

Universidade de Évora
Curso de Mestrado em Matemática Aplicada 1995/97

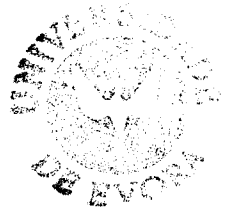
**TÉCNICAS
DE
AMOSTRAGEM
E
SUB-AMOSTRAGEM**

**Dissertação realizada por:
Carla Maria Lopes da Silva Afonso dos Santos**

**ÉVORA
1997**

Universidade de Évora
Curso de Mestrado em Matemática Aplicada 1995/97

**TÉCNICAS
DE
AMOSTRAGEM
E
SUB-AMOSTRAGEM**



87262

Dissertação realizada por:
Carla Maria Lopes da Silva Afonso dos Santos

ÉVORA
1997

Autora: Carla Maria Lopes da Silva Afonso dos Santos

Orientador: Professor Doutor João Tiago Mexia

Co-orientador: Professor Doutor Carlos Braumann

1. INTRODUÇÃO	1
2. AMOSTRAGEM ALEATÓRIA SIMPLES.....	2
2.1. Amostragem com probabilidades iguais, sem reposição, para características numéricas	2
2.2. Amostragem com probabilidades iguais, sem reposição, para atributos	9
2.3. Amostragem com probabilidades iguais, com reposição, para características numéricas ...	10
2.4. Amostragem com probabilidades iguais, com reposição, para atributos.....	12
2.5. Pré-amostragem	13
2.6. Utilização de tabelas de números aleatórios	16
3. AMOSTRAGEM SIMPLES COM PROBABILIDADES VARIÁVEIS	17
3.1. Sem reposição	17
3.2. Com reposição	23
4. AMOSTRAGEM ESTRATIFICADA	30
5. AMOSTRAGEM PARA AGREGADOS.....	36
6. SUB-AMOSTRAGEM.....	41
6.1. Resultados prévios	42
6.2. Notações.....	42
6.3. Sub-amostragem com probabilidades iguais, sem reposição	46
6.3.1. Construção do estimador.....	47
6.3.2. Variância do estimador.....	48
6.4. Sub-amostragem com probabilidades proporcionais.....	53
6.4.1. Construção do estimador.....	55
6.4.2. Variância do estimador.....	55
6.5. Custo Total	59

1. INTRODUÇÃO

Apesar de o principal objectivo desta dissertação ser a apresentação de resultados novos no domínio da sub-amostragem, pretende-se igualmente apresentar, de uma forma consistente, as principais técnicas de amostragem.

Começa-se pelo estudo da amostragem simples, com e sem reposição, analisando, separadamente, as modalidades em que a probabilidade de escolha, dos elementos da população, é igual, e é variável.

Em seguida considera-se a amostragem estratificada e a amostragem para agregados, assim como um caso especial desta última, a amostragem sistemática.

Finalmente, dado que a teoria, referente às técnicas de sub-amostragem, estava apenas desenvolvida para, quando muito, três níveis de fraccionamento, generalizámos a mesma para qualquer número (k) de níveis, abordando os casos em que as probabilidades de escolha são iguais e proporcionais. Para tal, utilizou-se uma notação que permite referenciar, sem ambiguidades, o nível de estrutura a que pertence determinada sub-população e qual a posição que ocupa, assim como, considerando as decomposições sucessivas da população inicial, saber, exactamente, quais as sub-populações que deram origem a essa sub-população.

2. AMOSTRAGEM ALEATÓRIA SIMPLES

A amostragem aleatória simples, apesar de na prática ser pouco usada, é um método de amostragem de grande importância pois funciona como ponto de partida para o estudo de outros métodos.

Neste método a população é considerada como um todo e as unidades estatísticas que vão pertencer à amostra são seleccionadas ao acaso de entre todos os elementos da população. Cada unidade estatística tem uma probabilidade conhecida e não nula de figurar na amostra, podendo esta probabilidade ser igual para todos os elementos ou variar entre eles.

A amostragem aleatória simples, com probabilidades iguais, pressupõe um total desconhecimento dos caracteres a estudar na população e como consequência todas as unidades estatísticas têm a mesma probabilidade, de serem escolhidas para a amostra.

