

The Estimate of variance and the Precision of Sampling

C. Hayashi

M.D. Ichida

In this paper it is considered to decide the sample size estimating the variance step by step of a sampling process. First drawing n random samples, we estimate the variance and calculate the sampling variance using it.

Then if this variance is over a criterion, we draw more samples and repeat the same process. As the variance is estimated, the sampling variance is not that of population.

This is in question. But if we can properly make the criterion, concerning with the estimated sampling variance the last result will be obtained under the desired level of confidence. Thus we continue the sampling following the criterion until we secure the desired precision from the point of view of the theory of probability.

In order to secure the desired sampling precision in the sense of level of confidence, the criterion must be decided.

The descision of this criterion will be described as follows and the example for the population normally distributed is given.

④

分散の推定とサンプリンクの精度

林 知 己 夫

石 田 正 次

母集団平均推定のためにサンプリンク調査を行ふ場合分散の判明して
ゐないのが通例である。此の様な時あらかじめ準備調査を行つて分散
を推定してゆくやり方があるが、準備調査としてどの位のサンプル数を
とつて分散を推定し、サンプル数を決定してゆくならば、本当の調査の
時散する丈の精度をどの位の信頼度を以て確保することが出来るであら
うか。

此の問題を考へる場合当然推定分散の精度を考慮に入れなければなら
ない。

こゝではあらかじめ準備調査を行ふと云ふ方法ではなく逐次サンプル
をとりながら、そのサンプルから母集団分散を推定し、此の推定分散を
用いて推定平均（サンプル平均）の信頼巾（信頼度一定）を求め此につ
いてサンプル数は十分か否か、少ければさらにサンプルを増加する等の
事を行ふ方法を考へてみよう。

此がある値より小になるならば散する丈の精度（母集団分散を用いた
時の信頼巾）をある信頼度の下に得てゐる筈である。その値より大なら
ば散する精度が得られてゐると言へぬからさらに又サンプルをとり、こ
こで同様の手続きをくりかへしてゆく様なことになる。こゝで又なぐり以上
の様なサンプリンクの方法をとるときある値は如何にきめられるべきで
あらうか。此の問題を少しく考へてゆく事にしよう。

実際問題でこの様な場合はよくあるものである。一例をあげよう。手許に莫大な資料がある。ここからいくつかのサンプルをとり此をしらべある標識についての総平均を一定の精度で推定したいと言ふ様な場合(その事についての分散不明)がさうである。

この様なとき莫大な資料から母集団を構成し、まづ少しのサンプルをとり此より分散を推定し、サンプル平均の精度を計算する。不十分なら又サンプルをとり又精度を計算し又不十分なら-----。

本論にうつらう。

調査対象の標識を x とする。此の各に等しい抽出確率をあてて母集団を構成する。此の大小は N とする。しかし N は十分大であり

$$\frac{N-1}{N-n} \doteq 1 \quad \text{の程度であるとする。}$$

n はサンプル数である。母集団平均を \bar{X} 、分散を σ^2

$$\frac{N_4}{\sigma^4} \text{ を } \beta_2 \text{ とする。ここは } \mu_4 = \sum_{i=1}^N (X_i - \bar{X})^4$$

である。即ち平均のまはりの4次のモーメントである。

母集団から n 個のサンプルを抽出し、母集団平均の推定のためのサンプル平均 \bar{x} をつくる時、 \bar{x} の信頼巾(信頼度一定)は

$$\alpha \sqrt{\frac{\sigma^2}{n}} = \gamma$$

によつてあてられる。

$$\frac{\gamma}{\alpha} = \varepsilon \quad \text{とおくと}$$

$$\frac{\sigma^2}{n} = \varepsilon \quad \text{となる。}$$

ε を相対精度とカリに名づけよう。この ε は欲する精度をきめておけば一定のものである。

$\frac{\sigma^2}{n}$ が ε より小になる様に n を決めねばならぬのである。

さて今 n' 個のサンプルを先づ抽出したとしよう。此から σ^2 の

偏りのない推定値 $S_{n'}^2$ をつくろう。 此から $\bar{x}_{n'}$ (サンプル平均) の相対精度の推定 $\frac{S_{n'}^2}{n'}$ を考える。

