

なぜ不偏分散は $N-1$ で割るのか *

小杉考司 †

1 母集団と標本

1.1 サンプリング

推測統計学の基本は、全体の情報を知るために、全ての要素をチェックしていくのは大変だから、少ないサンプルをもとに全体的特徴を推し量ろう、という考え方である。選挙でどこの政党がもっとも票を集めるのか、を知るために、日本国民全員に聞いていくのでは選挙を一度やる苦労と変わりがない。そこで小数の電話調査などで、だいたいの当たりをつけるのである。

このとき、全体のことを母集団といい、(すでに出てきたが) そこから集められる少数のデータのことをサンプルあるいは標本という。

問題は、サンプルが母集団の特徴をきちんと反映しているかどうか、である。母集団からサンプルを集める方法をサンプリングという。普通、サンプリングはランダム・サンプリング(無作為抽出)がなされる。ランダムとは、無作為、つまり調査者の意図が入っていないということである。調査者が自分に都合の良いようなサンプルを集めたら(例えば自民黨員ばかりに支持政党調査を行ったら)、母集団をうまく反映しない結果が出るのは容易に想像できるだろう。

サンプリングの方法については、様々なものがあるので他書に譲る。

1.2 なぜサンプルで正しいことがわかるか

1.2.1 経験的な説明

さて、では母集団からサンプルを取り出すと、何が分かるのだろうか。例を挙げてみてみよう。

$N(50, 10)$ の正規分布に従う乱数を用いて、100 個のデータを作ってみた (1)。100 人の学生さんによるテストの点数だとでも思ってくれればよい*¹。この 100 人のデータから、10 人分ずつサンプルを取ったとしよう。ここでは行の 10 人、あるいは列の 10 人を取ったとして、それぞれの平均値を算出してみた。サンプルの取り方によっては、57 点のような取り方もあれば、40 点になるような取り方もある。しかし、これら「サンプルの平均」の平均は 49.6 とほぼ 50 に一致していることがわかるだろう。

ところで、サンプルの平均の散らばりを表す、サンプル平均の標準偏差のことを、特に「標準誤差」とよぶ。

* written on 2008.02.06 / revised on 2015.02.26

† correspondence to kosugi@yamaguchi-u.ac.jp or kosugitti@gmail.com

*¹ これは Excel の分析ツールにある、「乱数の発生」を使っている。便利な世の中になったものだ。

表1 サンプル平均と母平均

47	24	45	61	43	57	49	55	47	51	→	47.84
37	64	70	38	33	53	45	50	46	58	→	49.48
52	37	59	34	32	41	57	39	63	59	→	47.31
63	43	74	57	40	48	46	32	49	44	→	49.61
62	58	43	56	42	51	58	58	48	41	→	51.77
67	55	67	72	29	56	36	54	50	47	→	53.22
28	59	34	64	44	51	42	56	47	42	→	46.70
48	56	55	63	46	41	35	52	72	42	→	50.95
61	36	59	51	51	69	46	40	33	46	→	49.20
39	39	69	50	46	55	50	62	43	45	→	49.86
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓		
50.46	47.14	57.45	54.77	40.70	52.12	46.23	49.91	49.79	47.35		

サンプルを繰り返せば、母集団の情報に近づいていく - これがサンプリングをするメリットである。特に、母集団がやたらと大きい場合は、母集団の悉皆調査が実際的に不可能であったとしても、サンプルを数回取り出すことぐらいなら、何とか実現可能な範囲内にあるはずだからである。

1.2.2 理論的な説明

上の例を、理論的にきちんとフォローしておこう。

サンプルを取ってくると、その平均値 (標本平均) は簡単に求められる。すなわち、

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n}(x_1 + x_2 + \dots + x_n) \quad (1)$$

