# 標本データの抽出、標本平均

## 1 有限母集団と無限母集団

母集団の構成要素、これから調べようとする対象の個数が有限個か無限個によって、有限母集団、無限母集団とよぶ。ただし実際には無限母集団というときには、標本の抽出を繰り返して何度でもできる場合が無限母集団の議論で、繰り返しによって、調査の対象の状況が変わってくる場合が有限母集団である。しかし理論的ではなく、近似的に大きな数、たとえば N が大きく 1/N がゼロと近いとみなせる場合には、標本の抽出の前後でによって変化しないと考えられるから、有限母集団と無限母集団は同じと考えられる。これに対して、抽出の前後で母集団が変わってくるとき、復元抽出と非復元抽出の区別をする必要がでてくる。いま母集団には2種類の属性 ( $\heartsuit$ と $\clubsuit$ ) がそれぞれ M,N 個あるものとしてこれから標本を取り出す。番号をつけて区別するとして、起こりえるあらゆる場合を表にまとめてみる。表の行は最初の取り出して、横の列は2回

目の取り出しの結果を表すとしよう。まず母集団は $\clubsuit1 \clubsuit2 \cdots \clubsuitM$   $\heartsuit1 \heartsuit2 \cdots \heartsuitN$   $\clubsuit k, k = 1, 2, \cdots M$  :  $\heartsuit j, j = 1, 2, \cdots N$  であり、2 個の標本をとりだしたときから考える。

[復元抽出]

	<b>4</b> 1	<b>4</b> 2	$\cdots$ $\clubsuit M$	$\heartsuit 1$	 $\lozenge N$
<b>\$</b> 1	<b>\$</b> 1 <b>\$</b> 1	<b>4</b> 1 <b>4</b> 2	<b>♣</b> 1 <b>♣</b> <i>M</i>	<b>♣</b> 1♡1	1 % N
<b>\$</b> 2	<b>♣</b> 2 <b>♣</b> 1	<b>♣</b> 2 <b>♣</b> 2	<b>♣</b> 2 <b>♣</b> <i>M</i>	<b>♣</b> 2♡1	-2 % N
:					
$\clubsuit M$	<b>♣</b> M <b>♣</b> 1	<b>♣</b> M <b>♣</b> 2	AMAM	$AM \heartsuit 1$	$\clubsuit M \heartsuit N$
$\bigcirc$ 1	♡1♣1	♡1♣2	♡1 <b>♣</b> M	♡1♡1	$\heartsuit 1 \heartsuit N$
:					
$\heartsuit N$	♡ <b>N</b> ♣1	$\heartsuit N \clubsuit 2$	$\heartsuit N \clubsuit M$	$\lozenge N \lozenge 1$	$\lozenge N \lozenge N$

## 要素の個数

2回目	<b>♣</b> 1 <b>♣</b> 2 ··· <b>♣</b> M	$\bigcirc$ 1 ··· $\bigcirc$ N
♣ 1 ♣ 2 ⋮	$M^2$	MN
♡ 1 : ♡ N	MN	$N^2$

2個の結果を属性でまとめて合計してみると

$$\left. \begin{array}{l} \clubsuit \clubsuit \\ \clubsuit \heartsuit + \heartsuit \clubsuit \\ \vdots \\ 2MN \\ \heartsuit \heartsuit \\ \vdots \\ N^2 \end{array} \right\} = (M+N)^2 = (M+N)(M+N)$$

つぎに非復元では1回目と2回目に同じものはないから、つまり非復元の場合には対角線の部分を除く。 [非復元抽出] 表中の $\times$  は起こらないことを意味している

	<b>\$</b> 1	<b>4</b> 2	$\cdots$ $\clubsuit M$	♡1	$\cdots$ $\heartsuit N$
<b>4</b> 1	×	<b>\$</b> 1 <b>\$</b> 2	<b>♣</b> 1 <b>♣</b> <i>M</i>	<b>♣</b> 1♡1	<b>♣</b> 1♡N
<b>4</b> 2	<b>♣</b> 2 <b>♣</b> 1	×	<b>♣</b> 2 <b>♣</b> <i>M</i>	<b>♣</b> 2♡1	-2 % N
:					
$\clubsuit M$	<b>♣</b> M <b>♣</b> 1	<b>♣</b> M <b>♣</b> 2	×	$AM \heartsuit 1$	$AM \heartsuit N$
♡1	♡1♣1	♡1♣2	♡1 <b>♣</b> M	×	$\heartsuit 1 \heartsuit N$
:					
$\lozenge N$	$\heartsuit N - 1$	$\heartsuit N \clubsuit 2$	$\heartsuit N \clubsuit M$	$\heartsuit N \heartsuit 1$	×

### したがって要素の個数は

	<b>4</b> 1	<b>4</b> 2	 $\clubsuit M$	♡1		$\lozenge N$
<b>4</b> 1						
<b>4</b> 2		$(M)_2$			MN	
:						
$\clubsuit M$						
$\heartsuit 1$						
:		MN			$(N)_2$	
$\bigcirc N$						