此がどの位小であるならば $\frac{\sigma^2}{n} < \varepsilon$ なることが言はれるであらうか。

今 $\frac{S_{n'}^2}{n'}$ が $\varepsilon(n')$ より小ならばある一度の信頼度の下に $\frac{\sigma^2}{n} < \varepsilon$ が言はれるものとしよう。

それならば $\varepsilon(n')$ は如何に定められるであらうか。

まづ $S_{n'}^2$ の分散を考へてみよう。

此が

$$\sigma_{S_{n'}^2}^2 = \frac{\sigma^4}{n} \left(\beta_2 - \frac{n'-3}{n'-1} \right) \quad \text{であることは}$$

よく知られてゐる所である。

$$\sigma_{S_{n'}^2} = \sigma^2 \sqrt{\frac{1}{n} \left(\beta_2 - \frac{n'-3}{n'-1} \right)}$$

さて、 σ^2 の推定のための信頼区間の問題を考へてみよう。此れには單なるシェアアイシェアの定理でなく高能化されたものを用ひよう。

つまり

$$\Pr \{ |y - E(y)| < k \sigma_y \} \geq 1 - \frac{\mu_{2\lambda}}{k^{2\lambda} \sigma_y^{2\lambda}}$$

λ は正の整数 σ_y^2 , $\mu_{2\lambda}$ は y の分散平均のまはりの 2 次のモーメント

$$\text{特に } y = \frac{y_1 + \dots + y_n}{n} \quad \text{と考へるならば}$$

一般に n がさう小でない時近似的に

$$\Pr \{ |y - E(y)| < k \sigma_y \} \geq 1 - \frac{1.3 \dots (2\lambda-1)}{k^{2\lambda}}$$

が言へる。此よりみるときは信頼度 95% ならば $k=2.4$ 程度
 でよいことが言へる。安全をみても $k=3$ とすれば十分信頼出
 来る。なほ、一般の場合でも $k=3$ と考へれば相当信頼ある結
 論である事が了解せられる。

我々の場にもどらう。

又として、 $S_{n'}^2$ を考へ $\sigma_y^2 = \sigma^2$ と考へよう。
 こうするならば サンプルなる $S_{n'}^2$ が

$$\sigma^2 \pm k\sigma^2 \sqrt{\frac{1}{n'} \left(\beta_2 - \frac{n'-3}{n'-1} \right)} \quad (k=2.5 \sim 3)$$

なる巾を逸脱する確率は十分小になつてくるのである。

ましてや此の下限 $\underline{S}_{n'}^2$

$$\underline{S}_{n'}^2 = \sigma^2 \left(1 - k \sqrt{\frac{1}{n'} \left(\beta_2 - \frac{n'-3}{n'-1} \right)} \right) \quad \text{よりも}$$

サンプルの $S_{n'}^2$ が、 k となつて現れる確率はますます小となつてき
 て「此より小なることはない」と言ふ Proposition は十分安全
 なものとなつてくるであらう。

此の下限 $\underline{S}_{n'}^2$ を用いて考へをすすめよう。

$$\frac{\underline{S}_{n'}^2}{n'} = \frac{\sigma^2}{n'} \left(1 - k \sqrt{\frac{1}{n'} \left(\beta_2 - \frac{n'-3}{n'-1} \right)} \right)$$

を考へる。 n' をもつて欲する精度のサンプルが得られるためには

即ち $\frac{\sigma^2}{n'} = \varepsilon$ であるためにはその条件の下で

$$\frac{\underline{S}_{n'}^2}{n'} = \varepsilon(n') \quad \text{ときめてあればよいであらう。}$$

つまり我々のサンプルからつくつた $\frac{\underline{S}_{n'}^2}{n'}$ が此の $\varepsilon(n')$ より

り大でなければサンプリグをやめてよいのである。即ち

$$\varepsilon(n') = \frac{\sigma^2}{n'} \left(1 - k \sqrt{\frac{1}{n'} \left(\beta_2 - \frac{n'-3}{n'-1} \right)} \right)$$

$$= \varepsilon \left(1 - k \sqrt{\frac{1}{n'} \left(\beta_2 - \frac{n'-3}{n'-1} \right)} \right)$$

