

統計分析シリーズ(Ⅱ)

茨城大学教授 所 一 夫

Ⅲ 平均と比率の推定 (大標本の場合)

1. 正規分布

前号で標本平均の分布は、標本の大きさ n が大であれば中心極限定理によりほとんど正規分布と見てよい事を示したが、本号ではこの大標本の場合を考えて、正規分布について概説しよう。

正規分布というのはそれに属する単位数は無限であり(無限母集団)、数学的に理想化されたモデルで、現実の問題ではこれに近いと思われる母集団が考えられるものである。そして正規分布に近いと見られる母集団ではその母平均を M 、標準偏差を S 、各単位の持つ数値を x とすれば、 $x=M$ のとき $x \pm d$ (d は小さい正数と考えてよい) の範囲内にはいる母集団の単位数が最大で、 x が M を離れるにしたがってこの数は小さく、全体の分布は $x=M$ を中心として左右対称になっている。

くわしくは平均を M 、標準偏差を S とした正規分布では(これを $N(M, S^2)$ と書く) x が M から、 $M+aS$ の範囲内にある単位数 ($M-aS$ の範囲内にある単位数も同じ) の全体に対する割合を $f(a)$ とするとき、 a の種々の値に対する $f(a)$ の値が表示されている。この一部分を次に示す。

a	0	0.5	1.0	1.5	2.0	2.5
$f(a)$	0.000	0.1915	0.3413	0.4332	0.4772	0.4938
	3.0	3.5	4.0			
	0.4987	0.4998	0.5			

この表によれば正規分布では $M-2S$ と $M+2S$ の間には全体の $0.4772 \times 2S = 0.9544$ すなわち約 95% がはいって居り、 $M-3S$ と $M+3S$ の間には全体の $0.4987 \times 2 = 0.9974$ すなわち約 99.7% がはいっていることがわかる。なおくわしい正規分布表があれば、一般に b と c の間に全体の何%がはいって居るかは、上表の a に対応する値が $(c-M)/S$ などであることから

$$f\left(\frac{c-M}{S}\right) - f\left(\frac{b-M}{S}\right)$$

として上表から算出できる。(ただし $f(-a) = -f(a)$ とする。)

2. 母平均の推定

(i) 母平均の区間推定

大きさ N の母集団から大きさ n のランダムサンプルを抽出しその標本平均 \bar{x} を算出した場合に (a) \bar{x} は母平均 M にはならない、(b) \bar{x} の平均は M である、(c) \bar{x} の標準偏差 $s(\bar{x})$ は $\sqrt{(N-n)/(N-1)} S / \sqrt{n}$ 、(d) n が大きい場合には \bar{x} の分布は正規分布に近い事を前回述べた。

そうすると母平均 M 、母集団の標準偏差が S なる母集

団から大きさ n のランダムサンプルを抽出した場合に、 \bar{x} は種々の値をとるが、正規分布の性質から、これらの値が $M \pm 2s(\bar{x})$ の範囲内にはいっている場合は約 95% で、 $M \pm 3s(\bar{x})$ の範囲内にはいっている場合は約 99.7% と考えられる。このことは \bar{x} が $M \pm 2s(\bar{x})$ の範囲内にはいる確率が 95%、 $M \pm 3s(\bar{x})$ の範囲内にはいる確率が 99.7% であることを示している。

すなわち \bar{x} は M に等しくはないが確率 95% で $M \pm 2s(\bar{x})$ 、確率 99.7% で $M \pm 3s(\bar{x})$ の範囲内にあると判断ができる事を示して居り、この事実を他の面から見ると、知ろうとする数値 M は \bar{x} に等しいとはいえないが、 M が $\bar{x} \pm 2s(\bar{x})$ の範囲内にある確率が 95% で、 $\bar{x} \pm 3s(\bar{x})$ の範囲内にある確率が 99.7% と判断できることを示している。

