

WILTON DE OLIVEIRA BUSSAB

Departamento de Estatística - IME

U.S.P.

SÔBRE A ESTIMAÇÃO DA MÉDIA OU TOTAL EM SUBPOPULAÇÕES

TRABALHO APRESENTADO AO INSTITUTO
DE MATEMÁTICA E ESTATÍSTICA DA U-
NIVERSIDADE DE SÃO PAULO, PARA MES-
TRADO EM ESTATÍSTICA APLICADA.

ABRIL 1971

P R E F Á C I O

A teoria da amostragem baseia a grande maioria das suas conclusões na fixação a priori do tamanho da amostra a ser colhida. Mas, em alguns esquemas êsse tamanho não pode ser préfixado, sendo uma variável aleatória.

A maioria dos livros não trazem um tratamento sistemático do assunto. O nosso propósito foi fazer um levantamento da literatura sobre êsse assunto, procurando sistematizar o tratamento, bem como desenvolver sugestões dadas pelos livros e não desenvolvidas. Em outros casos procuramos verificar a possibilidade de aplicação de certos estimadores desenvolvidos em casos gerais, ao nosso específico problema.

Dêsse modo dividimos o trabalho em três partes:

CAPÍTULO I - Introdução ao problema, às notações, e abordando o caso de amostragem casual simples. Inicialmente vemos algumas maneiras na estimação da média, procurando comparar os efeitos de cada estimador, e em seguida tratamos da estimativa do total;

CAPÍTULO II - Fazemos a mesma abordagem para a amostragem estratificada;

CAPÍTULO III - Uma rápida apresentação para a amostragem em dois estágios.

É nosso interêsse pesquisar no futuro, outros estimadores subpopulacionais bem como em outros esquemas de amostragem.

Finalmente queremos agradecer à tódos aqueles que diretamente ou indiretamente nos ajudaram na elaboração dêste trabalho. Principalmente ao Professor Doutor Lindo Fava, nosso orientador, sugerindo o tema dêste trabalho bem como nos elucidando várias dúvidas; e ao senhor João Baptista Esteves de Oliveira, pelo trabalho de datilografia.

São Paulo, abril de 1971

CAPÍTULO I -

1-INTRODUÇÃO

Consideremos uma população N , com N unidades de amostragem, e em relação a qual está associada uma variável Y . Consideremos ainda imersa em N uma subpopulação M , isto é, algumas unidades da população possuem em comum determinada característica que as faz pertencer a essa subpopulação. O que acontece frequentemente é estarmos interessados na estimação de parâmetros desta subpopulação. Vejamos alguns exemplos:

- (i) Quer-se estimar a renda média em domicílios que tenham mais de três adultos. Geralmente não dispomos da listagem dos domicílios com esta particular especificação, "mais de três adultos", mas a listagem de todos os domicílios. A identificação *a priori* das unidades amostrais da subpopulação pode ser muito cara ou difícil, o que nos impede de trabalhar apenas com a enumeração das unidades para a subpopulação, e temos que usar a listagem toda no processo de seleção da amostra.
- (ii) Numa amostragem realizada estimou-se determinado parâmetro para uma população de pessoas adultas. Usando a mesma amostra quer-se estimar o mesmo parâmetro, mas agora restrita apenas à subpopulação masculina.
- (iii) Uma grande empresa tem arquivado os canchotos de cheques emitidos no último ano, deseja-se estimar através de uma amostra o total em cruzeiros pagos com cheques de valor inferior a determinada quantia. Difícilmente dispõe-se da enumeração dos cheques nestas condições, e precisamos usar novamente a listagem de todos os cheques para colher a amostra.

Dos exemplos acima, constatamos que a grande dificuldade na estimação de parâmetros da subpopulação, é o desconhecimento da enumeração das unidades amostrais da mesma. Este fato sugere-nos dois caminhos a seguir para determinação da amostra;

- 1º) Vamos selecionando elementos da população até encontrarmos um número fixado de elementos da subpopulação, abandonando os elementos não pertencentes a ela, e que foram selecionados.
- 2º) Sorteamos uma amostra de tamanho fixado da população N , e desta amostra consideramos apenas os elementos pertencentes à subpopulação.

O primeiro procedimento leva-nos a uma amostragem casual simples de tamanho fixo. Mas podemos estar sujeito a colher um número

muito grande de elementos da população até atingirmos o número pré-fixado de elementos da subpopulação.

O segundo procedimento embora fixando o número de elementos de N que comporão a amostra, traz o inconveniente de que o número de elementos de M encontrados naquela amostra é uma variável aleatória.

Neste trabalho iremos tratar, principalmente, do uso do segundo processo, apresentado acima, para a estimação de parâmetros da subpopulação.

Embora usemos a palavra subpopulação, convém lembrar que o mesmo conceito aparece na literatura especializada sob outros títulos como: domínio, setor, subclasse, etc.

2. DEFINIÇÕES E NOTAÇÕES

Indicaremos a população N , do seguinte modo:

$$N = \{U_1, U_2, \dots, U_N\}$$

Os possíveis valores da variável Y , associada a esta população, serão indicados por:

$$Y_1, Y_2, \dots, Y_N$$

Para identificarmos os valores de Y correspondentes aos elementos da subpopulação, definiremos a variável auxiliar W do seguinte modo:

$$W_i = \begin{cases} Y_i & \text{se } U_i \in M \\ 0 & \text{se } U_i \notin M \end{cases} \quad (2.1)$$

Então, os principais parâmetros da subpopulação, serão:

- valor médio da subpopulação

$$\bar{W} = \frac{1}{M} \sum_{i=1}^N W_i \quad (2.2)$$

- total da subpopulação

$$W = \sum_{i=1}^N W_i = M\bar{W} \quad (2.3)$$

onde M indica o número de elementos da subpopulação.

Colhida uma amostra de n elementos da população queremos usar estes elementos para estimar (2.2) e (2.3). A nossa amostra

$$y_1, y_2, \dots, y_n$$

terá um número m de elementos da subpopulação, que iremos identificar do seguinte modo:

$$w_i = \begin{cases} y_i & \text{se } y_i \in M \\ 0 & \text{se } y_i \notin M \end{cases} \quad (2.4)$$

Nestas condições, se quisermos estimar a média \bar{W} de M , através de:

$$\bar{w} = \frac{1}{m} \sum_{i=1}^n w_i \quad (2.5)$$

não podemos usar os teoremas básicos de amostragem (Cochran, teoremas 2.1, 2.2 e 2.4), porque m não é um número fixo mas varia de amostra para amostra, ou seja, m é uma variável aleatória discreta assumindo valores de zero até n .

Neste trabalho, iremos ver algumas técnicas usadas na estimação de \bar{W} , e do total W da subpopulação. Além disso, quando nada se disser ao contrário, estamos considerando que a amostra foi tomada mediante amostragem casual simples, isto é, a probabilidade de qualquer unidade amostral ser incluída na amostra é constante e igual a $1/N$.

3. ESTIMAÇÃO COM m FIXO

Kish (pag.38) sugere como a primeira abordagem do problema, que em vez de especificarmos o tamanho n da amostra para a população, fixemos o tamanho m para a subpopulação. Em seguida vamos sorteando elementos de N , até completar os m diferentes elementos de M , e desprezando aqueles que não pertencem a esta subpopulação. Esse procedimento, leva-nos a uma amostragem casual simples de m elementos, retirados de uma população de tamanho M .

De acordo com os teoremas usuais deste tipo de amostragem, tiramos as seguintes conclusões:

- o estimador definido em (2.5)

$$\bar{w} = \frac{1}{m} \sum_{i=1}^n w_i = \frac{1}{m} \sum_{i=1}^m w_i \quad (*) \quad (3.1)$$

é um estimador não viesado, e cuja variância é

$$V(\bar{w}) = (1-f) \frac{S_w^2}{m} \quad (3.2)$$

onde

$$S_w^2 = \frac{1}{M-1} \sum_{i=1}^M (w_i - \bar{w})^2 \quad (3.3)$$

(*) Usaremos o símbolo $\sum_{i=1}^m w_i$, para indicar a soma estendida apenas aos valores não nulos da variável w_i na amostra.

$\sum_{i=1}^M W_i$, terá o mesmo significado para a população.

e

$$f' = \frac{m}{M} \quad (3.4)$$

e onde podemos usar como estimador de S_w^2 , a variância da amostra:

$$s_w^2 = \frac{1}{m-1} \sum (w_i - \bar{w})^2 \quad (3.5)$$

Tal procedimento levou-nos à seguinte preocupação: qual o número médio de elementos sorteados até completar a amostra de tamanho desejado m ?

Vamos dar uma resposta aproximada a este problema, definindo inicialmente, a seguinte variável:

S_i : - número de elementos selecionados de N , após a escolha do i -ésimo elemento de M até a escolha do $(i + 1)$ -ésimo elemento de M . (Para facilitar a linguagem, diremos que o processo está na i -ésima etapa).

