})\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (91)$$

の形で書くと

$$\mathbf{G}(\mathbf{x}) = [\mathbf{1}, \psi(a_1(\mathbf{x} - b_1)), \psi(a_2(\mathbf{x} - b_2)), \dots, \psi(a_Q(\mathbf{x} - b_Q))] \quad (92)$$

となる。

スプライン回帰や多項式と似ている点は、両者とも、本来の説明変数を  $G$  という形に展開し、それらの線形結合で  $y$  を予測しようと試みることである。

両者の違いは、スプラインの場合は、展開に際して用いられる接点  $c_\ell, \ell = 1, 2, \dots, K$  は定数であるのに対し、ニューラルネットワークにおいては、その係数  $(a_j, b_j), j = 1, 2, \dots, Q$  はパラメタとして推定されるものである点である。

また、多項式回帰の場合、中間層を増やすことは、結局、最初に展開された説明変数の線形回帰を作ることになり、意味がないが、スプラインや、特にニューラルネットワークの場合は、意味がある。

## 7.7 デザイン行列

水準数が  $A$  の 1 元配置の分散分析モデルは、各水準での観測値数を  $n_j, j = 1, 2, \dots, A$  とすると、

$$y_{ij} = \hat{y}_{ij} + \varepsilon_{ij}, \quad \hat{y}_{ij} = \mu_j, \quad i = 1, 2, \dots, n_j, \quad j = 1, 2, \dots, A \quad (93)$$

である。これを線形モデル  $\hat{y} = X\beta$  の形で書くと

$$\begin{pmatrix} \hat{y}_{11} \\ \hat{y}_{21} \\ \vdots \\ \hat{y}_{n_1 1} \\ \hat{y}_{12} \\ \hat{y}_{22} \\ \vdots \\ \hat{y}_{1A} \\ \vdots \\ \hat{y}_{n_A A} \end{pmatrix} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 1 \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} \times \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_A \end{pmatrix} \quad (94)$$

となる。ただし、 $\mu_j$  は、第  $j$  群で観測される  $y_{ij}, i = 1, 2, \dots, n_j$ 、の母平均である。(これをセル平均と呼ぶことがある。)

この同じモデルを、以下のように書き表すことも出来る。ただし、 $\mu$  は全水準の平均 (grand mean) を表す母数、 $\alpha_j$  は第  $j$  水準の効果 (effect) を表す母数である。

$$y_{ij} = \hat{y}_{ij} + \varepsilon_{ij}, \quad \hat{y}_{ij} = \mu + \alpha_j, \quad i = 1, 2, \dots, n_j, \quad j = 1, 2, \dots, A \quad (95)$$

である。これを線形モデル  $\hat{y} = X\beta$  の形で書くと

$$\begin{pmatrix} \hat{y}_{11} \\ \hat{y}_{21} \\ \vdots \\ \hat{y}_{n_1 1} \\ \hat{y}_{12} \\ \hat{y}_{22} \\ \vdots \\ \hat{y}_{1A} \\ \vdots \\ \hat{y}_{n_A A} \end{pmatrix} = \begin{bmatrix} 1 & 1 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & 0 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & 0 & 0 & \dots & 1 \end{bmatrix} \times \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_A \end{pmatrix} \quad (96)$$

しかし、このデザイン行列は、その第 2 列目から最後までを加えると第 1 列目の単位ベクトルとなる。このように、行列のある列が、残りの列の線形結合として表されるような場合にその行列の各列は 1 次独立ではない (1 次従属である) と言う。そして、その場合、 $(X'X)^{-1}$  は存在しない。従って、パラメタ  $\beta$  の推定値は一意には定まらない。

そこで、通常は、効果の和が 1 であるという制約条件

$$\sum_{j=1}^A \alpha_j = 0 \quad (97)$$

をもちいて、母数の数を一つ減らすことを考える。すると、 $\alpha_A = -\sum_{j=1}^{A-1} \alpha_j$  であるから、同じモデルを

$$\begin{pmatrix} \hat{y}_{11} \\ \hat{y}_{21} \\ \vdots \\ \hat{y}_{n_11} \\ \hat{y}_{12} \\ \hat{y}_{22} \\ \vdots \\ \hat{y}_{1A} \\ \vdots \\ \hat{y}_{n_{AA}} \end{pmatrix} = \begin{bmatrix} 1 & 1 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & -1 & -1 & \dots & -1 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & -1 & -1 & \dots & -1 \end{bmatrix} \times \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_{A-1} \end{pmatrix} \quad (98)$$

