

7章 標本と統計的推測

統計的推測の話の始まり

確率部分

確率変数

確率密度関数, 分布関数に基づき,

平均と分散などの統計量,

変数変換

統計では

目の前にあるものに基づいて

もとの確率変数のことを推測する

9.1 標本とパラメータ

用語と記号 (確率と統計では同じこととを少し違う用語で表現するところ)

サイコロ投げを例に、

母集団
(標本空間)母集団: サイコロの出目の全体 $\Leftarrow ?$
(一般的に) 考える対象とする集団

母集団分布:

サイコロの出目の現れ方は確率的に決まる。

母集団の確率的な挙動を表す確率分布

サイコロを母集団とした場合、

母集団分布

母集団分布は多項分布

$$\vec{X} = (X_1, X_2, \dots, X_n)$$

n回のサイコロ投げ

標本値
 x

$$\left\{ \begin{array}{l} X_1: \text{1回のサイコロ投げの値} \quad 1, 2, 3, 4, 5, 6 \text{ の出目} \\ P_1 = P_2 = \dots = P_6 = 1/6. \\ \text{具体的に現れる値 } x \end{array} \right.$$

標本

 x の背後に想定される確率変数 X を標本という

本書独特な捉え方。

本来の標本には、確率的な考え方は不必要。
(母集団があって、そこから取り出したものが標本)標本から母集団の様子を、うまく推測する
ための本書での考え方、および具体的に、

無作為抽出と無作為標本

サイコロを n 回振って得られると想定される
標本 X_1, \dots, X_n があつたとして

無作為抽出

random

sampling

無作為標本

それぞれの標本は無作為に得られる
このように標本を作る方法

得られた標本

確率変数の言葉では

確率変数 $X_i, i=1:n$ が

母集団の確率的挙動を表す確率変数 X と
同一の分布に従っていて、独立

$$X_i, i=1:n \underset{i.i.d.}{\sim} X$$

パラメータ

母集団を特徴づけるパラメータ (母数)

 $B(\theta)$ $N(\mu, \sigma^2)$

など

密度関数とパラメータ

必要に応じて

母集団のパラメータを記すことがある

$$f(x; \theta)$$

$$E_{\theta}[g(x)] = \int g(x) \cdot f(x; \theta) dx$$

パラメータの真値

?

想定する 集団が具体的なとき、
母集団を特徴づけている θ には、
何らかの真値 θ^* が存在している。

7.2 統計的推測

確率部分の流れ.

↑
統計の流れ

$$X \sim N(\mu, \sigma^2) \text{ としよう}$$

X の平均と分散が μ と σ^2 であり

$$Y = aX + b \sim N(a\mu + b, a^2\sigma^2)$$

$$X \sim \dots, X_{i=1:n} \rightarrow \bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

$$\bar{X} \xrightarrow{d} \mu$$

統計的推測 (statistical inference)

標本に基づいて、母の確率変数 (つまり母集団) のことを推測する

統計的推測の例

歪められたコインの表の出る確率 (真推定)

- 表が出る $X=1$, 裏が出る $X=0$

- 表が出る確率をパラメータ θ で表わす
(θ は不明)

- θ を推し量りたい

- 無作為標本 $X_{i=1:n}$ に対し \bar{X} から θ を推し量る

- コインを 100 回振って、標本値 $X_{i=1:100}$

- 表の数が 36 であった. $\bar{X} = 36$, $\theta = 0.36$

7.3 標本平均と標本分散

無作為標本 $X_{i=1:n}$ に対し

標本平均	$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$	} 二つの統計量
" 分散	$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$	

統計量 標本だけによって定義される量
(statistic)