#### 合計は

$$\left. \begin{array}{ll} \clubsuit \clubsuit & : (M)_2 \\ \clubsuit \heartsuit + \heartsuit \clubsuit & : 2MN \\ \heartsuit \heartsuit & : (N)_2 \end{array} \right\} = (M+N)_2 = (M+N)(M+N-1)$$

さらに3回くり返しを考えると、つぎの表の行には2回目までの属性による合計結果を表し、横の列は3回目の結果で分けてみる。

### [復元抽出] での要素の数とその合計は

3回目2回目まで	*	$\Diamond$	*+*+*	$: M^3$	
$\clubsuit \clubsuit : M^2$	$M^3$	$M^2N$	<b>♣</b> + ♣ + ♡ <b>♣</b> + ♡ + ♡	: 31VI IV	$= (M+N)^3$
$\clubsuit \heartsuit : 2MN$	$2M^2N$	$2MN^2$			
$\heartsuit \ \heartsuit : N^2$	$MN^2$	$N^3$	$\triangle + \triangle + \triangle$	$: N^3$	J

## [非復元抽出]での要素の数とその合計は

3回目2回目まで	*	$\Diamond$	<b>*</b> + <b>*</b> + <b>*</b>		
$\clubsuit \clubsuit : (M)_2$	$(M)_3$	$(M)_2N$	♣ + ♣ + ♡ ♣ + ♡ + ♡	3M(N)	$= (M+N)_3$
$\clubsuit \heartsuit : 2MN$	$2(M)_2N$	$2M(N)_2$	\times + \times + \times	` '	
$\heartsuit \ \heartsuit : (N)_2$	$M(N)_2$	$(N)_3$		. (1.)3	,

以上の計算から、一般にn回くり返した場合には「復元抽出」

「非復元抽出」

確率変数 X を抽出における標本での n 個のうち  $\P$  の個数を表すとすると  $k=0,1,2\cdots,n$  に対して

復元抽出はパラメータ  $n,p=\frac{M}{M+N}$  の 2 項分布  $Binom\left(n,\frac{M}{M+N}\right)$  であり、

$$f(k) = \binom{n}{k} p^k (1-p)^{n-k}, k = 0, 1, 2, \dots, n$$

また非復元抽出の場合には超幾何分布 Hyper(M+N,M,n) と呼ばれ、

$$f(k) = \binom{n}{k} \frac{(M)_k(N)_{n-k}}{(M+N)_n} = \frac{\binom{M}{k} \binom{N}{n-k}}{\binom{M+N}{n}}, k = 0, 1, 2, \dots, n$$

と表されることが多い。

これを比較すればわかるよう、べき乗の数と組合せの数が対応している。復元抽出では  $M^n=\overbrace{M\cdot M\cdots M}^n$  に対して非復元抽出では  $(M)_n=\overbrace{M(M-1)\cdots (M-n+1)}^n$  となっていることに注意する。

# 2 標本平均の平均と分散

標本データから求められる基本的な統計量としては、標本平均がある。単純に算術平均を求めることである。いま母集団には 2 種類の属性があるときに、標本平均に対する平均(期待値)と分散を求める。属性が 2 種類しかないから、標本平均というより、標本比率あるいは標本割合ということが多い。つまり母集団の属性が 3, 2 として、3 番目 3 は 4 が抽出されたとき 4 の 4 が加出されたとき

とおくと、これらの和  $X=X_1+X_2+\cdots+X_n$  の分布は復元抽出ならば、2 項分布  $Binom\left(n,\frac{M}{M+N}\right)$  であって、もし非復元抽出ならば、超幾何分布 Hyper(M+N,M,n) にしたがう。したがってこれを利用すれば、標本平均(標本比率)

$$\overline{X}_n = \frac{X}{n} = \frac{1}{n} \left( X_1 + X_2 + \dots + X_n \right)$$

の平均と分散が求められる。

復元抽出: 
$$E(\overline{X}_n) = \frac{M}{M+N}$$
,  $V(\overline{X}_n) = \frac{1}{n} \frac{M}{M+N} \frac{N}{M+N}$  非復元抽出:  $E(\overline{X}_n) = \frac{M}{M+N}$ ,  $V(\overline{X}_n) = \frac{1}{n} \frac{M}{M+N} \frac{N}{M+N} \frac{M+N-n}{M+N-1}$ 

非復元抽出の分散で  $\frac{M+N-n}{M+N-1}$  は有限補正項とよばれるが、標本数 n に比べて母集団サイズが大きく M+N が大とすると、この値は 1 に近いとみなせるから、復元と非復元の違いはなくなってくることが分かる。

非復元抽出で平均  $E(X_i)=p$ , 分散  $V(X_j)=p(1-p)$  ただし p=M/(M+N) を示し、また共分散  $Cov(X_i,X_j)=P(X_j=1)P(X_j=1|X_i=1)-p^2$  を計算せよ。