と書くこともできるが、このデザイン行列の各列は1次独立となる。

この二つのセル平均の表現

$$\hat{y}_{ij} = \mu_j \quad \text{and} \quad \hat{y}_{ij} = \mu + \alpha_j, \quad \sum_{j=1}^A \alpha_j = 0 \quad (99)$$

を比べると、

$$\mu = \frac{1}{A} \sum_{j=1}^A \mu_j \quad (\text{母平均の平均}) \quad (100)$$

$$\alpha_j = \mu_j - \mu \quad (\text{母平均の } \mu \text{ からのズレ}) \quad (101)$$

という関係があることがわかる。

2元配置の場合のモデルは

$$y_{ijk} = \hat{y}_{ijk} + \varepsilon_{ijk}, \quad \hat{y}_{ijk} = \mu + \alpha_j + \beta_k, \quad i = 1, 2, \dots, n_{ij}, \quad j = 1, 2, \dots, A \quad k = 1, 2, \dots, B \quad (102)$$

であるが、全ての  $n_{ij}$  を 1 とした場合のデザイン行列は

$$\begin{pmatrix} \hat{y}_{111} \\ \hat{y}_{112} \\ \vdots \\ \hat{y}_{11B} \\ \hat{y}_{121} \\ \hat{y}_{122} \\ \vdots \\ \hat{y}_{1A1} \\ \vdots \\ \hat{y}_{1AB} \end{pmatrix} = \begin{bmatrix} 1 & 1 & 0 & \dots & 0 & 1 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & 1 & 0 & \dots & 0 & 0 & 0 & \dots & 1 \\ 1 & 0 & 1 & \dots & 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & 0 & 0 & \dots & 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & 0 & 0 & \dots & 1 & 0 & 0 & \dots & 1 \end{bmatrix} \times \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_A \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_B \end{pmatrix} \quad (103)$$

となる。このデザイン行列の場合、要因 A に対応する列の和、並びに要因 B に対応する列の和がそれぞれ単位ベクトルになり、二つの線形従属関係がある。従って、パラメタの推定値は一意には定まらない。なお、このデザイン行列を  $X = \{1, A, B\}$  と分割し、すべてのデータの分だけ生成すると、二つの要因の水準間の同時度数行列  $F$  は

$$F = A'B \quad (104)$$

という  $AxB$  の行列で与えられる。

通常の制約条件  $\sum_{j=1}^A \alpha_j = 0, \sum_{k=1}^B \beta_k = 0$  を入れると、 $\beta_B = -\sum_{k=1}^{B-1} \beta_k$  であるから、同じモデルを

$$\begin{pmatrix} \hat{y}_{111} \\ \hat{y}_{112} \\ \vdots \\ \hat{y}_{11B} \\ \hat{y}_{121} \\ \hat{y}_{122} \\ \vdots \\ \hat{y}_{1A1} \\ \vdots \\ \hat{y}_{1AB} \end{pmatrix} = \begin{bmatrix} 1 & 1 & 0 & \dots & 0 & 1 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & 1 & 0 & \dots & 0 & -1 & -1 & \dots & -1 \\ 1 & 0 & 1 & \dots & 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & -1 & -1 & \dots & -1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & -1 & -1 & \dots & -1 & -1 & -1 & \dots & -1 \end{bmatrix} \times \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_{A-1} \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_{B-1} \end{pmatrix} \quad (105)$$

と書くこともできる。このデザイン行列の各列は1次独立である。

なお、モデルに交互作用が組み込まれている場合には、交互作用のパラメタ  $\alpha\beta_{jk}$  に対応するデザイン行列の列を、 $\alpha_j$  と  $\beta_k$  に対応する  $X$  の2つの列  $x_{(1+j)}, x_{(1+A+k)}$  の要素毎の積を要素として持つベクトルとして定義していけばよい。