Então, o número de elementos selecionados até o final da i -ésima etapa, é a variável aleatória

$$T_i = S_0 + S_1 + \dots + S_{i-1} \quad (3.6)$$

Determinamos a distribuição de S_i . Ao iniciar-se a i -ésima etapa, já foram removidos T_{i-1} elementos de M , e teremos na população e subpopulação, respectivamente, os seguintes números de elementos:

$$\begin{aligned} A &= N - T_{i-1} \\ B &= M - i \end{aligned} \quad (3.7)$$

A probabilidade^(*) de que na etapa i , sejam removidos k elementos, é dada pela probabilidade de escolhermos os primeiros $(k - 1)$ elementos não pertencentes a M , vezes a probabilidade do k -ésimo elemento pertencer a M . Então:

$$\begin{aligned} P(S_i = k) &= \frac{(A-B)}{A} \cdot \frac{(A-B-1)}{A-1} \dots \frac{(A-B-k+2)}{A-k+2} \cdot \frac{B}{A-k+1} = \frac{(A-B)!}{(A-B-k+1)!} \cdot \frac{(A-k)!}{A!} \cdot B = \\ &= \frac{(A-B)!}{(A-B-k+1)!} \cdot \frac{(A-k)!}{A!} \cdot \frac{B!}{(B-1)!} = \frac{(A-k)!}{(A-k-B+1)! (B-1)!} \cdot \frac{(A-B)! B!}{A!} \\ P(S_i = k) &= \frac{\binom{A-k}{B-1}}{\binom{A}{B}} \end{aligned} \quad (3.8)$$

(*) Essa probabilidade é calculada, condicionada a um particular valor da variável T_{i-1} , isto é, supondo T_{i-1} fixo.

onde, $k = 1, 2, 3, \dots, A - B + 1$. Como

$$1 = \sum_{k=1}^{A-B+1} P(S_i=k) = \sum \frac{\binom{A-k}{B-1}}{\binom{A}{B}} = \frac{1}{\binom{A}{B}} \sum \binom{A-k}{B-1}$$

vem imediatamente que:

$$\sum_{k=1}^{A-B+1} \binom{A-k}{B-1} = \binom{A}{B} \tag{3.9}$$

Mas

$$E(S_i) = \sum_{k=1}^{A-B+1} k \cdot P(S_i=k) = \frac{1}{\binom{A}{B}} \sum k \binom{A-k}{B-1} \tag{3.10}$$

usando a igualdade $k = (A-B+1) - (A-B-k+1)$, em (3.10), vem

$$\begin{aligned} E(S_i) &= (A-B+1) - \frac{1}{\binom{A}{B}} \sum_{k=1}^{A-B} (A-B-k+1) \frac{(A-k)!}{(A-B-k+1)!(B-1)!} = \\ &= (A-B+1) - \frac{B}{\binom{A}{B}} \sum_{k=1}^{A-B} \binom{A-k}{B} = (A-B+1) - \frac{B}{\binom{A}{B}} \binom{A}{B+1} \end{aligned}$$

esta última igualdade decorre da expressão (3.9). Simplificando os binomiais e reduzindo ao mesmo denominador chegamos a

$$E(S_i) = \frac{A+1}{B+1}$$

ou substituindo por (3.7)

$$E(S_i) = \frac{N - T_{i-1} + 1}{M - i + 1} \tag{3.11}$$

O maior valor para $E(S_i)$, é aquele obtido quando nas etapas anteriores $0, 1, 2, \dots, (i-1)$, o primeiro elemento sorteado já pertence a M , isto é, $S_0 = 1, S_1 = 1, \dots, S_{i-1} = 1$. Acontecendo isto na etapa i , teremos a maior quantidade possível de elementos não pertencentes a M , e conseqüentemente $T_i = i$. Então,

$$E(S_i) = \frac{N - T_{i-1} + 1}{M - i + 1} \leq \frac{N - i + 1}{M - i + 1} \tag{3.12}$$

Queremos determinar o número de elementos removidos até completar a amostra de tamanho m desejado, que é dado por

$$T_m = S_0 + S_1 + \dots + S_{m-1} \tag{3.13}$$

logo:

$$E(T_m) = E(S_0) + E(S_1) + \dots + E(S_{m-1}) \tag{3.14}$$

das considerações anteriores obtemos:

$$\begin{aligned}
 E(T_m) &\leq \frac{N+1}{M+1} + \frac{N}{M} + \frac{N-1}{M-1} + \dots + \frac{N-m+2}{M-m+2} = \\
 &= \frac{M+D+1}{M+1} + \frac{M+D}{M} + \dots + \frac{M+D-m+2}{M-m+2} = \\
 &= m + D \frac{1}{M+1} + \frac{1}{M} + \frac{1}{M-1} + \dots + \frac{1}{M-m+2} < \\
 &< m + D \frac{1}{M-m+2} + \frac{1}{M-m+2} + \dots + \frac{1}{M-m+2} = \\
 &= m + \frac{m \cdot D}{M-m+2} = m \cdot 1 + \frac{D}{M-m+2}
 \end{aligned}$$

onde $D = N - M$. Substituindo e simplificando

$$E(T_m) < m \frac{N-m+2}{M-m+2}$$

por outro lado o limite inferior de $E(T_m)$ é o valor \underline{m} , logo

$$m < E(T_m) < m \frac{N-m+2}{M-m+2} = m \cdot 1 + \frac{Q}{(1-f')P} \quad (3.15)$$

onde

$$P = \frac{M}{N} \quad \text{e} \quad Q = 1 - P \quad (3.16)$$

A expressão (3.15) nos informa que quanto menor for P , maior será o extremo superior do número médio esperado de unidades removidas.

Como o custo de uma amostragem é fundamental na escolha do esquema a ser usado, vamos determinar uma expressão que nos informe sobre o custo total do esquema sugerido nesta seção.

Suponhamos que para investigar uma unidade amostrada gasta-se a cruzeiros. Caso a unidade pertença a subpopulação, gasta-se mais b cruzeiros para completar a investigação. Assim a variável C , custo total, é dada por

$$C = C_0 + mb + aT_m$$

onde C_0 é o custo fixo. Então,

$$E(C) = C_0 + mb + a E(T_m)$$

Esta última expressão combinada com (3.15), nos dá subsídios para analisar o custo esperado da pesquisa.

4. AMOSTRA "CONDICIONAL" DE m ELEMENTOS

Cochran (pag.33), quando aborda inicialmente o problema de subpopulações, e Kish (pag.42), sugerem um método para estimar \bar{w} . Consideram a distribuição de \bar{w} sobre as amostras de n elementos da população, nas quais aparecem exatamente m elementos da subpopulação. Isto é, dentre todas as possíveis amostras que podemos formar

de tamanho n , determinamos a distribuição condicional de \bar{w} restrita àquelas que possuam exatamente m elementos de M .

O número de amostras nas condições acima é

$$\binom{N - M}{n - m} \cdot \binom{M}{m} \quad (4.1)$$

pois temos n lugares dos quais m devem ser preenchidos pelos M elementos da subpopulação, e os $n-m$ restantes lugares pelos $N-M$ elementos de M' (complementar de M em relação ao universo N). Dentre estas amostras o número daquelas que apresentam uma específica amostra de m unidades da subpopulação é

$$\binom{N - M}{n - m} \quad (4.2)$$

pois colocada a específica configuração, sobram $n-m$ posições a serem preenchidas por $N-M$ elementos. Então, a probabilidade de que seja escolhida uma particular amostra de m elementos de M é dada por

$$\frac{\binom{N - M}{n - m}}{\binom{N - M}{n - m} \binom{M}{m}} = \frac{1}{\binom{M}{m}}$$

O que equivale a afirmar que estamos fazendo uma amostragem casual simples, de m elementos de uma população de tamanho M . Se queremos estimar \bar{W} , podemos usar o estimador \bar{w} definido em (2.5), e que está nas mesmas condições da secção anterior, é um estimador não viesado, cuja variância é dada por (3.2)

$$V(\bar{w}) = \frac{1 - f'}{m} S_w^2 \quad (4.3)$$

Frequentemente, desconhecemos a grandeza M , para calcularmos f' . Uma aproximação é substituir f' por

$$f = \frac{n}{N} \quad (4.4)$$

sendo esta uma aproximação razoável já que

$$E(f') = E \frac{m}{M} = \frac{n}{N} = f \quad (4.5)$$

Para demonstrar esta última afirmação introduziremos a variável auxiliar X , do seguinte modo:

$$X_i = \begin{cases} 1 & \text{se o elemento pertence a } M \\ 0 & \text{nos demais casos} \end{cases} \quad (4.6)$$

Usaremos x_i para indicar a mesma variável mas agora em relação a amostra. Assim

$$\sum_{i=1}^N X_i = M \quad \text{e} \quad \sum_{i=1}^n x_i = m \quad (4.7)$$

então:

$$E\left(\frac{m}{M}\right) = \frac{1}{M} E\left(\sum_{i=1}^n x_i\right) = \frac{1}{M} \sum_{i=1}^n E(x_i) \quad (4.8)$$

mas, $E(x_i) = \bar{X} = \frac{\sum_{i=1}^n x_i}{N} = \frac{M}{N}$ (4.9)

substituindo em (4.8), vem:

$$E\left(\frac{m}{M}\right) = \frac{1}{M} \sum_{i=1}^n \frac{M}{N} = \frac{1}{M} \cdot m \cdot \frac{M}{N} = \frac{n}{N} \quad (4.10)$$

o que demonstra (4.5).

