

④④ 観測値の組分けについて

青 山 博 次 郎

§ 1. 緒 言 R.V. Mises¹⁾ は先に観測値の組分けについて次のような問題を考えた。即ち一つの集団があつて、この各々に試行を施すとその結果一定の数 x (実数) を示し、又集団に属する各々のものはすべて n 個の class の中の一つに属しているものと仮定する。この n 個の class は一定の x に対して n 個の確率密度 $f_1(x), f_2(x), \dots, f_n(x)$ によつて特色づけられているとする。今観測の結果 x を得た時、これはどの class に属しているかを決定したいというのがその問題である。

この応用として林知己夫氏の依沢放についての興味ある論文がある。
2) このでは Mises の立場と若干異なる立場から同様な結果が得られることを示そう。

§ 2. class が二つの場合

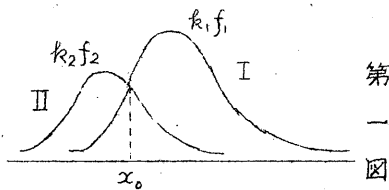
我々は一定の無限母集団を考え、この中には二つの密度函数 $f_1(x), f_2(x)$ によつて特色づけられる classes が混合しているものとする。

即ち

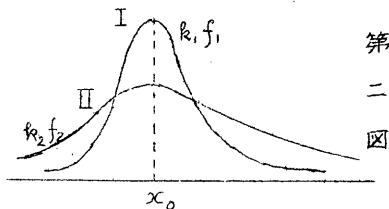
$$\int_{-\infty}^{\infty} f_1(x) dx = 1, \quad \int_{-\infty}^{\infty} f_2(x) dx = 1 \quad (1)$$

$$\int_{-\infty}^{\infty} \{k_1 f_1(x) + k_2 f_2(x)\} dx = 1 \quad (2)$$

こゝに k_1, k_2 は $k_1 + k_2 = 1$ を満足する定数であつて我々には予め分つてゐるものとは限らない。このとき一つの観測値 x を得たとして、一定数 x_0 をとり、 $x > x_0$ ならば f_1 の class に、 $x < x_0$ ならば f_2 の class に属するものと判断することになると、信頼度はいかなる x_0 をえらぶときに最大となるかが問題になる。 $f_1(x), f_2(x)$ は共に unimodal とすると、図の如き場合が考えられよう。



第一図



第二図

信頼度 P を正しい判断をした率で表わすことにすれば

$$P = \int_{x_0}^{\infty} k_1 f_1(x) dx + \int_{-\infty}^{x_0} k_2 f_2(x) dx \quad (3)$$

とおくことができる。

先づ f_1, f_2 共に正規分布でその平均はそれぞれ m_1, m_2 標準偏差はそれぞれ σ_1, σ_2 とし、 k_1, k_2 が既知の場合を考へ

てみよう。信頼度 P を最大とするためには $\frac{\partial P}{\partial x_0} = 0$ より x_0 は

$$\frac{k_1}{\sigma_1} e^{-\frac{(x_0 - m_1)^2}{2\sigma_1^2}} = \frac{k_2}{\sigma_2} e^{-\frac{(x_0 - m_2)^2}{2\sigma_2^2}}$$

を満足することが分る。即ち

$$(\sigma_2^2 - \sigma_1^2)x_0^2 - 2x_0(m_1\sigma_2^2 - m_2\sigma_1^2) + (m_1^2\sigma_2^2 - m_2^2\sigma_1^2 - 2\sigma_1^2\sigma_2^2 \log \frac{k_1\sigma_2}{k_2\sigma_1}) = 0 \quad (4)$$

もし $m_1 > m_2$ ならば、この二根のうち $x_0 > m_2$ なるものをえらぶ。特に $\sigma_1 = \sigma_2 = \sigma$ のときは

$$x_0 = \frac{1}{2} \left\{ (m_1 + m_2) - \frac{2\sigma^2}{m_1 - m_2} \log \frac{k_1}{k_2} \right\} \quad (5)$$

これらの何れの場合に於ても、信頼度は(3)から計算できる。
 以上は k_1, k_2 が既知の場合であるが、もし未知の場合ならば(3)を最大にする如き x, k_1 を定めることにする。即ち

$$P = \frac{k_2}{\sqrt{2\pi}\sigma_2} \int_{-\infty}^x e^{-\frac{(x-m_2)^2}{2\sigma_2^2}} dx + \frac{k_1}{\sqrt{2\pi}\sigma_1} \int_x^{\infty} e^{-\frac{(x-m_1)^2}{2\sigma_1^2}} dx \quad (6)$$

$$\text{但し } k_1 + k_2 = 1$$

$$\frac{\partial P}{\partial k_1} = 0 \quad \text{より}$$

$$\frac{1}{\sigma_1} \int_{x_0}^{\infty} e^{-\frac{(x-m_1)^2}{2\sigma_1^2}} dx = \frac{1}{\sigma_2} \int_{-\infty}^{x_0} e^{-\frac{(x-m_2)^2}{2\sigma_2^2}} dx \quad (7)$$

$$\text{このとき } \frac{\partial P}{\partial x} = 0 \quad \text{より}$$

$$\frac{f_1(x_0)}{f_2(x_0)} = \frac{k_2}{k_1} \quad (8)$$

が得られる。

一般の密度函数については(7)の代りに

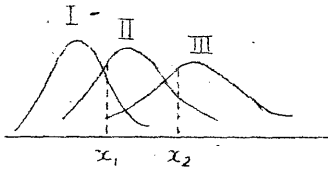
$$\int_{x_0}^{\infty} f_1(x) dx = \int_{-\infty}^{x_0} f_2(x) dx \quad (9)$$