As unidades estatísticas podem ser retiradas, da população, com ou sem reposição.

A selecção é considerada com reposição se retiramos um elemento da população, e após a sua observação o devolvemos à população. Assim, a população não é alterada, de extracção para extracção, e portanto, cada unidade estatística pode aparecer mais do que uma vez na amostra.

As extracções são independentes, visto que, a extracção de uma determinada unidade estatística não depende das extracções anteriores.

A selecção considera-se sem reposição se retiramos um elemento da população, observamo-lo e não o devolvemos à população. A população é alterada de extracção para extracção, ou seja, a extracção da segunda unidade estatística é feita numa população que difere da população inicial porque já não contém a unidade estatística retirada em primeiro lugar. Neste caso, cada elemento só pode aparecer uma vez na amostra. A extracção de uma determinada unidade estatística depende das extracções anteriores, pelo que, as extracções são acontecimentos dependentes.

2.1. AMOSTRAGEM COM PROBABILIDADES IGUAIS, SEM REPOSIÇÃO, PARA CARACTERÍSTICAS NUMÉRICAS

Este é um método de selecção de n elementos de entre N existentes na população, no qual em cada tiragem qualquer elemento da população tem igual probabilidade de ser escolhido. Para além disso cada elemento que é extraído é removido da população para as tiragens subseqüentes.

Os n elementos escolhidos constituem a amostra.

Sejam y_i , $i=1,2,\dots,N$, os valores da característica numérica, para a qual se está a amostrar, os quais indicam o número de ordem dos diferentes elementos da população e Y_j , $j=1,2,\dots,n$, as variáveis aleatórias que dão os resultados obtidos nas várias tiragens. Estas variáveis tomam como valores, os valores da característica numérica nos elementos que vão sendo seleccionados.

Ao considerar que as probabilidades de selecção são iguais para todos os elementos da população temos

$$(2.1) \quad P(Y_1 = y_i) = \frac{1}{N}$$

e para $j > 1$ temos

$$(2.2) \quad \begin{cases} P \left[Y_j = y_i \mid \bigcap_{j'=1}^{j-1} (Y_{j'} \neq y_i) \right] = \frac{1}{N-j+1} \\ P \left[Y_j = y_i \mid \bigcup_{j'=1}^{j-1} (Y_{j'} = y_i) \right] = 0 \end{cases}$$

Para provar que a probabilidade de selecção é igual em todas as extracções, ou seja, que o j -ésimo elemento (y_j) tem a mesma probabilidade de ser seleccionado em qualquer extracção, consideremos, em primeiro lugar, a probabilidade deste elemento ser seleccionado na primeira extracção, que é:

$$(2.3) \quad P_j^1 = P(Y_1 = y_j) = \frac{1}{N} \quad j = 1, \dots, N$$

A probabilidade de que y_j seja seleccionado na segunda extracção é:

$$(2.4) \quad \begin{aligned} P_j^2 &= P(Y_2 = y_j) = \sum_{k \neq j} P(Y_1 = y_k, Y_2 = y_j) = \\ &= \sum_{k \neq j} P(Y_2 = y_j \mid Y_1 = y_k) \cdot P(Y_1 = y_k) = \\ &= \sum_{k \neq j} \frac{1}{N-1} \cdot \frac{1}{N} = (N-1) \frac{1}{N-1} \cdot \frac{1}{N} = \frac{1}{N} \quad j = 1, \dots, N \end{aligned}$$

A probabilidade de que y_j seja seleccionado na terceira extracção é:

$$(2.5) \quad \begin{aligned} P_j^3 &= P(Y_3 = y_j) = \sum_{j \neq k \neq m} P(Y_1 = y_k, Y_2 = y_m, Y_3 = y_j) = \\ &= \sum_{j \neq k \neq m} P[Y_3 = y_j \mid (Y_1 = y_k, Y_2 = y_m)] \cdot P(Y_2 = y_m \mid Y_1 = y_k) \cdot P(Y_1 = y_k) = \\ &= \sum_{j \neq k \neq m} P[Y_3 = y_j \mid (Y_1 = y_k, Y_2 = y_m)] \cdot \sum_{k \neq m} \frac{1}{N} \cdot \frac{1}{N-1} = \\ &= \sum_{j \neq k \neq m} \frac{1}{N-2} \cdot (N-1) \cdot \frac{1}{N-1} \cdot \frac{1}{N} = \\ &= (N-2) \cdot \frac{1}{N-2} \cdot (N-1) \cdot \frac{1}{N-1} \cdot \frac{1}{N} = \frac{1}{N} \quad j = 1, \dots, N \end{aligned}$$

