

医学情報処理演習第4回「標本統計量と母数推定」

中澤 港 (nminato@med.gunma-u.ac.jp)

2004年11月1日

我々が統計解析で相手にするのは標本データの場合が多く、得られる統計量も標本統計量であるが、本当に知りたいのは母集団の統計量（母数）であるため、標本統計量から母数を推定しなければならない。その仕組みを考えるため、まず標本抽出をシミュレートしてみる。

標本抽出

100万人の成人男性の身長からなる母集団が、次の表と図のようになっていたとしよう*1。

表1. 成人男子の身長からなる母集団と母平均 μ , 母分散 σ^2 の計算

身長 (cm) (x)	度数	相対度数 ($p(x)$)	$x \cdot p(x)$	$(x - \mu)^2 \cdot p(x)$
144	2	0.000002	0.000288	1.151979e-03
145	3	0.000003	0.000435	1.586970e-03
146	6	0.000006	0.000876	2.903942e-03
147	17	0.000017	0.002499	7.496844e-03
148	19	0.000019	0.002812	7.599834e-03
149	56	0.000056	0.008344	2.021553e-02
150	125	0.000125	0.018750	4.049901e-02
151	219	0.000219	0.033069	6.328937e-02
152	463	0.000463	0.070376	1.185248e-01
153	915	0.000915	0.139995	2.058690e-01
154	1609	0.001609	0.247786	3.153541e-01
155	2649	0.002649	0.410595	4.476659e-01
156	4550	0.004550	0.709800	6.551761e-01
157	7214	0.007214	1.132598	8.728592e-01
158	11005	0.011005	1.738790	1.100452e+00
159	16081	0.016081	2.556879	1.302498e+00
160	22098	0.022098	3.535680	1.414195e+00
161	29903	0.029903	4.814383	1.465155e+00
162	39048	0.039048	6.325776	1.405625e+00
163	48312	0.048312	7.874856	1.207694e+00
164	57703	0.057703	9.463292	9.231469e-01
165	66639	0.066639	10.995435	5.996634e-01
166	73332	0.073332	12.173112	2.932638e-01
167	78051	0.078051	13.034517	7.801682e-02
168	79829	0.079829	13.411272	3.828679e-02
169	77866	0.077866	13.159354	7.790011e-02
170	73767	0.073767	12.540390	2.951326e-02
171	66321	0.066321	11.340891	5.969761e-03
172	57993	0.057993	9.974796	9.279896e-03
173	48410	0.048410	8.374930	1.210356e+00
174	39081	0.039081	6.800094	1.407019e+00
175	29967	0.029967	5.244225	1.468475e+00
176	22055	0.022055	3.881680	1.411597e+00
177	15810	0.015810	2.798370	1.280672e+00
178	10875	0.010875	1.935750	1.087548e+00
179	7309	0.007309	1.308311	8.844242e-01
180	4596	0.004596	0.827280	6.618482e-01
181	2726	0.002726	0.493406	4.607095e-01
182	1519	0.001519	0.276458	2.977333e-01
183	939	0.000939	0.171837	2.112812e-01
184	462	0.000462	0.085008	1.182752e-01
185	224	0.000224	0.041440	6.473767e-02
186	128	0.000128	0.023808	4.147301e-02
187	50	0.000050	0.009350	1.805042e-02
188	31	0.000031	0.005828	1.240027e-02
189	14	0.000014	0.002646	6.174129e-03
190	5	0.000005	0.000950	2.420048e-03
191	4	0.000004	0.000764	2.116040e-03
合計	1000000	1.00	$\mu = 167.9998$	$\sigma^2 = 25.09521$ $\sigma = 5.009512$

*1 こんなにきれいな正規分布に近い表が現実のデータとして存在するはずがなく、実は、これは R の正規乱数を使って得た人工データである。母平均 168, 母標準偏差 5 の正規乱数を 100 万個発生させるためには、`X <- rnorm(1000000, 168, 5)` とすればよい。実行するたびに異なるデータになってしまうのを避けるためには、このコマンドを実行する前に乱数の初期値を `RNGkind("Mersenne-Twister", normal.kind="Inversion"); set.seed(1)` などとして固定すればよい。