が成立する。

§ 3. class が n 個の場合

n 個の classes の密度函数をそれぞれ $f_1(x), f_2(x), \dots, f_n(x)$ とし、これらがすべて k_1, k_2, \dots, k_n の割合で母集団中に含まれているものとする。

前と同様にして信頼度を



第三圖

$$P = k_1 \int_{-\infty}^{x_1} f_1(x) dx + k_2 \int_{x_1}^{x_2} f_2(x) dx + \dots + k_n \int_{x_{n-1}}^{\infty} f_n(x) dx \quad (10)$$

$$\text{但し} \quad k_1 + k_2 + \dots + k_n = 1 \quad (11)$$

で表わすとき、この P を最大にする如く $k_1, k_2, \dots, k_{n-1}, x_1, \dots, x_{n-1}$ を求めれば、これらの分点によつて

$$x < x_1, \quad x_1 < x < x_2, \quad x_2 < x < x_3, \dots, \quad x_{n-1} < x$$

に属する x をそれぞれ第 1, 第 2, \dots , 第 n 組に属するものと考えればよい。

$$\frac{\partial P}{\partial k_i} = 0, \quad \frac{\partial P}{\partial x_i} = 0 \quad (i = 1, 2, \dots, n-1) \text{ を計算して}$$

$$\int_{-\infty}^{x_1} f_1(x) dx = \int_{x_1}^{x_2} f_2(x) dx = \dots = \int_{x_{n-1}}^{\infty} f_n(x) dx \quad (12)$$

を満足する如く x_1, x_2, \dots, x_{n-1} を定めればよいことが分る。このとき

$$\frac{k_1}{f_2(x_1)} = \frac{k_2}{f_1(x_1)}, \quad \frac{k_2}{f_3(x_2)} = \frac{k_3}{f_2(x_2)}, \dots, \quad \frac{k_{n-1}}{f_n(x_{n-1})} = \frac{k_n}{f_{n-1}(x_{n-1})} \quad (13)$$

が成立する。即ち分点 x_1, x_2, \dots, x_{n-1} に於て

$$\frac{f_v(x)}{f_{v+1}(x)} = \frac{k_{v+1}}{k_v} \quad (v = 1, 2, \dots, n-1) \quad (14)$$

が成立する。

もし k_1, k_2, \dots, k_n が既知のときは (14) を満足する如く x_1, x_2, \dots, x_{n-1} をえらべば, その信頼度は (10) によつて計算される。

§ 4 多変数の場合

n 個の classes が密度函数 $f_1(x_1, x_2, \dots, x_m)$, $f_2(x_1, x_2, \dots, x_m), \dots, f_n(x_1, x_2, \dots, x_m)$ をもつて特長づけられているときも同様にして, 信頼度は

$$P = k_1 \int_{R_1} f_1 dR_1 + k_2 \int_{R_2} f_2 dR_2 + \dots + k_n \int_{R_n} f_n dR_n \quad (15)$$

但し $k_1 + k_2 + \dots + k_n = 1$

で與えられる。こゝに R_v は m 次元のある領域を示し, $R_1 + R_2 + \dots + R_n$ は全 m 次元空間と一致するものとする。

こゝで $\frac{\partial P}{\partial k_v} = 0$ より

$$\int_{R_1} f_1 dR_1 = \int_{R_2} f_2 dR_2 = \dots = \int_{R_n} f_n dR_n \quad (16)$$

なることは直ちに分る。また R_v を種々に変へるとき (15) の最大値を求めるとは相隣る領域 R_v, R_{v+1} にそれぞれ微小な増分及び減分 ΔR を與へるとき $\Delta P = 0$ なるべきことより

$$\frac{f_v(x_1, x_2, \dots, x_m)}{f_{vH}(x_1, x_2, \dots, x_m)} = \frac{k_{v+1}}{k_v} \quad (17)$$

が境界上で成立すること示される。

(註) 1) R.V. Mises: On the classification of observation data into distinct groups, Annals of Math. Statist. Vol. XVI No. 1, 1945

2) 林 知 巳 夫 : , Parole Prediction に 於 ける 統 計 的
方 法 の 一 応 用 に つ い て , 再 犯 調 査 の 基 礎 , ケ ー ス ワ ー ク 研 究 会