同様に、サンプルの分散 (標本分散) も次式で得られる。

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (2)$$

これは以下のように変形できる。通常の計算はこちらの方が楽だろう。

$$\begin{aligned} s^2 &= \frac{1}{n} \{(x_1^2 - 2x_1\bar{x} + \bar{x}^2) + \dots + (x_n^2 - 2x_n\bar{x} + \bar{x}^2)\} \\ &= \frac{1}{n} \{(x_1^2 + x_2^2 + \dots + x_n^2) - 2\bar{x}(x_1 + x_2 + \dots + x_n) + n\bar{x}^2\} \\ &= \frac{1}{n} \{\sum x_i^2 - 2\bar{x}(n\bar{x}) + n\bar{x}^2\} \\ &= \frac{1}{n} \sum x_i^2 - \bar{x}^2 \end{aligned} \quad (3)$$

余談だが、共分散も同様に

$$s_{xy} = \frac{1}{N} \sum x_i y_i - \bar{x}\bar{y} \quad (4)$$

で求められる。

さて、推測統計においては、サンプル x は母集団分布に従う確率変数 X が x という形に具現化したもの、と考える。ここで、標本の大きさが n の時、その期待値は

$$E[\bar{X}] = E\left[\frac{1}{n}(X_1 + X_2 + \dots + X_n)\right] = \frac{1}{n}\{E[X_1] + E[X_2] + \dots + E[X_n]\}$$

であるから、

$$E[\bar{X}] = \frac{1}{n}(\mu + \mu + \dots + \mu) = \mu \quad (5)$$

である。これは、標本平均の期待値が母平均に等しいことを意味しており、サンプルを何度も取り、その平均値の平均値は母平均に等しくなることが理論的に示される。

ところで、後々必要になってくるから、ここで標本を繰り返したときの、平均値の散らばりについて考えておく。標本平均 \bar{X} の散らばりの期待値だから、

$$\begin{aligned} E[(\bar{X} - \mu)^2] &= E\left[\left\{\frac{1}{n}(X_1 + X_2 + \cdots + X_n - n\mu)\right\}^2\right] \\ &= E\left[\left\{\frac{1}{n}\{(X_1 - \mu) + (X_2 - \mu) + \cdots + (X_n - \mu)\}\right\}^2\right] \\ &= \frac{1}{n^2} \sum_{i=1}^n E[(X_i - \mu)^2] \\ &= \frac{1}{n^2} \underbrace{(\sigma^2 + \sigma^2 + \cdots + \sigma^2)}_{n \text{ 個}} \end{aligned}$$

となり、

$$= \frac{1}{n} \sigma^2 \quad (6)$$

であることがわかる。

では次に標本分散を見てみよう。この例の全体の分散は 116.68 である。先ほどと同じように、各列を $n=10$ のサンプルだと考えて計算してみると 10 サンプル得られるから、これの平均を出してみると、96.951 となり、ずいぶんと外れたものになってしまっている。ただのサンプリングミスだろうか？

では元の式に戻って、考え直してみよう。標本の大きさが n のとき、標本分散 S^2 の期待値を母分散で表したい。

$$S^2 = \frac{1}{n} \{(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \cdots + (X_n - \bar{X})^2\}$$

で、 $\sigma^2 = E[(X - \mu)^2]$ である。なんとかして $X - \mu$ を式の中に入れたいので、

$$S^2 = \frac{1}{n} \{(X_1 - \mu - \bar{X} + \mu)^2 + \cdots + (X_n - \mu - \bar{X} + \mu)^2\}$$

と置こう。これを展開すると、

$$\begin{aligned} &= \frac{1}{n} \{[(X_1 - \mu) - (\bar{X} - \mu)]^2 + \cdots + [(X_n - \mu) - (\bar{X} - \mu)]^2\} \\ &= \frac{1}{n} \{(X_1 - \mu)^2 - 2(X_1 - \mu)(\bar{X} - \mu) + (\bar{X} - \mu)^2 + \cdots + (X_n - \mu)^2 - 2(X_n - \mu)(\bar{X} - \mu) + (\bar{X} - \mu)^2\} \\ &= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - 2 \frac{1}{n} \sum_{j=1}^n (X_j - \mu)(\bar{X} - \mu) + (\bar{X} - \mu)^2 \end{aligned}$$