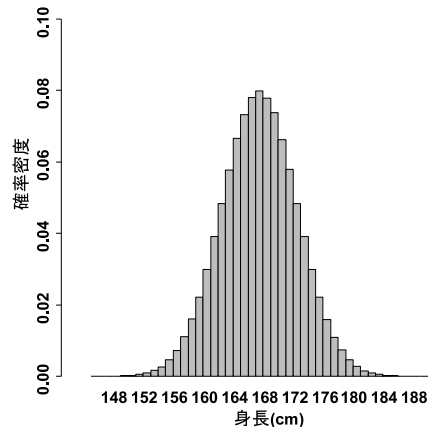


表 1 のように計算すれば， X の平均 $\mu = \sum x \cdot p(x) = 167.9998$ ， X の分散 $\sigma^2 = 25.09521$ である (X は成人男性の身長からなる元の母集団を代表している)。ここから 10 個の標本を抽出することを考える。なるべく母集団全体の情報を偏りなく代表させるためには，身長が書かれた 100 万枚のカードを袋に入れてよくかきまぜ，10 枚取り出すことが考えられる。無作為に抽出される最初のカードが x である確率は， $p(x)$ である。この 1 回目の抽出を X_1 であらわす。2 枚目のカードについても同じである (厳密に言えば，1 枚としては元に戻して再びよくかきまぜて抽出する「復元抽出」でないと，1 枚目に何が出たかに依存して相対度数が変化するため，条件付き確率はまったく同じにはならないが，100 万くらいの大きな母集団からサイズ 10 の標本を取り出すのであれば，「非復元抽出」(1 度サンプルを取り出したら元に戻さない) であっても， $p(x)$ はほとんど変化しないし，1 枚目に何がでるかという場合を総当りして合計すれば，2 枚目に何が出るかという確率も厳密に $p(x)$ となる)。このように，無作為に選ばれるカードは，すべて母集団における相対度数によって与えられる確率で，母集団の値のどれかをとりうる，確率変数である。

サイズ 10 の標本の標本平均 \bar{X} は， $\bar{X} = \frac{1}{10}(X_1 + X_2 + \dots + X_{10})$ なので，その期待値は，

$$E(\bar{X}) = \frac{1}{10}(E(X_1) + E(X_2) + \dots + E(X_{10})) = \frac{1}{10}(\mu + \mu + \dots + \mu) = \frac{1}{10} \cdot 10 \cdot \mu = \mu$$

となって母平均に一致する。標本平均の分散は

$$\begin{aligned} \text{var} \bar{X} &= \text{var} \left\{ \frac{1}{10}(X_1 + X_2 + \dots + X_{10}) \right\} \\ &= \left(\frac{1}{10} \right)^2 \{ \text{var}(X_1) + \text{var}(X_2) + \dots + \text{var}(X_{10}) \} \\ &= \left(\frac{1}{10} \right)^2 (\sigma^2 + \sigma^2 + \dots + \sigma^2) = \left(\frac{1}{10} \right)^2 \cdot 10 \cdot \sigma^2 \\ &= \frac{\sigma^2}{10} \end{aligned}$$

したがって，標本平均の標準偏差は， $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{10}}$ となる。これは，標本から推定される母集団の平均値の信頼性を示す値なので，平均の標準誤差と呼ばれる。

R で標本抽出をする関数は，`sample()` である。`replace=オプション` で復元，非復元を選択できる (デフォルトは非復元)。

```

X <- c(rep(144,2),rep(145,3),rep(146,6),rep(147,17),rep(148,19),rep(149,56),
rep(159,125),rep(151,219),rep(152,463),rep(153,915),rep(154,1609),rep(155,2649),
rep(156,4550),rep(157,7214),rep(158,11005),rep(159,16081),rep(160,22098),
rep(161,29903),rep(162,39048),rep(163,48312),rep(164,57703),rep(165,66639),
rep(166,73332),rep(167,78051),rep(168,79829),rep(169,77866),rep(170,73767),
rep(171,66321),rep(172,57993),rep(173,48410),rep(174,39081),rep(175,29967),
rep(176,22055),rep(177,15810),rep(178,10875),rep(179,7309),rep(180,4596),
rep(181,2726),rep(182,1519),rep(183,939),rep(184,462),rep(185,224),
rep(186,128),rep(187,50),rep(188,31),rep(189,14),rep(190,5),rep(191,4))
par(mfrow=c(2,2))
hist(s10 <- sample(X,10))
hist(s100 <- sample(X,100))
hist(s1000 <- sample(X,1000))
hist(s10000 <- sample(X,10000))

```

サンプルサイズが大きくなるにつれて、標本分布は母集団の分布に近づく。しかし、標本平均の分布については、もっと多くのことがいえる。

中心極限定理

ほとんどの任意の母集団から抽出された無作為標本の平均 \bar{X} の分布は、標本の大きさが増大するにつれて、平均 μ 、標準偏差 σ/\sqrt{n} の正規分布に近づくことがわかっている。これを中心極限定理という。R で試してみよう。

```

X <- c(runif(5000,0,100),rnorm(5000,100,10))
par(mfrow=c(2,2))
hist(X,xlim=c(0,140))
Z5 <- rep(0,1000)
for (i in 1:1000) { Z5[i] <- mean(sample(X,5)) }
# サンプルサイズ 5 ではまだ変な形
hist(Z5,xlim=c(0,140))
Z10 <- rep(0,1000)
for (i in 1:1000) { Z10[i] <- mean(sample(X,10)) }
hist(Z10,xlim=c(0,140))
Z30 <- rep(0,1000)
for (i in 1:1000) { Z30[i] <- mean(sample(X,30)) }
hist(Z30,xlim=c(0,140))

```

サンプルサイズが大きくなるほど、正規分布に近づくと同時に、標本平均の標準偏差（つまり標準誤差）が小さくなっていき、ヒストグラムの幅が狭くなっていくのが一目瞭然である。