が、サムアリンクの続行、中止を決定する規準となるのである。

なんとならば

$$\frac{S_{n'}^2}{n'} < \varepsilon \left(1 - k \sqrt{\frac{1}{n'} \left(\beta_2 - \frac{n'-3}{n'-1} \right)} \right)$$

とすれば、 $S_{n'}^2$ は十分安全に (信頼度高く)

$$S_{n'}^2 \geq \sigma^2 \left(1 - k \sqrt{\frac{1}{n'} \left(\beta_2 - \frac{n'-3}{n'-1} \right)} \right)$$

であるから

$$\frac{\sigma^2}{n'} \left(1 - k \sqrt{\frac{1}{n'} \left(\beta_2 - \frac{n'-3}{n'-1} \right)} \right) < \varepsilon \left(1 - k \sqrt{\frac{1}{n'} \left(\beta_2 - \frac{n'-3}{n'-1} \right)} \right)$$

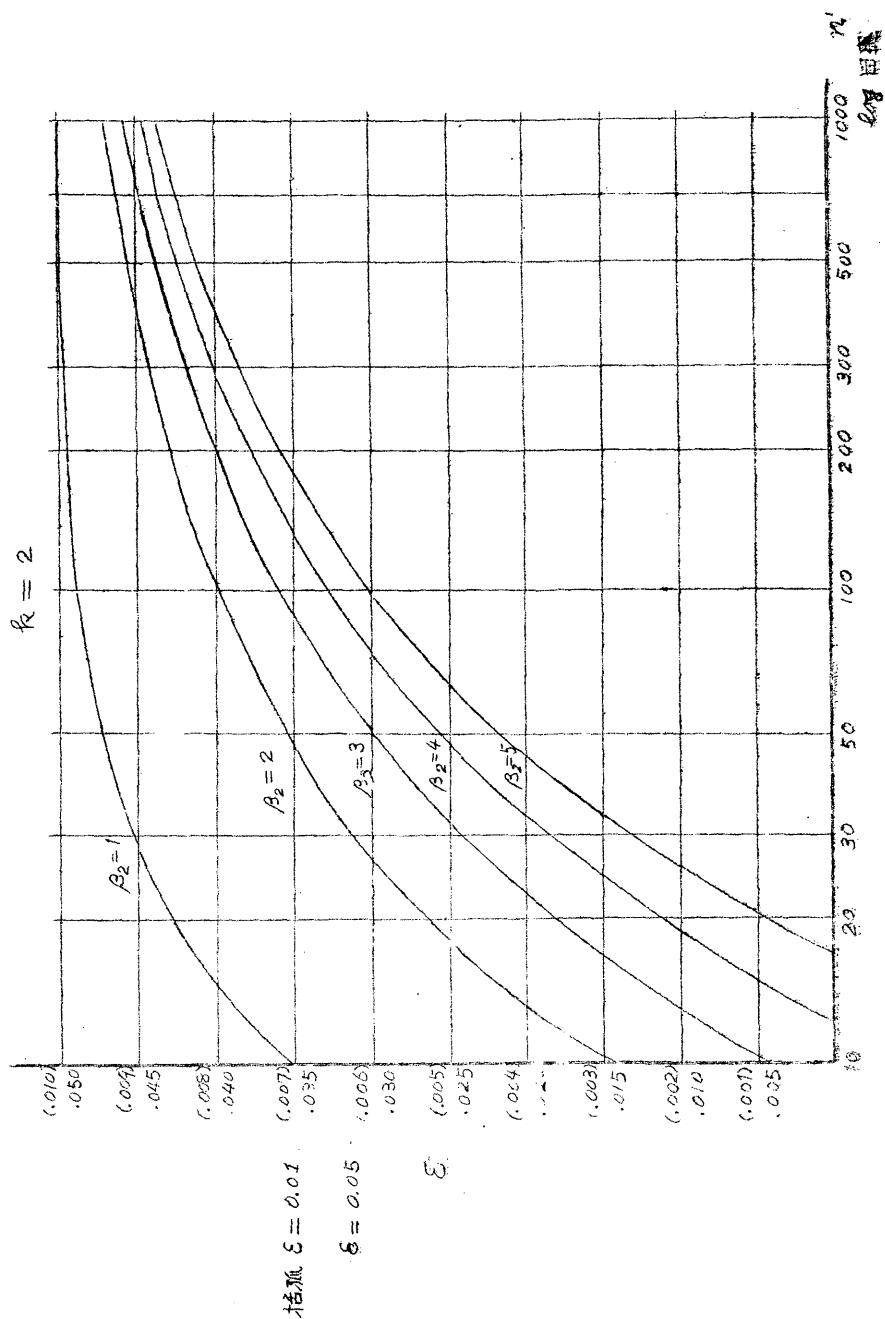
となり、 $\frac{\sigma^2}{n'} < \varepsilon$ となるからである。