Por recorrência, conclui-se que,

$$(2.6) \quad P_j^k = P(Y_k = y_j) = \frac{1}{N} \quad ; \quad j=1, \dots, N; \quad 1 \leq k \leq n; \quad n \leq N$$

é a probabilidade de y_j ser seleccionado na k -ésima extracção.

Pode ainda provar-se que, no método de amostragem aleatória simples com probabilidades iguais e sem reposição, a probabilidade de que um qualquer elemento pertença a uma amostra, de dimensão n , é $\frac{n}{N}$.

Para que o elemento y_j pertença a uma amostra, é necessário que tenha sido escolhido uma e uma só vez, visto tratar-se de extracções sem reposição.

Seja

$$(2.7) \quad P_j = P(y_j \in \text{amostra}) = 1 - P(y_j \notin \text{amostra})$$

Mas a probabilidade de y_j não pertencer à amostra é igual à probabilidade dos $N-1$ elementos y_k , com $k \neq j$, poderem pertencer à amostra.

Assim

$$(2.8) \quad P(y_j \notin \text{amostra}) = \frac{{}^{N-1}C_n}{{}^N C_n}$$

e portanto

$$(2.9) \quad P(y_j \in \text{amostra}) = 1 - \frac{{}^{N-1}C_n}{{}^N C_n} = \frac{{}^N C_n - {}^{N-1}C_n}{{}^N C_n} = \frac{{}^{N-1}C_{n-1}}{{}^N C_n} = \frac{n}{N},$$

pois ${}^{N-1}C_n + {}^{N-1}C_{n-1} = {}^N C_n$.

Para determinar quantas amostras diferentes, de n elementos, é possível obter com os N elementos da população há que ter em consideração se a ordem pela qual os elementos são seleccionados tem ou não interesse para a constituição da amostra, e também que, neste caso, a amostragem é feita sem reposição.

Se ignorarmos a ordem pela qual os elementos vão sendo obtidos, isto é, se considerarmos iguais as amostras que contêm os mesmos elementos, obtidos por diferentes ordens, o número de amostras possíveis é:

$$(2.10) \quad {}^N C_n = \frac{N!}{(N-n)!n!}.$$

Se tiver interesse a ordem pela qual se vão obtendo os elementos que irão compor a amostra, isto é, se considerarmos distintas as amostras compostas pelos mesmos elementos, mas obtidos por ordens diferentes, o número de amostras possíveis será:

$$(2.11) \quad {}^N A_n = \frac{N!}{(N-n)!}.$$

É fácil ver que todas as amostras distintas têm a mesma probabilidade de serem seleccionadas.

Consideremos um conjunto de n elementos específicos $\{y_1, \dots, y_n\}$, que constituem uma determinada amostra. Os valores dessa amostra correspondem, tendo em consideração que a amostragem é sem reposição, a elementos distintos da população.

A probabilidade de esta amostra ser colhida é:

(2.12)

$$P\left[\bigcap_{j=1}^n (Y_j = y_j)\right] = P\left[Y_n = y_n \mid \bigcap_{j=1}^{n-1} (Y_j = y_j)\right] \cdot P\left[\bigcap_{j=1}^{n-1} (Y_j = y_j)\right] = \dots = P\left[Y_n = y_n \mid \bigcap_{j=1}^{n-1} (Y_j = y_j)\right] \cdot P\left[Y_{n-1} = y_{n-1} \mid \bigcap_{j=1}^{n-2} (Y_j = y_j)\right] \cdot P\left[Y_2 = y_2 \mid Y_1 = y_1\right] \cdot P(Y_1 = y_1).$$

Como todos os elementos têm igual probabilidade de serem escolhidos, a probabilidade de ser obtida uma amostra específica, de n elementos com uma determinada ordem, é:

