

⑧ 実験データの棄却について

阪大理学部数学教室

小川潤次郎

1組の実験データがあるとき、その中に異質的なものが混入しているならば、これを摘出するとか、又既に得られている1組のデータに対して新たに得られたデータが等質的なものと見られるかどうかを判定する場合に所謂“棄却検定法”(Rejection Test)が用いられる。

この棄却検定法については幾多の文献があるが⁽¹⁾、その中には概念の混乱も見受けられるようである。⁽²⁾ この棄却検定法従って棄却限界の概念を明確にするには矢張り統計的仮説検定論(Theory of Testing Statistical Hypotheses)⁽³⁾の立場に立つのが最も良いが、一般には複合仮説の検定となるので、対立仮説が複雑で仰々ろまく数学的取扱いが出来ないのは E. S. Pearson 及び C. Chandra Sekar⁽⁴⁾の指摘する通りである。

ここでは先づ最も単純な Situation について考えよう。未知なる真値 m をもつ量 X を n 回測定して、 n 回の実験値 x_1, x_2, \dots, x_n を得たとする。但しこの測定には Systematic Error がないものとする。そうすれば、この実験値 x_1, x_2, \dots, x_n は真値 m に夫々偶然誤差 $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ が加わったものと考えてよいであろう。

このとき $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ は互に独立に平均 0 なる正規分布に従うものとする。

このような Situation の下で、例えば x_1 が残餘の x_2, \dots, x_n に対して異質的であるかどうかという問題は、一応次のように Formulate されるであろう。

問題: x_1 は正規母集団 $N(m, \lambda^2 \sigma^2)$ から抽出された Random

Sample, x_1, \dots, x_n の同一の正規母集団 $N(m, \sigma^2)$ から抽出された Size $(n-1)$ の Random Sample であり, 従つて x_1, x_2, \dots, x_n の尤度函数が

$$p(x_1, x_2, \dots, x_n; \lambda, m, \sigma) = \frac{1}{\lambda} \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \left(\frac{(x_1-m)^2}{\lambda^2} + (x_2-m)^2 + \dots + (x_n-m)^2 \right) \right\}$$

であるとき, パラメーター λ に関する統計的仮説

$$H: \lambda = 1$$

を検定する Optimum な Critical Region を求めよ。

この場合には, パラメーターは三つあって, その admissible values を

$$-\infty < m < +\infty, \quad 0 < \sigma < +\infty, \quad 0 < \lambda < +\infty$$

であるから, 仮説 H を明かに複合仮説である。

吾々以下に於ては, λ の admissible values として

$$\lambda \geq 1$$

のみ, 従つて H の対立仮説の指定する λ の値として

$$\lambda > 1$$

のみ考へることとするが, 棄却検定の意味を考へるとき, これは不都合でわれないであらう。

先づ初めに, σ が既知の場合を考へる。これは森口氏⁽⁵⁾が取扱つた場合である。このとき $\sigma = 1$ として一般性を失われない。

J. Neyman & E.S. Pearson の一般的方法⁽⁶⁾を適用する爲に

$$p_0(x_1, x_2, \dots, x_n) = p(x_1, \dots, x_n; \mu, \sigma^2)$$

$$= \left(\frac{1}{\sqrt{2\pi}} \right)^n \exp \left\{ -\frac{1}{2} \left((x_1 - \mu)^2 + \dots + (x_n - \mu)^2 \right) \right\}$$

とかくと,

$$\phi = \frac{\partial \log p_0}{\partial \mu} = n(\bar{x} - \mu), \quad \phi' = \frac{\partial \phi}{\partial \mu} = -n$$

であるから, Neyman-Pearson の条件⁽⁷⁾

$$\phi' = A + B\phi$$

が満たれる。このとき

$$\phi = \text{Const.}$$

なる Hyper Surface を, 超平面

$$\bar{x} = \text{Const.}$$

となる。今求める Critical Region を w とし, k と上の超平面との交わりを $w(\bar{x})$ とすると, k は

$$p_t(x_1, \dots, x_n) \geq k(\bar{x}) \cdot p_0(x_1, \dots, x_n)$$

で定められる。但し,

$$p_t(x_1, \dots, x_n) = p(x_1, \dots, x_n; \lambda_t, \mu_t, \sigma^2) \quad \lambda_t > 1,$$