Substituindo-se (4.4) e (3.5) em (4.3), temos que a variância de \bar{w} "condicionada" ao encontro de m elementos da subpopulação é dada por:

$$V(\bar{w}) = \frac{1-f}{m} \cdot s_w^2 = \left(1 - \frac{n}{N}\right) \cdot \frac{1}{m} \cdot \frac{1}{m-1} \sum_{i=1}^m (w_i - \bar{w})^2 \quad (4.11)$$

5. \bar{w} COMO ESTIMADOR RAZÃO

Usando as variáveis W e X , definidas por (2.1) e (4.6), respectivamente verificamos as seguintes relações:

$$\bar{w} = \frac{1}{M} \sum_{i=1}^N W_i = \frac{\sum_{i=1}^N W_i}{\sum_{i=1}^N X_i} = \frac{\bar{w}_0}{\bar{x}} = R \quad (5.1)$$

onde:

$$\bar{w}_0 = \frac{1}{N} \sum_{i=1}^N W_i \quad (5.2)$$

e $\bar{x} = \frac{1}{N} \sum_{i=1}^N X_i = \frac{M}{N}$ (5.3)

bem como $\bar{w} = \frac{1}{m} \sum_{i=1}^n w_i = \frac{\sum_{i=1}^n w_i}{\sum_{i=1}^n x_i} = \frac{\bar{w}_0}{\bar{x}} = r$ (5.4)

onde $\bar{w}_0 = \frac{1}{n} \sum_{i=1}^n w_i$ (5.5)

e $\bar{x} = \sum x_i = \frac{m}{n}$ (5.6)

Para determinarmos as propriedades deste estimador, vamos usar a técnica usada por P.V.Sukahtme (sec. 4a.3 pag.139), para estimador razão geral, e particularizemos para a nossa situação.

Façamos:

$$w_i = \bar{w}_0 + E_i \quad \text{o que resulta} \quad \bar{w}_0 = \bar{w}_0 + \bar{E} \quad (5.7)$$

onde

$$E(\bar{E}) = 0 \quad \text{e} \quad V(\bar{E}) = E(\bar{E}^2) = (1-f) \cdot \frac{S_0^2}{n} \quad (5.8)$$

onde

$$S_0^2 = \frac{1}{N-1} \sum_{i=1}^N (w_i - \bar{w}_0)^2 \quad (5.9)$$

Do mesmo modo:

$$x_i = \bar{x} + E'_i \quad \bar{x} = \bar{x} + \bar{E}' \quad (5.10)$$

com

$$E(\bar{E}') = 0 \quad \text{e} \quad V(\bar{E}') = (1-f) \frac{S_x^2}{n} \quad (5.11)$$

onde

$$S_x^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 = \frac{M}{N-1} \left(1 - \frac{M}{N}\right) \quad (5.12)$$

Substituindo (5.7) e (5.10) em (5.4), teremos:

$$n = \frac{\bar{w}_0}{\bar{x}} = \frac{\bar{w}_0 + \bar{E}}{\bar{x} + \bar{E}'} = \frac{\bar{w}_0}{\bar{x}} \cdot \frac{1 + \frac{\bar{E}}{\bar{w}_0}}{1 + \frac{\bar{E}'}{\bar{x}}} \quad (5.13)$$

Como n , é um quociente de duas variáveis aleatórias, vamos calcular a esperança aproximada, desenvolvendo:

$$\left(1 + \frac{\bar{E}'}{\bar{x}}\right)^{-1}$$

em série, mas para que a série seja convergente, precisamos ter

$$\left|\frac{\bar{E}'}{\bar{x}}\right| < 1$$

isto é

$$\left|\frac{\bar{x} - \bar{x}}{\bar{x}}\right| < 1 \rightarrow \left|\frac{\bar{x}}{\bar{x}} - 1\right| < 1 \rightarrow 0 < \frac{\bar{x}}{\bar{x}} < 2$$

ou substituindo

$$0 < \frac{\frac{m}{n}}{\frac{M}{N}} < 2$$

que resulta em:

$$\frac{m}{n} > 0 \quad \text{e} \quad \frac{m}{n} < 2 \frac{M}{N}$$

A primeira desigualdade, quase sempre está satisfeita, porque o número de amostras de tamanho n que não contenha nenhum elemento da subpopulação, pode ser consideraconsiderado como uma raridade, e quanto maior for n , mais difícil fica não aparecer elementos de M . A segunda desigualdade já é mais delicada. Necessita que a fração da subpopulação em relação a população toda na amostra, seja infe-

rior ao dõbro da mesma fração no todo. Como frequentemente as subpopulações são frações muito pequenas, esta desigualdade estará satisfeita num número reduzido de vêzes. Mas no caso em que subpopulação é grande mais que a metade da população, essa desigualdade se verifica sempre.

$$\left(1 + \frac{\bar{E}'}{\bar{X}}\right)^{-1} = 1 - \frac{\bar{E}'}{\bar{X}} + \left(\frac{\bar{E}'}{\bar{X}}\right)^2 - \dots \quad (5.14)$$

Substituindo êste valôr em (5.13), teremos:

$$\begin{aligned} n &= \frac{\bar{W}'_0}{\bar{X}} \left(1 + \frac{\bar{E}}{\bar{W}'_0}\right) \left\{1 - \frac{\bar{E}}{\bar{X}} + \left(\frac{\bar{E}'}{\bar{X}}\right)^2 + \dots\right\} \\ &= \frac{\bar{W}'_0}{\bar{X}} \left\{1 + \frac{\bar{E}}{\bar{W}'_0} - \frac{\bar{E}'}{\bar{X}} - \frac{\bar{E}'\bar{E}}{\bar{X}\bar{W}'_0} + \left(\frac{\bar{E}'}{\bar{X}}\right)^2 + \left(\frac{\bar{E}'}{\bar{X}}\right)^2 \cdot \frac{\bar{E}}{\bar{W}'_0} + \dots\right\} \quad (5.15) \end{aligned}$$

Consideremos negligenciáveis a contribuição dos tãermos envolvendo \bar{E}' e \bar{E} , com potências superior ao segundo grau, e calculando a esperança, teremos:

$$\begin{aligned} E(r) &= E\left(R \left[1 + \frac{\bar{E}}{\bar{W}'_0} - \frac{\bar{E}'}{\bar{X}'} - \frac{\bar{E}'\bar{E}}{\bar{X}\bar{W}'_0} + \left(\frac{\bar{E}'}{\bar{X}}\right)^2\right]\right) = \\ &= R \left\{1 + \frac{E(\bar{E})}{\bar{W}'_0} - \frac{E(\bar{E}')}{\bar{X}'} - \frac{E(\bar{E}'\bar{E})}{\bar{X}\bar{W}'_0} + \frac{E(\bar{E}'^2)}{\bar{X}^2}\right\} \end{aligned}$$

mas

$$\begin{aligned} E(\bar{E}) &= 0, & E(\bar{E}') &= 0 \\ E(\bar{E}'^2) &= (1-f) \frac{S_x^2}{n} & \text{segundo (5.11)} \\ E(\bar{E}'\bar{E}) &= \frac{N-n}{N} \frac{1}{n} \rho S_n \cdot S_0 \end{aligned}$$

esta última expressão encontra-se desenvolvida em Sukahtme.

Substituindo êstes valôres em $E(r)$, encontramos:

$$\begin{aligned} E(r) &= R \cdot \left\{1 + \frac{(1-f)}{\bar{X}^2} \frac{S_x^2}{n} - \frac{(1-f)}{\bar{X}\bar{W}'_0} \cdot \frac{\rho S_x S_0}{n}\right\} = \\ &= R \cdot \left\{1 + \frac{(1-f)}{n} \frac{S_x^2}{\bar{X}^2} - \frac{\rho S_x S_0}{\bar{X}\bar{W}'_0}\right\} \quad (5.16) \end{aligned}$$

Para que seja um estimador não viesado, é preciso que a expressão entre parêntesis seja igual a 0 (zero), isto é,:

$$\frac{S_x^2}{\bar{X}^2} - \frac{\rho S_x S_0}{\bar{X}\bar{W}'_0} = 0$$

$$\frac{\rho S_o}{\bar{w}_o} = \frac{S_x}{\bar{x}}$$

$$\bar{w}_o = \rho \frac{S_o}{S_x} \bar{x} \quad (5.17)$$

Esta última expressão significa que a reta de regressão de W em função de X, passa pela origem. Para o nosso problema este fato sempre se verifica, porque a todo elemento não pertencente a subpopulação M, está associado a observação (0,0). O que nos mostra que o estimador assim definido, é não viesado quando se usa a aproximação até segunda ordem. Então,

$$E(\bar{w}) = E(r) = R = \bar{w} \quad (5.18)$$

Algumas vezes é bastante importante o conhecimento da magnitude do verdadeiro viés, e isso pode ser calculado usando a sugestão dada por Hartley e Ross, como encontramos em Fava (pág.9). Indicando por ϵ o verdadeiro viés de \bar{w} , calculado como estimador razão, teremos:

$$\begin{aligned} \epsilon &= E(\bar{w}) - \bar{w} = - \{ \bar{w} - E(\bar{w}) \} = - \left[\frac{\bar{w}_o}{\bar{x}} - E \left(\frac{\bar{w}_o}{\bar{x}} \right) \right] = \\ &= - \left[\frac{\bar{w}_o}{\bar{x}} - E \left(\frac{\bar{w}_o}{\bar{x}} \right) \right] = - \frac{1}{\bar{x}} \left[\bar{w}_o - \bar{x} E \left(\frac{\bar{w}_o}{\bar{x}} \right) \right] = \\ &= - \frac{1}{\bar{x}} E \left\{ \left[\frac{\bar{w}_o}{\bar{x}} - E \left(\frac{\bar{w}_o}{\bar{x}} \right) \right] \left[\bar{x} - E(\bar{x}) \right] \right\} = - \frac{1}{\bar{x}} \text{Cov} \left(\frac{\bar{w}_o}{\bar{x}}, \bar{x} \right) \end{aligned} \quad (5.19)$$

onde $\text{Cov} \left(\frac{\bar{w}_o}{\bar{x}}, \bar{x} \right)$ é a covariância entre as variáveis $\frac{\bar{w}_o}{\bar{x}}$ e \bar{x} . Então,

$$|\epsilon| = \frac{1}{\bar{x}} \left| \text{Cov} \left(\frac{\bar{w}_o}{\bar{x}}, \bar{x} \right) \right| \leq \frac{1}{\bar{x}} S \left(\frac{\bar{w}_o}{\bar{x}} \right) S_{\bar{x}} \quad (5.20)$$

pois $|\rho| \leq 1$.