(2.13)

$$P\left[\bigcap_{j=1}^n (Y_j = y_j)\right] = \frac{1}{N-n+1} \cdot \frac{1}{N-n+2} \cdot \dots \cdot \frac{1}{N-1} \cdot \frac{1}{N} = \frac{(N-n)!}{N!}$$

Se, pelo contrário, a ordem pela qual são obtidos os elementos que irão compor a amostra for ignorada, a probabilidade de se obter essa amostra passa a ser:

(2.14)

$$P\left[\bigcap_{j=1}^n (Y_j = y_j)\right] = \frac{n!(N-n)!}{N!} = \frac{1}{{}^N C_n}$$

onde $n!$ indica por quantas ordens diferentes podem ser obtidos os n elementos.

Consideremos agora os acontecimentos, que se verificam quando o i -ésimo elemento é escolhido na j -ésima extracção,

(2.15)

$$A_{j,i} = (Y_j = y_i) ; j = 1, \dots, n, i = 1, \dots, N,$$

que como já vimos têm todos probabilidade $\frac{1}{N}$.

Assim as variáveis aleatórias Y_j , $j = 1, \dots, n$, são identicamente distribuídas, correspondendo-lhe o esquema

(2.16)

$$Y_j \left\{ \begin{array}{l} y_1, \dots, y_N \\ \frac{1}{N}, \dots, \frac{1}{N} \end{array} ; j = 1, \dots, n. \right.$$

Considerando

(2.17)

$$T = \sum_{i=1}^N y_i$$

e

$$(2.18) \quad S = \sum_{i=1}^N y_i^2 \quad ,$$

obtém-se o valor médio da característica (o parâmetro que se pretende estimar)

$$(2.19) \quad \mu = \mu(y_j) = \frac{T}{N} \quad , \quad j = 1, \dots, n$$

e a variância da característica

$$(2.20) \quad \sigma^2 = \sigma^2(Y_j) = \mu(Y_j^2) - \mu^2(Y_j) = \frac{1}{N} \left(S - \frac{T^2}{N} \right) \quad ; \quad j = 1, \dots, n,$$

com

$$(2.21) \quad \mu(Y_j^2) = \frac{S}{N} \quad ; \quad j = 1, \dots, n .$$

Para calcular $P\left[(Y_j = y_i) \cap ((Y_{j'} = y_{i'}))\right]$, $i \neq i'$, consideremos as amostras que é possível obter fixados dois elementos. Para construir essas amostras, há que obter amostras de dimensão $n-2$ da população original, uma vez excluídos os dois elementos dados.

Se considerarmos, com importância, a ordem pela qual vão sendo obtidos os elementos da amostra, existem

$$(2.22) \quad \frac{(N-2)!}{[(N-2)! - (n-2)!]} = \frac{(N-2)!}{(N-n)!}$$

amostras de dimensão $n-2$.

Se ignorarmos essa ordem, existirão

$$(2.23) \quad \frac{(N-2)!}{(n-2)! [(N-2)! - (n-2)!]} = \frac{(N-2)!}{(n-2)! (N-n)!}$$

Obtemos então, para o primeiro caso:

$$(2.24) \quad P\left[(Y_j = y_i) \cap (Y_{j'} = y_{i'})\right] = \frac{(N-2)! (N-n)!}{(N-n)! N!} = \frac{1}{N(N-1)}$$

e para o segundo:

$$(2.25) \quad P\left[(Y_j = y_i) \cap (Y_{j'} = y_{i'})\right] = \frac{(N-2)!}{(n-2)! (N-n)!} \frac{n!(N-n)!}{N!} = \frac{n(n-1)}{N(N-1)}$$

De modo a não haver uma duplicação desnecessária de cálculos, devido à grande semelhança dos dois casos, de agora em diante apenas se efectuará o estudo do caso em que a ordem de selecção dos elementos é tida em consideração.