で定数 $k(\bar{x})$ を, 1st kind の Error の確率を ε (即ち Critical region の Size ε) とすれば

$$\int \int_{w(\bar{x})} p_0(x_1, \dots, x_n) dx_1, \dots, dx_n$$

$$= E \int \int_{\bar{x}} p_0(x_1, \dots, x_n) dx_1, \dots, dx_n$$

より定められる。領分領域 \bar{x} というのわ、考える超平面 $\bar{x} = \text{Const.}$ の全体の意味である。

ところで $w(\bar{x})$ を書き直すと

$$(x_i - m_i)^2 \geq C'(\bar{x})$$

又は

$$|x_i| \geq C(\bar{x})$$

であつて、 x_i と \bar{x} の同時分布の密度は、仮説 H の下では

$$\text{Const. } e^{-\frac{1}{2}(x_i - \bar{x})^2 - \frac{n}{2}\bar{x}^2}$$

であるから、 \bar{x} given のときの x_i の Conditional frequency function は

$$\text{const. } e^{-\frac{1}{2}(x_i - \bar{x})^2}$$

である、よつて $w(\bar{x})$ は、

$$|x_i - \bar{x}| \geq \lambda_p$$

但し

$$1 - 2 \int_0^{\lambda_p} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt = p/100. \equiv \varepsilon.$$

によつて定められる。

よつて、この場合の最良棄域は、

$$|x_i - \bar{x}| \geq \lambda_p \sigma$$

である。従って棄却限界は

$$\bar{x} - \lambda_p \sigma, \quad \bar{x} + \lambda_p \sigma$$

である。これは森口氏が何人ともなく作ったものである。

次には σ を λ known とする。この場合は

$$p_0(x_1, \dots, x_n) = p(x_1, \dots, x_n; \mu, \sigma)$$

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

$$\log p_0 = \text{const.} - n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\phi_1 = \frac{\partial \log p_0}{\partial \mu} = \frac{n}{\sigma^2} (\bar{x} - \mu), \quad \phi_1' = \frac{\partial \phi_1}{\partial \mu} = -\frac{n}{\sigma^2}$$

$$\phi_2 = \frac{\partial \log p_0}{\partial \sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum (x_i - \mu)^2$$

$$\phi_2' = \frac{\partial \phi_2}{\partial \sigma} = \frac{n}{\sigma^2} - \frac{3}{\sigma^4} \sum (x_i - \mu)^2$$

故から、Neyman-Pearson の Condition ⁽⁵⁾

$$\phi_1' = A_1 + B_1 \phi_1$$

$$\phi_2' = A_2 + B_2 \phi_2$$

を満足す。そしてこのときは

$$\phi_1 = \text{const.} \quad \phi_2 = \text{const.}$$

となる。よって仮説 H の Best Critical Region とこれらの Hyper surfaces との交わりを $w(\bar{x}, s^2)$ とすれば、それは、

$$p_t(x_1, \dots, x_n) \geq k(\bar{x}, s^2) p_0(x_1, \dots, x_n)$$

で定められる。但し、定数 $k(\bar{x}, s^2)$ は、

$$\iint_{\omega(\bar{x}, s^2)} p_0(x_1, \dots, x_n) dx_1, \dots, dx_n = \varepsilon \iint_{\bar{x}, s^2} p_0(x_1, \dots, x_n) dx_1, \dots, dx_n$$

によつて定められる。従つて

$$\omega(\bar{x}, s^2) : (x_i - m)^2 \leq c(\bar{x}, s^2).$$

今仮説 H が真なるときの x_1, \dots, x_n の分布は

$$\begin{aligned} & \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum (x_i - m)^2 \right\} dx_1, \dots, dx_n \\ &= \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{n}{2\sigma^2} (s^2 + (\bar{x} - m)^2) \right\} dx_1, \dots, dx_n \end{aligned}$$

$$y_1 = \sqrt{n} \bar{x} = \frac{1}{\sqrt{n}} x_1 + \frac{1}{\sqrt{n}} x_2 + \dots + \frac{1}{\sqrt{n}} x_n$$

$$y_2 = \sqrt{\frac{n}{n-1}} (x_1 - \bar{x}) = \sqrt{\frac{n-1}{n}} x_1 - \frac{1}{\sqrt{n(n-1)}} x_2 - \dots - \frac{1}{\sqrt{n(n-1)}} x_n$$