$$\left| \frac{\text{Cov} \left(\frac{\bar{w}_o}{\bar{x}}, \bar{x} \right)}{S \left(\frac{\bar{w}_o}{\bar{x}} \right) S_{\bar{x}}} \right| \leq 1 \rightarrow \left| \text{Cov} \left(\frac{\bar{w}_o}{\bar{x}}, \bar{x} \right) \right| \leq S \left(\frac{\bar{w}_o}{\bar{x}} \right) S_{\bar{x}}$$

onde ρ é o coeficiente de correlação entre as duas variáveis. Finalmente

$$|\epsilon| \leq \frac{1}{\bar{x}} S_{\bar{w}} S_{\bar{x}} \quad \text{ou} \quad \frac{|\epsilon|}{S_{\bar{w}}} \leq \frac{S_{\bar{x}}}{\bar{x}} \quad (5.21)$$

mas como $S_{\bar{x}}$ é igual a S_x dividido por \sqrt{n} , e como já vimos em (5.12)

que $S_x = \sqrt{\frac{M}{N-1} \left(1 - \frac{M}{N}\right)} \approx \sqrt{P \cdot Q}$

onde $P = \frac{M}{N}$ e $Q = 1 - P$, substituindo em (5.21) encontramos:

$$\frac{|e|}{S_x} < \frac{1}{n} \frac{P}{Q} \tag{5.22}$$

e esta expressão nos fornece qual a intensidade do viés. Se a fração P for pequena, vemos que para manter o viés reduzido teremos necessidade de que a amostra seja bem grande. Vejamos um exemplo numérico. Suponhamos que queremos manter um viés relativo inferior à 0,1. Então devemos ter:

$$\sqrt{\frac{1}{n} \frac{Q}{P}} = 0,1 \longrightarrow \frac{1}{n} \frac{Q}{P} = 0,01$$

$$n = \frac{Q}{0,01P}$$

Suponhamos inicialmente que a fração da subpopulação é 80% da população, isto é, $P = 0,8$, então necessitamos $n = 25$ elementos, se $P = 0,5$ então o tamanho da amostra será 100 elementos, e se P é apenas 20% necessitaremos de 400 elementos na amostra. O que nos mostra que necessitamos de amostras muito grandes quando a proporção da subpopulação é pequena.

6. CALCULO DA VARIÂNCIA DE \bar{w} , COMO ESTIMADOR RAZÃO

Para calcular a esperança de \bar{w} , tivemos que fazer aproximações, no cálculo da variância também vamos fazê-las, mas em duas aproximações. A primeira é usando o desenvolvimento de r obtido em (5.15), apenas com potências até o primeiro grau. Esta é a técnica usada por Sukahtme (pág.146), para estimador razão no caso geral, e que particularizaremos para este problema. Assim temos:

$$r = R \left\{ 1 - \frac{\bar{E}}{\bar{W}_0} + \frac{\bar{E}'}{\bar{X}} \right\} \tag{6.1}$$

Mas do fato de $E(r) = R$, obtido na secção anterior, teremos:

$$V_1(\bar{w}) = V(r) = E\{(r - E(r))^2\} = E\{(r - R)^2\} =$$

$$= R^2 \left\{ \frac{E(\bar{E}'^2)}{\bar{X}^2} + \frac{E(\bar{E}^2)}{\bar{W}_0^2} - \frac{2E(\bar{E}'\bar{E})}{\bar{X}\bar{W}_0} \right\}$$

Calculando as esperanças, segundo sugestão da secção anterior,

$$V_1(\bar{w}) = R^2 \left\{ \frac{(1-f)}{\bar{X}^2 n} S_x^2 + \frac{(1-f)}{\bar{W}_0^2 n} S_0^2 - \frac{2(1-f)}{n} \frac{\rho S_x S_0}{\bar{X}\bar{W}_0} \right\} =$$

$$= R^2 \left\{ \frac{(1-f)}{n} \frac{S_x^2}{\bar{X}^2} + \frac{S_0^2}{\bar{W}_0^2} - \frac{2\rho S_x S_0}{\bar{X}\bar{W}_0} \right\}$$

mas de (5.17); tiramos imediatamente que:

$$V_1(\bar{w}) = R^2 \frac{(1-f)}{n} \left[\frac{S_O^2}{\bar{w}^2} + \frac{S_X^2}{\bar{x}^2} \right] \quad (6.2)$$

Vamos tornar esta expressão mais simples, lembrando que:

$$R = \bar{w} \quad \text{de (5.1)}$$

$$S_X^2 = \frac{M}{N-1} \left(1 - \frac{M}{N} \right) \quad \text{de (5.12)}$$

$$\bar{x} = \frac{M}{N} \quad \text{de (5.3)}$$

e determinando as novas relações

$$\begin{aligned} \bar{w}_O &= \frac{1}{N} \sum^N w_i = \frac{M}{N} \frac{1}{N} \sum^M w_i = \frac{M}{N} \bar{w} = P\bar{w} \quad (6.3) \\ (N-1) S_O^2 &\equiv \sum^N (w_i - \bar{w}_O)^2 = \sum^N w_i^2 - N\bar{w}_O^2 = \sum^N w_i^2 - M\bar{w}^2 + M\bar{w}^2 - N\bar{w}_O^2 = \\ &= \sum^M w_i^2 - M\bar{w}^2 + M\bar{w}^2 - N \frac{M^2}{N^2} \bar{w}^2 = (M-1) S_w^2 + M\bar{w}^2 Q \end{aligned}$$

$$\text{isto é } S_O^2 \equiv P S_w^2 + PQ \bar{w}^2 \quad (6.4)$$

Substituindo estes últimos resultados em (6.2), e simplificando temos:

$$V_1(\bar{w}) = \frac{(1-f)}{n} \frac{N}{M} S_w^2 = (1-f) \frac{S_w^2}{fM} \quad (6.5)$$

Usa-se como estimador de S_w^2 a variância amostral

$$s_w^2 = \frac{1}{m-1} \sum^m (w_i - \bar{w})^2 \quad (6.6)$$

Como segunda aproximação usaremos o processo desenvolvido por Hansen (pág. 114, 2º vol.),

$$V_2(\bar{w}) = V(r) \approx E\{(r - R)^2\}$$

e que através da propriedade de esperança condicional, pode ser escrita da seguinte maneira:

$$V_2(\bar{w}) = E \left\{ E_m (r - R)^2 \right\} \quad (6.7)$$

onde E_m significa a esperança condicionada à ocorrência de m indivíduos da subpopulação.

$$E_m \left\{ (r - R)^2 \right\} = (1-f') \frac{S_w^2}{m} = \left(\frac{M-m}{M} \right) \frac{1}{m} S_w^2 = \left(\frac{1}{m} - \frac{1}{M} \right) S_w^2 \quad (6.8)$$

Substituindo êsse resultado em (6.7), vem

$$V_2(\bar{w}) = E \left\{ \left(\frac{1}{m} - \frac{1}{M} \right) S_w^2 \right\} = S_w^2 \left(E \left(\frac{1}{m} \right) - \frac{1}{M} \right) \quad (6.9)$$

De acôrdo com a definição de \bar{E}' dada em (5.10), podemos es-
crever:

$$m = n\bar{x} = n(\bar{X} + \bar{E}') = n\bar{X} \left(1 + \frac{\bar{E}'}{\bar{X}}\right) \quad (6.10)$$

mas $\left(1 + \frac{\bar{E}'}{\bar{X}}\right)^{-1} = 1 - \frac{\bar{E}'}{\bar{X}} + \left(\frac{\bar{E}'}{\bar{X}}\right)^2 + \dots$

logo
$$E\left(\frac{1}{m}\right) = \frac{1}{n\bar{X}} E\left(1 - \frac{\bar{E}'}{\bar{X}} + \left(\frac{\bar{E}'}{\bar{X}}\right)^2 + \dots\right) \approx \frac{1}{n\bar{X}} \left(1 + \frac{E(\bar{E}'^2)}{\bar{X}^2}\right) =$$

$$= \frac{N}{nM} \left\{ 1 + \frac{(1-f) \cdot \frac{1}{N} \frac{M}{N-1} \left(1 - \frac{M}{N}\right)}{\frac{M^2}{N^2}} \right\} \quad (6.11)$$

Substituindo em (6.9), vamos encontrar:

$$V_2(\bar{w}) = \left\{ \frac{N}{nM} \left[1 + \frac{(1-f)}{n} \frac{N}{M} \left(1 - \frac{M}{N}\right) \frac{N}{N-1} \right] - \frac{1}{M} \right\} S_w^2 =$$

$$\approx \frac{N}{nM} \left\{ (1-f) \left[1 + \frac{1}{n} \frac{Q}{P} \left(\frac{N}{N-1}\right) \right] \right\} S_w^2 =$$

$$\approx \frac{1-f}{fM} \left(1 + \frac{1}{n} \frac{Q}{P}\right) S_w^2 \quad (6.12)$$

Na prática usamos s_w^2 definido em (6.6) como estimador de S_w^2 ; m como es-
timador de fM e $\frac{m}{n}$ como estimador de P .

