

## 統計解析手法 (初回除く)

8/6

## (2) 連続型

累積分布関数

 $F(x)$  が右から左からも連続なとき  $(\lim_{s \rightarrow 0} F(x+s) = \lim_{s \rightarrow 0} F(x-s) = F(x))$ 

連続

連続確率変数

・  $F(x)$  が連続 ・  $F(x)$  の導関数  $F'(x) = \frac{dF(x)}{dx} = f(x)$  が存在

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(t) dt$$

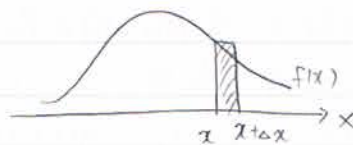
 $f(x)$ : 確率密度関数  $f(x) \geq 0, \int_{-\infty}^{\infty} f(x) dx = 1$ 

(離散型と確率分布と区別)

ここで微小な  $\Delta x > 0$  に対して

$$P(x \leq X \leq x + \Delta x) \approx f(x) \cdot \Delta x$$

と近似できる。図のように連続では

確率・面積, 離散は  $f(x_i) = P(X = x_i)$  $f(x)$  は  $P(X=x)$  を表さない (例:  $P(X=a) = 0, P(a < X < b) = F(b) - F(a) = \int_a^b f(x) dx$ )

## 2.2 期待値と分散

## (1) 定義と性質

$$\text{期待値 } E(X) = \begin{cases} \sum_i x_i f(x_i) & \text{離散} \\ \int_{-\infty}^{\infty} x f(x) dx & \text{連続} \end{cases}$$

母集団分布の期待値:  $\mu$  で表現

$$\text{分散 } \text{var}(X) = E(X - \mu)^2 = \begin{cases} \sum_i (x_i - \mu)^2 f(x_i) & \text{離散} \\ \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx & \text{連続} \end{cases}$$

母集団分布の分散:  $\sigma^2$  で表現標準偏差  $\sigma = \sqrt{\text{var}(X)}$  ( $E(X)$  と決まると合わせて)

## ・ 性質

①  $E(a + bX) = a + bE(X)$

②  $E(a) = a$

③  $E(X + Y) = E(X) + E(Y)$

④  $\text{var}(a + bX) = b^2 \text{var}(X)$

⑤  $\text{var}(a) = 0$

無次元量 (異種のデータや、同種の多数のデータを比較するに役立つ)

## (2) 標準化

標準化変数  $z = \frac{x - \mu}{\sigma}$  ( $a = -\frac{\mu}{\sigma}$ ,  $b = \frac{1}{\sigma}$  としたときの線形変換  $a + bx$ )

$$E(z) = E(a + bx) = a + b\mu = 0.$$

$$Var(z) = Var(a + bx) = b^2 Var(x) = b^2 \sigma^2 = 1$$

平均は0, 分散は1

## 2.3. 多変量分布

同時確率分布と周辺確率分布

ex) 1枚のコインを3回投げ、表H, 裏Tとする。

標本空間  $\Omega = \{HHH, HHT, HTH, THH, HTT, THT, TTH, TTT\}$

確率変数  $X$ : 最初2回の表の回数

$Y$ : 最後の "

$(X=x) \cap (Y=y)$  の確率:

$f(x, y)$  同時確率

HHH なら  $x=y=2$ , HHT なら  $x=2, y=1$

$$f(0, 1) = P[(X=0) \cap (Y=1)] = P(TTH) = 1/8.$$

$$f(1, 1) = P[(X=1) \cap (Y=1)] = P(HTH) + P(THT) = 2/8.$$

$X \backslash Y$	0	1	2	
0	1/8	1/8	0	2/8
1	1/8	2/8	1/8	4/8
2	0	1/8	1/8	2/8
	2/8	4/8	2/8	1.

$$P(X=1) = g(1) = f(1, 0) + f(1, 1) + f(1, 2) = 4/8.$$

$$P(Y=2) = h(2) = \dots = 2/8.$$

$g(x), h(y)$  は周辺確率分布.

離散確率変数:  $X, Y$

$$f(x_i, y_j) \geq 0, \sum_i \sum_j f(x_i, y_j) = 1, g(x) = \sum_j f(x, y_j), h(y) = \sum_i f(x_i, y)$$

連続確率変数:  $X, Y \rightarrow f(x, y)$ : 同時確率密度関数

$$f(x, y) \geq 0, \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1, g(x) = \int_{-\infty}^{\infty} f(x, y) dy, h(y) = \int_{-\infty}^{\infty} f(x, y) dx.$$

周辺確率密度関数.

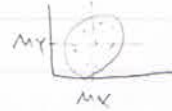
## 2.4 共分散・相関係数

## (1) 共分散

$X, Y$  がそれぞれの期待値  $\mu_X, \mu_Y$  から互いに関連しあっている程度

$$\text{cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

$\text{cov}(X, Y) > 0 \rightarrow X, Y$  は大抵同傾向  
 $< 0 \rightarrow$  " 反対傾向  
 $= 0 \rightarrow$  上記の関係なし


 $\text{cov} > 0$ 

 $\text{cov} = 0$ 

 $\text{cov} < 0$ 

## (2) 相関係数

・共分散は、変化の傾向を表すので単位、異なる2つの共分散は比較できない。

$\rightarrow$  標準化  $Z_X = \frac{X - \mu_X}{\sigma_X}, Z_Y = \frac{Y - \mu_Y}{\sigma_Y}$

・相関係数  $\rho$

$$\begin{aligned} \rho &= \text{cov}(Z_X, Z_Y) = E[(Z_X - \mu_{Z_X})(Z_Y - \mu_{Z_Y})] \\ &= E(Z_X, Z_Y) = E\left[\left(\frac{X - \mu_X}{\sigma_X}\right)\left(\frac{Y - \mu_Y}{\sigma_Y}\right)\right] = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} \end{aligned}$$

・ $\rho$  は  $-1 \leq \rho \leq 1$  の値をとる。

1に近いほど強い正の相関, -1に近いほど強い負の相関, 0:無相関

・もし、 $\rho = \pm 1$  なら  $X$  と  $Y$  は線形関係にある ( $Y = a + bX$  で表せる)

・相関係数は  $X$  と  $Y$  の1次結合の強さを測る尺度

## 2.5 条件付確率分布

## (1) 条件付確率分布

3回コインを投げた例 (cf 2.3.)

$$P(X=0 | Y=1) = \frac{P(X=0, Y=1)}{P(Y=1)} = \frac{1/8}{4/8} = 1/4$$

$$g(0|1) = \frac{f(0,1)}{h(1)}$$

$$\text{条件付確率分布 } g(x|y) = \frac{f(x,y)}{h(y)}, \quad h(y) \neq 0, \quad g(x|y) \neq 0$$

(2) 独立

・ 事象Aの確率が事象Bに影響しない。つまり  $P(A|B) = P(A)$ , 逆に  $P(B|A) = P(B)$  も成立するならば、 $P(A \cap B) = P(A)P(B)$  となる。事象AとBは独立であるという。  
↳ 従属

・ 確率分布の場合

$f(x, y)$  のある  $x, y$  について  $f(x, y) = g(x) \cdot h(y)$  が成立 = 独立

一般化すると、 $f(x_1, x_2, \dots, x_n) = f_1(x_1) \cdot f_2(x_2) \cdot \dots \cdot f_n(x_n) = \prod_{i=1}^n f_i(x_i)$   
添え字?

・ 共分散と独立

独立(離散):  $f(x_i, y_j) = g(x_i)h(y_j)$  となる。

$$E(XY) = \sum_i \sum_j x_i y_j f(x_i, y_j) = \sum_i \sum_j x_i y_j g(x_i)h(y_j) \\ = \sum_i x_i g(x_i) \cdot \sum_j y_j h(y_j) = E(X)E(Y)$$

よって共分散  $cov(X, Y) = E(XY) - E(X)E(Y) = 0$ 。  
↑ やってやれ。

$X, Y$  が独立  $\Rightarrow$  共分散 = 0, 相関係数 = 0       $\times$  逆は成り立たない。

独立: 変数間に関連がない  $\rightarrow$  確率分布そのものが決定

共分散・相関係数: 平均的な性質  $\rightarrow$  確率分布から決まる量決定

1/23

§3. 確率分布

- ・ 観測データ = 標本  $\leftarrow$  母集団: 確率分布によって表現
- ・ 確率分布を特徴づける定数 = パラメータ

3.1. 離散確率変数の確率分布

(1) ベルヌーイ試行  $\leftarrow$  また分布の性質

条件

- ① 結果が2種類 (ex. 「成功(S)」, 「失敗(F)」)
- ② 確率  $p = P(S)$  は一定  $\rightarrow P(F) = 1 - p$
- ③ 試行は独立



いま、 $S \rightarrow x=1$ ,  $F \rightarrow x=0$  とする。

$$\text{ベルヌーイ分布} \quad f(x) = p^x (1-p)^{1-x} \quad (x=0, 1)$$

$$E(x) = p, \quad \text{var}(x) = p(1-p)$$

これは1回の試行のみ。n回おこなうと

## (2) 二項分布

ベルヌーイ試行をn回行う。

$S: x$ 回,  $F: n-x$ 回起こる。とすると,  $(x=0, 1, \dots, n)$

$$\text{二項分布: } f(x) = n C_x p^x (1-p)^{n-x} \quad (x=0, 1, \dots, n)$$

パラメータ:  $n, p \xrightarrow{\text{明記}} B(n, p)$  と表すこともある

$$\sum f(x) = (p + 1 - p)^n = 1 \quad \leftarrow \text{二項定理より}$$

$$E(x) = np, \quad \text{var}(x) = np(1-p)$$

∴ n回の独立な試行。→ x回の成功確率  $p^x$ ,  $(n-x)$ 回の失敗確率  $(1-p)^{n-x}$

$$\text{最初のx回の成功なら, } p(\underbrace{S, S, \dots, S}_{x\text{回}}, \underbrace{F, \dots, F}_{n-x\text{回}}) = p^x (1-p)^{n-x}$$

x回の成功はどこで起してもよい = 組み合わせ  $n C_x$  をかける。

5.2.2 応用

## (3) ポアソン分布

二項分布の  
期待値

$\lambda, \lambda = \text{起こり}$

二項分布において  $np = \lambda$  (一定) とおく。  $n \rightarrow \infty$ ,  $p \rightarrow 0$  (稀小現象) とおけば、

$$\begin{aligned} n C_x p^x (1-p)^{n-x} &= \frac{n(n-1)\dots(n-x+1)}{x!} \cdot \left(\frac{\lambda}{n}\right)^x \cdot \left(1 - \frac{\lambda}{n}\right)^{n-x} \\ &= \frac{1}{x!} \lambda^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \cdot \frac{n(n-1)\dots(n-x+1)}{n^x} \\ &= \frac{1}{x!} \lambda^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \cdot \underbrace{1 \cdot \left(1 - \frac{\lambda}{n}\right) \cdot \dots \cdot \left(1 - \frac{\lambda-x+1}{n}\right)}_{n \rightarrow \infty \text{ で } 1} \end{aligned}$$

$$\frac{\lambda}{n} = \frac{1}{2} \text{ とおくと, } n = 2\lambda. \quad \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^{n-x} = \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n = \lim_{z \rightarrow \infty} \left\{ \left(1 - \frac{1}{2}\right)^{-2} \right\}^{-\lambda} = e^{-\lambda}$$

∴  $n \rightarrow \infty$ ,  $p \rightarrow 0$  の極限で,  $n C_x p^x (1-p)^{n-x} = \frac{\lambda^x e^{-\lambda}}{x!}$  (ポアソンの分布の法則)

$$\text{ポアソン分布: } f(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

パラメータ:  $\lambda \rightarrow \mu(\lambda)$  と表すこともある。  $E(x) = \lambda, \quad \text{var}(x) = \lambda$