信頼区間

標本平均の期待値は、母平均に一致するので、標本平均は、母平均のよい点推定量になっている。しかし、必ず一致するとは限らない。そこで、それがどれくらい確からしい推定かということを示すことを考える。通常、標本から計算される、ある区間の中に、正しく μ を含む割合が 95% であるような区間を推定したい（つまり、95% の信頼度をもって区間推定をしたい）と考える。

標本が大きい場合は、 \bar{X} の正規分布において、ちょうど 95% の確率を囲むような最短の範囲を選択すればよい。そのために、正規分布の両側の裾 2.5% 分を除く中央部分をとることができる。この部分は、母平均 μ を中心として、大きい方と小さい方に、 \bar{X} の標準偏差（つまり標準誤差）の 1.96 倍（1.96 は標準正規分布の 97.5% 点である）だけ動かした範囲を含む。つまり、

$$Pr\left(\mu - 1.96 \frac{\sigma}{\sqrt{n}} < \bar{X} < \mu + 1.96 \frac{\sigma}{\sqrt{n}}\right) = 0.95$$

である。この式を変形すると、

$$Pr\left(\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}}\right) = 0.95$$

となる。 σ は未知なので、通常、標本の不偏標準偏差 $sd(X)$ を用いる。つまり、 μ の 95% 信頼区間は、 $\bar{X} - 1.96sd(X)/\sqrt{n}$ から $\bar{X} + 1.96sd(X)/\sqrt{n}$ までとなる。

標本が小さい場合は、未知の σ に対して $sd(X)$ を代入することが無視できない誤差の原因となる。そこで、 $sd(X)$ を使って定義される t という統計量、

$$t = \frac{\bar{X} - \mu}{sd(X)/\sqrt{n}}$$

が自由度 $n - 1$ の t 分布に従うことから、 t 分布の 2.5% 点から 97.5% 点までを 95% 信頼区間とすることで、この誤差を回避することができる。自由度が無限大になれば t 分布は正規分布に一致するので、結局常にこちらで計算すればよいことになる。

すなわち、標本サイズ n 、標本平均 \bar{X} 、標本の不偏標準偏差 $sd(X)$ のとき、母平均の 95% 信頼区間は、

$$\bar{X} - t_{0.025}sd(X)/\sqrt{n}$$

から

$$\bar{X} + t_{0.025}sd(X)/\sqrt{n}$$

までになる。R で自由度 $n - 1$ の t 分布の 97.5% 点を与える関数は `qt(0.975, n-1)` なので、R のプログラム上で、例えば上で定義した `s100` から母平均の 95% 信頼区間を推定するには、次のようにすればよい。

```
barX <- mean(s100)
sdX <- sd(s100)
t975 <- qt(0.975, length(s100)-1)
rootn <- sqrt(length(s100))
print(barX - t975*sdX/rootn)
print(barX + t975*sdX/rootn)
```

この結果から、95%信頼区間を表示するときは、[166.96, 168.96] のように記載するのが普通である。ただし、もちろん、R にはもっと楽にこの計算をしてくれる関数が用意されていて、標本データが付値されている数値型変数 X について、 $t.test(X)$ とすれば、母平均がゼロであるという帰無仮説の検定結果（検定については後日説明する）とともに、母平均の推定値と母平均の 95%信頼区間が表示される。

自由度

自由度について直感的に理解するためには、次のように考えると良い。 n 個の数値からなる標本には、はじめ n の自由度がある。しかし、自由度の 1 つは、 \bar{X} を計算するときに使われてしまい、 $sd(X)$ を計算するための偏差 $X_i - \bar{X}$ に対しては $n - 1$ の自由度しか残らない。逆に見れば、サイズ n の標本について、はじめの $n - 1$ 個の偏差は自由だが、最後の偏差は、全偏差の合計がゼロにならねばならない（そのように平均という量はとられている）ために決定済みであり、結局自由に決められる個数は $n - 1$ となる。

一般には、データの数から母数の数を引いた値が自由度になる。

課題

2001 年に厚生科学研究で行われた「少子化の見通しに関する専門家調査」の結果の一部を見てみる。この調査は、「人口学、経済学、家族社会学、公衆衛生学を中心とした専門家を対象として少子化研究会のメンバーが対象候補者を抽出し、回答者の偏りや不足等について検討を加えた上で、748 名を対象として調査を実施した」もので、回収率は 44 % であった。この調査では、2025 年の合計出生率がいくつになるかという推定値が尋ねられていて、生データを見ると、

1.38 1.50 1.30 1.40 1.40 1.15 1.31 1.37 1.50 1.55 1.55 1.56
1.50 1.56 1.50 1.38 1.50 1.20 1.20 1.50 1.25 1.25 1.22 1.40
1.80 1.37 1.35 1.70 1.35 1.50 ... (後略)

のようになっていて、回答数は 311、平均値は 1.385、不偏分散は 0.0252 であった。以上の情報から、専門家母集団における、2025 年の合計出生率の推定値の 95%信頼区間を求めよ。

結果は配布する紙に学籍番号、氏名と共に自筆して提出すること。結果の提出をもって出席確認とする。