此の $\varepsilon(n')$ の曲線を n' の値についてあらかじめつくっておき、一サンプルを n' 個とり出して $\frac{S_{n'}^2}{n'}$ をつくり此の値が上の曲線の上にあるか下にあるかによつてサムアリンクの逐次試行を考へればよいであらう。一般に β_2 は同様未知であるがサンプルの値から一応推定されると考へられるから実際のサムアリンクの用に供す爲め次に k の値、 β_2 の値、 ε の値について $\varepsilon(n')$ の曲線を描いてみよう。

実際の時は k 、 ε の値はきまつてゐるから β_2 の値を推定し其に応ずる $\varepsilon(n')$ をつかへばよいのである。

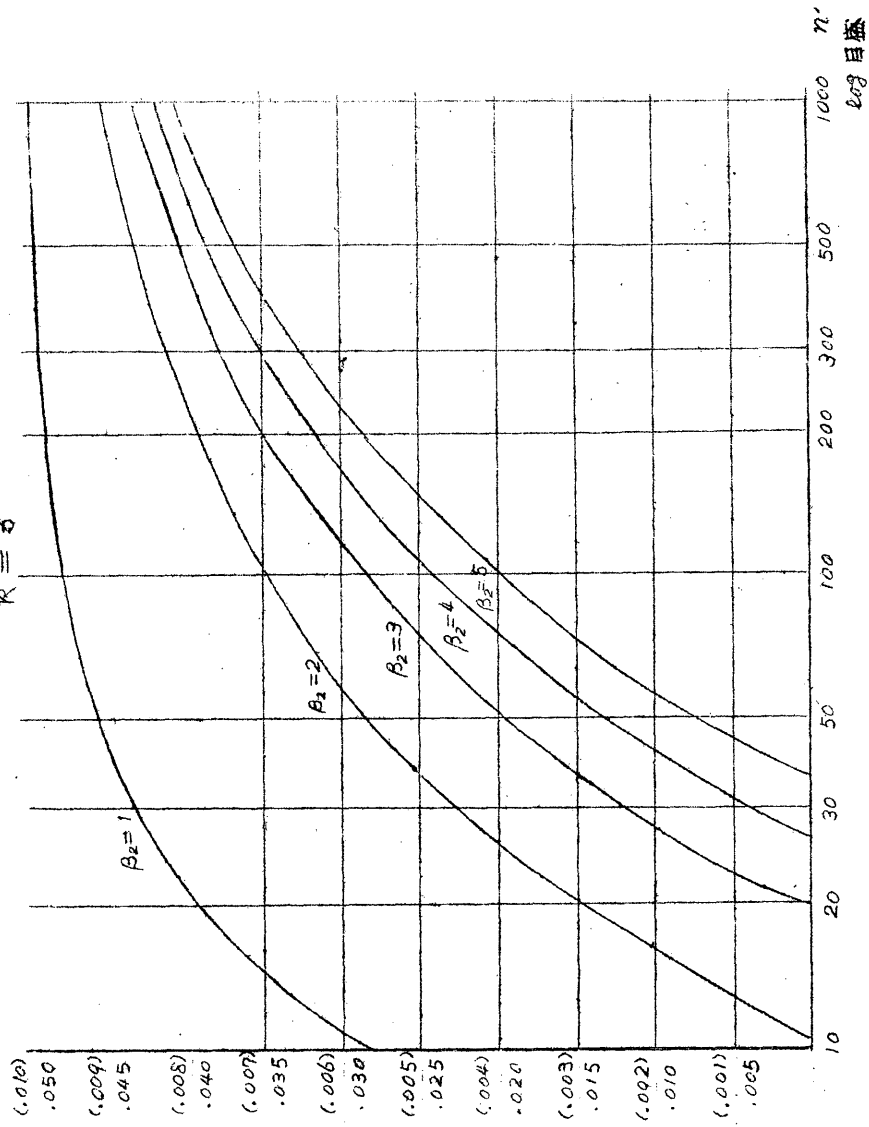
なほ、グラフは n' について \log 目盛によつて描いておいた。

$$\varepsilon(n) = \varepsilon \left(1 - k \sqrt{\frac{1}{n} \left(\beta_2 - \frac{n-3}{n-1} \right)} \right)$$