Obtemos então:

$$\begin{aligned}
 \mu(Y_j; Y_{j'}) &= \frac{1}{N(N-1)} \sum_{i=1}^N \sum_{i'=1}^N Y_i Y_{i'} = \frac{1}{N(N-1)} \sum_{i=1}^N Y_i (T - Y_i) = \\
 (2.26) \quad &= \frac{1}{N(N-1)} \left(T \sum_{i=1}^N Y_i - S \right) = \frac{T^2 - S}{N(N-1)}
 \end{aligned}$$

já que, devido à amostragem ser sem reposição, $P\left[(Y_j = y_i) \cap (Y_j = y_{i'})\right] = 0$, e

$$\begin{aligned}
 \sigma(Y_j; Y_{j'}) &= \mu(Y_j; Y_{j'}) - \mu(Y_j)\mu(Y_{j'}) = \frac{T^2 - S}{N(N-1)} - \frac{T^2}{N^2} = \\
 (2.27) \quad &= -\frac{1}{N(N-1)} \left(S - \frac{T^2}{N} \right) = -\frac{\sigma^2}{N-1}
 \end{aligned}$$

Assim dado o vector aleatório

$$(2.28) \quad Y' = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}$$

tem-se

$$(2.29) \quad \mu'(Y') = \begin{bmatrix} \mu \\ \vdots \\ \mu \end{bmatrix}$$

bem como

$$(2.30) \quad \sigma^2(Y') = \sigma^2 \begin{bmatrix} 1 & \dots & -\frac{1}{N-1} \\ \vdots & & \vdots \\ -\frac{1}{N-1} & \dots & 1 \end{bmatrix}$$

Pretendemos agora obter o estimador linear, centrado com variância mínima, de μ .

Dado

$$(2.31) \quad \mu^* = k'^T Y' = \sum_{j=1}^n k_j Y_j$$

tem-se

$$(2.32) \quad \mu(\mu^*) = k'^T \mu'(Y') = \sum_{j=1}^n k_j \mu = \mu \sum_{j=1}^n k_j.$$

Para que μ^* seja centrado é necessário que

$$(2.33) \quad \sum_{j=1}^n k_j = 1.$$

Por outro lado

$$(2.34) \quad \begin{aligned} \sigma^2(\mu^*) &= k'^T \sigma^2(Y') k' = \\ &= [k_1 \quad \dots \quad k_n] \sigma^2 \begin{pmatrix} 1 & \dots & -\frac{1}{N-1} \\ \vdots & & \vdots \\ -\frac{1}{N-1} & \dots & 1 \end{pmatrix} \begin{bmatrix} k_1 \\ \vdots \\ k_n \end{bmatrix} = \\ &= \sigma^2 \left(\sum_{j=1}^n k_j^2 - \frac{1}{N-1} \sum_{j=1}^n \sum_{j'=1}^n k_j k_{j'} \right) = \\ &= \sigma^2 \left(\frac{N}{N-1} \sum_{j=1}^n k_j^2 - \frac{1}{N-1} \sum_{j=1}^n \sum_{j'=1}^n k_j k_{j'} \right) = \\ &= \sigma^2 \left(\frac{N}{N-1} \sum_{j=1}^n k_j^2 - \frac{1}{N-1} \left(\sum_{j=1}^n k_j \right) \left(\sum_{j'=1}^n k_{j'} \right) \right). \end{aligned}$$

Para estimadores centrados que satisfazem a condição (2.33), a expressão (2.34) reduz-se a

$$(2.35) \quad \sigma^2(\mu^*) = \frac{\sigma^2}{N-1} \left(N \sum_{j=1}^n k_j^2 - 1 \right)$$

O problema limita-se então a minimizar $\sum_{j=1}^n k_j^2$, sob a condição (2.33).

Utilizando multiplicadores de Lagrange, com uma função auxiliar do tipo

$$(2.36) \quad g(k_1, \dots, k_n, \lambda) = \sum_{j=1}^n k_j^2 - \lambda \left(\sum_{j=1}^n k_j - 1 \right)$$

procuramos a solução do sistema

$$(2.37) \quad \left\{ \begin{array}{l} \frac{\partial g}{\partial k_j}(k_1, \dots, k_n, \lambda) = 2k_j - \lambda = 0 \quad ; j = 1, \dots, n \\ \frac{\partial g}{\partial \lambda}(k_1, \dots, k_n, \lambda) = \sum_{j=1}^n k_j - 1 = 0 \end{array} \right.$$

que é

$$(2.38) \quad k_j = \frac{1}{n} ; j = 1, \dots, n.$$

Substituindo (2.38) na expressão (2.31), obtemos então

$$(2.39) \quad \mu^* = \sum_{j=1}^n k_j Y_j = \frac{\sum_{j=1}^n Y_j}{n} = \bar{Y}_n$$

pelo que se conclui que, o estimador centrado, com variância mínima, é a média amostral.