下記の説明に主に使われる。

① 稀小現象の生起回数

② 到着数

e.g) 交通事故の1年間で起こる回数

1年間を8760時間に分割 → 事故が起こるかをベルヌーイ試行

1時間内に複数回起こる → 分割を増やしベルヌーイ試行

↓ 分割で限りなく増やす

$n \rightarrow \infty$   $n$  = 項分布 = ポアソン分布

### 3.2 連続確率変数の確率分布

#### (1) 指数分布

ポアソン分布: 稀小現象のある基準時間内に起こる回数

↓

ある稀小現象の起こるまでの時間(待ち時間)をとする (ポアソン分布で  $\lambda \rightarrow \lambda t$ )

= つまり 1日とも現象が起きていない:  $x=0$

累積確率密度関数  $F(t) = P(X \leq t) = 1 - e^{-\lambda t}$

指数分布:  $f(t) = \frac{dF(t)}{dt} = \lambda e^{-\lambda t}$

パラメータ:  $\lambda$   $E(X) = \frac{1}{\lambda}$ ,  $\text{Var}(X) = \frac{1}{\lambda^2}$

#### (2) 正規分布

二項分布の  $n \rightarrow \infty$  の極限分布

正規分布:  $f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

パラメータ:  $\mu, \sigma^2 \rightarrow X \sim N(\mu, \sigma^2)$   $E(X) = \mu$ ,  $\text{Var}(X) = \sigma^2$

標準化  $Z = \frac{X-\mu}{\sigma} \sim N(0, 1)$  - 標準正規分布

$E(X) = E(\mu + \sigma Z) = \mu + \sigma E(Z) = \mu$ ,  $\text{Var}(X) = \text{Var}(\mu + \sigma Z) = \sigma^2 \text{Var}(Z) = \sigma^2$

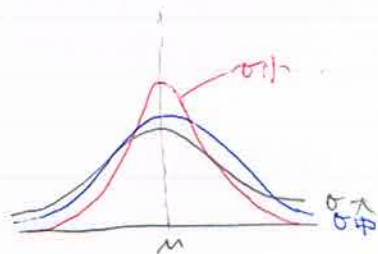
μを中心に対称

対称性から

$$\textcircled{1} P(\mu+a \leq X) = P(X \leq \mu-a)$$

$$\textcircled{2} P(X \leq \mu+a) = P(\mu-a \leq X)$$

$$\textcircled{3} P(\mu-a \leq X \leq \mu) = P(\mu \leq X \leq \mu+a)$$



$$Y = a + bX \sim N(a + b\mu, b^2\sigma^2) \quad (T.F.L \quad X \sim N(\mu, \sigma^2))$$

$$\therefore X = \frac{Y-a}{b}, \quad \left| \frac{dy}{dx} \right| = \frac{1}{|b|}$$

$$f_Y(y) = f_X\left(\frac{y-a}{b}\right) = \frac{1}{\sqrt{2\pi}|b|\sigma} \exp\left\{-\frac{(y-a-b\mu)^2}{2b^2\sigma^2}\right\}$$

$$\cdot \text{いま } P(a \leq X \leq b) \text{ を求めたい} \dots P(b) - P(a)$$

$$\text{標準正規分布を利用} \quad \frac{a-\mu}{\sigma} \leq \frac{X-\mu}{\sigma} \leq \frac{b-\mu}{\sigma} \quad Z \sim N(0, 1)$$

$$f(z) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{z^2}{2}\right\}$$

$$\text{よって } Z \text{ から与えられる } Y \text{ は } X = \sigma Z + \mu$$

$$\begin{cases} Z = \pm 1 \rightarrow X = \mu \pm \sigma & \text{"1}\sigma\text{範囲"} & P(\mu - \sigma \leq X \leq \mu + \sigma) = 0.68 \left( \approx \frac{2}{3} \right) \\ Z = \pm 2 \rightarrow X = \mu \pm 2\sigma & \text{"2}\sigma\text{"} & P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) = 0.95 \\ Z = \pm 3 \rightarrow X = \mu \pm 3\sigma & \text{"3}\sigma\text{"} & P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) = 0.99 \end{cases}$$

cf). 偏差値:  $\mu = 50, \sigma^2 = 10^2$  に調整したもの.

### (3) 多変量正規分布

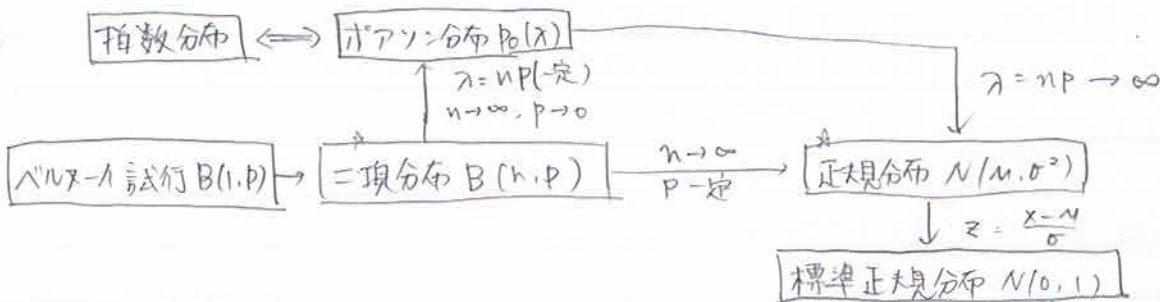
$n$  次元のデータ

$$f(x) = \frac{1}{(\sqrt{2\pi})^n \sqrt{|\Sigma|}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\} \sim N(\mu, \Sigma)$$

$$x = (x_1, x_2, \dots, x_n)^T, \quad \mu = (\mu_1, \mu_2, \dots, \mu_n)^T: \text{平均ベクトル}$$

$$\Sigma = \begin{pmatrix} \sigma_{11} & \dots & \sigma_{1n} \\ \vdots & & \vdots \\ \sigma_{n1} & \dots & \sigma_{nn} \end{pmatrix}: \text{分散共分散行列}$$

### 3.3 主要な確率分布



4/30

§4. 大数の法則と中心極限定理

4.1 チェビシェフの不等式

確率分布の期待値と分散しかわかっていないとき、確率変数のとりうる値(区間)に対する確率の値の見当をつけたい。

確率分布に依存しない

チェビシェフ・不等式  $P(|X - \mu| \leq k\sigma) \geq 1 - \frac{1}{k^2}$

下限値を見積もる

ex)  $k=2$   $P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \geq 1 - \frac{1}{2^2} = \frac{3}{4} = 0.75$

$k=3$   $P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \geq 1 - \frac{1}{3^2} = \frac{8}{9} = 0.89$

証明 (X連続なとき)

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx \geq \int_{|x - \mu| \geq k\sigma} (x - \mu)^2 f(x) dx \geq \int_{|x - \mu| \geq k\sigma} (k\sigma)^2 f(x) dx = (k\sigma)^2 P(|X - \mu| \geq k\sigma)$$

$\therefore \frac{1}{k^2} \geq P(|X - \mu| \geq k\sigma)$  余事として、 $P(|X - \mu| \leq k\sigma) \geq 1 - \frac{1}{k^2}$

意義

- ①  $\mu$  と  $\sigma$  によって、Xの分布の状態が示される。
- \*② Xがどのような確率分布であるかと成立し、 $|X - \mu| \leq k\sigma$ の確率の下限と与える。
- ③ 大数の(弱)法則を導くことができる。

標準母集団を推定するとき利用

4.2 大数の法則

$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$   $k\sigma = c$  とおくと、 $P(|X - \mu| \geq c) \leq \frac{\sigma^2}{c^2}$

Xの標本  $\{x_1, \dots, x_n\}$  があり、この平均  $\bar{x}$  を考える (計算平均)

$E(X) = \mu$ ,  $\text{var}(X) = \sigma^2$  とき、 $E(\bar{x}) = \mu$ ,  $\text{var}(\bar{x}) = \frac{\sigma^2}{n}$  となる。

$P(|\bar{x} - \mu| \geq c) \leq \frac{(\sigma^2/n)}{c^2} \xrightarrow{n \rightarrow \infty} 0$

大数の法則  $\lim_{n \rightarrow \infty} P(|\bar{x} - \mu| \geq c) = 0$

$\bar{x}$  は  $\mu$  に確率収束する。

- $\Rightarrow$  十分な大きさの標本を調べれば、母集団の特性を知ることができる。
- $\Rightarrow$  統計的推測の理論。



直接試験には出にくい  
理解する

NO.

9

DATE

平均したものを

Cauchy分布など

和をすれば正規分布を仮定できる

$E(X)$ と $\text{Var}(X)$ の分らない分布  
には適用できない

### 4.3. 中心極限定理

$X$  の標本  $\{X_1, \dots, X_n\}$  があり.  $E(X) = \mu$ ,  $\text{Var}(X) = \sigma^2$

$$\lim_{n \rightarrow \infty} \sum X_i \sim N(n\mu, n\sigma^2)$$

厳密に表す.  $\lim_{n \rightarrow \infty} P(a \leq \frac{\sum X_i - n\mu}{\sqrt{n}\sigma} \leq b) = \int_a^b \frac{1}{\sqrt{2\pi}} \exp(-\frac{x^2}{2}) dx$

(標準化)

和の期待値は期待値の和

証明

$E(\sum X_i) = n\mu$ ,  $\text{Var}(\sum X_i) = n\sigma^2$  標準化する. ( $Z_i = \frac{X_i - \mu}{\sigma} \rightarrow E(Z_i) = 0, \text{Var}(Z_i) = 1$ )

$$\frac{\sum X_i - n\mu}{\sqrt{n}\sigma} = \frac{\sum Z_i}{\sqrt{n}}$$

∴  $e^x$  の展開式  $e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$   $x$  に  $tX$  を代入.

$e^{tX} = 1 + tX + \frac{(tX)^2}{2!} + \frac{(tX)^3}{3!} + \dots$  両辺の期待値をとる.

$M_X(t) = E(e^{tX})$ : モーメント母関数. ← 左辺

$$= 1 + tE(X) + \frac{t^2 E(X^2)}{2!} + \dots$$

(標本を考慮すると)  
 $X_1$  と  $X_2$  は独立

$E(Z_i) = 0$ ,  $E(Z_i^2) = \text{Var}(Z_i) - \{E(Z_i)\}^2 = 1$  だから.

$$M_{Z_i}(t) = 1 + \frac{t^2}{2} + \dots$$

$$\begin{aligned} \text{また, } M_{X_1+X_2}(t) &= E(e^{t(X_1+X_2)}) = E(e^{tX_1} \cdot e^{tX_2}) = E(e^{tX_1}) E(e^{tX_2}) \\ &= M_{X_1}(t) \cdot M_{X_2}(t) \end{aligned}$$

$$\therefore M_{\sum Z_i}(t) = \{M_{Z_i}(t)\}^n = (1 + \frac{t^2}{2} + \dots)^n$$

$t \rightarrow \frac{t}{\sqrt{n}}$  に置きかえる.

$$M_{\sum Z_i}(\frac{t}{\sqrt{n}}) = \{M_{Z_i}(\frac{t}{\sqrt{n}})\}^n = (1 + \frac{t^2}{2n} + \dots)^n = \exp(\frac{t^2}{2})$$

標準正規分布のモーメント母関数

∴ 分布とモーメント母関数は1対1に対応する

∴  $\frac{\sum Z_i}{\sqrt{n}}$  は標準正規分布

⇒  $\sum X_i \sim N(n\mu, n\sigma^2)$  となる

$$(X \sim N(\mu, \sigma) \rightarrow M_X(t) = \exp\{\mu t + \frac{\sigma^2 t^2}{2}\})$$

$eX$  = 二項分布 → 正規分布

$$B(n, p) : f_X = n(x)^n (1-p)^{n-x}$$

$$X \sim B(1, p)$$

成功回数:  $\sum X_i$ ,  $E(\sum X_i) = np$ ,  $\text{Var}(\sum X_i) = np(1-p)$  ← 前回やった.