ここで第二項は、

$$\begin{aligned} -2 \frac{1}{n} \sum_{j=1}^n (X_j - \mu)(\bar{X} - \mu) &= -2(\bar{X} - \mu) \frac{1}{n} \sum_{j=1}^n (X_j - \mu) \\ &= -2(\bar{X} - \mu)(\bar{X} - \mu) \\ &= -2(\bar{X} - \mu)^2 \end{aligned}$$

だから、元の式は

$$\begin{aligned}
 &= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - 2(\bar{X} - \mu)^2 + (\bar{X} - \mu)^2 \\
 &= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - (\bar{X} - \mu)^2
 \end{aligned}$$

となる。これの期待値、 $E[S^2]$ は

$$\begin{aligned}
 E[S^2] &= \frac{1}{n} \sum_{i=1}^n E[(X_i - \mu)^2] - E[(\bar{X} - \mu)^2] \\
 &= \sigma^2 - E[(\bar{X} - \mu)^2]
 \end{aligned}$$

である。これの第二項は、標本平均の分散であり、 σ^2/n で得られるのだったから (式 6)、 $E[S^2]$ は、

$$\begin{aligned}
 E[S^2] &= \sigma^2 - \frac{1}{n}\sigma^2 \\
 &= \frac{n-1}{n}\sigma^2
 \end{aligned}$$

となる。さて、標本分散の期待値が $\frac{n-1}{n}\sigma^2$ であるから、 $\frac{1}{n} \sum (X_i - \bar{X})^2$ で求めた標本分散とずれていることがわかるだろう。標本分散は

$$\frac{1}{n} \sum (X_i - \bar{X})^2 \times \frac{n}{n-1} = \frac{1}{n-1} \sum (X_i - \bar{X})^2$$

で求めるべきであり、このようにして求める分散のことを特に不偏分散という。

分散だけ $n-1$ で割るのは、どうも不公平な感じがする、という方がいるかもしれない。試しに、実際のデータで見てみよう。先ほどの表 1 の例では、母分散が 116.68 であるのに、標本分散が 96.951 であった。では、不偏分散の平均値を取ってみよう。表 2 にあるように、不偏分散を

表 2 標本分散と不偏分散

	47	24	45	61	43	57	49	55	47	51	
	37	64	70	38	33	53	45	50	46	58	
	52	37	59	34	32	41	57	39	63	59	
	63	43	74	57	40	48	46	32	49	44	
	62	58	43	56	42	51	58	58	48	41	
	67	55	67	72	29	56	36	54	50	47	
	28	59	34	64	44	51	42	56	47	42	
	48	56	55	63	46	41	35	52	72	42	
	61	36	59	51	51	69	46	40	33	46	
	39	39	69	50	46	55	50	62	43	45	分散の期待値
nで割る	150.327	150.276	155.649	124.092	47.747	60.431	52.904	83.666	105.838	38.581	96.951
n-1で割る	167.030	166.973	172.944	137.880	53.053	67.145	58.782	92.962	117.598	42.868	107.724

10 サンプル取ったときの平均値は 107.724 と、幾分母分散に近づいている。

これは 10 サンプルの例であるが、サンプル数を増やせばもっとわかりやすい。図 1.2.2 は、サンプル数を増やしていったときの標本分散と不偏分散の平均値の推移である。この図の描き方は以下の通りである。まず、 $N(50, 10)$ の乱数を 5000 ケース発生させ、10 ケースで 1 サンプルとして 4990 サンプルを得、標本分散と不偏分散を求めた*2。次に、2 サンプル、3 サンプル・・・4990