例えば、 $A = 2, B = 3$  の場合には以下ようになる。

$$\mathbf{y} = \begin{pmatrix} y \\ y_{111} \\ y_{211} \\ \vdots \\ y_{N_{11}11} \\ y_{112} \\ y_{212} \\ \vdots \\ y_{N_{12}12} \\ y_{113} \\ y_{213} \\ \vdots \\ y_{N_{13}13} \\ y_{121} \\ y_{221} \\ \vdots \\ y_{N_{22}21} \\ y_{122} \\ y_{222} \\ \vdots \\ y_{N_{22}22} \\ y_{123} \\ y_{223} \\ \vdots \\ y_{N_{23}23} \end{pmatrix} \quad \mathbf{X} = \begin{bmatrix} \mu & \alpha_1 & \beta_1 & \beta_2 & \alpha\beta_{11} & \alpha\beta_{12} \\ 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 0 & 1 & 0 & 1 \\ 1 & 1 & -1 & -1 & -1 & -1 \\ 1 & 1 & -1 & -1 & -1 & -1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & -1 & -1 & -1 & -1 \\ 1 & -1 & 1 & 0 & -1 & 0 \\ 1 & -1 & 1 & 0 & -1 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & -1 & 1 & 0 & -1 & 0 \\ 1 & -1 & 0 & 1 & 0 & -1 \\ 1 & -1 & 0 & 1 & 0 & -1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & -1 & 0 & 1 & 0 & -1 \\ 1 & -1 & -1 & -1 & 1 & 1 \\ 1 & -1 & -1 & -1 & 1 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & -1 & -1 & -1 & 1 & 1 \end{bmatrix} \quad \boldsymbol{\beta} = \begin{pmatrix} \mu \\ \alpha_1 \\ \beta_1 \\ \beta_2 \\ \alpha\beta_{11} \\ \alpha\beta_{12} \end{pmatrix} \quad (106)$$

ただし、たとえば  $y_{113}$  の要素のところ入っている  $-1$  は、本来は  $y_{113} = \mu + \alpha_1 + \beta_3 + \alpha\beta_{13} + e_{113}$  であるが、そこに、制約条件を使って  $\beta_3 = -\beta_1 - \beta_2$  を代入したためである。同様にまた、 $y_{121}$  の要素のところ入っている  $-1$  は、 $\alpha_2 = -\alpha_1$  を代入したためである。また、交互作用  $\alpha\beta_{ij}$  に対応するデザイン行列の列は  $\alpha_1$  と  $\beta_j, j = 1, 2$  に対応するデザイン要列の列の積として定義されている。

2 要因モデルの場合のセル平均と効果を表す母数の関係は、交互作用の効果に関しても

$$\sum_{j=1}^A \alpha\beta_{jk} = \sum_{k=1}^B \alpha\beta_{jk} = 0 \quad (107)$$

を仮定すると、

$$\mu = \frac{1}{AB} \sum_{j=1}^A \sum_{k=1}^B \mu_{jk} \quad (\text{母平均の平均}) \quad (108)$$

$$\alpha_j = \mu_{.j} - \mu \quad (\text{行平均の } \mu \text{ からのズレ}) \quad (109)$$

$$\beta_k = \mu_{.k} - \mu \quad (\text{列平均の } \mu \text{ からのズレ}) \quad (110)$$

$$\alpha\beta_{jk} = \mu_{jk} - (\mu + \alpha_j + \beta_k) \quad (\text{加算モデルからのズレ}) \quad (111)$$

となる。ただし、

$$\mu_{j.} = \frac{1}{B} \sum_{k=1}^B \mu_{jk}, \quad (\text{行平均}) \quad \mu_{.k} = \frac{1}{A} \sum_{j=1}^A \mu_{jk} \quad (\text{列平均}) \quad (112)$$