Substituindo os $k_j, j=1, \dots, n$, por $\frac{1}{n}$, na expressão (2.35), vem

$$(2.40) \quad \sigma^2(\bar{Y}_n) = \frac{\sigma^2}{N-1} \left(N \sum_{j=1}^n \frac{1}{n^2} - 1 \right) = \frac{\sigma^2}{N-1} \left(\frac{N}{n} - 1 \right) = \frac{\sigma^2}{N-1} \frac{N-n}{n} = \frac{\sigma^2}{n} \frac{N-n}{N-1}.$$

2.2. AMOSTRAGEM COM PROBABILIDADES IGUAIS, SEM REPOSIÇÃO, PARA ATRIBUTOS

Admitamos agora que, a amostragem é feita com o objectivo de estimar a probabilidade, p , de um elemento, extraído ao acaso da população, possuir determinado atributo. A percentagem dos elementos da população com esse atributo será $100p$.

De modo a podermos aplicar os resultados obtidos para características numéricas, consideramos uma característica numérica X que toma os valores $\underline{1}$ para os elementos da população que possuem o atributo e $\underline{0}$ para os elementos que não possuem. Isto é, a variável aleatória $X_i = 1$ se o elemento y_i possui o atributo e $X_i = 0$ se o elemento y_i não possui o atributo.

Seja M o número de elementos da população que possuem o atributo, então temos

$$(2.41) \quad p = \frac{M}{N}$$

e

$$(2.42) \quad T = S = M$$

logo

$$(2.43) \quad \mu = \frac{M}{N} = p$$

e

$$(2.44) \quad \sigma^2 = \frac{1}{N} \left(M - \frac{M^2}{N} \right) = p(1-p).$$

Numa amostra de dimensão n , onde m elementos possuem o atributo tem-se

$$(2.45) \quad p^* = \bar{Y}_n = \frac{m}{n}.$$

Como já vimos anteriormente, no caso das características numéricas

$$(2.46) \quad \mu(\bar{Y}_n) = \mu$$

e

$$(2.47) \quad \sigma^2(\bar{Y}_n) = \frac{N-n}{N-1} \frac{\sigma^2}{n}$$

logo

$$(2.48) \quad \mu(p^*) = p$$

e

$$(2.49) \quad \sigma^2(p^*) = \frac{N-n}{N-1} \frac{p(1-p)}{n},$$

portanto p^* é estimador centrado de p .

2.3. AMOSTRAGEM COM PROBABILIDADES IGUAIS, COM REPOSIÇÃO, PARA CARACTERÍSTICAS NUMÉRICAS

Neste caso, o elemento que é escolhido, em cada extracção, é devolvido à população, pelo que esta mantém-se igual em todas as extracções. Todos os elementos têm igual probabilidade de escolha, em todas as extracções, quer tenha ou não sido escolhido anteriormente.

A probabilidade de um qualquer elemento y_i ser seleccionado na j -ésima extracção é então:

$$(2.50) \quad P_i^j = P(Y_j = y_i) = \frac{1}{N} \quad ; j = 1, \dots, n \quad ; i = 1, \dots, N.$$

As variáveis aleatórias Y_j , $j=1, \dots, n$ são independentes e identicamente distribuídas, correspondendo-lhes o esquema

$$(2.51) \quad Y_j \begin{cases} y_1, \dots, y_N \\ \frac{1}{N}, \dots, \frac{1}{N} \end{cases}; j = 1, \dots, n$$

Podemos provar agora que, no método de amostragem aleatória simples com probabilidades iguais e com reposição, a probabilidade P_i de que o elemento y_i , $i=1, \dots, N$, pertença a uma qualquer amostra de dimensão n , é $1 - \left(1 - \frac{1}{N}\right)^n$.