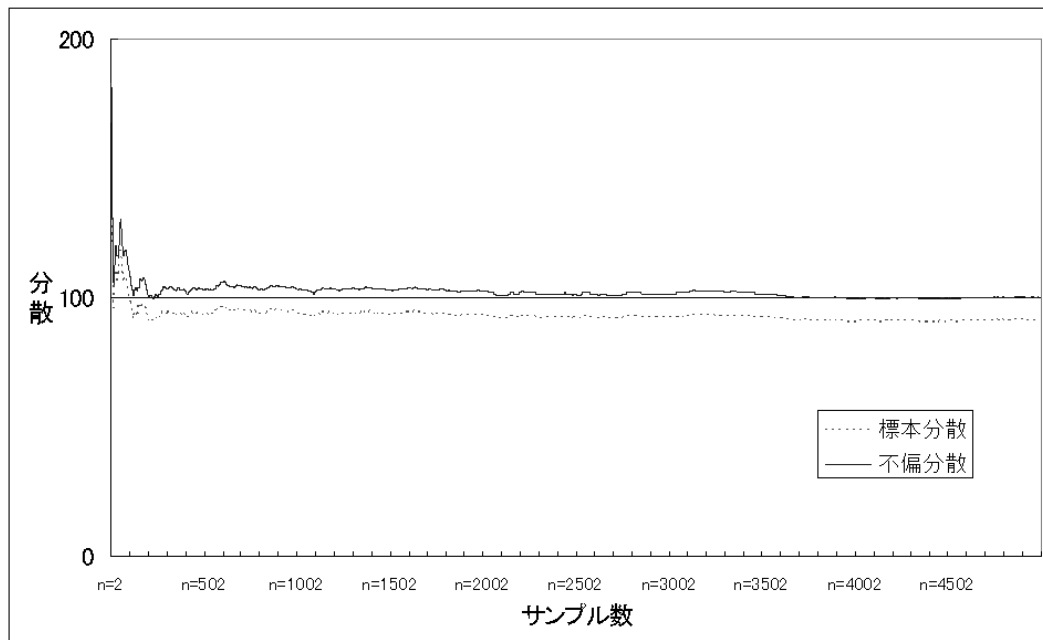


図1 標本分散と不偏分散の期待値

サンプルの平均値を求め、グラフ化する。図 1.2.2 にあるように、標本分散は母分散 100 から外れたところを推移し、決して母分散に近づくことはない。最終的に、標本分散の期待値は 91.26、不偏分散の期待値は 100.38 と、後者がより母分散に近いことは明らかだ^{*3}。

分散を推定値と見なす場合は、以後 $\frac{1}{n-1}$ で除したものをを用いる。

1.3 必要なサンプル数の求め方

1.4 推定のひみつ

さて、サンプリングからいくつかの基本的なカラクリが見えてくる。

まず重要なのは、サンプルの平均値 \bar{x} の期待値は、母集団の平均値 \bar{X} と一致するということ(式 5)。繰り返しになるが、何度もサンプルを繰り返して、サンプル平均の平均をとると、母平均になるということ。これを仮にサンプリングの第一定理と呼ぼう。

次に、サンプルの平均値の分散は、母分散を標本の大きさに割ったものに等しくなるということ(式 6)。同じく本書ではこれを第二定理と呼ぶ。この第二の定理からは、もうひとつ重要なことがわかる。サンプル平均の標準偏差、つまり何度も繰り返して取られるサンプルの、平均の散らばり具合は、 $\sqrt{\sigma/n}$ に一致するとのこと。この式の分母に入っているのは、サンプルの大きさ

^{*2} Excel では、標本分散を求める関数が *VARP*、不偏分散を求める関数が *VAR* である。普通 *VAR* と入れてしまいうようになるが、これは Microsoft が「分散を求めるようなヤツは推測統計をやるもの、と決めてかかっているからか？」

^{*3} 母分散はこのとき 100.16 である。乱数なので、母分散が 100 というのは理論値ではない。

n (の平方根)である。分母が大きくなればなるほど、ここで得られる値は小さくなる。つまり、サンプルのサイズが大きくなればなるほど、精度の高い母平均の推定が出来ることを意味している。これも一緒に頭の中に入れておこう。