中心極限定理より.  $\lim_{n \rightarrow \infty} \frac{\sum X_i - np}{\sqrt{np(1-p)}} \sim N(0, 1)$

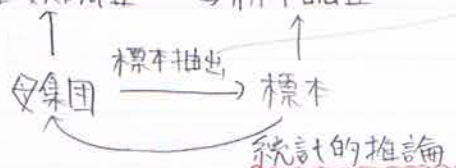
↑ これは統計的推測の準備

## §5 標本分布

## 5.1 母集団と標本 (cf 1.2)

## (1) 母集団と標本 (総括調査)

調査方法 ① 全数調査 ② 標本調査 無作為抽出

母集団: 母集団分布, 母数 (パラメータ) (母平均  $\mu$ , 母分散  $\sigma^2$  など)

標本: 標本分布, 統計量 (標本平均, 標本分散 など)

母数推定, 重要な手がかり (cf. 大数の法則)

標本を要約し, 母数の推定に使われるもの

標本  $(x_1, \dots, x_n)$  の関数  $t(x_1, \dots, x_n)$ 

未知パラメータは含まれない

観測値  $(x_1, \dots, x_n)$  の  $t(x_1, \dots, x_n)$ : 統計値 (統計量は計算方法)

5/13

## 12) 統計量

代表的な母数:  $\mu, \sigma^2$  → 標本平均  $\bar{x}$ , 標本分散  $s^2$ 以下, 標本  $(x_1, \dots, x_n)$  は 母集団分布  $(\mu, \sigma^2)$  に従う独立な確率変数

## ○ 標本平均

$$\bar{x} = \frac{x_1 + \dots + x_n}{n}, \quad E(\bar{x}) = E\left(\frac{x_1 + \dots + x_n}{n}\right) = \frac{n\mu}{n} = \mu$$

$$\text{var}(\bar{x}) = \text{var}\left(\frac{x_1 + \dots + x_n}{n}\right) = \frac{1}{n^2} \text{var}(x_1 + \dots + x_n) = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$$

## ○ 標本分散

$$s^2 = \frac{1}{n-1} \{ (x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \}, \quad E(s^2) = \sigma^2$$

期待値が母数のものと, 不偏推定量という  
 $\therefore s^2$ : 不偏分散

 $n-1$ : 自由度 (自由に動かせる変数の数)

$$(x_1 - \bar{x}) + (x_2 - \bar{x}) + \dots + (x_n - \bar{x}) = 0 \text{ より, } x_n = g(x_1, \dots, x_{n-1})$$

よって自由度は  $n-1$

## (3) 標本比率

2値確率変数  $\{0, 1\}$  の場合で、1 とした回数が  $X$  回であったとする。

→ 二項分布 (cf. 3.1.(2))

標本比率 (確率):  $p = \frac{x}{n}$  を考える

$$E(p) = E\left(\frac{X}{n}\right) = \frac{1}{n} E(X) = \frac{np}{n} = p$$

$$\text{var}(p) = \text{var}\left(\frac{X}{n}\right) = \frac{1}{n^2} \text{var}(X) = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}$$

★

## (4) 統計量の標本分布

以下、母集団が正規分布に従うとして、次の分布を扱う。

- ・ 標本平均 ( $\sigma^2$  が既知) の分布 → 正規分布
- ・ 標本分散の分布 →  $\chi^2$  分布 ★  $\sigma^2$  は大抵未知 (だから推定する)
- ・  $\sigma^2$  が未知ときの標本平均の分布 →  $t$  分布
- ・ 標本分散比の分布 →  $F$  分布

以下これらを詳しくみていく。

5.2. 標本平均の分布 ( $\sigma^2$  既知)

$$X \sim N(\mu, \sigma^2)$$

- ・ 正規分布・再生性について

$$C_i \text{ を定数列とすれば } \sum_i C_i X_i \sim N\left(\sum C_i \mu, \sum (C_i \sigma)^2\right)$$

(cf. 3.2.(2) の  $Y = aX + b$  の一般化)

✓ (推定値, 精度) がある

- ・  $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$  ★  $\sigma^2$  既知なら正規分布

( $\bar{X}$  は  $X_1, \dots, X_n$  の標本平均)

- ・ ここでの分散は精度を表す →  $n$  が増加すれば、 $\bar{X}$  は  $\mu$  より正確な推定値。
- ・ 推定精度 (分散) は、標本数に対して  $1/n$  のオーダーでの減少になる。
- ・  $\bar{X}$  も正規分布 → 単独の  $X_i$  よりもすぐれた推定値

✓ 標本平均の  
精度は高い

## 5.3. 標本分散の分布

$\sigma^2$  は不明の場合が多い → 分散の推定値の精度を知りたい。

$$S^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2 \quad \therefore (n-1) S^2 = \sum (X_i - \bar{X})^2 \quad \dots (*)$$

ここで

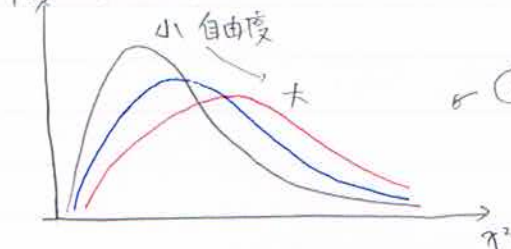


•  $\chi^2$  分布について

標準化  $z_i = \frac{x_i - \mu}{\sigma} \sim N(0, 1)$  について考える.  $\chi^2 = z_1^2 + z_2^2 + \dots + z_k^2$  とする.

確率変数  $\chi^2 \sim$  自由度  $k$  の  $\chi^2$  分布  $\chi^2(k)$  (と定義),  $E(\chi^2) = k$ ,  $\text{var}(\chi^2) = 2k$

pdf と 確率密度関数



$$E(z_i^2) = 1 + 1, \quad f(x) = \frac{1}{2^{\frac{k}{2}} \Gamma(\frac{k}{2})} x^{\frac{k}{2}-1} e^{-\frac{x}{2}} \quad (x > 0)$$

$$(P(z) = \int_0^\infty e^{-x} x^{k-1} dx \quad (z > 0))$$

自由度が大きいほど

自由度が大きいほど

$$\chi^2 = \sum z_i^2 = \sum \frac{(x_i - \mu)^2}{\sigma^2} \quad (*) \text{ 対し, } \frac{(n-1)S^2}{\sigma^2} = \sum \frac{(x_i - \bar{x})^2}{\sigma^2} = \chi^2$$

$$\therefore \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$$

$\mu \rightarrow \bar{x}$  におきかえれば,  $(x_1 - \bar{x}) + \dots + (x_n - \bar{x})$  の自由度 1 減る:  $n-1$

$$\therefore E(S^2) = \frac{\sigma^2}{n-1} E(\chi^2) = \frac{\sigma^2}{n-1} \cdot (n-1) = \sigma^2 \quad \text{不偏推定量} \Rightarrow \text{母分散の推定に利用}$$

※

5.4  $\sigma^2$  未知のときの標本平均の分布

自然対数

$\sim N(0, 1)$

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n}) \quad \text{標準化 } z_i = \frac{\bar{x} - \mu}{\sqrt{\sigma^2/n}}$$

$\sigma^2$  は不明の場合...  $\rightarrow$  標本分散  $S^2$  で代替したのが現実的

• t 統計量

$$t = \frac{\bar{x} - \mu}{\sqrt{S^2/n}} \quad \text{※ } t \text{ は } N(0, 1) \text{ に似ていない}$$

• 標本平均  $\times$  標本標準偏差の比の分布

$$t = \frac{\bar{x} - \mu}{\sqrt{S^2/n}} = \frac{\frac{\bar{x} - \mu}{\sqrt{\sigma^2/n}}}{\sqrt{\frac{(n-1)S^2}{\sigma^2} / (n-1)}}$$

③:  $\bar{X}$  と  $S^2$  は独立

$\chi^2(n-1)$  ... ②

ここで

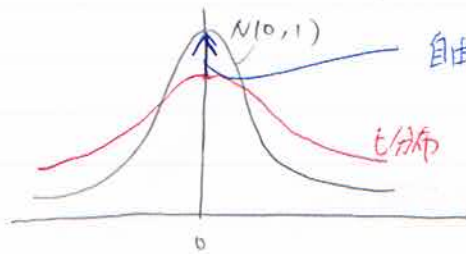
t 分布

一般に ①  $Z \sim N(0, 1)$  ②  $Y \sim \chi^2(k)$  ③  $Z, Y$  は独立 とき

$$t = \frac{Z}{\sqrt{Y/k}} \sim \text{自由度 } k \text{ の } t \text{ 分布 } t(k) \quad E(t) = 0 \quad (k > 1), \quad \text{var}(t) = \frac{k}{k-2} \quad (k > 2)$$



$$\text{pdf: } f(t) = \frac{\Gamma\left(\frac{k+1}{2}\right)}{\sqrt{k\pi} \Gamma\left(\frac{k}{2}\right)} \cdot \frac{1}{\left\{\frac{k}{2} + 1\right\}^{\frac{k+1}{2}}} \quad (-\infty < t < \infty) \quad \triangle \text{ だーでーもい...}$$



自由度  $\rightarrow \infty$  で  $t$  分布は  $N(0,1)$  に近づく。

$$\therefore n \rightarrow \infty \Rightarrow S^2 \hat{=} \sigma^2$$

(外れた値をとり  
確率が大きめ)

$E(t) = 0$  で対称、分布、両方が  $N(0,1)$  より広い。

★  $t$  分布は小標本の厳密な標本分布

以上より、
$$t = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t(n-1)$$

★ 標本平均の分布のまとめ

母集団		正規分布	非正規分布
$\sigma^2$	$n$		
既知	大	①	⑤
既知	小	②	⑥
未知	大	③	⑦
未知	小	④	⑧

①, ②  $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ ,  $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$

⑤  $\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$ ,  $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \approx N(0,1)$  ( $\approx$ : 中心極限定理による)

③, ⑦  $\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$ ,  $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \approx N(0,1)$

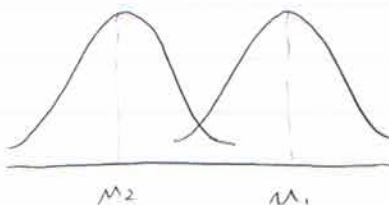
④  $t = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t(n-1)$

⑥, ⑧  $\rightarrow \bar{X}$  の分布は分からない

5.2 おり

### 3.5. 2標本問題

F分布



第1標本  $\{X_1, \dots, X_{n_1}\} \sim N(\mu_1, \sigma_1^2)$

第2標本  $\{Y_1, \dots, Y_{n_2}\} \sim N(\mu_2, \sigma_2^2)$

$X_1, \dots, X_{n_1}, Y_1, \dots, Y_{n_2}$  は独立

標本分散比の分布

$$s_1^2 = \frac{1}{n_1-1} \sum (x_i - \bar{x})^2, s_2^2 = \frac{1}{n_2-1} \sum (y_j - \bar{y})^2$$

F統計量  $F = \frac{s_1^2/s_1^2}{s_2^2/s_2^2}$

$$F = \frac{\frac{(n_1-1)s_1^2}{s_1^2}}{\frac{(n_2-1)s_2^2}{s_2^2}} = \frac{\chi^2(n_1-1) \textcircled{1}}{\chi^2(n_2-1) \textcircled{2}}$$

③:  $s_1^2, s_2^2$  は独立.