である。

## 8 実験計画・分散分析

### 8.1 用語

- 要因 (因子) factor  
実験結果に影響を与えると考えるもの。(名義尺度の説明変数)
- 水準 level  
要因として用いられた変数の値
- 処理 treatment  
要因の各水準の組み合わせ
- 効果 effect 各要因の各水準の効果を数値化したもの
- 実験計画 design of experiment, n 元配置デザイン  
n は要因の数
- 被験者間要因 between subject factor、被験者内要因 within subject factor  
その要因の各水準が異なる被験者に割り当てられるか否か。  
被験者内要因は繰り返し測定、反復測定、対応のある測定、対比較とも呼ばれる。
- 繰り返し  
各処理の下での測定数。
- 平方和 Sum of Squares (SS)、平均平方等 Mean Square (MS)  
F 統計量を計算するために使われる統計量で、平均平方は平方和を対応する自由度で割ったもの。  
F 統計量は平均平方の比となる。
- 自由度 degrees of freedom  
平均平方の分母で、主効果の自由度は要因の水準の数 - 1。
- 平方和、自由度の分解  
恒等式による平方和等の加算的分解。
- 主効果 main effect、交互作用 interaction effect
- 固定効果 (母数モデル) fixed effect, 変量効果 (変量モデル) random effect
- 線形モデル linear model, 加算モデル additive model
- crossed, nested factors
- 多重比較 multiple comparison
- フィッシャーの3原則 (反復、無作為化、局所管理)
- 制御因子、標示因子、ブロック因子

## 8.2 1元配置分散分析

要因 A				
水準 1	水準 2	水準 3		水準 A
$\mu_1$	$\mu_2$	$\mu_3$		$\mu_A$
				$\mu$

### 8.2.1 モデル

$$Y_{ij} = \mu_j + \varepsilon_{ij} \quad (113)$$

$$= \mu + \alpha_j + \varepsilon_{ij} \quad (114)$$

$$\varepsilon_{ij} \sim N(0, \sigma_e^2) \quad (115)$$

ただし

$$\mu = \frac{1}{A} \sum_{j=1}^A \mu_j \quad \text{or} \quad \mu = \frac{1}{\sum_{j=1}^A n_j} \sum_{j=1}^A n_j \mu_j \quad (116)$$

$$\alpha_j = \mu_j - \mu \quad (117)$$

### 8.2.2 検定仮説

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_A \quad (118)$$

$$H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_A = 0 \quad (119)$$

### 8.2.3 データ、自由度、平方和の分解

$$Y_{ij} - \mu = \mu_j - \mu + Y_{ij} - \mu_j \quad (120)$$

$$Y_{ij} - \bar{Y}_{..} = \bar{Y}_{.j} - \bar{Y}_{..} + Y_{ij} - \bar{Y}_{.j} \quad (121)$$

$$\sum_{j=1}^A \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{..})^2 = \sum_{j=1}^A n_j (\bar{Y}_{.j} - \bar{Y}_{..})^2 + \sum_{j=1}^A \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_{.j})^2 \quad (122)$$

$$SS_{total} = SS_A + SS_{error} \quad (123)$$

$$= SS_{between} + SS_{within} \quad (124)$$

$$\sum_{j=1}^A n_j - 1 = A - 1 + \sum_{j=1}^A n_j - A \quad (125)$$

$$df_{total} = df_A + df_{error} \quad (126)$$

$$= df_{between} + df_{within} \quad (127)$$

### 8.2.4 分散分析表

source	SS	df	MS	EMS	F
A	$SS_A$	$df_A$	$\frac{SS_A}{df_A}$	$\sigma_e^2 + \frac{\sum_{j=1}^A n_j \alpha_j^2}{A-1}$	$\frac{MS_A}{MS_e}$
error	$SS_{error}$	$df_{error}$	$\frac{SS_{error}}{df_{error}}$	$\sigma_e^2$	
total	$SS_{total}$	$df_{total}$			

### 8.3 2元配置分散分析 (全ての $n_{jk}$ が等しい、バランスが取れている、直交している場合)

		要因 B ( $k$ )				
		水準 1	水準 2	水準 3	水準 B	
要因 A ( $j$ )	水準 1	$\mu_{11}$	$\mu_{12}$	$\mu_{13}$	$\mu_{1B}$	$\mu_{1.}$
	水準 2	$\mu_{21}$	$\mu_{22}$	$\mu_{23}$	$\mu_{2B}$	$\mu_{2.}$
	水準 A	$\mu_{A1}$	$\mu_{A2}$	$\mu_{A3}$	$\mu_{AB}$	$\mu_{A.}$
		$\mu_{.1}$	$\mu_{.2}$	$\mu_{.3}$	$\mu_{.B}$	$\mu$