第三に、サンプルの平均値 \bar{x} は、母集団の平均値 \bar{X} を中心にした正規分布に従うこと(第三定理)。

これら三つの特徴から、サンプルの平均値 \bar{x} がわかれば、母平均は $N(\bar{x}, \sigma^2/n)$ の正規分布のどこかにあることがわかる。正規分布の確率密度関数は、理論的に明らかであるから、??ページで触れたように、 $\pm 1\sigma$ の範囲に全体の何 % が含まれるか、ということが算出できる。例えば、 $\bar{x} \pm \sqrt{\sigma/n}$ の範囲内に母平均が入る確率が、67.3% あるということだ。

この辺りに推測統計学のカラクリが潜んでいる。

母集団から、サンプルを取ってきたとする。サンプル平均 \bar{x} とその標準偏差 σ が得られる。さて、ここで母平均を求めたいとする。サンプルを何度も取り、その平均値をどんどん均していけば、上の第一定理により母平均に一致するはずだ。でもサンプルを何度も取るのは大変。だから、一回のサンプルを信じて、サンプル平均 \bar{x} が母平均と一緒に考えよう(これが推測の第一歩。手抜きの手続きでもある)。ただどさずがに、一回のサンプル平均が母平均とぴったり一致するなんてコトはありそうにない。だから、あんまり「これ(サンプル平均)が母平均なんです」と言い切るのはよろしくない。そこで、第二定理を使う。母分散がわかっているならば、 $\sqrt{\sigma/n}$ の散らばりをもつ正規分布に従うのだから、「一回のサンプル平均 \bar{x} の、周囲 $\bar{x} \pm \sqrt{\sigma/n}$ の範囲内に母平均が入る確率が、67.3% ある。 $\pm 2\sqrt{\sigma/n}$ の範囲内には 95.3% の確率で入ってくる。」と言えるだろう。しかし、ここにも問題がある。普通、母分散なんてわからないのだ。だから、その時はサンプルの分散 S^2 が母分散の推定値だと考えて(推測の第二歩、さらに足下が怪しくなる)、サンプルの分散を使いながら、母平均の入りそうな範囲を確率と共に推定する。^{*4}

1.5 有限母集団からのサンプリング

ところで、サンプルの分散を、そのまま推定値として使うのには、ちょっと問題がある。今までのサンプルと母集団の関係式(式6など)は、母集団が無限であることを前提とした式になっている。しかし、社会調査のシーンで使われる母集団は有限(人類全体とか)であるから、有限を前提とした式になるように、修正が必要である。少し横道に逸れるが、以下にそのプロセスを辿ってみよう。

まず、サンプル平均の期待値の算出について^{*5}。

母集団の要素が X_1, X_2, \dots, X_n である(有限ですね)とき、ここから大きさ n のサンプルを取り出す取り出し方は、 ${}_N C_n$ だから $N!/(N-n)!n!$ 通りある。この中から、たまたま X_1 がサンプルされて取り出された、としよう。このような取り出され方は何回あるだろうか？

n 個のサンプルのうち、ひとつが X_1 で、残り $n-1$ 個は X_1 以外の何であっても良い。 X_1 以外を取り出すのは、 ${}_{N-1} C_{n-1}$ 通りである。このことに注意した上で、サンプル平均の期待値を求めよう。全ての平均値の平均値だから、

^{*4} このように、推測の幅を持たせることを区間推定といい、そうでない点推定とは区別する。

^{*5} 母分散の推定値についての項目なのだが、計算の途中でサンプル平均が関係してくるので、まずそれをやっつけておくのです。

$$E(\bar{x}) = \frac{1}{N C_n} \left\{ \frac{1}{n} (X_1 + X_2 + \cdots + X_n) + \frac{1}{n} (X_1 + X_2 + \cdots + X_{n-1} + X_{n+1}) \right. \\ \left. + \cdots + \frac{1}{n} (X_{N-n+1} + X_{N-n+2} + \cdots + X_N) \right\}$$