7. COMPARAÇÃO DAS VARIÂNCIAS OBTIDAS PARA \bar{w}

Vamos procurar comparar as variâncias obtidas na secção ante-
rior, com a variância de \bar{w} obtida quando se supõe a escolha de uma amos-
tra casual simples, de m elementos de uma população com M elementos, e
como já foi visto em (3.2) é dada por

$$V(\bar{w}) = (1-f') \frac{S_w^2}{m} \quad (7.1)$$

Se compararmos com a variância do estimador razão, obtidas em
(6.5)

$$V_1(\bar{w}) = (1-f) \frac{S_w^2}{fM} \quad (7.2)$$

vemos que as diferenças são mínimas, visto que, podemos supor $f \approx f'$,
e $E(m) = fM$, e sabemos que m pode ser confundido com a sua esperança,
o que nos leva a concluir que as diferenças nas duas formas são despre-
zíveis. Para aplicação, ambas as formas devem sofrer alteração quando
se desconhece M . Para (7.1) devemos substituir f' por f , e em (7.2) fM
por m .

Comparamos agora com a variância obtida em (6.12), ou seja, como estimador razão em segunda aproximação,

$$V_2(\bar{w}) = (1-f) \left(1 + \frac{1}{n} \frac{Q}{P} \right) \frac{S_w^2}{fM} \quad (7.3)$$

A variância neste caso é maior que a variância obtida pela fórmula (7.1), esse acréscimo é devido à aleatoriedade de m na amostra. Analisemos um pouco mais detidamente o acréscimo

$$\frac{Q}{nP} - \frac{S_w^2}{fM} \quad (7.4)$$

cuja interpretação é bastante intuitiva:

- i) Esse acréscimo diminui quando o número de unidades amostradas aumenta;
- ii) Quanto menor a fração da subpopulação em relação ao tamanho da população, maior é o acréscimo observado na variância.

8. ESTIMATIVA DO VALOR TOTAL W DAS SUBPOPULAÇÕES

Suponhamos agora que estamos interessados em estimar o total

$$W = \sum_{i=1}^N W_i = \sum_{i=1}^M W_i = M\bar{w}$$

Para isso vamos dividir o tratamento em três casos:

1. O total Y da população é conhecido;
2. O tamanho M da subpopulação é conhecido;
3. Y e M desconhecidos.

8.1-O total Y é conhecido

Existem algumas pesquisas em subpopulação, que por algum motivo, conhecemos o valor total da população. Como vimos no início deste trabalho, tínhamos uma população na qual observávamos uma variável Y, e através de um artifício criamos uma variável auxiliar W. Assim para cada unidade da população temos associado o par (W_i, Y_i) , o que nos permite usar um estimador razão, na estimativa do total W. Analisando a relação entre a variável W e Y podemos nos decidir por qualquer um dos estimadores razão.

8.2-O tamanho M da subpopulação é conhecido

Neste caso a estimativa do total pode ser dado através de:

$$\tilde{w} = M\bar{w} \quad (8.1)$$

Na qual podemos usar \tilde{w} , como qualquer um dos estimadores propostos nas secções anteriores. Usando o estimador definido na secção 3, ou no caso da amostragem "condicional", visto na secção 4, \tilde{w} é um estimador não viesado com a seguinte variância:

$$V(\tilde{w}) = \frac{1-f'}{m} M^2 S_w^2 \quad (8.2)$$

de acôrdo com as conclusões das citadas secções.

Se usarmos \tilde{w} como o estimador razão da secção 6, e em concordância com as conclusões lá obtidas, temos em \tilde{w} um estimador não viesado, até segunda ordem de aproximação, e cuja variância é dada, em primeira aproximação por (6.5)

$$V_1(\tilde{w}) = \frac{1-f}{f} M S_w^2 \quad (8.3)$$

e em segunda aproximação, conforme (6.12)

$$V_2(\tilde{w}) = \frac{1-f}{f} M \left(1 + \frac{1}{n} \frac{Q}{P} \right) S_w^2 \quad (8.4)$$

8.3-Y e M desconhecidos

Este é o caso mais geral e mais frequente. Usaremos aqui o artifício encontrado em Cochran (pág.34). O estimador de W proposto é o obtido pela multiplicação do total observado na amostra pelo fator de crescimento $\frac{N}{n}$:

$$\hat{w} = \frac{N}{n} w \quad (8.5)$$

onde

$$w = \sum_{i=1}^n w_i$$

então

$$\hat{w} = \frac{N}{n} \sum_{i=1}^n w_i = N\bar{w}_0 \quad (8.6)$$

$$E(\hat{w}) = N E(\bar{w}_0) = N \bar{w}_0 = N \frac{\sum w_i}{N} = \sum w_i = w \quad (8.7)$$

ou seja, este é um estimador não viesado.

Da expressão (8.6) podemos calcular a variância rapidamente lembrando que \bar{w}_0 , representa a média de uma amostra de n elementos, colhida de uma população com N elementos, e pelas conclusões da amostragem casual simples:

$$V(\bar{w}_0) = (1-f) \frac{S_0^2}{n} \quad (8.8)$$

Substituindo em (8.6)

$$V(\hat{w}) = V(N\bar{w}_O) = N^2 V(\bar{w}_O) = N^2 (1-f) \frac{S_O^2}{n} \quad (8.9)$$

onde usamos

$$S_O^2 = \frac{1}{n-1} \{ \sum (w_i - \bar{w}_O)^2 \}$$

como estimador de S_O^2 . O que pode causar um pouco de estranheza na última expressão, é o fato da variância ser calculada inclusive para os zeros da população, o que a torna bem maior do que aquela obtida dos valores não nulos.

Algumas vezes, é conveniente exprimir (8.9) em função de S_w^2 . Usando o resultado (6.4) e (6.13), podemos escrever:

$$S_O^2 = \frac{M-1}{N-1} S_w^2 + \frac{M}{N-1} \left(1 - \frac{M}{N}\right) \bar{w}^2 \cong PS_w^2 + PQ\bar{w}^2 \quad (8.10)$$

substituindo em (8.9)

$$V(\hat{w}) = \frac{N^2 (1-f)}{n} \{PS_w^2 + PQ\bar{w}^2\} \quad (8.11)$$

Como frequentemente desconhecemos S_w^2 , P , Q e \bar{w} , um estimador de (8.11) é

$$v(\hat{w}) = (1-f) \frac{N^2}{n} \frac{m}{n} S_w^2 + \frac{m}{n} \left(1 - \frac{m}{n}\right) \bar{w}^2$$

8.4-Comparação das variâncias

Vamos comparar as variâncias obtidas quando se desconhece todas as informações da subpopulação com aquela obtida do conhecimento do tamanho da mesma. (Cochran - pág. 36).

Da expressão (8.2) temos que:

$$V(\tilde{w}) = (1-f') \frac{M^2}{m} S_w^2 \quad (8.12)$$

onde m é uma variável aleatória, mas já vimos que $E(m) = nP$. Então, supondo $m = nP$, teremos:

$$f' = \frac{m}{M} = \frac{nP}{NP} = f$$

substituindo em (8.12)

$$V(\tilde{w}) = (1-f) \frac{N^2 P^2}{nP} S_w^2 = (1-f) \frac{N^2 P}{n} S_w^2 \quad (8.13)$$

A variância quando se desconhece o tamanho da subpopulação é dada por (8.11):

$$V(\hat{w}) = (1-f) \frac{N^2}{n} P \{S_w^2 + Q\bar{w}^2\} \quad (8.14)$$

então

$$\frac{V(\hat{w})}{V(\bar{w})} = \frac{S_w^2 + Q\bar{w}^2}{S_w^2} = 1 + \frac{Q}{C_w^2} = \frac{C_w^2 + Q}{C_w^2} \quad (8.15)$$

onde

$C_w = \frac{S_w}{\bar{w}}$ é o coeficiente de variação da subpopulação. Eviden-

temente a variância é maior quando não se tem informações sobre o tamanho da subpopulação. Em outras palavras, o lucro que se tem pelo conhecimento desse tamanho aumenta, quanto menor for a fração P,

8.5-Outra aplicação

As conclusões obtidas na secção anterior são de grande utilidade no seguinte problema (Kish, pág. 435). Suponhamos que queremos estimar o total Y, de uma população em que muito dos seus valores são zeros. A pergunta que se coloca é a seguinte: De quanto podemos reduzir a variância se restringirmos nossa amostragem apenas aos valores não nulos da população?