ここで、

F分布

一般に ①  $U \sim \chi^2(k_1)$  ②  $V \sim \chi^2(k_2)$  ③  $U, V$  は独立 のとき

$F = \frac{U/k_1}{V/k_2} \sim$  自由度  $(k_1, k_2)$  の F分布:  $F(k_1, k_2)$

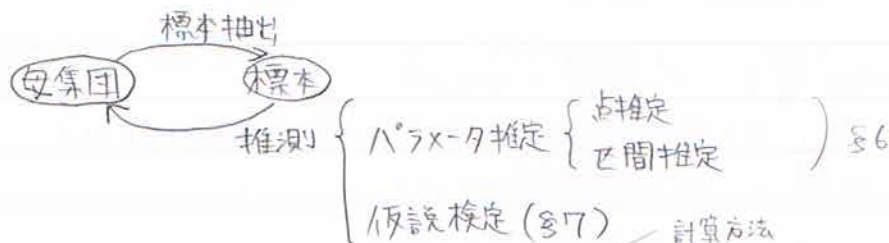
$$E(F) = \frac{k_2}{k_2-2} \quad (k_2 > 2), \quad \text{Var}(F) = \frac{2k_2^2(k_1+k_2-2)}{k_1(k_2-2)^2(k_2-4)} \quad (k_2 > 4)$$

$E(F)$  と  $\text{Var}(F)$  は参考で.

以上より、 $F = \frac{s_1^2/s_1^2}{s_2^2/s_2^2} \sim F(n_1-1, n_2-1)$

5/28

§6 パラメータの推定



パラメータ推定のための標本から求めた統計量: 推定量 - 確率変数

計算値 → 統計値: 推定値 - 実現値

点推定: パラメータを1つの値で定める方法

区間推定: パラメータが存在する区間を推定する方法

↑  
仮説検定

一般にパラメータ

## 6.1 点推定・推定法

推定量には、 $\wedge$  (ハット) をつける。 ex)  $\hat{\mu}$ ,  $\hat{\sigma}^2$

## (1) モーメント法

$k$  次モーメント  $= E(X^k)$  により推定する方法

期待値:  $\mu = E(X)$ ,  $\bar{x} = \frac{1}{n} \sum x_i$  により推定  $\rightarrow \hat{\mu} = \bar{x} = \frac{1}{n} \sum x_i$

分散:  $\sigma^2 = E(X^2) - \{E(X)\}^2$ ,  $\frac{1}{n} \sum x_i^2$  により推定  $\rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum x_i^2 - (\bar{x})^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$

\*  $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$ : 不偏分散より,  $\hat{\sigma}^2$  は真の分散を過小推定している。

い、 $\hat{\sigma}^2$  は不偏ではない。

\* この手法は確率分布によるもの。

★ (1) より精度↑

## (2) 最尤法

母分布を必要とする。パラメータ  $\theta$  (定数) と表現する。

同時確率分布  $f(x_1, x_2, \dots, x_n; \theta)$

(パラメータ  $\theta$  固定したときに  $x_1, \dots, x_n$  の値をとる確率)

ここで、

$(x_1, \dots, x_n)$ : 観測ごとに異なる値をとる確率変数

いす、1組の観測値  $(x_1, \dots, x_n)$  が得られた (所与)

↓

$(x_1, \dots, x_n)$  はもう確率変数ではなく、その実現値、

$(x_1, \dots, x_n)$  をもたせし  $\theta$  は様々な値が考えられる。

likelihood  $\Rightarrow \theta$  を変数と考えると、

$L(\theta; x_1, \dots, x_n)$  or 単に  $L(\theta)$

↑ 変数 定数

得られた観測値  $(x_1, \dots, x_n)$  の確率  $\propto \theta$  によってどう変化するか示す関数

$\theta = \theta_0$  のとき (5% 同値をとる確率分布ではない、値では大小が重要になる)

$L(\theta_0; x_1, \dots, x_n) = f(x_1, \dots, x_n; \theta_0)$

$L(\theta)$  の大小関係が重要

同じ観測結果に対し  $L(\theta_1) > L(\theta_2)$  なら、 $\theta = \theta_1$  から得られた標本を認める方が

も、ともし。

maximum likelihood: ML

↓

最尤法: 尤度関数を最大にする  $\theta$  を  $\hat{\theta}$  とする。

最尤推定量 (推定値)

$L(\theta)$  の最大化  $\frac{dL(\theta)}{d\theta} = 0$

パラメータ  $\theta_1, \dots, \theta_k$  の場合  $\frac{\partial L(\theta)}{\partial \theta_j} = 0 \quad (j=1, \dots, k) \rightarrow$  連立方程式

あるいは (計算方法を77=13)

対数変換: 単調増加変数

$$\max L(\theta) \rightarrow \max \log L(\theta)$$

対数尤度関数を用いてもよい

ex) 正規分布

$\theta_1, \theta_2: \mu, \sigma^2 \quad X_i \sim N(\mu, \sigma^2), (x_1, \dots, x_n)$  所与

$$L(\mu, \sigma^2) = \prod_{i=1}^n \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(x_i - \mu)^2}{2\sigma^2} \right\} \right]$$

確率密度関数  $X_i \sim \tau_i$  実現値

$$= (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\}$$

(x1, ..., xn) の実現する確率

$$\log L(\mu, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

(対数尤度の計算にもなる)

$$\frac{\partial \log L(\mu, \sigma^2)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0$$

$$\frac{\partial \log L(\mu, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

(σ²で1階微分)

$$\therefore \hat{\mu} = \bar{x}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

← モーメント法と同じ (正規分布のときは)

## 6.2 望ましい推定量の特性

μ の推定量の候補: 標本平均, メディアン, モード...

σ², S² の相違

⇒ 推定量の望ましい特性が必要

小標本特性: n が大ききに関わりなくも、ているべき特性

(不偏性, 有効性, 最良線形不偏性)

大標本特性: n が大きくな、ていくとも、ているべき特性 (漸近的特性)

(漸近的不偏性, 一致性, 漸近的有效性)

### (1) 不偏性

θ が 過大・過小推定でない  $\therefore E(\hat{\theta}) = \theta$  — 不偏推定量

ex)  $E(\bar{x}) = \mu$  (cf 5-1(2))

$$E(\hat{\sigma}^2) = \frac{n-1}{n} \sigma^2 = \sigma^2 - \frac{1}{n} \sigma^2$$

← 不偏推定量じゃない

過小推定





ラグランジュの未定乗数法

$$\phi = \sum w_i^2 - 2\lambda (\sum w_i - 1)$$
  
ラグランジュ関数  
都合上  $\lambda$  で  $2\lambda$  とした

$$\begin{cases} \frac{\partial \phi}{\partial w_i} = 2w_i - 2\lambda = 0 & (\text{for } \forall i) \\ \frac{\partial \phi}{\partial \lambda} = -2(\sum w_i - 1) = 0 \end{cases} \Rightarrow w_i = \lambda, \lambda = \frac{1}{n}$$

$\therefore w_1 = w_2 = \dots = \frac{1}{n} \Rightarrow \text{よって } \bar{X} \text{ は BLUE}$

大標本特性

(1) 一貫性

大きさ  $n$  の標本からの推定量:  $\hat{\theta}_n$

$$\forall \epsilon > 0 \text{ に対して, } \lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| \geq \epsilon) = 0 \quad \text{: 一致推定量}$$



(十分条件:  $\lim_{n \rightarrow \infty} E(\hat{\theta}_n) = \theta$  (漸近的 unbiasedness)  $\wedge \lim_{n \rightarrow \infty} \text{var}(\hat{\theta}_n) = 0$ )

e.g.  $E(\bar{X}) = \mu, \text{var}(\bar{X}) = \frac{\sigma^2}{n} \xrightarrow{n \rightarrow \infty} 0$   $\therefore \bar{X}$  は  $\mu$  の一致推定量

$\lim_{n \rightarrow \infty} (\hat{\sigma}^2) = \sigma^2, \lim_{n \rightarrow \infty} \text{var}(\hat{\sigma}^2) = \lim_{n \rightarrow \infty} \frac{2(n-1)\sigma^4}{n^2} = 0$

$\therefore \hat{\sigma}^2$  は  $\sigma^2$  の一致推定量,  $S^2$  も  $\sigma^2$  の一致推定量

6/4 仮説検定とセクトに33と合のりやろ

6.3. 区間推定

パラメータが存在すると予想される区間を確率的に推定



「 $\theta$  が区間  $[a_1, a_2]$  に存在する確率が  $1-\alpha$ 」という推定

$$\therefore P(a_1 \leq \theta \leq a_2) = 1-\alpha$$

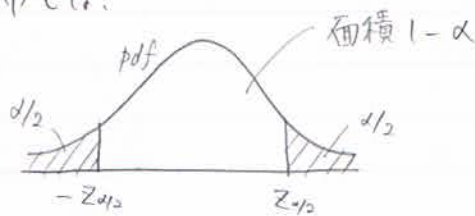
信頼係数 (有意水準)  $\checkmark$  通常 0.95, 0.99 が多い  
 $\uparrow \quad \uparrow$   
2σ      3σ

区間  $[a_1, a_2] = 100(1-\alpha)\%$  信頼区間

$\rightarrow \theta$  が区間内に含まれる割合が  $1-\alpha$

$\alpha$   $\left( \begin{array}{l} * \theta_1, \theta_2 \text{ は実験で毎回変わる} \\ \text{その区間に含まれる確率が } 1-\alpha \\ \text{(定数区間に入る確率じゃない)} \end{array} \right)$

標準正規分布では、



( $Z_{\alpha/2}$ :  $Z$  の点より上側の確率が  $100 \times \alpha/2$  % となる  $Z$  = パーセント点)

(1) 平均に関する推定

① 標本平均の分布 (付 5.4: 標本平均の分布のまとめの表)

①: ( $\sigma^2$  既知, 母集団分布: 正規分布,  $n$  = 大)

②: ( " " " "  $n$  = 小) では、

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

・ 区間推定

$$P(-Z_{\alpha/2} \leq Z \leq Z_{\alpha/2}) = 1 - \alpha$$

限界誤差

$$P(\bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

上下限の差: 区間幅  $2Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

区間幅が小さいほど  $\bar{X}$  の推定精度が上がる

・  $\sigma$  が小さいほど

・  $n$  が大きいほど (一限界誤差  $D = Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$  が所与ならば  $n = (\frac{Z_{\alpha/2} \sigma}{D})^2$ )

・  $Z_{\alpha/2}$  が小さいほど (信頼度は低下)

③, ④ は  $\sigma \rightarrow S$  だけ。

④  $t$  分布  $\sigma \rightarrow S, Z_{\alpha/2} \rightarrow t_{\alpha/2}(n-1)$

⑦. ⑧ ex)  $X_1, \dots, X_5 \sim N(\mu, 0.2^2)$

$\bar{x} = 19.5$  を観測  $\mu$  の 99% 信頼区間?