$$\varepsilon(n') = \varepsilon \left(1 - k \sqrt{\frac{1}{n'} \left(\beta_2 - \frac{n'-3}{n'-1} \right)} \right)$$

$k = 3$



恒取 $\varepsilon = 0.01$

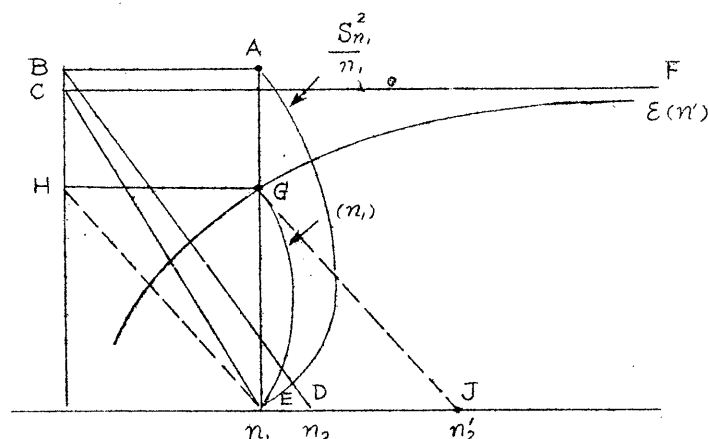
$\varepsilon = 0.05$

ε

n' 個のサンプルをとり $\frac{S_{n'}^2}{n'}$ をつくり此が $\varepsilon(n')$ の曲線より
上にあつたとする時次に何個のサンプルをとるべきかの目安のつけ方
を示しておかう。

まづ n_1 個のサンプルをとる。

A 点は $\frac{S_{n_1}^2}{n_1}$ をあらはす。



$AB \parallel FC \parallel DE \parallel GH$ とする。

此のとき次にとるべきサンプル数は $(n_2 - n_1)$ と目安がつけられ
る。ここに $EC \parallel DB$ である。何とならば、近似的にみて
次の様なことが考えられるからである。

$$\frac{\sigma^2}{n_1} \div \frac{S_{n_1}^2}{n_1}$$

$$\frac{\sigma^2}{n_2} = \varepsilon$$

$$n_2 : n_1 = \frac{S_{n_1}^2}{n_1} : \varepsilon$$

此の方法によつて逐次求めてゆくのであるが或は $HE \parallel GJ$ によ
つて $(n_2' - n_1)$ 個とる目安をつけることもよい。後者の方がよ
り安全な仕組みと考えられる。

次に此の様にしてサンプル数 n' をつくつてゆくとき、 σ^2 を知つてゐる場合にくらべて同一の精度をうるためにどの位と言ふ形でサンプル数は増大するものであらうか、又 n' の分布 (n' が確率変数となることに注意) はどうなるであらうか。

一つの簡単な実例について此の様態をしめしてみよう。

母集団の分布が $\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$ であたへられておると考える。

この時は $\sum_{i=1}^n x_i^2$ は χ^2 分布にしたがふ事に注意されたい。

此の時 $S^2 = \frac{1}{n'} \sum_{i=1}^{n'} x_i^2$ が母集団分散の偏りのない推定値となる。

$$\sigma_{S^2}^2 = \frac{\sigma^4}{n'} (\beta_2 - 1) = 2 \frac{\sigma^4}{n'} \quad \text{となる。}$$

此の様な値をつかひ $\varepsilon(n')$ が決定されるのである。

数値的には χ^2 分布の表を用ひて容易に決定されるのである。

1個サンプルをとり出し $x_1^2 = y_1$ をつくる。

次に、又一個とり出す。云々、 n' 個をとり出し此でサンプリングのおはる確率如何、

$$y_i \text{ の分布は } \frac{1}{\sqrt{2\pi}} (y)^{-\frac{1}{2}} e^{-\frac{y}{2}} dy = p(y_i) dy_i$$