右辺第二の項は、母集団の X_n 番目の要素がないサンプルである。それ以降の項は同様に、 n 個のサンプルにおいて、欠けている要素がひとつずつズレながら、 X_N 番目の要素まで続いている。この中で、 X_1 が入っている項は ${}_{N-1}C_{n-1}$ 個、 X_2, X_3, \dots, X_N もそれぞれ、 ${}_{N-1}C_{n-1}$ 個入っている。とすると、この式の $\{$ の内側は、 ${}_{N-1}C_{n-1}$ 個の X_i を全部足して、 n で割っていることになる。つまり、右辺は

$$= \frac{1}{n} \frac{1}{N C_n} {}_{N-1}C_{n-1} (X_1 + X_2 + \cdots + X_N)$$

である。これを紐解くと、

$$= \frac{1}{n} \frac{1}{N! / (N-n)! n!} \frac{(N-1)!}{((N-1) - (n-1))! (n-1)!} (X_1 + X_2 + \cdots + X_N) \\ = \frac{1}{n} \frac{n! (N-1)! (N-n)!}{N! (N-1)! (N-n)!} (X_1 + X_2 + \cdots + X_N) \\ = \frac{1}{n} \frac{n}{N} (X_1 + X_2 + \cdots + X_N) \\ = \frac{1}{N} (X_1 + X_2 + \cdots + X_N) \\ = \bar{X}$$

となる。有限であれ、無限であれ、サンプル平均の期待値は母平均に一致するというわけだ。t ところがサンプル平均の分散の場合は少し話が違ってくる。サンプル平均の分散の期待値、 $E(\bar{x}^2)$ はどうなるかというと、

$$E(\bar{x}^2) = \frac{1}{N C_n} \left[\left\{ \frac{1}{n} (X_1 + X_2 + \cdots + X_n) - E(\bar{x}) \right\}^2 \right. \\ \left. + \left\{ \frac{1}{n} (X_1 + X_2 + \cdots + X_{n+1}) - E(\bar{x}) \right\}^2 + \cdots \right. \\ \left. + \left\{ \frac{1}{n} (X_{N-n+1} + X_{N-n+2} + \cdots + X_N) - E(\bar{x}) \right\}^2 \right]$$

となる。

これをそのまま計算すると非常に面倒なので、 $(a-b)^2 = a^2 - 2ab + b^2$ の公式に従って分解してみよう。もちろん a に当たるのがサンプルで、 b に当たるのが $E(\bar{x})$ である。

第一の項はこんな感じである。

$$\text{第一項} = \frac{1}{n^2} \left\{ (X_1 + X_2 + \cdots + X_n)^2 + (X_1 + X_2 + \cdots + X_{n-1} + X_{n+1})^2 + \cdots + (X_{N-n+1} + \cdots + X_N)^2 \right\}$$

. 第二項は、

$$\text{第二項} = -2E(\bar{x}) \left\{ \frac{1}{n} (X_1 + X_2 + \cdots + X_n) + \frac{1}{n} (X_1 + X_2 + \cdots + X_{n-1} + X_{n+1}) + \cdots + \frac{1}{n} (X_{N-n+1} + \cdots + X_N) \right\}$$

第三項は、

$$\text{第三項} = [{}_N C_n \{E(\bar{x})\}^2] = \{E(\bar{x})\}^2$$

これら全てに $\frac{1}{{}_N C_n}$ がかかっていることを忘れずに。

さてまず第一項。これは () 内の n 個の要素を足して、二乗するのだが、記号だけでやるので得てしてわけが分からなくなる。でも振り落とされずについて来て下さい。

簡単な例から考えてみよう。 $(a+b+c+d)^2 = a^2+b^2+c^2+d^2+2ab+2ac+2ad+2bc+2bd+2cd$ である。ということは、第一項のなかの最初の項からは、 $(X_1^2+X_2^2+\dots+X_n^2+X_1X_2+X_1X_3+\dots+X_{n-1}X_n)$ が得られることになる。つまり X_i^2 と、 X_jX_k の組み合わせがつつらと。特に後者は、 $j < k$ という制限を付けると、 $2 \sum_j \sum_k X_jX_k$ と書ける。これが $(X_{N-n+1} + X_{N-n+2} + \dots + X_N)$ の項まで同様に続くわけだ。では結局どれだけの数が含まれるのだろうか。