$\alpha = 0.01, Z_{\alpha/2} = 2.576, n = 5, \sigma = 0.2$

$$\bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = 19.5 \pm 2.576 \cdot \frac{0.2}{\sqrt{5}} = 19.5 \pm 0.2304$$

95% 信頼区間  
 $Z_{\alpha/2} = 1.96$

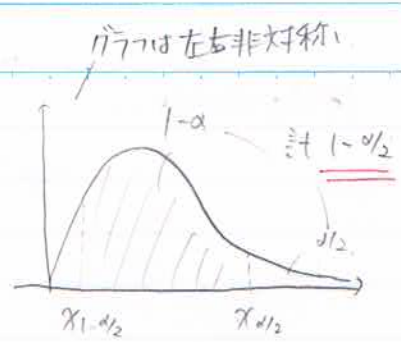
(2) 分散に関する推定 (正規分布)

$$\chi^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} = \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1) \text{ とおく}$$

$$\chi^2_{1-\alpha/2}(n-1) \text{ と } \chi^2_{\alpha/2}(n-1) \text{ とおく}$$

$$P(\chi^2_{1-\alpha/2} \leq \chi^2 \leq \chi^2_{\alpha/2}) = 1-\alpha$$

② 右図



$$\therefore P\left(\frac{(n-1)S^2}{\chi^2_{\alpha/2}} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi^2_{1-\alpha/2}}\right) = 1-\alpha$$

(3) 比率の区間推定 (カイ二乗)

結果が 2 事象の n 個の標本のうち、着目する事象の標本が k 個であったとする

標本比率  $\hat{p} = \frac{k}{n}$

標本数 n が十分大きい場合 (二項 → 正規)  $\hat{p} \sim N\left(p, \frac{p(1-p)}{n}\right)$

標準化  $z = \frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \sim N(0, 1)$

分母の p を  $\hat{p}$  で近似すれば、母比率 p の区間推定

$$P\left(\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right) = 1-\alpha$$

大抵

87 仮説検定

9.1 検定の考え方

母集団について仮定した命題を標本に基づいて検証

母集団 → 標本

- 仮説 H
- ・ H が著しくはズれる → H を否定する (棄却する)
  - ・ H が著しくはズれない → H を認める (採択する)
- 有意水準  $\alpha$  により判断

例) 仮説: 「地球温暖化が進行」

観測: 「18年間連続で異常高温が続いている」

ここでの異常: 「たかだか 30 年に 1 度の頻度」

対立する仮説: 「地球平均気温は定常」

→ ある年に「異常気温」が生じる確率  $1/30$

各年の気温は独立と仮定



⇒  $p = 1/30$ ,  $n = 18$  → 二項分布

$$f(x) = n C_x p^x (1-p)^{n-x} = {}_{18}C_8 (1/30)^8 (1-1/30)^{10} = 2.58 \times 10^{-27}$$

解釈

- ① 観測対象は稀ではなく、誤った前提により稀な確率が計算された。
- ② 非常に稀だが全く起こらないわけではない事実が起こった。

↑ ex) に当て

## 仮定の有意性の検証 = 仮説検定

検定の流れ

- (1) 仮説の設定 <sup>(1) ①</sup> この点に注意
- (2) 検定統計量とその分布の決定
- (3) 有意水準と棄却域の決定
- (4) 検定の実施

新しいもの

ex) 母分散既知の正規分布からの小標本に対する検定 (cf 5.4 ②) <sup>表の</sup>

新製品の寿命のメーカー公称値: 9hour

旧製品:  $\mu = 8\text{hour}$ ,  $\sigma = 1\text{hour}$  ~ 正規分布

(新製品の旧製品の  
寿命を比較する)

新製品で10回実験 →  $n = 10$ ,  $\bar{x} = 8.8$

(1) 仮説の設定

知りたいことと仮定

新製品も旧製品  $\mu = 8$  と同じではないか? → これを検証

$H_0: \mu = 8$  帰無仮説

データが  $H_0$  から著しく外れたときのみ  $H_0$  を棄却

$H_1: \mu > 8$  対立仮説

	$H_0$ が正しい	$H_0$ は誤り
$H_0$ を棄却しない	①	③
$H_0$ を棄却する	②	④

①, ④ は正しい検証

②  $H_0$  は正しい →  $H_0$  棄却: 第I種の過誤

③  $H_0$  は誤り →  $H_0$  採択: 第II種の過誤

☆ ~~正しいという目的~~ (証明済)  
☆ 一 別のこともあって間違っている  
検定 あるいは  
考へず

上記では、③お④の過誤を大きな誤りと考える。

↑ 背理法を考へ (棄却することの目的)

$H_0$  採採でも、積極的に支持したわけではなく「矛盾はない」だけ。

(2) 検定統計量とその分布の決定

$\mu$  の良い推定量:  $\bar{x}$  → 検定にも  $\bar{x}$  を用いる: 検定統計量

$n$ : 小,  $\sigma^2$  既知  $\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$  (of 6.3 (1))  
( $\sigma^2$  未知な分布)

4/1

④ 検定の流れ

- (1) 仮説の設定
- (2) 検定統計量とその分布の決定
- (3) 有意水準と棄却域の決定
- (4) 検定の実施

↓ 今のはこのへり

(1) 仮説の設定

帰無仮説  $H_0: \mu = 8$ , 対立仮説:  $\mu > 8$

これを棄却したい。

(2) 検定統計量とその分布

検定統計量:  $\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$  例の8.1.9-2で分類

(3) 有意水準と棄却域の決定

$\bar{x}$  の値が 8 を超え大きいほど、 $H_0: \mu = 8$  を棄却する証拠

$H_0$  を棄却すべき  $\bar{x}$  の領域: 棄却域  $R$ .

||

$\bar{x}$  がある値  $c$  を超える領域

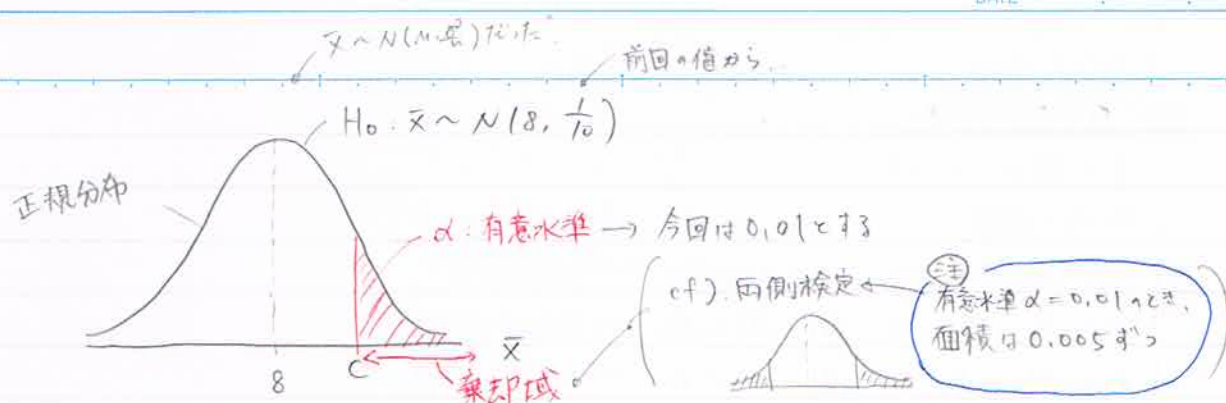
$R = \{ \bar{x}; \bar{x} > c \}$  ← 片側検定

( of  $R = \{ \bar{x}; |\bar{x}| > c \}$  : 両側検定 ( $H_1: \mu \neq 8$ ) )

$H_0$  を棄却 ← データが  $H_0$  から著しくはずれる. (0.05 くらいある)

それだけ小さな確率  $\alpha$  でしか生じない。

有意水準



前提  $\alpha$  は依存

$$P(\bar{x} > c | H_0 \text{ 正しい}) = 0.01 \quad \text{となる } c \text{ を計算}$$

$$Z = \frac{\bar{x} - 8}{1/\sqrt{10}} \sim N(0, 1) \quad \therefore P(Z > \frac{c-8}{1/\sqrt{10}}) = 0.01 \quad (Z > Z_{0.01})$$

$\left( \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \right)$  (p.15)

標準正規分布の右片側 0.01 の確率を考慮する。  $Z_{0.01} = 2.326$

$$Z = \frac{c-8}{1/\sqrt{10}} = 2.326 \quad \therefore c = 8.74$$

$\therefore$  棄却域  $R = \{\bar{x} ; \bar{x} > 8.74\}$   $\rightarrow$  棄却域が決まった!!

(4) 検定の実施  $\leftarrow$  実際のデータから行う

$$\begin{cases} \bar{x} \in R \Rightarrow H_0 \text{ 棄却 (} H_1 \text{ 採択)} \\ \bar{x} \notin R \Rightarrow H_0 \text{ 棄却 しない} \end{cases}$$

今回は、 $\bar{x} = 8.8 \in R \setminus \bar{x} ; \bar{x} > 8.74 \Rightarrow$  有意水準 1% で  $H_0: \mu = 8$  は棄却される。

$$P(Z = \frac{8.8-8}{1/\sqrt{10}}) = 0.0059 \quad \text{P値}$$

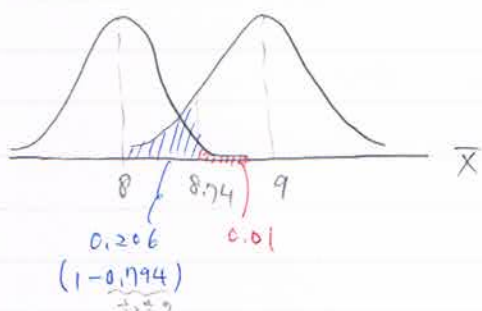
(おまけ) 前提を変えてみる。  $H_1$  を帰無仮説に小さくした。

$$P(\bar{x} > 8.74 | H_1 \text{ 正しい (例えば } \mu = 9 \text{ とする)}) = P(Z > \frac{8.74-9}{1/\sqrt{10}}) = 0.794$$

(-1 が出てきた)

。過誤について

第I種の過誤の確率



$$\begin{cases} P(I) = P(\bar{x} > 8.74 | \mu = 8) = 0.01 = \alpha \\ P(II) = P(\bar{x} \leq 8.74 | \mu = 9) = 0.206 = \beta \end{cases}$$

$\alpha$  を小さくしたい  $\Rightarrow \beta$  は大きくなる。

( $\alpha$  と  $\beta$  は同時に小さくできない)

$\therefore$  第I種の過誤の危険  $\rightarrow$  小

$\Rightarrow$  第II種の過誤の危険  $\rightarrow$  大

$\alpha$ : 有意水準 = 危険率

$1-\beta$ : 検出力

$H_0$  が偽っているときにちゃんと棄却できる確率。

## 7.2 母平均の検定

7.1 の通りにやってみよう

### (1) 仮説の設定

- (A)  $H_0: \mu = \mu_0, H_1: \mu \neq \mu_0 \rightarrow$  両側検定  
 (B)  $H_0: \mu = \mu_0, H_1: \mu > \mu_0 \rightarrow$  右片側検定  
 (C)  $H_0: \mu = \mu_0, H_1: \mu < \mu_0 \rightarrow$  左片側検定

### (2) 検定統計量と分布

1-1 p.13

5.4 「標本平均の分布のまとめ」表で分類

(a) ①, ②, ⑤  $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$  (③, ④では  $\sigma \rightarrow s$ )

(b) ④  $t = \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t(n-1) \leftarrow$  t検定 平均の検定

★

### (3) 棄却域

	a	b
A	$ Z  > Z_{\alpha/2}$	$ t  > t_{\alpha/2}(n-1)$
B	$Z > Z_{\alpha}$	$t > t_{\alpha}(n-1)$
C	$Z < -Z_{\alpha}$	$t < -t_{\alpha}(n-1)$

t分布は正規分布と同様に左右対称  
(ただし  $n$  が小さいと重い)

## 7.3 正規分布の母分散の検定

### (1) 仮説

- (A)  $H_0: \sigma^2 = \sigma_0^2, H_1: \sigma^2 \neq \sigma_0^2 \rightarrow$  両側検定  $\sigma_a^2 > \sigma_b^2$   
 (B)  $H_0: \sigma^2 = \sigma_0^2, H_1: \sigma^2 > \sigma_0^2$  or  $\sigma^2 = \sigma_a^2 \rightarrow$  右片側検定  
 (C)  $H_0: \sigma^2 = \sigma_0^2, H_1: \sigma^2 < \sigma_0^2$  or  $\sigma^2 = \sigma_b^2 \rightarrow$  左片側検定  $\sigma_b^2 < \sigma_a^2$

### (2) 検定統計量と分布

$\chi^2 = \sum \frac{(x_i - \bar{x})^2}{\sigma_0^2} = \frac{(n-1)s^2}{\sigma_0^2} \sim \chi^2(n-1) \leftarrow \chi^2$  検定.