A variável Y da população N, pode ser separada em duas subpopulações a saber: M dos elementos em que Y não é nulo e M' dos zeros da população. Usaremos aqui novamente, a variável W definida em (2.1). Assim este problema, é o mesmo dos tratados nas secções anteriores, ou seja, queremos comparar a variância entre dois tipos de estimadores,

$$\hat{w} = N \bar{w}_0 = \frac{N}{n} w$$

e

$$\tilde{w} = M \bar{w} = \frac{M}{m} w$$

Aqui, ao contrário da secção (8.4), é conveniente exprimir as variâncias em função de S_0 , que é a variância da população toda. Assim de (8.10), temos que

$$S_w^2 = \frac{S_0^2}{P} - Q\bar{w}^2 \quad (8.16)$$

Substituindo em (8.2) teremos:

$$V(\tilde{w}) = \frac{(1-f')}{m} M^2 \left(\frac{S_0^2}{P} - Q\bar{w}^2 \right)$$

ou usando transformações anteriores

$$V(\tilde{w}) = \frac{(1-f)}{n P} N^2 P^2 \left\{ \frac{S_0^2 - PQ\bar{w}^2}{P} \right\} = \frac{(1-f)}{n} N^2 \{S_0^2 - PQ\bar{w}^2\}$$

por outro lado, de (8.9)

$$V(\hat{w}) = \frac{(1-f)}{n} N^2 S_0^2$$

comparando os dois últimos resultados teremos:

$$\frac{V(\tilde{w})}{V(\hat{w})} = \frac{S_o^2 - PQ\bar{w}^2}{S_o^2} = 1 - PQ \frac{\bar{w}^2}{S_o^2}$$

porém, lembrando que

$$\bar{w} = \frac{N}{M} \bar{w}_o = \frac{\bar{w}_o}{P} \text{ e fazendo } C_o = \frac{S_o}{\bar{w}_o}, \text{ vem}$$

$$\frac{V(\tilde{w})}{V(\hat{w})} = 1 - \frac{Q}{PC_o^2}$$

onde o fator $\frac{Q}{PC_o^2}$ indica o quanto decresce a variância quando se res

tringe a pesquisa apenas à parte não nula da população. Notamos que ês se decréscimo é maior, quando a proporção Q de nulos é grande. Ou quando o coeficiente de variação C_o é pequeno.

Evidentemente êste lucro na variância deve ser balanceado com o custo necessário para a eliminação, *a priori*, dos elementos nulos da população.

CAPÍTULO 11

1. INTRODUÇÃO

Nêste capítulo iremos tratar da amostragem estratificada. Suponhamos nossa população N , dividida em k estratos mutuamente exclusivos $N_1, N_2, N_3, \dots, N_k$, contendo cada estrato respectivamente $N_1, N_2, N_3, \dots, N_k$ elementos. Logo o número total N de elementos da população é:

$$N = \sum_{i=1}^k N_i \quad (1.1)$$

A subpopulação M interseccionará os estratos, em k partes mutuamente exclusivas, $M_1, M_2, M_3, \dots, M_k$, cada uma contendo $M_1, M_2, M_3, \dots, M_k$ elementos. O número total M de elementos da subpopulação é dado por

$$M = \sum_{i=1}^k M_i \quad (1.2)$$

Os N valôres da variável Y , serão indicados por dois índices. Assim Y_{ij} estará indicando a observação feita na j -ésima unidade do i -ésimo estrato. De acôrdo com o exposto acima Y_{ij} os índices terão os seguintes domínios: $i = 1, 2, \dots, k$ e $j = 1, 2, \dots, N_i$.

Como foi introduzido na primeira parte dêste trabalho, vamos definir a variável auxiliar:

$$W_{ij} = \begin{cases} Y_{ij} & \text{se } Y_{ij} \in M \\ 0 & \text{se } Y_{ij} \notin M \end{cases} \quad (1.3)$$

O nosso problema consiste em estimar, para a subpopulação, a média \bar{W} ou o total W . Com:

$$\bar{W} = \frac{1}{M} \sum_{i=1}^k \sum_{j=1}^{N_i} W_{ij} = \frac{1}{M} \sum_{i=1}^k N_i \bar{W}_{io} = \frac{1}{M} \sum_{i=1}^k M_i \bar{W}_i \quad (1.4)$$

e

$$W = \sum_{i=1}^k \sum_{j=1}^{N_i} W_{ij} = \sum_{i=1}^k M_i \bar{W}_i = \sum_{i=1}^k N_i \bar{W}_{io} \quad (1.5)$$

onde

$$\bar{W}_i = \frac{1}{M_i} \sum_{j=1}^{N_i} W_{ij} \quad (1.6) \quad \text{e} \quad \bar{W}_{io} = \frac{1}{N_i} \sum_{j=1}^{N_i} W_{ij} \quad (1.7)$$

Retiremos agora dentro de cada estrato, uma amostra casual simples, contendo n_1, n_2, \dots, n_k elementos, respectivamente de cada

estrato. O número total n da amostra será:

$$n = \sum_{i=1}^k n_i \quad (1.8)$$

Essa amostra conterà m_1 elementos de M_1 , m_2 elementos de M_2 , etc. Dêsse modo o número m de elementos da subpopulação, que fazem parte da amostra é

$$m = \sum_{i=1}^k m_i \quad (1.9)$$

Através das mesmas considerações feitas na primeira parte do trabalho, notamos que fixado o tamanho n da amostra, o tamanho m é uma variável aleatória, bem como serão m_1, m_2, \dots, m_k .

Ainda de acôrdo com a nomenclatura usada na primeira parte, onde as letras minúsculas eram usadas na representação da amostra, teremos:

$$\bar{w} = \frac{1}{m} \sum_{i=1}^k \sum_{j=1}^{n_i} w_{ij} = \frac{1}{m} \sum_{i=1}^k m_i \bar{w}_i = \frac{1}{m} \sum_{i=1}^k n_i \bar{w}_{i0} \quad (1.10)$$

onde

$$\bar{w}_i = \frac{1}{m_i} \sum_{j=1}^{n_i} w_{ij} \quad e \quad \bar{w}_{i0} = \frac{1}{n_i} \sum_{j=1}^{n_i} w_{ij} \quad (1.11)$$

ou para indicar o total da amostra:

$$w = \sum_{i=1}^k \sum_{j=1}^{n_i} w_{ij} = \sum_{i=1}^k m_i \bar{w}_i = \sum_{i=1}^k n_i \bar{w}_{i0} \quad (1.12)$$

Usaremos ainda neste capítulo, as notações:
a variável auxiliar:

$$x_{ij} = \begin{cases} 1 & \text{se } y_{ij} \in M \\ 0 & \text{se } y_{ij} \notin M \end{cases} \quad (1.13)$$

logo

$$M = \sum_{i=1}^k \sum_{j=1}^{N_i} x_{ij} = \sum_{i=1}^k N_i \bar{x}_i = \sum_{i=1}^k M_i \quad (1.14)$$

$$m = \sum_{i=1}^k \sum_{j=1}^{n_i} x_{ij} = \sum_{i=1}^k n_i \bar{x}_i = \sum_{i=1}^k m_i$$

as proporções :

$$P_i = \frac{M_i}{N_i} \quad e \quad Q_i = 1 - P_i = 1 - \frac{M_i}{N_i} \quad (1.15)$$

então $\bar{X}_i = P_i$

as frações amostrais

$$f_i = \frac{n_i}{N_i} \quad \text{e} \quad f'_i = \frac{m_i}{M_i} \quad (1.16)$$

2. ESTIMATIVA DE \bar{W} COM OS M_i CONHECIDOS

Quando conhecem-se os valores de M_i , Cochran (pág.147) sugere um processo que é a extensão das suposições feitas na secção I.4. Usando como estimador \bar{w} , \bar{w}' definido do seguinte modo:

$$\bar{w}' = \frac{1}{M} \sum_{i=1}^k M_i \bar{w}_i \quad (2.1)$$

onde de acôrdo com as conclusões daquela secção, \bar{w}_i é um estimador não viesado de \bar{W}_i , logo

$$E(\bar{w}') = \frac{1}{M} \sum_{i=1}^k M_i E(\bar{w}_i) = \frac{1}{M} \sum_{i=1}^k M_i \bar{W}_i = \bar{W} \quad (2.2)$$

e ainda de acôrdo com (I.4.3)

$$V(\bar{w}_i) = (1-f'_i) \frac{S_{wi}^2}{m_i} \quad (2.3)$$

onde S_{wi}^2 indica a variância da variável W no i -ésimo estrato,

$$S_{wi}^2 = \frac{1}{M_i - 1} \sum_j^{M_i} (W_{ij} - \bar{W}_i)^2 \quad (2.4)$$

onde $\sum_j^{M_i}$ indica a somatória estendida apenas aos valores não nulos de W_{ij} , no i -ésimo estrato.