#### 8.3.1 モデル

$$Y_{ijk} = \mu_{jk} + \varepsilon_{ijk} \quad (128)$$

$$= \mu + \alpha_j + \beta_k + \alpha\beta_{jk} + \varepsilon_{ijk} \quad (129)$$

$$\varepsilon_{ijk} \sim N(0, \sigma_e^2) \quad (130)$$

ただし

$$\mu = \frac{1}{AB} \sum_{j=1}^A \sum_{k=1}^B \mu_{jk} \quad (131)$$

$$\alpha_j = \mu_{j.} - \mu \quad (132)$$

$$\beta_k = \mu_{.k} - \mu \quad (133)$$

$$\alpha\beta_{jk} = \mu_{jk} - (\mu + \alpha_j + \beta_k) = \mu_{jk} - \mu_{j.} - \mu_{.k} + \mu \quad (134)$$

(交互作用は加算モデルからのズレを表している。)

#### 8.3.2 検定仮説

$$H_{0A} : \alpha_1 = \alpha_2 = \dots = \alpha_A = 0 \quad (\mu_{1.} = \mu_{2.} = \dots = \mu_{A.}) \quad (135)$$

$$H_{0B} : \beta_1 = \beta_2 = \dots = \beta_B = 0 \quad (\mu_{.1} = \mu_{.2} = \dots = \mu_{.B}) \quad (136)$$

$$H_{0AB} : \alpha\beta_{11} = \alpha\beta_{12} = \dots = \alpha\beta_{AB} = 0 \quad (137)$$

#### 8.3.3 データ、自由度、平方和の分解

$$Y_{ijk} - \mu = \mu_{jk} - \mu + Y_{ijk} - \mu_{jk} \quad (138)$$

$$= \mu_{j.} - \mu + \mu_{.k} - \mu + \mu_{jk} - \mu_{j.} - \mu_{.k} + \mu + Y_{ijk} - \mu_{jk}$$

$$Y_{ijk} - \bar{Y}_{...} = \bar{Y}_{.jk} - \bar{Y}_{...} + Y_{ijk} - \bar{Y}_{.jk} \quad (139)$$

$$= \bar{Y}_{.j.} - \bar{Y}_{...} + \bar{Y}_{.k.} - \bar{Y}_{...} + \bar{Y}_{.jk} - \bar{Y}_{.j.} - \bar{Y}_{.k.} + \bar{Y}_{...} + Y_{ijk} - \bar{Y}_{.jk}$$

$$\sum_{j=1}^A \sum_{k=1}^B \sum_{i=1}^n (Y_{ijk} - \bar{Y}_{...})^2 = nB \sum_{j=1}^A (\bar{Y}_{.j.} - \bar{Y}_{...})^2 + nA \sum_{k=1}^B (\bar{Y}_{.k.} - \bar{Y}_{...})^2 + \quad (140)$$

$$n \sum_{j=1}^A \sum_{k=1}^B (\bar{Y}_{.jk} - \bar{Y}_{.j.} - \bar{Y}_{.k.} + \bar{Y}_{...})^2 + \sum_{i=1}^n \sum_{j=1}^A \sum_{k=1}^B (Y_{ijk} - \bar{Y}_{.jk})^2 \quad (141)$$

$$SS_{total} = SS_A + SS_B + SS_{AB} + SS_{error} \quad (142)$$

$$nAB - 1 = A - 1 + B - 1 + (A - 1)(B - 1) + AB(n - 1) \quad (143)$$

$$df_{total} = df_A + df_B + df_{AB} + df_{error} \quad (144)$$

### 8.3.4 分散分析表

source	SS	df	MS	EMS	F
A	$SS_A$	$df_A$	$\frac{SS_A}{df_A}$	$\sigma_e^2 + \frac{nB \sum_{j=1}^A \alpha_j^2}{A-1}$	$\frac{MS_A}{MS_e}$
B	$SS_B$	$df_B$	$\frac{SS_B}{df_B}$	$\sigma_e^2 + \frac{nA \sum_{k=1}^B \beta_k^2}{B-1}$	$\frac{MS_B}{MS_e}$
AB	$SS_{AB}$	$df_{AB}$	$\frac{SS_{AB}}{df_{AB}}$	$\sigma_e^2 + \frac{n \sum_{j=1}^A \sum_{k=1}^B \alpha\beta_{kj}^2}{(A-1)(B-1)}$	$\frac{MS_{AB}}{MS_e}$
error	$SS_{error}$	$df_{error}$	$\frac{SS_{error}}{df_{error}}$	$\sigma_e^2$	
total	$SS_{total}$	$df_{total}$			