$\chi^2$  のグラフは左右対称じゃないでこんな書き方になる

### (3) 棄却域 (cf. 6.3 (2))

A:  $\chi^2 < \chi^2_{1-\alpha/2}, \chi^2 > \chi^2_{\alpha/2}$

B:  $\chi^2 > \chi^2_{\alpha}$

C:  $\chi^2 < \chi^2_{1-\alpha}$



ex). 正規分布と仮定したものが本当に正規分布か。

\*  $\chi^2$ 検定は適合度検定にも用いられる。

\* 母分散の比の検定  $\rightarrow$  F検定

母平均の差の検定  $\rightarrow$  t検定

おまけ

## ★ §8. 回帰分析

空間情報学でも使う。

### 8.1. 回帰モデル

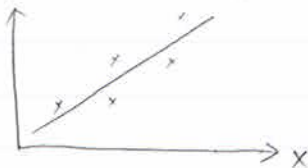
2変数:  $X, Y$

$Y$ の変動を $X$ で説明:  $X$ と $Y$ の定量的な関係(モデル)を知りたい。

$X$ : 説明変数, 独立変数  $\leftarrow$  横軸が多い

$Y$ : 被説明変数, 従属変数  $\leftarrow$  縦軸 "

ex).  $Y$ : ばねの長さ,  $X$ : おもりの重さ



$$Y = \beta_0 + \beta_1 X + \varepsilon$$

← 誤差項  
← 回帰式

回帰式が直線の時: 線形回帰

$\beta_0, \beta_1$ : 回帰係数,  $\varepsilon$ : 誤差項

非線形回帰モデル

③ にての線形は,

$\beta_0, \beta_1$ に拘りていない

$(Y = \beta_0 + \beta_1^2 X \leftarrow$  非線形)

$(Y = \beta_0 + \beta_1 X^2 \leftarrow$  線形)

説明変数1つ: 単回帰モデル

" 複数: 重回帰モデル

○ (線形)単回帰モデル  $\leftarrow$  最も簡単なもの

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad (\beta_0, \beta_1: \text{パラメータ}) \quad (i: \text{サンプル})$$

★ ★ ★ 仮定

仮定

<  $X$ への仮定 >

①  $X$ は非確率変数。

<  $\varepsilon$ への仮定 >

②  $E(\varepsilon_i) = 0 \quad (i = 1, 2, \dots, n)$

③ 分散一定:  $\text{var}(\varepsilon_i) = E(\varepsilon_i^2) = \sigma^2 \quad (i = 1, 2, \dots, n)$

④ 異なる誤差項は無相関(共分散0):  $\text{cov}(\varepsilon_i, \varepsilon_j) = E(\varepsilon_i \varepsilon_j) = 0 \quad (i \neq j, i, j = 1, 2, \dots, n)$

⑤  $\varepsilon_i$ は正規分布

誤差の諸要因:  $\varepsilon_{i1}, \varepsilon_{i2}, \dots \rightarrow \varepsilon_i = \varepsilon_{i1} + \varepsilon_{i2} + \dots$

中心極限定理より,  $\varepsilon_i$ は正規分布

② ~ ⑤ より,  $\varepsilon_i \sim N(0, \sigma^2)$

cf). コロケーション ( $X$ も確率変数)

所与のデータ

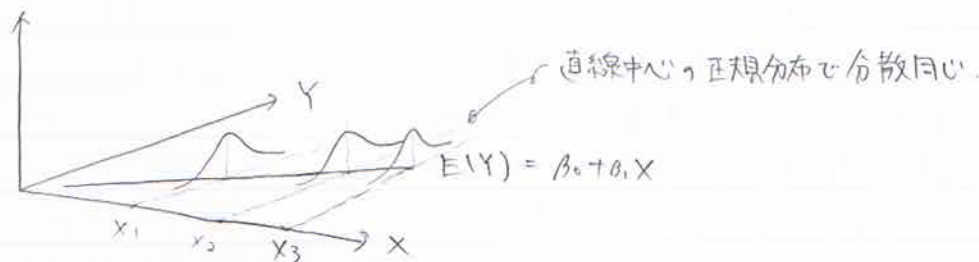
ここで、 $E(Y_i) = \beta_0 + \beta_1 X_i$  (①, ②より)

$\therefore Y_i = E(Y_i) + \varepsilon_i$  ←  $Y_i$  = 平均的な大きさ + 誤差

$$\text{Var}(Y_i) = E[\underbrace{Y_i - E(Y_i)}_{\varepsilon_i}]^2 = E(\varepsilon_i^2) = \text{Var}(\varepsilon_i) = \sigma^2 \quad (\text{③より})$$

$$\text{Cov}(Y_i, Y_j) = E[(Y_i - E(Y_i))(Y_j - E(Y_j))] = \varepsilon_i \varepsilon_j = 0 \quad (\text{④より})$$

以上から、仮定の1x-2j



・目的

① 未知パラメータ  $\beta_0, \beta_1, \sigma^2$  を推定

② モデルの説明力の確認

③ 推定量の統計的推測

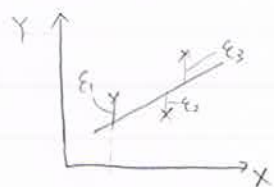
最小二乗法

→ 次回

4/18

## 8.2. 最小二乗法

(1) 回帰係数の推定



$$\varepsilon_i = Y_i - (\beta_0 + \beta_1 X_i) \quad \text{全部正しく}$$

誤差の二乗和を最小にしたい

cf). GLS

$$\sum \varepsilon_i^2 = \sum \{Y_i - (\beta_0 + \beta_1 X_i)\}^2 \rightarrow \min : \text{最小二乗法 (OLS)}$$

$$\begin{cases} \frac{\partial \sum \varepsilon_i^2}{\partial \beta_0} = -2 \sum (Y_i - \beta_0 - \beta_1 X_i) = 0 \\ \frac{\partial \sum \varepsilon_i^2}{\partial \beta_1} = -2 \sum (Y_i - \beta_0 - \beta_1 X_i) \cdot X_i = 0 \end{cases}$$

$$\Leftrightarrow \begin{cases} n\beta_0 + \beta_1 \sum X_i = \sum Y_i \\ \beta_0 \sum X_i + \beta_1 \sum X_i^2 = \sum X_i Y_i \end{cases} \quad \text{正規見直し式}$$

$$\bar{X} = \frac{1}{n} \sum X_i, \quad \bar{Y} = \frac{1}{n} \sum Y_i \text{ とおくと, } \beta_0 = \bar{Y} - \beta_1 \bar{X} \text{ とおいて代入}$$

$$\beta_1 (\sum X_i^2 - \bar{X} \sum X_i) = \sum X_i Y_i - \bar{Y} \sum X_i$$

$$\Leftrightarrow \beta_1 \sum (X_i - \bar{X})^2 = \sum (X_i - \bar{X})(Y_i - \bar{Y})$$

→  $(\bar{x}, \bar{y})$  は必ず通る

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

最小二乗推定量

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i \quad \checkmark \quad \text{データ } x_i \text{ から } Y_i \text{ を推測する式}$$

$$e_i = Y_i - \hat{Y}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i \quad \text{残差} \neq \text{誤差 } \varepsilon_i$$

$$\sum e_i = 0, \quad \sum e_i X_i = 0$$

← (多次元  $e_i \cdot X_i = 0$ )

$$E(\hat{\beta}_0) = \beta_0, \quad E(\hat{\beta}_1) = \beta_1 \rightarrow \text{不偏推定量}$$

理由は省略

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum X_i^2}{n \sum (x_i - \bar{x})^2}, \quad \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2} \quad \leftarrow \text{最小分散 (有効性)}$$

※ 最小二乗推定量は BLUE (cf. 6.2) (計算過程 19)

P.25 で仮定  $\varepsilon_i \sim N(0, \sigma^2)$  誤差の分散を残差の分散から推定する

## (2) 分散の推定

$$\sigma^2 \text{ の推定量 } s^2 = \frac{\sum e_i^2}{n-2} \quad \leftarrow \text{パラメータ } \beta_0, \beta_1 \text{ の2つより自由度は } n-2$$

$$\textcircled{1} \sum e_i = \sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) = 0 \quad \leftarrow \frac{\partial \sum e_i^2}{\partial \beta_0} = 0 \text{ から導かれる}$$

$$\sum e_i X_i = \sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) X_i = 0 \quad \leftarrow \frac{\partial \sum e_i^2}{\partial \beta_1} = 0 \quad //$$

上記2つの制約1=0より、自由度は2減る

一般的には、データ数  $n$ , 回帰係数 (パラメータ) の数  $p \rightarrow$  自由度  $n-p$

(cf. 最大法)

## (3) 最大推定量との関係

P.25

仮定より  $\varepsilon_i \sim N(0, \sigma^2)$ ,  $Y_i$  は  $\varepsilon_i$  の線形関数, 均一分散

$\therefore Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2) \quad \leftarrow (P.24 \text{ の図})$

$Y_1, \dots, Y_n$  の pdf (確率密度関数) は、 $n$  が独立

$$(2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum (Y_i - \beta_0 - \beta_1 X_i)^2\right\} = L(\beta_0, \beta_1, \sigma^2) \quad \leftarrow \text{尤度関数}$$

対数として考えて

$$\left[ \frac{\partial \log L}{\partial \beta_0} = \frac{1}{\sigma^2} \sum (Y_i - \beta_0 - \beta_1 X_i) = 0 \right]$$

$$\frac{\partial \log L}{\partial \beta_1} = \frac{1}{\sigma^2} \sum (Y_i - \beta_0 - \beta_1 X_i) X_i = 0$$

$$\frac{\partial \log L}{\partial \sigma^2} = -\frac{1}{2\sigma^2} + \frac{1}{2\sigma^4} \sum (Y_i - \beta_0 - \beta_1 X_i)^2 = 0$$

最小二乗法と同じ式!!

(誤差項に正規分布を仮定しているため)

$$\therefore \hat{\beta}_{0u} = \hat{\beta}_0, \quad \hat{\beta}_{1u} = \hat{\beta}_1, \quad \hat{\sigma}_u^2 = \frac{1}{n} \sum e_i^2$$

(最大法を用いて)

→ 不偏でない (cf. p.16)

$\varepsilon_i$  が正規分布という仮定のもとでは、 $\beta_0, \beta_1$  の最大推定量と最小二乗推定量

は一致 (分散は異なる)

## 8.2 モデルの説明力

$X$  が  $Y$  をどの程度説明できるか  $\rightarrow$  回帰モデルの当てはまりよさ

ここで、残差  $e_i = Y_i - \hat{Y}_i \rightarrow Y_i = \hat{Y}_i + e_i$  両辺から  $\bar{Y}$  を引く

$$Y_i - \bar{Y} = \hat{Y}_i - \bar{Y} + e_i \quad \text{この2乗和を3つ}$$

$$\sum (Y_i - \bar{Y})^2 = \sum (\hat{Y}_i - \bar{Y})^2 + \sum e_i^2 + 2 \sum (\hat{Y}_i - \bar{Y}) e_i$$

$$(\text{右辺第3項}) = \sum (\hat{Y}_i - \bar{Y}) e_i = \sum \hat{Y}_i e_i \quad (\because \sum \bar{Y} e_i = \bar{Y} \sum e_i = 0)$$

$$= \sum (\beta_0 + \beta_1 X_i) e_i = 0 \quad (\because \sum e_i = 0, \sum X_i e_i = 0)$$

$$\therefore \sum (Y_i - \bar{Y})^2 = \sum (\hat{Y}_i - \bar{Y})^2 + \sum e_i^2$$

全変動 モデルによる説明 残差平方和

全変動をモデルでどれだけ説明できるか

$$\text{説明力の尺度} = \text{決定係数 } r^2 = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2} = 1 - \frac{\sum e_i^2}{\sum (Y_i - \bar{Y})^2} = r_{Y\hat{Y}}^2 \quad \left( \begin{array}{l} \text{相関係数} \\ \text{の2乗} \end{array} \right)$$

$$0 \leq r^2 \leq 1$$

$$r^2 = 1 \text{ のとき, } \sum e_i^2 = 0 \Rightarrow \hat{Y}_i = Y_i$$

$$r^2 = 0 \text{ のとき, } \sum (\hat{Y}_i - \bar{Y})^2 = 0 \Rightarrow \hat{Y}_i = \bar{Y}$$

各  $X_i$  が  $Y$  の  $X_i$  の力からどれだけ説明できるか

## 8.4 統計的推測

$\varepsilon_i \sim N(0, \sigma^2)$   $\beta_0, \beta_1: \varepsilon_i$  の線形関数  $\rightarrow$  正規分布 (再生性)

$$E(\hat{\beta}_0) = \beta_0, E(\hat{\beta}_1) = \beta_1, \text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum X_i^2}{n \sum (X_i - \bar{X})^2}, \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{n \sum (X_i - \bar{X})^2}$$

$$\therefore \hat{\beta}_0 \sim N(\beta_0, \text{Var}(\hat{\beta}_0)), \hat{\beta}_1 \sim N(\beta_1, \text{Var}(\hat{\beta}_1))$$

$$\sigma^2 \text{ 未知} \rightarrow S^2 = \frac{\sum e_i^2}{n-2} \text{ で置き換える}$$

$\hat{\beta}_0, \hat{\beta}_1$  の分散の不偏推定量は

$$S_{\hat{\beta}_0}^2 = \frac{S^2 \sum X_i^2}{n \sum (X_i - \bar{X})^2}, S_{\hat{\beta}_1}^2 = \frac{S^2}{n \sum (X_i - \bar{X})^2}$$

母集団正規分布,  $\sigma^2$  未知より

(cf. P.113)

$\hat{\beta}_0, \hat{\beta}_1$  を標準化したものは、自由度  $n-2$  の  $t$  分布になる。

$$t_0 = \frac{\hat{\beta}_0 - \beta_0}{S_{\hat{\beta}_0}} \sim t(n-2), t_1 = \frac{\hat{\beta}_1 - \beta_1}{S_{\hat{\beta}_1}} \sim t(n-2)$$

( $H_0: \beta_0 = 0$  と仮定したとき...)