 \bar{X} と S^2 の性質

X

母集団の平均と分散: μ と σ^2 \bar{X}

標本平均 \bar{X} の 平均 μ (⇒ 2.6 節)
分散 σ^2/n

 S^2

$$\begin{aligned}
 S^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \\
 &= \frac{n}{n-1} \cdot \frac{1}{n} \sum \{(X_i - \mu) - (\bar{X} - \mu)\}^2 \\
 &= \frac{n}{n-1} \cdot \frac{1}{n} \sum \{(X_i - \mu)^2 + (\bar{X} - \mu)^2 - 2(X_i - \mu)(\bar{X} - \mu)\} \\
 &= \frac{n}{n-1} \cdot \frac{1}{n} \left\{ \sum \{(X_i - \mu)^2 + (\bar{X} - \mu)^2\} - \frac{2(\bar{X} - \mu) \sum (X_i - \mu)}{2n(\bar{X} - \mu)^2} \right\} \\
 &= \frac{n}{n-1} \left\{ \frac{1}{n} \sum (X_i - \mu)^2 - (\bar{X} - \mu)^2 \right\}
 \end{aligned}$$

$E[S^2]$

$$\begin{aligned}
 E[S^2] &= \frac{n}{n-1} \left\{ \frac{1}{n} \sum_{i=1}^n E[(X_i - \mu)^2] - E[(\bar{X} - \mu)^2] \right\} \\
 &= \frac{n}{n-1} \cdot \frac{1}{n} \left\{ \sigma^2 - \frac{\sigma^2}{n} \right\} \\
 &= \sigma^2
 \end{aligned}$$

 \bar{X} と S^2 の 確率収束大数の
法則より

$$X, \mu = E[X], \sigma^2 = V[X]$$

$$X_{i=1:n} \sim_{i.i.d} X \Rightarrow \bar{X} = \frac{1}{n} \sum X_i \xrightarrow{P} \mu$$

算術平均

 $S^2 \xrightarrow{P} \sigma^2$ の証明

$$S^2 = \frac{n}{n-1} \left\{ \frac{1}{n} \sum (X_i - \mu)^2 - (\bar{X} - \mu)^2 \right\}$$

↓ μ 大数

 $\lim_{n \rightarrow \infty} \text{「大数」} = 0$

より従って、

$$\Rightarrow 1 \left\{ \underbrace{\frac{1}{n} \sum (X_i - \mu)^2}_{E[(X - \mu)^2]} - 0 \right\}$$

$$\xrightarrow{P} \sigma^2$$

標本分布

統計量が内在している確率分布のこと

$$\text{母集団分布} \sim N(\mu, \sigma^2)$$

のとき

$$\text{標本平均 } \bar{X} \sim N(\mu, \sigma^2/n)$$

つまり

$$\bar{X} \text{ の標本分布は } N(\mu, \sigma^2/n) \text{ である} \text{ といふ}$$

標本分散の
標本分布は?

$$\begin{cases} (n-1)S^2/\sigma^2 \sim \chi_{n-1}^2 \\ \bar{X} \text{ と } S^2 \text{ は独立である} \end{cases}$$

証明.

まず準備

$$\vec{j} = (1, \dots, 1)^t, \quad \vec{a} = \vec{j}/\sqrt{n} \quad (\text{長さ } 1 \text{ のベクトル})$$

標本平均

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \vec{j}^t \vec{X} = \frac{1}{\sqrt{n}} \vec{a}^t \vec{X}$$

行列の準備

$$A = (\vec{a}, A_2)$$

(n次)
Aは直交行列

$$A^t A = A \cdot A^t = I_n$$

$$\vec{0} = (0, \dots, 0)^t$$

$$A_2^t \vec{a} = \vec{0}$$

新しい確率変数

$$\vec{Y} = A^t \vec{X}$$

$$\vec{X} \sim N_n(\mu \cdot \vec{j}, \sigma^2 I_n)$$

S^2 の標本分布の続き

準備

$$A^t \cdot \vec{j} = \begin{pmatrix} \vec{a}^t \cdot \vec{j} \\ A_2^t \cdot \vec{j} \end{pmatrix} = \begin{pmatrix} \vec{a}^t \cdot \sqrt{n} \vec{a} \\ A_2^t \cdot \sqrt{n} \vec{a} \end{pmatrix} = \begin{pmatrix} \sqrt{n} \\ \vec{0} \end{pmatrix}$$

\vec{Y} は \vec{X} の線形変換なので

$$\begin{aligned} \vec{Y} = A^t \cdot \vec{X} &\sim N_n(\mu A^t \vec{j}, \sigma^2 A^t A) \\ &\stackrel{d}{=} N_n(\mu A^t \vec{j}, \sigma^2 \vec{I}_n) \\ &= N_n((\sqrt{n}\mu, 0, \dots, 0)^t, \sigma^2 \vec{I}_n) \end{aligned}$$