この場合の範囲をその区間推定に対する信頼区間と

言いその確率を信頼係数と云っている。

以上より標本から算出した平均を \bar{x} 、要求された信頼係数を 95%、母集団の大きさを N 、その標準偏差を S とするとき M の信頼区間は

$$\left(\bar{x} - 2\sqrt{\frac{N-n}{N-1}} \frac{S}{\sqrt{n}}, \bar{x} + 2\sqrt{\frac{N-n}{N-1}} \frac{S}{\sqrt{n}} \right)$$

であり信頼係数を 99.7% まで要求された場合には上の -2 、 $+2$ の代りに -3 、 $+3$ とすればよい。

上の話ではまだ未解決の問題が残っている。それは母集団の標準偏差 S が未知の場合が多いからである。この解決策として実際に用いられている方法は

(a) 過去の資料、実験、または同様な資料について調査した結果から得られた標準偏差 S' をもって S の代用とする。

(b) n が大である場合を考えているので、この場合には抽出した標本の n 個の値から標本標準偏差 s を算出してこれを母集団標準偏差 S の代用とすることである。

例 世帯数 5,000 の町で、大きさ 100 のランダムサンプルを抽出して家計調査を行なった結果、1ヶ月の平均支出は 95,000 円、標準偏差は 10,000 円であった。この結果より全町での平均支出はいくらと推定されるか、信頼係数 95% でその信頼区間を求めよ。

$$\begin{aligned} \text{解 } s(\bar{x}) &= \sqrt{\frac{N-n}{N-1}} \frac{s}{\sqrt{n}} = \sqrt{\frac{5,000-100}{5,000-1}} \\ & \frac{10,000}{\sqrt{100}} = 990 \end{aligned}$$

したがって信頼係数 95% の信頼区間は

$$(95,000 - 2 \times 990, 95,000 + 2 \times 990)$$

すなわち (93,020, 96,980) である。

実際問題では n を大にすれば $s(\bar{x})$ は小さくなり、上

例では次の(ii)で示すまでに、 $n=400$ とすれば $s=10,000$ ならば $2s(\bar{x})=990$ で信頼区間は $95,000$ 円 $\pm 1,000$ 円以内となり、これを表示する場合に $1,000$ 円単位で ± 990 円を略して $95,000$ 円とする場合が多い。

(ii) 標本の大きさ n の決定

上例の場合、 $\bar{x} \pm 2s(\bar{x})$ によつて示された $2s(\bar{x})$ の値を定められた値以下にするためには n をどの程度にすればよいかを調べるのが調査者にとって問題である。

これは結局 $s(\bar{x}) = \sqrt{(N-n)/(N-1)} S / \sqrt{n}$ を与えられた b 以下におさえる問題で

$$\sqrt{\frac{N-n}{N-1}} \frac{S}{\sqrt{n}} < \frac{S}{\sqrt{n}} \leq d \quad \text{より} \quad \frac{S}{d} \leq \sqrt{n},$$

$$\frac{S^2}{d^2} \leq n$$

とすればよい。ここに d は与えられた値であり、 S は未知であるが、前に行なった調査または同様な他の調査より推定するか、または予備調査を行なって少し過大に見積って置いて、この b と S (またはその代用値) より上式によつて n を求めればよい。

上例で $\pm 2s(\bar{x})$ を $\pm 1,000$ とするように n を求めて見ると

$$s=10,000, \quad 2d=1,000 \quad \text{より} \quad d=500$$

$$s^2/d^2 = (10,000)^2 / (500)^2 = 400 \leq n$$

すなわち $S=10,000$ と考えられるならば $n=400$ でよいことがわかる。

3. 比率の推定

(i) 母集団比率の信頼区間

前号で示したように、母集団比率はその属性を持った単位は 1 とし他の単位を 0 とした場合の平均 M である。したがって N 個の中でその属性を持ったものを R 個とすると比率 P は $P=M=R/N$ であり、この場合の標準偏差 S は