$$y_3 = \dots \dots \dots \sqrt{\frac{n-2}{n-1}} x_2 - \dots - \frac{1}{\sqrt{(n-1)(n-2)}} x_n$$

$$y_n = \dots \dots \dots \sqrt{\frac{1}{2}} x_{n-1} - \frac{1}{\sqrt{2}} x_n$$

なる変数変換をすれば、

$$\left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} (y_2^2 + \dots + y_n^2 + (y_1 - \sqrt{n}m)^2) \right\} dy_1 dy_2 \dots dy_n$$

$$n s^2 = y_1^2 + \dots + y_n^2$$

であるが、こゝで、 \bar{x} , s^2 given の conditional prob を考える爲に

$$y_1 = \sqrt{n} s Z_1, \dots, y_n = \sqrt{n} s Z_n$$

とおくと

$$Z_1^2 + \dots + Z_n^2 = 1$$

こゝで変数を \bar{x} , s , Z_1, \dots, Z_{n-2} に変換すると

$$\frac{2n^{\frac{n}{2}}}{(2\pi\sigma^2)^{\frac{n}{2}}} e^{-\frac{n}{2\sigma^2}(s^2 + (\bar{x}-m)^2)} \frac{s^{n-2}}{\sqrt{1-Z_1^2-\dots-Z_{n-2}^2}} d\bar{x} ds dz_1 \dots dz_{n-2}$$

こゝで、 \bar{x} , s given のときの Z_1 の分布を考える。それ結局

$$\begin{aligned} & \text{Const.} \int \int \dots \int \frac{dz_2 \dots dz_{n-2}}{\sqrt{1-Z_1^2-Z_2^2-\dots-Z_{n-2}^2}} dz_1 \\ & = \text{Const.} \times (1-Z_1^2)^{\frac{n-4}{2}} dz_1 \\ & = \text{Const.} \left(1 - \frac{\tau^2}{n-1}\right)^{\frac{n-4}{2}} d\tau \end{aligned}$$

但
$$\tau = \frac{x_1 - \bar{x}}{s}$$

これはつまり、Thompson の棄却検定である。

この様な簡単な場合には Thompson の棄却限界は Optimum であることが証明され反駁である。

もっと複雑なときには、Thompson の方法は必ずしも良くない。これについては、E.S. Pearson 及 C. Chandra Sekar の前出の論文参

照

Thompson Test の Efficiency については別に述べる積りである。

E. S. Pearson 及 C. Chandra Sekar の前記論文の指摘に従って、Order を考慮した。棄却検法が Grubb の最近の論文⁽⁹⁾で扱われていることを指摘し度い。

しかし、いづれも Power という見地からの研究はしていない。

参 考 文 献

- (1) W.R.Thompson: On a criterion for the rejection of observations and the distribution of the ratio of deviation to sample standard deviation, Ann. of Math. Statist. Vol. 6 (1935), p. 214.
E. S. Pearson and C. Chandra Sekar: The efficiency of statistical tools and a criterion for rejection of outlying observations: Biometrika. Vol. 28 (1936), p. 308.
森口 繁一: 実験データの棄却について, 数学, Vol. 2, No. 1 (1949), p. 65
- (2) 森口, loc. cit. 参照、尤も森口氏の論文もはっきりと棄却域の概念がつかめているとは思われない点がある。
- (3) J. Neyman and E. S. Pearson: On the problem of the most efficient tests of statistical hypotheses. Phil. Trans. of the Roy. Society, London, A. 231. (1933), p. 289
- (4) E. S. Pearson and C. Chandra Sekar. loc. cit. § 2. (2) 参照
- (5) 森口, loc. cit.

- (6) J. Neyman and E. S. Pearson, loc. cit.
J. Neyman : On a statistical problem arising in routine analyses and in sampling inspection in man production : Ann. of Math. Statist. Vol. 12. (1941)
- (7) J. Neyman and E. S. Pearson, loc. cit.
- (8) J. Neyman and E. S. Pearson, loc. cit.
- (9) Frank E. Grubb: Sample Criteria For Testing Outlying Observations Ann. of Math. Statist. Vol. 21. No. 1. March. (1950) pp. 27-58.

(受付 1950. 10. 10)