$\stackrel{d}{=}$ は両辺の確率分布が等しいという意味,

$\vec{Y} = A^t \cdot \vec{X} \sim N_n((\sqrt{n}\mu, 0, \dots, 0)^t, \sigma^2 \vec{I}_n)$ が
意味するのは $\vec{Y} = (Y_1, \dots, Y_n)^t$ とし

$$\left\{ \begin{array}{l} Y_1 \sim N(\sqrt{n}\mu, \sigma^2) \\ Y_{i=2:n} \sim N(0, \sigma^2) \\ Y_1, \dots, Y_n \text{ は互いに独立} \end{array} \right\} \text{ が導かれた。}$$

\bar{X} と S^2 は \vec{Y} で表現

$$\bar{X} = Y_1 / \sqrt{n}$$

$$\begin{aligned} (n-1)S^2 &= \sum_{i=1}^n (X_i - \bar{X})^2 & \sum X_i^2 - \underbrace{2X_i \bar{X}}_{-2\bar{X} \cdot n\bar{X}} + \underbrace{\bar{X}^2}_{n\bar{X}^2} \\ &= \sum X_i^2 - n\bar{X}^2 \\ &= \vec{X}^T \vec{X} - n\bar{X}^2 \end{aligned}$$

$$\begin{aligned} \vec{X} &= A^T \cdot \vec{Y} \\ (\vec{X})^T &= (A^T \cdot \vec{Y})^T = A \cdot \vec{Y}^T \end{aligned} \quad \begin{matrix} [A_{ji}] \\ A_1 \end{matrix} \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

$$= \vec{Y}^T \cdot \vec{Y} - Y_1^2$$

$$= Y_2^2 + \dots + Y_n^2$$

$\therefore \bar{X}$ と S^2 は互いに独立

$$\bar{X} = \bar{X}(Y_1), \quad S^2 = S^2(Y_2, \dots, Y_n)$$

$$\bar{X} \sim N(\sqrt{n}\mu, \sigma^2)/\sqrt{n} = N(\mu, \sigma^2/n) \quad \text{正規分布}$$

$$\frac{n-1}{\sigma^2} S^2 \sim \chi_{n-1}^2, \quad \frac{Y_2}{\sigma}, \dots, \frac{Y_n}{\sigma} \sim N(0, 1) \quad \text{標準正規}$$

7.4 標準化とステューデント化

大前提

母集団 (μ, σ^2) からの無作為標本 $X_{i=1:n}$

ステューデント化

確率変数の標準化 (確率の話が登場)
 { ステューデント化 (統計 ")
 もいはい

標本平均の
標準化

$$Z_n = \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}}$$

標本平均

標本平均の
ステューデント化

$$T_n = \frac{\bar{X} - \mu}{\sqrt{S^2/n}}$$

σ^2 に確率収束する
 確率変数 X で置き換える。

母集団 $\sim N(0, 1)$ の時

...

$$Z_n \sim N(0, 1)$$

$$T_n \sim \text{自由度 } n-1 \text{ の } t\text{-分布}$$

$$T_n = \frac{\bar{X} - \mu}{\sqrt{S^2/n}} = \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \cdot \frac{1}{\sqrt{S^2/\sigma^2}}$$

$$= \frac{Z_n \sim N(0,1)}{\sqrt{\underbrace{(n-1)S^2/\sigma^2}_{\chi^2_{n-1}}/(n-1)}}$$

$$\stackrel{d}{=} \frac{N(0,1)}{\sqrt{\chi^2_{n-1}/(n-1)}}$$

$$\stackrel{d}{=} t_{n-1}$$

3.2.5

 χ^2 分布

t-分布

母集団分布が不明のとき

 Z_n は漸近的に $\sim N(0,1)$

中心極限定理より

 T_n も漸近的に $\sim N(0,1)$

$$T_n = \frac{\bar{X} - \mu}{\sqrt{S^2/n}}$$

$$= \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n} \cdot \sqrt{S^2/\sigma^2}}$$

$$\xrightarrow{d} N(0,1) : 1$$