である。一般に

$$\frac{S_{n'}^2}{n'} \leq \varepsilon(n') \quad \text{即ち} \quad S_{n'}^2 < n' \varepsilon(n')$$

ならばサンプリングがおはる。

$$\text{即ち} \quad x_1^2 + \dots + x_{n'}^2 \leq n' \varepsilon(n') = \tau(n')$$

$$y_1 + \dots + y_{n'} \leq \tau(n')$$

ならおはるのである。

したがつて n' 個サンプルを抽出しサンプリングのおはる確率 g_n

は

$$q_{n'} = \int \int \cdots \int_{\substack{\infty > y_1 \geq \tau(1) \\ \infty > y_1 + y_2 \geq \tau(2) \\ \vdots \\ \infty > y_1 + \cdots + y_{n'-1} \geq \tau(n'-1) \\ \tau(n) \leq y_1 + \cdots + y_{n'} \leq 0}} p_1(y_1) p_2(y_2) \cdots p_n(y_n) dy_1 dy_2 \cdots dy_n$$

である。

$$\begin{aligned} \text{今} \quad & \int \int \cdots \int_{\substack{\infty > y_1 \geq \tau(1) \\ \vdots \\ \infty > y_1 + \cdots + y_{n'-1} \geq \tau(n'-1) \\ \infty > y_1 + \cdots + y_{n'-1} + y_{n'} \geq \tau(n')}} p_1 \cdots p_n dy_1 \cdots dy_n = \tilde{r}_{n'} \end{aligned}$$

とおくと $q_{n'} = \tilde{r}_{n'-1} - \tilde{r}_{n'}$ が成立する。 ($n' \geq 1, \tilde{r}_0 = 1$)
さて、此の \tilde{r}_n を、ちりくれ変換によって変数の変換をやれば

$$\tilde{r}_{n'} = \int_{\tau(n')}^{\infty} p(x_n^2) dx_n^2$$

となる。此を用いて $q_{n'}$ を出し n' のモーメントを出さう。

$$E(n) = 1 + \sum_{i=1}^{\infty} \int_{\tau(i)}^{\infty} p(x_i^2) dx_i^2$$

$$E(n^2) = 1 + \sum_{i=1}^{\infty} ((i+1)^2 - i^2) \tilde{r}_i$$

$$E(n^3) = 1 + \sum_{i=1}^{\infty} ((i+1)^3 - i^3) \tilde{r}_i$$

$$E(n^4) = 1 + \sum_{i=1}^{\infty} ((i+1)^4 - i^4) \gamma_i$$

て (n') が求められ定められてゐるから、 γ_n は χ^2 分布の表から計算することができしたがつてモーメントを数値的に計算することができる。此の計算は相当面倒であり且つ χ^2 分布の表が階差をとつて行つてみるとすこぶる数値にあやしい所が多いのできはめて困難を感じた。(χ^2 分布の表の計算を正しくしなほし、且つ詳細なものにする必要が痛感せられる。) したがつて我々の場合 n' のモーメントを決定するのに相当近似的な計算を行つた。

今假りに $\varepsilon = 0.2$ ($\sigma^2 = 1$ であるから $n = 5$ の場合) とつて試みに計算を行つてみると次の様になつた。

平 均	11.5
分 散 σ^2	12.5
三次のモーメント μ_3	-16.8
四次のモーメント μ_4	512.7

を得た。

$$\text{此より} \quad \beta_1 = \frac{(\mu_3)^2}{(\sigma^2)^3} = 0.14$$

$$\beta_2 = \frac{\mu_4}{\sigma^4} = 3.28$$

を得た。 n' の分布の形をピアソン系によつて考へてみるに判別條件によると

$$\text{IV 型} \quad y = y_0 \left(1 + \frac{x^2}{a^2} \right)^m e^{-v \tan^{-1} \frac{x}{a}}$$

と言ふことになる。しかしピアソン系のあてはめには相当疑向があると考へられるから實際に実験せざるかぎり n' の分布の形がこれの型であると言ふのはきはめて危険である。しかし平均、分散 β_1, β_2 よりみて n' の様子はほぼ察せられることであらう。