$$S = \sqrt{\frac{(1-P)^2 R + (0-P)^2 (N-R)}{N}} =$$

$$\sqrt{(1-P)^2 P + P^2 (1-P)} = \sqrt{P(1-P)}$$
 である。

いま大きさ n のランダムサンプルについての標本比率を p とすれば $p = \bar{x}$ であり標本比率 p についての標準偏差 $S(p)$ は $S(\bar{x})$ であるから

$$s(p) = s(\bar{x}) = \sqrt{\frac{N-n}{N-1}} \frac{S}{\sqrt{n}} = \sqrt{\frac{N-n}{N-1}} \sqrt{\frac{P(1-P)}{n}}$$

したがって標本より求めた比率が p であるとき母集団における比率は信頼係数を 95% とすれば

$$\left(p - 2\sqrt{\frac{N-n}{N-1}} \sqrt{\frac{P(1-P)}{n}}, p + 2\sqrt{\frac{N-n}{N-1}} \sqrt{\frac{P(1-P)}{n}} \right)$$

である。しかしここでも P は未知であるので、 n が大である場合には標本から求めた p によつて P の代用をさせることが多い。すなわち上の信頼区間を

$$\left(p - 2\sqrt{\frac{N-n}{N-1}} \sqrt{\frac{p(1-p)}{n}}, p + 2\sqrt{\frac{N-n}{N-1}} \sqrt{\frac{p(1-p)}{n}} \right)$$

として実用に供するのである。

また P はどのような値であっても常に $\sqrt{P(1-P)} \leq 0.5$ 、 $\sqrt{(N-n)/(N-1)} < 1$ であることより少し範囲を広くとって上の信頼区間を

$$\left(p - 2\sqrt{\frac{0.5}{n}}, p + 2\sqrt{\frac{0.5}{n}} \right) \text{すなわち} \left(p - \frac{1}{\sqrt{n}}, p + \frac{1}{\sqrt{n}} \right)$$

とすれば真の P がこの区間内に存任する確率は 95% 以上となる。(信頼係数は 95% より大とする。)

例 農林省統計情報部がまとめた「46年度農家の形態別にみた農家経済」では全国約 517 万戸の農家より $11,100$ 戸を抽出して諸調査をした結果、稲作経営の場合の所得の農業依存度は 28% と発表された。この結果より信頼係数 95% で全国稲作農家についての農業依存度の信頼区間を求めて見よう。

解 N は大きいので $\sqrt{(N-n)/(N-1)} = 1$ と見る、
 $n=11,100$ $p=0.28$ より

$$2\sqrt{\frac{p(1-p)}{n}} = 2\sqrt{\frac{0.28 \times 0.72}{11,100}} = 2 \times 0.00426 = 0.0085 < 0.01$$

したがって信頼係数 95% の信頼区間は $(0.28 - 0.0085, 0.28 + 0.0085)$ すなわち $(0.272, 0.289)$ となる。この結果より依存度を 28% と発表された意味もわかる。

(ii) 標本の大きさ、の決定

上述より信頼係数 95% で考える場合には $\pm 1/\sqrt{n}$ を与えられた d より小さくしたい場合に n をどの程度に

すればよいかの計算が直ちにできる。すなわち $1/\sqrt{n}$

$\leq d$ として $1/n \leq d^2$ 、 $1/d^2 \leq n$ とすればよい。すなわち \pm の部分を 1% 以下におさえようとすれば上式から $n \geq 10,000$ とすれば十分であることがわかる。

統計ニュース

— 9 月 の 行 事 —

9月 3日～5日 地方統計職員業務研修(専門研修)

4日～6日 刊行物ブロック会議

13日～14日 労働力特調新設集団住宅ブロック会議

14日 統計グラフコンクール

19日～20日 漁業センサス本調査ブロック会議

30日 工作機械調査日