X_i^2 の項は、 ${}_N C_n$ 個ある中括弧 {} の中に、全部で $n \cdot {}_N C_n$ 個入っている。さらに、 X_1, X_2, \dots, X_n のどの項も同じ数だけ入っているはずだから、 X_1^2 は $n \cdot {}_N C_n / n$ 個あるはずだ。 X_2 や X_3 など、全ての項が同じだけあるはずなので、ここは $n \cdot {}_N C_n / n (X_1^2 + X_2^2 + \dots + X_n^2)$ と書けるだろう。

では X_jX_k の項はどうなるか。これはともかく、二つの項の組み合わせである。 N 個の中から 2 つの項を取り出す取り出し方は、 ${}_N C_2$ 通り。ひとつのサンプルから 2 つの項を取り出す取り出し方は、 ${}_n C_2$ 通りある。全てのサンプリング方法、 ${}_N C_n$ 通りの中には、 ${}_n C_2 / {}_N C_2$ の割合で X_jX_k が入っているに違いない。ということは全部で $2 \frac{{}_n C_2 {}_N C_n}{{}_N C_2}$ 個あるはずなのだ。

さて少し式の展開をしてみよう。第一項を更に展開すると以下のようになる。

$$\text{第一項} = \frac{1}{{}_N C_n} \left[\frac{1}{n^2} \left\{ (X_1 + X_2 + \dots + X_n)^2 + (X_1 + X_2 + \dots + X_{n-1} + X_{n+1})^2 + \dots + (X_{N-n+1} + \dots + X_N)^2 \right\} \right]$$

$$\begin{aligned} &= \frac{1}{{}_N C_n} \left[\frac{1}{n^2} \left\{ \frac{{}_N C_n}{N} \sum X_i^2 + 2 \frac{{}_n C_2 {}_N C_n}{{}_N C_2} \sum_j \sum_k X_j X_k \right\} \right] \\ &= \frac{1}{Nn} \sum X_i^2 + 2 \frac{1}{{}_n C_2} \sum_j \sum_k X_j X_k \\ &= \frac{1}{Nn} \sum X_i^2 + 2 \frac{1}{n^2} \frac{\frac{n!}{(n-2)!2!}}{\frac{N!}{(N-2)!2!}} \sum_j \sum_k X_j X_k \end{aligned}$$

ここで後者の項が少し面倒になってきたので、整理しよう。

$$\begin{aligned} \frac{1}{{}_n C_2} \frac{{}_n C_2}{{}_N C_2} &= \frac{1}{n^2} \frac{\frac{n!}{(n-2)!2!}}{\frac{N!}{(N-2)!2!}} = \frac{1}{n^2} \frac{n!(N-2)!2!}{(n-2)!2!N!} \\ &= \frac{1}{n^2} \frac{(n \times n - 1 \times n - 2 \times \dots \times 1) \times (N-2 \times N-3 \times \dots \times 1) \times 2 \times 1}{(n-2 \times n-3 \times \dots \times 1) \times 2 \times 1 \times (N \times N-1 \times N-2 \times \dots \times 1)} \end{aligned}$$

分子と分母に同じものがあるので、それらはキャンセルしあって、

$$= \frac{n-1}{Nn(N-1)}$$

これだけになる。さあこれを元の式に組み込もう。

$$\text{第一項} = \frac{1}{Nn} \sum X_i^2 + 2 \frac{n-1}{Nn(N-1)} \sum_j \sum_k X_j X_k$$

やっと見やすい形になった。しかしこれはまだ一つ目、第二項もやっつけよう。第二項は全部書くと、

$$-\frac{1}{nC_n} 2E(\bar{x}) \left\{ \frac{1}{n}(X_1 + X_2 + \cdots + X_n) + \cdots + \frac{1}{n}(X_{N-n+1} + X_{N-n+2} + \cdots + X_N) \right\}$$