Usando (2.3) a variância de \bar{w}' reduz-se à

$$V(\bar{w}') = \frac{1}{M^2} \sum_{i=1}^k M_i^2 (1-f'_i) \frac{S_{wi}^2}{m_i} \quad (2.5)$$

No caso particular de partilha proporcional ao tamanho dos estratos, e de variância constante em cada estrato, (2.5) se reduz a:

$$V(\bar{w}') = \frac{1}{Mf'} (1-f') S_w^2 = (1-f') \frac{S_w^2}{E(m)} \quad (2.6)$$

3. \bar{w} COMO ESTIMADOR RAZÃO

Como desejamos estimar a média \bar{W} da população, definida em

(1.7), uma sugestão é determinar um estimador não viesado para o numerador e outro para o denominador. Assim um estimador não viesado para:

$$W = \sum_{i=1}^k \sum_{j=1}^{N_i} W_{ij} = \sum_{i=1}^k N_i \bar{w}_{i0} \quad (3.1)$$

e

$$w = \sum_{i=1}^k N_i \bar{w}_{i0} = \sum_{i=1}^k \frac{N_i}{n_i} \sum_{j=1}^{n_i} w_{ij} \quad (3.2)$$

E um estimador não viesado para:

$$M = \sum_{i=1}^k \sum_{j=1}^{N_i} W_{ij} = \sum_{i=1}^k M_i = \sum_{i=1}^k N_i \bar{x}_i \quad (3.3)$$

e

$$x = \sum_{i=1}^k N_i \bar{x}_i \quad (3.4)$$

onde $\bar{x}_i = \frac{m_i}{n_i}$, mas conforme vimos na secção I.4

$$E(\bar{x}_i) = \frac{M_i}{N_i} \quad (3.5)$$

então, calculando a esperança de x, temos imediatamente que

$$E(x) = M \quad (3.6)$$

Vamos definir como estimador de \bar{W} , o quociente das expressões obtidas em (3.2) e (3.4), assim:

$$\bar{w}_e = \frac{\sum N_i \bar{w}_{i0}}{\sum N_i \bar{x}_i} = \frac{\frac{1}{N} \sum N_i \bar{w}_{i0}}{\frac{1}{N} \sum N_i \bar{x}_i} = \frac{\bar{w}_{es}}{\bar{x}_{es}} \quad (3.7)$$

onde

$$\bar{w}_{es} = \frac{1}{N} \sum N_i \bar{w}_{i0} \quad (3.8)$$

$$e \quad \bar{x}_{es} = \frac{1}{N} \sum N_i \bar{x}_i \quad (3.9)$$

Este estimador, que é o estimador razão combinado devido a Hansen, Hurwitz e Gurney, é proposto por Cochran (pág.148) e por Durbin (pág. 114), para estimar \bar{W} .

Fava (pág.29) calculando a esperança matemática desse estimador, no caso geral, chegou em primeira aproximação ao seguinte resultado:

$$E(\bar{w}_e) = \bar{W} + \epsilon \quad (3.10)$$

onde ϵ , o viés, é dado pela seguinte expressão:

$$\epsilon = \frac{1}{\bar{X}^2} \sum_i \frac{N_i(N_i - n_i)}{N n_i} (\bar{W} S_{xi}^2 - S_{wxi}) \quad (3.11)$$

particularizando para a nossa variável,

$$S_{xi}^2 = \frac{1}{N_i - 1} \left\{ \sum_{j=1}^{N_i} (X_{ij} - \bar{X}_i)^2 \right\} = \frac{M_i}{N_i - 1} \left(1 - \frac{M_i}{N_i} \right) \cong P_i Q_i \quad (3.12)$$

$$e \quad S_{wxi} = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (W_{ij} - \bar{W}_{i0}) (X_{ij} - \bar{X}) =$$

$$= \frac{N_i \bar{W}_{i0}}{N_i - 1} \left(1 - \frac{M_i}{N_i} \right) \cong P_i Q_i \bar{W}_i \quad (3.13)$$

substituindo estes dois resultados em (3.11), vem :

$$\epsilon = \frac{N^2}{M^2} \sum_i \frac{N_i(N_i - n_i)}{N n_i} P_i Q_i (\bar{W} - \bar{W}_i) = \frac{N}{M^2} \sum_i \frac{N_i^2}{n_i} (1 - f_i) P_i Q_i (\bar{W} - \bar{W}_i) \quad (3.14)$$

Caso a estratificação tenha sido feita com partilha proporcional, isto é,

$$f_i = \frac{n_i}{N_i} = \frac{n}{N} \quad \text{para } i = 1, 2, \dots, k$$

então

$$\epsilon = \frac{N^2}{M^2} \frac{(1 - f)}{n} \sum_i N_i P_i Q_i (\bar{W} - \bar{W}_i) \quad (3.15)$$

ou lembrando que $P_i Q_i \leq \frac{1}{4}$, sempre, vem:

$$\epsilon \leq \frac{1}{4} \frac{N^2}{M^2} \frac{(1 - f)}{n} \sum_i N_i (\bar{W} - \bar{W}_i)$$

o que nos permite ver mais claramente, a magnitude do viés em função de seus fatores componentes.

Se a amostragem é feita com partilha proporcional e as variâncias e covariâncias são constantes de extrato para extrato é fácil mostrar usando raciocínio visto no Capítulo I que o viés é zero.

4. CÁLCULO DA VARIÂNCIA DE \bar{w}_i .

Queremos determinar :

$$V(\bar{w}_e) = E(\bar{w}_e - \bar{W})^2$$

mas

$$\bar{w}_e - \bar{W} = \frac{\bar{w}_{es}}{\bar{x}_{es}} - \bar{W} = \frac{\bar{w}_{es} - \bar{W}\bar{x}_{es}}{\bar{x}_{es}}$$

substituíamos a variável \bar{x}_{es} pela sua esperança \bar{X} . Assim:

$$\bar{w}_e - \bar{W} = \frac{1}{\bar{X}} (\bar{w}_{es} - \bar{W}\bar{x}_{es}) = \frac{1}{\bar{X}} \bar{u}_{es}$$

onde $\bar{u}_{es} = \bar{w}_{es} - \bar{W}\bar{x}_{es}$

então
$$V(\bar{w}_e) = \frac{1}{\bar{X}^2} V(\bar{u}_{es}) \quad (4.1)$$

e \bar{u}_{es} representa a média de uma amostra da variável aleatória

$w_{ij} - \bar{W}x_{ij}$, colhida através de amostragem estratificada. Assim:

$$V(\bar{u}_{es}) = \frac{1}{N^2} \sum_i \frac{N_i(N_i - n_i)}{n_i} S_{ui}^2 \quad (4.2)$$

onde

$$\begin{aligned} S_{ui}^2 &= \frac{1}{N_i - 1} \sum_j^{N_i} \left\{ (w_{ij} - \bar{W}x_{ij}) - (\bar{w}_{io} - \bar{W}\bar{X}_i) \right\}^2 = \\ &= \frac{1}{N_i - 1} \sum_j^{N_i} \left\{ (w_{ij} - \bar{W}x_{ij})^2 - N_i (\bar{w}_{io} - \bar{W}\bar{X}_i)^2 \right\} = \\ &= \frac{1}{N_i - 1} \sum_j^{M_i} \left\{ (w_{ij} - \bar{w}_i)^2 + M_i (\bar{w}_i - \bar{W})^2 - N_i P_i^2 (\bar{w}_i - \bar{W})^2 \right\} = \\ &= \frac{1}{N_i - 1} \sum_j^{M_i} \left\{ (w_{ij} - \bar{w}_i)^2 + M_i Q_i (\bar{w}_i - \bar{W})^2 \right\} \quad (4.3) \end{aligned}$$

substituindo (4.3) em (4.2), e o resultado em (4.1). vem:

$$V(\bar{w}_e) = \frac{1}{M^2} \sum_i \frac{N_i^2(1 - f_i)}{n_i(N_i - 1)} \left\{ \sum_j^{M_i} (w_{ij} - \bar{w}_i)^2 + M_i Q_i (\bar{w}_i - \bar{W})^2 \right\} \quad (4.4)$$

que é a variância procurada. Cochran (pág. 149) sugere como estimador dessa variância a expressão:

$$v(\bar{w}_e) \cong \frac{1}{x^2} \sum_i \frac{N_i^2(1 - f_i)}{n_i(n_i - 1)} \left\{ \sum_j^{m_i} (w_{ij} - \bar{w}_i)^2 + m_i \left(1 - \frac{m_i}{n_i}\right) (\bar{w}_i - \bar{w}_e)^2 \right\}$$

onde x é dado por (3.4):

$$x = \sum_i^k N_i \bar{x}_i = \sum_i^k \frac{N_i}{n_i} m_i$$

Durbin (pág. 115) procurou medir qual o aumento sofrido pela variância devido a aleatoriedade de m . Para isso comparou a variância do estimador \bar{w}_e com \bar{w}' , obtido na secção 2. Mas para isso é necessário algumas aproximações. Suponhamos que a estratificação foi tomada, em ambos os casos, com partilha proporcional ao tamanho do estrato, e que $f' \approx f$. Nestas condições, de (4.4), vem:

$$V(\bar{w}_e) = \frac{1}{M^2} \frac{(1-f)}{f} \left\{ \sum_i^k \sum_j^i \frac{M_i}{j} (\bar{w}_{ij} - \bar{w}_i)^2 + \sum_i M_i Q_i (\bar{w}_i - \bar{w})^2 \right\} \quad (4.5)$$

e de (2.5)

$$V(\bar{w}') = \frac{1}{M^2} \frac{(1-f)}{f} \sum_i^k \sum_j^i \frac{M_i}{j} (\bar{w}_{ij} - \bar{w}_i)^2 \quad (4.6)$$

logo

$$\frac{V(\bar{w}_e)}{V(\bar{w}')} = 1 + \frac{\sum_i^k M_i Q_i (\bar{w}_i - \bar{w})^2}{\sum_i^k \sum_j^i \frac{M_i}{j} (\bar{w}_{ij} - \bar{w}_i)^2}$$

e o último quociente, é que mede o aumento da variância devido ao desconhecimento do tamanho dos estratos subpopulacionais.