$\Rightarrow$  パラメータに対する区間推定 (cf. 6.3) や仮説検定 (cf. 7) が可能



単回帰 ( $X$  が 1)

## 8.5. 重回帰モデル

線形回帰モデル

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + \varepsilon_i$$

↓ 行列に直す

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}, X = \begin{bmatrix} 1 & X_{11} & \dots & X_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & X_{n1} & \dots & X_{np} \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix}, \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} \text{ とする,}$$

$$Y = X\beta + \varepsilon$$

おくべき仮定は 8.1 と同じ。ただし、 $\varepsilon \sim N(0, \sigma^2 I)$ 最小二乗法  $\varepsilon' \varepsilon \rightarrow \min \therefore \frac{\partial \varepsilon' \varepsilon}{\partial \beta} = 0$  単位行列

$$\rightarrow \hat{\beta} = (X'X)^{-1} X'Y \quad \text{覚える} \quad , E(\hat{\beta}) = \beta, \text{Var}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$$

$$\text{残差: } e = Y - X\hat{\beta}$$

$$\text{不偏分散 } s^2 = \frac{e'e}{n-p} = \frac{(Y - X\hat{\beta})'(Y - X\hat{\beta})}{n-p} \quad (p = q+1, \quad p = \text{rank}(X))$$

一般に、重回帰式の方が決定係数高。(説明力が高まる)が、自由度が減少

そのため不偏分散が大きくなり回帰直線の汎化性が低下する。

Akaike Information Criterion

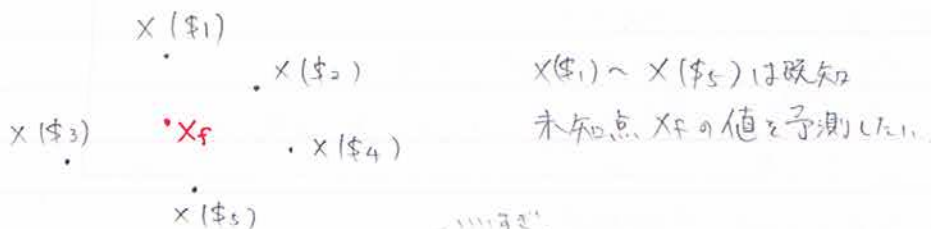
説明変数の数も含めたモデルの決定  $\Rightarrow$  情報量規準 (AIC)

乗回帰分析

6/15 8/9 時系列解析 (中西さん)  $\rightarrow$  フォント

## 8/10 空間データの統計解析

降雨量、標高など 2次元

通常は独立ではなく、近い所ほど近い値を持つ傾向 = 空間相関

空間相関を考慮し、よりよい予測値を求める。アプローチは2つ。

- ① 空間計量経済学 (10.2)
- ② 空間統計学, 地球統計学 (10.3, 10.4)

試験出るかも

27.770-47 共通

10.1. 一般化最小二乗法

重回帰分析 (cf. 8.5) : 通常最小二乗法 (ordinary least squares : OLS)

誤差項において空間相関を考慮 : 一般化最小二乗法 (generalized least squares : GLS)

線形回帰モデル

$$y = X\beta + \varepsilon$$

自由回帰と仮定がせや異なる

仮定

$$\textcircled{1} E(\varepsilon) = 0 \quad \leftarrow \text{平均ゼロ}$$

共分散なし

$$\textcircled{2} E(\varepsilon\varepsilon') = V = \sigma^2 \Omega \quad \leftarrow \text{共分散行列}$$

$$(cf) OLS: \varepsilon \sim N(0, \sigma^2 I)$$

$$\Omega: \text{正値定符号行列} \Leftrightarrow \forall x, x' \Omega x > 0$$

$$\textcircled{3} \varepsilon: \text{多変量正規分布} \quad \leftarrow \text{単位行列}$$

$$\Omega \text{ は正値定符号行列だから, } P\Omega P' = I \text{ なる } P \text{ が存在. } (\Omega^{-1} = P'P)$$

$$Py = Px\beta + P\varepsilon \quad \leftarrow P \text{ が行列}$$

$$\text{ここで } E(P\varepsilon) = 0, \text{ var}(P\varepsilon) = E(P\varepsilon\varepsilon'P) = \sigma^2 I \quad \leftarrow OLS \text{ 的だ!!}$$

 $P\varepsilon$  を被説明変数,  $PX$  を説明変数として OLS で推定可能

$$\begin{aligned} \hat{\beta}_{GLS} &= [(PX)'(PX)]^{-1} (PX)'Py = (X'P'PX)^{-1} X'P'Py \\ &\quad \leftarrow P'P \\ &= (X'\Omega^{-1}X)^{-1} X'\Omega^{-1}y = (X'\Omega^{-1}X)^{-1} X'\Omega^{-1}y \quad \leftarrow V = \sigma^2 \Omega \end{aligned}$$

平均と分散

$$\hat{\beta}_{GLS} = (X'\Omega^{-1}X)^{-1} X'\Omega^{-1}(X\beta + \varepsilon) = \beta + (X'\Omega^{-1}X)^{-1} X'\Omega^{-1}\varepsilon \quad \text{だから}$$

$$E(\hat{\beta}_{GLS}) = E(\beta) = \beta \quad (\because E(\varepsilon) = 0)$$

$$\text{var}(\hat{\beta}_{GLS}) = \sigma^2 (X'\Omega^{-1}X)^{-1} = (X'\Omega^{-1}X)^{-1} \quad \leftarrow \text{実数値行列}$$

$$S_G^2 = \frac{(y - X\hat{\beta}_{GLS})' \Omega^{-1} (y - X\hat{\beta}_{GLS})}{n - p} \quad OLS \text{ との違いをチェック}$$

 $V$  (or  $\Omega$ ) が定義できれば、空間相関を考慮したモデルがでる

 $\hookrightarrow$  この定義の仕方がアプロ-チの違いとなる。

10.2 ~ 10.4 はあまり試験に出ないよ

存在だけ知ってほしい

## 10.2 空間自己回帰モデル

誤差項をモデル化することにより、間接的に  $\nabla$  をモデル化

$$y = X\beta + u$$

空間関数

$$u = \lambda W u + \varepsilon \quad \lambda: \text{係数}, W: \text{空間重行列}$$

(自己で帰ってくる)

error component model

$$\text{Var}(u) = \nabla = \sigma^2 \Omega$$

$$\text{Var}(\varepsilon) = \sigma^2 I$$

) と仮定

$$u = (I - \lambda W)^{-1} \varepsilon$$

OLS と GLS あり

$$\therefore \nabla = E(uu') = \sigma^2 [(I - \lambda W)'(I - \lambda W)]^{-1}$$

これを用いて GLS によりパラメータ推定

## 10.3 確率場と定常性

位置に依存した確率変数

位置  $\phi \in \mathbb{R}^2$  に対して確率変数  $Z(\phi)$  が対応 $\{Z(\phi) = \phi \in \mathbb{R}^2\}$ : 確率場  $\Leftarrow$  cf. 確率過程 (89)統計解析のための  $Z(\phi)$  への仮定: 定常性・ 強定常  $\Leftarrow$  おまけ程度同時分布関数  $F_{\phi_1, \dots, \phi_n}(Z(\phi_1), \dots, Z(\phi_n))$  とする $\forall \phi_1, \dots, \phi_n, h \in \mathbb{R}^2$  に対して

$$F_{\phi_1, \dots, \phi_n}(Z(\phi_1), \dots, Z(\phi_n)) = F_{\phi_1+h, \dots, \phi_n+h}(Z(\phi_1), \dots, Z(\phi_n))$$

 $\Downarrow$ ・  $Z(\phi)$  の分布関数は任意の位置  $\phi$  で同じ・  $Z(\phi_1), \dots, Z(\phi_n)$  の同時分布関数は、位置差  $h$  のみに依存△  
シ  
ッ  
テ  
ハ  
ー  
リ  
ッ

強定常: 分布そのものにこの仮定

通常  $F_{\phi_1, \dots, \phi_n}$  は未知  $\rightarrow$  仮定をゆるめた (弱定常)  $\Leftarrow$  大事・ 2次定常性  $\checkmark$  こゝからまた大事

$$E[Z(\phi)] = \mu$$

$$\text{cov}[Z(\phi_i), Z(\phi_j)] = c(\phi_i - \phi_j) = c(h)$$

 $h$ : 点  $\phi_i, \phi_j$  の相対的位置 (ベクトル) $c(h)$ : 共分散関数 = コバリオグラム  $\Leftarrow$  空間重行列

さらに簡単にしたい

 $\rightarrow$  次のページ

★ 等方的2次定常性

$E[z(s)] = \mu$   
 $Cov[z(s_1), z(s_2)] = c(h)$  スカラー

$h$ : 点  $s_1, s_2$  の距離 (スカラー)

$c(h)$  を定義  $\rightarrow \nabla$  (or  $\Omega$ ) が定義できる  $\rightarrow$  次

任意の位置の埋め込みを調べたい

10.4 クリギング

(1) 一般的可予測



既知点:  $z_i (i=1, \dots, n)$

未知点:  $z_f$  forecast

重みつき平均 (内挿) 剛から未知点への推定

重みつき平均 (内挿)

$$\hat{z}_f = \lambda_1 z_1 + \dots + \lambda_n z_n = \lambda' z \quad (\lambda: \text{重み})$$

(2) 誤差項の予測

空間相関と誤差項で考える  $\rightarrow$  未知点の 残差 を既知点の残差で予測

重みつき平均

$$\hat{\varepsilon}_f = \lambda_1 \varepsilon_1 + \dots + \lambda_n \varepsilon_n = \lambda' \varepsilon$$

既知点の残差:  $\varepsilon = (y - X\beta_0)$

ここで、既知点間の分散共分散行列:  $V = \begin{bmatrix} \sigma_{11} & \dots & \sigma_{1n} \\ \vdots & & \vdots \\ \sigma_{n1} & \dots & \sigma_{nn} \end{bmatrix}$

既知点と未知点の共分散ベクトル:  $d = \begin{bmatrix} \sigma_{1f} \\ \vdots \\ \sigma_{nf} \end{bmatrix} = \begin{bmatrix} c(|s_f - s_1|) \\ \vdots \\ c(|s_f - s_n|) \end{bmatrix}$  キヨリ (コバリエーション) がないとノイズ

誤差と予測残差の差を最小とする重み  $E(\lambda' \varepsilon \varepsilon' \lambda)$

$$E[(\varepsilon_f - \hat{\varepsilon}_f)^2] = E(\varepsilon_f^2) - 2E(\varepsilon_f \hat{\varepsilon}_f) + E(\hat{\varepsilon}_f^2)$$

$$= \sigma^2 - 2\lambda' d + \lambda' V \lambda \rightarrow \min$$

とある:  $d$  は関係する  $\varepsilon_f$  の  $d$  の長さ

$$\frac{\partial E[\dots]}{\partial \lambda} = -2d + 2V\lambda = 0 \quad \therefore \lambda' = d'V^{-1}$$

1X上より

$$\hat{\varepsilon}_f = \lambda' \varepsilon = d'V^{-1} (y - X\beta_0)$$



点検定 (1971年 第2次試験)

数346

点検定 or 仮説検定

NO.