となる項である。ここで、 $2E(\bar{x})$ を除いた箇所、つまり

$$\frac{1}{nC_n} \left\{ \frac{1}{n}(X_1 + X_2 + \cdots + X_n) \cdots \right\}$$

は、結局のところ $E(\bar{x})$ の式と変わらないわけである。つまり、第二項は $-2E(\bar{x})E(\bar{x})$ である。

第二項と第三項を併せて考えてみよう。

$$\begin{aligned} \text{第二項} + \text{第三項} &= -2E(\bar{x})E(\bar{x}) + \{E(\bar{x})\}^2 = -\{E(\bar{x})\}^2 = -\bar{X}^2 \\ &= -\left\{ \frac{1}{N}(X_1 + X_2 + \cdots + X_N) \right\}^2 \\ &= -\frac{1}{N^2} \left\{ (X_1^2 + X_2^2 + \cdots + X_N^2) + 2(X_1X_2 + X_1X_3 + \cdots + X_{N-1}X_N) \right\} \\ &= -\frac{1}{N^2} \sum X_i^2 - 2\frac{1}{N^2} \sum \sum X_jX_k \end{aligned}$$

さてこれで、もともとの式に戻ってみると、

$$\begin{aligned} E(\bar{x}^2) &= \frac{1}{Nn} \sum X_i^2 + 2\frac{n-1}{Nn(N-1)} \sum \sum X_jX_k - \frac{1}{N^2} \sum X_i^2 - 2\frac{1}{N^2} \sum \sum X_jX_k \\ &= \left(\frac{1}{Nn} - \frac{1}{N^2} \right) \sum X_i^2 + 2 \left\{ \frac{n-1}{Nn(N-1)} - \frac{1}{N^2} \right\} \sum \sum X_jX_k \\ &= \frac{N-n}{N^2n} \sum X_i^2 + 2\frac{n-N}{N^2n(N-1)} \sum \sum X_jX_k \\ &= \frac{N-n}{n(N-1)} \left\{ \frac{N-1}{N^2} \sum X_i^2 - \frac{2}{N^2} \sum \sum X_jX_k \right\} \\ &= \frac{N-n}{n(N-1)} \left[\frac{N}{N^2} \sum X_i^2 - \frac{1}{N^2} \sum X_i^2 - \frac{2}{N^2} \sum \sum X_jX_k \right] \\ &= \frac{N-n}{n(N-1)} \left[\frac{N}{N^2} \sum X_i^2 - \frac{1}{N^2} \left\{ \sum X_i^2 + 2 \sum \sum X_jX_k \right\} \right] \\ &= \frac{N-n}{n(N-1)} \left[\frac{1}{N} \sum X_i^2 - \frac{1}{N^2} \left\{ \sum X_i \right\}^2 \right] \\ &= \frac{N-n}{n(N-1)} \left[\frac{1}{N} \sum X_i^2 - \left\{ \frac{1}{N} \sum X_i \right\}^2 \right] \\ &= \frac{N-n}{n(N-1)} \left\{ \frac{1}{N} \sum X_i^2 - \bar{X}^2 \right\} \end{aligned}$$

さてここで、最後に残った $1/N \sum X_i^2 - \bar{X}^2$ は分散の公式である。ここから以下の関係が得られる。

$$E(\bar{x}^2) = \frac{N-n}{N-1} \frac{\sigma^2}{n} \quad (7)$$

つまり、有限母集団を対象とする場合は、有限修正項 $N - n/N - 1$ をつけないと無限母集団の値と一致しないことがわかる。もっとも、多くのサンプリング調査の場合は、 n に比べて N が非常に大きい(10万人とか、500万人とか)ため、 $N - n/N - 1$ がほとんど 1.0 に近くなるから、標本誤差に大きな影響を与えることはない。