5. ESTIMATIVA DO TOTAL W DA SUBPOPULAÇÃO

5.1 Usando o fator de crescimento $\frac{N_i}{n_i}$

Quando desconhecem-se os valores M_i , podemos usar o fator de crescimento $\frac{N_i}{n_i}$, e aplicando as mesmas considerações da secção I.8.3, dentro de cada estrato, na estimativa do total W da subpopulação, podemos usar o estimador não viesado

$$\hat{w} = \sum_i^k \frac{N_i}{n_i} \sum_j^i w_{ij} \quad (5.1)$$

e de acordo com o que vimos naquela secção, a variância desse estimador será dado por:

$$V(\hat{w}) = \sum_i \frac{N_i^2 (1 - f_i)}{n_i} S_{w_{io}}^2 \quad (5.2)$$

com

$$S_{w_{io}}^2 = \frac{1}{N_i - 1} \sum_j^i (w_{ij} - \bar{w}_{io})^2 = \frac{1}{N_i - 1} \left\{ (N_i - 1) S_{w_i}^2 + M_i Q_i \bar{w}_i^2 \right\} \quad (5.3)$$

ou seja,

$$V(\hat{w}) = \sum_i^k \frac{N_i^2(1-f_i)}{n_i(N_i-1)} \{ (M_i-1) S_{wi}^2 + M_i Q_i \bar{w}_i^2 \} \quad (5.4)$$

Este estimador \hat{w} , está em Cochran (pág.147), que apresenta como estimador da variância (5.2), a expressão:

$$v(\hat{w}) = \sum_i^k \frac{N_i^2(1-f_i)}{n_i(n_i-1)} \left\{ \sum_j^m w_{ij}^2 - n_i \bar{w}_{io}^2 \right\} \quad (5.5)$$

5.2 Usando o estimador razão separado

Para estimar o total da subpopulação quando conhecem-se os valores M_i , podemos usar o estimador

$$w_s = \sum_i^k M_i \bar{w}_i = \sum_i^k M_i \frac{\bar{w}_{io}}{\bar{x}_i} \quad (5.6)$$

que identificamos imediatamente como estimador razão separado (conf. Fava, pág.28). Vamos determinar a esperança e a variância desse estimador, a partir dos resultados chegados por Fava (pág.30), para estimador razão separado no caso geral. Assim, em primeira aproximação

$$E(w_s) = W + \epsilon_s \quad (5.7)$$

onde

$$\epsilon_s = \sum_i^k \epsilon_i = \sum_i^k \frac{N_i - n_i}{n_i \bar{x}_i} (R_i S_{xi}^2 - S_{wxi}) \quad (5.8)$$

com

$$R_i = \frac{\bar{w}_{io}}{\bar{x}_i}$$

$$S_{xi}^2 = \frac{1}{N_i - 1} \sum_j^{N_i} (x_{ij} - \bar{x}_i)^2 = \frac{M_i}{N_i - 1} Q_i$$

$$S_{wi} = \frac{1}{N_i - 1} \sum_j^{N_i} (w_{ij} - \bar{w}_{io})(x_{ij} - \bar{x}_i) \quad (5.9)$$

Mas a expressão entre parêntesis da equação (5.8), é nula, pois dentro de cada estrato, a reta de regressão de W em função de X passa pela origem, e estamos na situação vista na secção I.5. Ou seja, w_s é em primeira aproximação, um estimador não viesado de W .

Para se ter uma idéia da magnitude do verdadeiro viés, podemos fazer as seguintes suposições:

- $\epsilon_i = \epsilon$, para todo $i = 1, 2, \dots, k$, isto é, temos aproximadamente o mesmo viés de estrato para estrato;
- A variância de W dentro de cada estrato, é a mesma;
- $C.V.(\bar{x}_i) = C.V.(\bar{x})$, isto é, também o coeficiente de variação das médias de x em cada estrato são aproximadamente iguais.

Nessas condições teremos

$$\frac{|\epsilon_s|}{\sqrt{V(w_s)}} \leq \frac{\sqrt{k} Q}{n} \quad (5.10)$$

sendo esta uma expressão para indicar a magnitude do verdadeiro viés.

A partir da expressão 2.3.2 de Fava (Pág.31), determinamos a variância de w_s

$$V(w_s) = \sum_i^k \frac{N_i^2 (1-f_i)}{n_i (N_i-1)} (M_i-1) S_{wi}^2 \quad (5.11)$$

Se compararmos este resultado com (5.4), teremos uma expressão para indicar em quanto diminui a variância pelo fato de conhecermos os valores M_i .

CAPÍTULO III

AMOSTRAGEM EM DOIS ESTÁGIOS

Vamos desenvolver aqui a mesma técnica dada por Durbin (pág.117). Consideremos inicialmente a amostragem em dois estágios, na qual k unidades são selecionadas, com reposição, na primeira fase de uma população de K unidades, com probabilidade de seleção proporcional aos números N_1, N_2, \dots, N_K , que são os tamanhos das unidades amostrais da primeira fase. Na segunda fase é selecionada, em cada unidade da primeira fase, um número fixo n de elementos. Do mesmo modo como foi apresentado em II.1, queremos estimar parâmetros de uma subpopulação M . Por exemplo a média da subpopulação

$$\bar{w} = \frac{1}{M} \sum_{i=1}^K \sum_{j=1}^{N_i} w_{ij} \quad (1.1)$$

e pode ser estimado por:

$$\bar{w}_1 = \frac{1}{m} \sum_{i=1}^k \sum_{j=1}^n w_{ij} \quad (1.2)$$

Numa primeira ordem de aproximação, podemos dizer que:

$$E(\bar{w}_1) = \frac{E\left(\sum_{i=1}^k \sum_{j=1}^n w_{ij}\right)}{E(m)} = \frac{\frac{kn}{N} \sum_{i=1}^K \sum_{j=1}^{N_i} w_{ij}}{\frac{kn}{N} M} = \bar{w} \quad (1.3)$$

ou seja, é um estimador aproximadamente não viesado. Para calcular a variância usaremos processo semelhante ao desenvolvido na secção II.4, assim

$$\begin{aligned} \bar{w}_1 - \bar{w} &= \frac{\sum_{i=1}^k \sum_{j=1}^n w_{ij}}{\sum_{i=1}^k \sum_{j=1}^n x_{ij}} - \bar{w} = \frac{\sum_{i=1}^k \sum_{j=1}^n (w_{ij} - \bar{w}x_{ij})}{m} = \\ &= \frac{k}{m} \frac{1}{k} \sum_{i=1}^k \sum_{j=1}^n u_{ij} \end{aligned} \quad (1.4)$$

onde $u_{ij} = w_{ij} - \bar{w}x_{ij}$. E a variância será dada por

$$v(\bar{w}_1) = \left(\frac{k}{m}\right)^2 v(\bar{u}) \quad (1.5)$$

onde \bar{u} é a média aritmética simples de k valores da forma:

$$u_i = \sum_{j=1}^n u_{ij} = w_i - \bar{w}m_i, \text{ e de acordo com a teoria da amostragem:}$$

$$v(\bar{u}) = \frac{1}{k} \frac{1}{k-1} \sum_i^k (u_i - \bar{u})^2 \quad (1.6)$$

$$\text{mas } \sum_i^k (u_i - \bar{u})^2 = \sum u_i^2 - k\bar{u}^2 = \sum (w_i - m_i\bar{w})^2 - \frac{m^2}{k} (\bar{w}_1 - \bar{w})^2 \quad (1.7)$$

Mas esta expressão contém o valor desconhecido \bar{w} , que não pode ser simplesmente substituído por \bar{w}_1 que é viesado. Entretanto, se os valores m_i forem iguais entre si, e portanto iguais a $\frac{m}{k}$, teremos

$$\sum_i^k (w_i - \frac{m}{k}\bar{w})^2 = \sum (w_i - \frac{m}{k}\bar{w}_1)^2 + \frac{m^2}{k} (\bar{w}_1 - \bar{w})^2$$

então, é razoável supor válida a relação:

$$\sum_i^k (w_i - m_i\bar{w})^2 \approx \sum (w_i - m_i\bar{w}_1)^2 + \frac{m^2}{k} (\bar{w}_1 - \bar{w})^2 \quad (1.8)$$

substituindo este último resultado em (1.7), vem

$$\sum_i^k (u_i - \bar{u})^2 = \sum (w_i - m_i\bar{w}_1)^2$$

ou seja a variância estimada em (1.5) reduz-se a

$$v(\bar{w}_1) = \frac{k}{m^2(k-1)} \sum_i^k (w_i - m_i\bar{w}_1)^2 \quad (1.9)$$

que é o estimador da variância procurada.

BIBLIOGRAFIA:

- Cochran, W.G - SAMPLING TECHNIQUES - second edition - John Wiley & Sons - 1965.
- Durbin, J. SAMPLING THOERY FOR ESTIMATES BASED ON FEWER INDIVIDUAL THAN THE MEMBER SELECTED - Bulletin of International Statistical Institute, Vol.36, pág. 113-119 - 1958.
- FAVA, L. - CONTRIBUIÇÃO PARA O ESTUDO DA ESTIMATIVA RAZÃO - Tese apresentada para Concurso de livre docência - F.F.C.L. - U.S.P. - 1959.
- Feller, W. - AN INTRODUCTION TO PROBABILITY THEORY AND APPLICATIONS John Wiley & Sons - Vol. 1 - 1950.
- Hansen, M.H. - SAMPLE SURVEY, METHODS AND THEORY (Vol. I e II) John Hurwitz, W.N. Wiley & Sons , 1953.
- Madow, W.G.
- Kish, L. - SURVEY SAMPLING - John Wilwy & Sons, 1967.
- Sukhatme, P.V. - SAMPLING THEORY OF SURVEY WITH APPLICATIONS - The Indian Society of Agricultural Statistics, New Delhi, 1954.
- YATES, F. - SAMPLING METHODS FOR CENSUSES AND SURVEYS - Charles Griffin and Co., London, 1953.