33

点検定 (通称「一般検定」) の主成分分析

DATE

1971

1971

1-27

ここで 線形回帰モデル

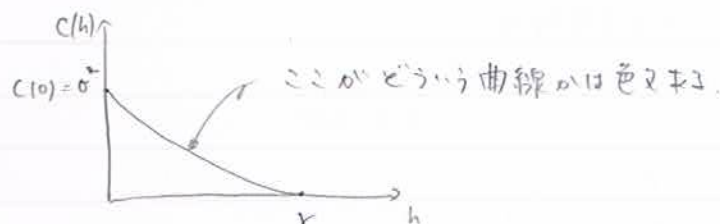
$x_f$ : 未知点の説明変数ベクトル  
 $y_f$ : " 被説明変数  
 $\varepsilon_f$ : " 誤差

とすると、

$$\hat{y}_f = x_f \hat{\beta} + \varepsilon_f = x_f \hat{\beta} + (y - X\hat{\beta})$$

モデル推定への残差

(3) コバリログラムの例



- ・ 指数型  $C(h) = \sigma^2 \exp(-\frac{h}{\nu})$
- ・ ガウス型  $C(h) = \sigma^2 \exp[-(\frac{h}{\nu})^2]$
- ・ 球状型  $C(h) = \sigma^2 [1 - \frac{3}{2}(\frac{h}{\nu}) + \frac{1}{2}(\frac{h}{\nu})^3]$

異なること多い

7/7

§11 多変量解析

principal component analysis

11.1 主成分分析 (PCA) の統計意味

ex) ある人の各科目の得点:  $x_1, \dots, x_p$ 総合得点:  $z = x_1 + \dots + x_p$  ← 本当に良し、総合評価か?

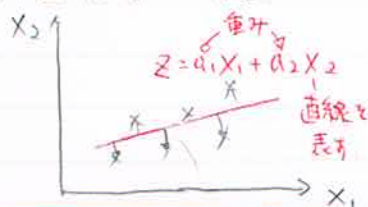
各科目の得点分布にもばらつきがある。これを考慮 (重みをつける)

各点に重みをつけて、総合得点を求める。より一般的には、

「各変数に重みをつけ、より少ない 新たな変数で総合化する方法

= 主成分分析」

(1) 2変数の場合

 $x_2$  だけで評価 → 差が大きい $x_1$  " → 総合化にならな

個人差が &gt; 総合得点

2次元 → 1次元

無意味

$$\lambda^2 - (S_{11} + S_{22})\lambda + S_{11}S_{22} - S_{12}^2 = 0$$

個人差をばらばらにする

軸を変換したとき、その軸上で値の分散が最大になる軸を求める。

$z = a_1 x_1 + a_2 x_2$  を考える。

$x_1, x_2$  の (標本) 分散共分散行列  $\begin{pmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{pmatrix}$  →  $z$  を考える

$\text{var}(z) = \frac{1}{n-1} \sum (z_i - \bar{z})^2$  ← 標本から  $x_1, x_2$  について自由度 1 がある

$$\begin{aligned} &= \frac{1}{n-1} \sum \{(a_1 x_{1i} + a_2 x_{2i}) - (a_1 \bar{x}_1 + a_2 \bar{x}_2)\}^2 \\ &= \frac{1}{n-1} \sum \{a_1^2 (x_{1i} - \bar{x}_1)^2 + 2a_1 a_2 (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2) + a_2^2 (x_{2i} - \bar{x}_2)^2\} \\ &= a_1^2 S_{11} + 2a_1 a_2 S_{12} + a_2^2 S_{22} \quad \rightarrow \text{最大化したい} \end{aligned}$$

制約条件:  $a_1^2 + a_2^2 = 1$  (式自体に 1 はないが、制約がないと  $\text{var}(z)$  無限に大きくなる)

制約条件つき最大化 → ラグランジュの未定乗数法

$$f(a_1, a_2, \lambda) = a_1^2 S_{11} + 2a_1 a_2 S_{12} + a_2^2 S_{22} - \lambda (a_1^2 + a_2^2 - 1)$$

制約条件

$$\begin{cases} \frac{\partial f}{\partial a_1} = 2a_1 S_{11} + 2a_2 S_{12} - 2a_1 \lambda = 0 \\ \frac{\partial f}{\partial a_2} = 2a_1 S_{12} + 2a_2 S_{22} - 2a_2 \lambda = 0 \\ \frac{\partial f}{\partial \lambda} = a_1^2 + a_2^2 - 1 = 0 \end{cases}$$

上の2式から、 $\begin{pmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \lambda \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}$  : 固有値問題になる。

$\lambda$  : 固有値,  $\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}$  : 固有ベクトルを表す。

$$\begin{pmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \begin{pmatrix} \lambda & 0 \\ 0 & \lambda \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} \Leftrightarrow \begin{pmatrix} S_{11} - \lambda & S_{12} \\ S_{21} & S_{22} - \lambda \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$a_1 = a_2 = 0$  以外の解を持つためには、

$$\begin{vmatrix} S_{11} - \lambda & S_{12} \\ S_{21} & S_{22} - \lambda \end{vmatrix} = 0 \quad \therefore \lambda = \frac{(S_{11} + S_{22}) \pm \sqrt{(S_{11} - S_{22})^2 + 4S_{12}^2}}{2}$$

$$\lambda_1 > \lambda_2 \text{ とすると、 } a_{1(1)} = \frac{S_{12}}{\sqrt{S_{12}^2 + (\lambda_1 - S_{11})^2}}, \quad a_{2(1)} = \frac{\lambda_1 - S_{11}}{\sqrt{S_{12}^2 + (\lambda_1 - S_{11})^2}}$$

$\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}$  は規格化  
 $\frac{\partial f}{\partial \lambda} = 0$  に代わり  
3 番目の式...

同様にして  $\lambda = \lambda_2 \rightarrow a_{1(2)}, a_{2(2)}$  も求まる。

新しい軸  $z_1$  : 第1主成分,  $z_2$  : 第2主成分 ← 固有ベクトルは直交 (幾何学的には回転)

$z_1 = a_{1(1)} x_1 + a_{2(1)} x_2$  : 主成分得点 ( $z_2$  は4行する)

固有値問題

対角化

$$\begin{pmatrix} s_{11} & s_{12} \\ s_{12} & s_{22} \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \end{pmatrix}$$

(cf) 2つの最小二乗法

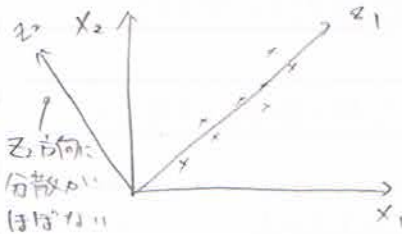
分散共分散行列の対角化  $\rightarrow z_1, z_2$  は無相関

処理を簡単にする

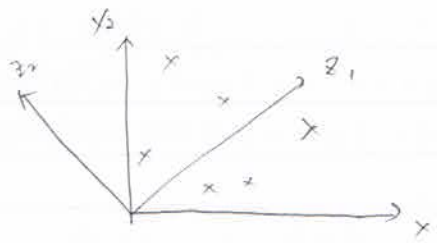
主成分分析: 相関のあるデータを無相関のデータに変換

・  $z_1$  で全体をどれだけ表現できるか $z_1$  軸上の分散:  $\lambda_1$ ,  $z_2$  軸上の分散:  $\lambda_2$ 寄与率 (第1主成分):  $\frac{\lambda_1}{\lambda_1 + \lambda_2}$ 

①: 元のデータの違いを十分に表現  
②: できていない



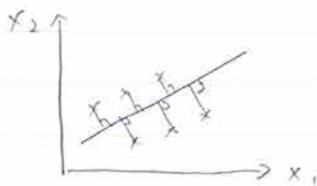
寄与率 ①

 $\Rightarrow$  次元圧縮

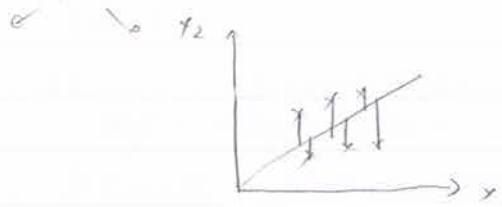
寄与率 ②

・ 回帰分析との相違

月データ



主成分

 $x_1$  と  $x_2$  同等 $\rightarrow$  主成分 1つだけ

回帰

 $x_1$  と  $y_2$  を説明 $\rightarrow y_2$  軸上での差

(2) 多変数の場合

$$Z = a_1 x_1 + \dots + a_p x_p$$

$$\text{Var}(Z) = \frac{1}{n-1} \sum \{a_1(x_{1i} - \bar{x}_1) + \dots + a_p(x_{pi} - \bar{x}_p)\}^2 = \sum_k \sum_k S_{kk} a_k a_k$$

$$\text{制約条件 } a_1^2 + \dots + a_p^2 = 1$$

ラグランジュの未定乗数法から、

$$\begin{pmatrix} S_{11} & \dots & S_{1p} \\ \vdots & \ddots & \vdots \\ S_{p1} & \dots & S_{pp} \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_p \end{pmatrix} = \lambda \begin{pmatrix} a_1 \\ \vdots \\ a_p \end{pmatrix} \quad : \text{固有値問題を解く (対角化)}$$

## 11.2. 因子分析 ✓ 心理学的由来、わかりやすい。

p変数の測定値、データ数: n

n &gt; p+1 ならば、最大限 p 個の主成分が得られる。

第1主成分の寄与率が著しく大きく、第2主成分以下の寄与率が低くなることを望まれる。

しかしこれはきわめてまれ → 主成分をいくつで止めたらよいのか？

そこで、寄与率の小さい主成分は誤差成分としてまとめ、できる限り少数の意味ある主成分を得ることを考える。

変数:  $x_1, \dots, x_p$ 、主成分:  $z_1, \dots, z_r$  として、次のモデルを考える。

$$\begin{cases} x_1 = a_{11}z_1 + \dots + a_{1r}z_r + \varepsilon_1 \\ \vdots \\ x_p = a_{p1}z_1 + \dots + a_{pr}z_r + \varepsilon_p \end{cases} \quad \begin{matrix} \text{XをZで表している} \\ \text{(主成分を全く逆!!)} \end{matrix}$$

主成分  $z \rightarrow$  因子分析では 共通因子 とよむ。以下  $f$  とする。

$$x_j = a_{j1}f_1 + \dots + a_{jr}f_r + \varepsilon_j \quad (j=1, \dots, p)$$

 $x_j$  の変動をいくつかの共通因子の加重和に分解。 $a_{jk}$ : 変数  $j$  の共通因子の影響度合 = 因子負荷量 $\varepsilon_j$ : 誤差 = 独自因子。

複数の変数からなるデータを少数の共通因子を使って表現し、データが現れる仕組みのモデルを見つけ出す方法、= 因子分析 (因子負荷量と共通因子の値を求めよう)



cf) 主成分: 与えられた変数を合成して全体の変異力になるべく少数の主成分で説明  
因子分析: ある変数の共通変動のみをいくつかの因子に分解