A probabilidade de y_i pertencer à amostra é a probabilidade de y_i ter sido escolhido, pelo menos uma vez, para pertencer à amostra.

Atendendo a (2.50) a probabilidade de y_i não ser escolhido na j -ésima extracção é:

$$(2.52) \quad P(Y_j \neq y_i) = 1 - \frac{1}{N}.$$

Tendo em conta que a dimensão da amostra é n , se y_i não foi escolhido para a amostra não pode ter sido escolhido em nenhuma das n extracções. Como as extracções são independentes, a probabilidade de y_i não ser escolhido em nenhuma extracção é:

$$(2.53) \quad 1 - P_i = \prod_{i=1}^n \left(1 - \frac{1}{N}\right) = \left(1 - \frac{1}{N}\right)^n$$

A probabilidade de y_i ser seleccionado pelo menos uma vez é então:

$$(2.54) \quad P_i = 1 - \left(1 - P_i\right) = 1 - \left(1 - \frac{1}{N}\right)^n \quad \text{c. q. d.}$$

Como se pode ver facilmente, também neste caso se verificam as expressões (2.19) e (2.20). Visto o estimador linear centrado, com variância mínima, do valor médio comum a n variáveis aleatórias independentes e identicamente distribuídas ser a respectiva média aritmética, tomamos

$$\mu^* = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i, \text{ provamos que ele é estimador centrado de } \mu \text{ e que a sua variância é } \frac{\sigma^2}{n}.$$

Para tal, calculemos $\mu(\mu^*)$. Aplicando as propriedades da esperança matemática, obtemos

$$(2.55) \quad \mu(\mu^*) = \mu\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \frac{1}{n} \sum_{i=1}^n \mu(Y_i).$$

Para calcular $\mu(Y_i)$, usamos a definição que diz que, o valor esperado de uma variável, V , é a soma dos produtos obtidos, multiplicando cada possível valor de V pela sua probabilidade, somados sobre todos os possíveis valores de V , e também (2.50). Assim temos:

$$(2.56) \quad \mu(Y_i) = \sum_{j=1}^N y_j \cdot P_j^i = \frac{1}{N} \sum_{j=1}^N y_j = \mu$$

então

$$(2.57) \quad \mu(\mu^*) = \frac{1}{n} \cdot n \cdot \mu = \mu,$$

o que prova que, μ^* é estimador centrado de μ .

Vamos então calcular a variância de μ^* . Por definição de μ^* e aplicando as propriedades da variância temos

$$(2.58) \quad \sigma^2(\mu^*) = \sigma^2\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \frac{1}{n^2} \sigma^2\left(\sum_{i=1}^n Y_i\right)$$

como as extracções são independentes, temos

$$(2.59) \quad \sigma^2(\mu^*) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2(Y_i),$$

mas por definição

$$(2.60) \quad \sigma^2(Y_i) = \mu \left[(Y_i - \mu(Y_i))^2 \right]$$

e como já vimos em (2.56), $\mu(Y_i) = \mu$, logo:

$$(2.61) \quad \sigma^2(Y_i) = \sum_{j=1}^N (y_j - \mu)^2 P_j^i = \frac{1}{N} \sum_{j=1}^N (y_j - \mu)^2 = \sigma^2$$

Substituindo (2.61) em (2.59) temos:

$$(2.62) \quad \sigma^2(\mu^*) = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n}.$$

Comparando as variâncias $\sigma^2(\mu^*)$ obtidas para a amostragem com e sem reposição, vemos que, utilizando amostras da mesma dimensão, a amostragem aleatória simples com probabilidades iguais e sem reposição conduz a estimadores com menor variância, pois $\frac{N-n}{N-1} < 1$.

2.4. AMOSTRAGEM COM PROBABILIDADES IGUAIS, COM REPOSIÇÃO, PARA ATRIBUTOS

Para ser possível aplicar os resultados obtidos para as características numéricas consideramos novamente uma característica numérica que toma os valores 1 ou 0.

Apesar de a amostragem ser agora com reposição, as expressões (2.19), (2.20) e (2.45), relativas à amostragem para atributos sem reposição, continuam a ser válidas, e portanto p^* é, também neste caso, um estimador centrado de p , ou seja, $\mu(p^*) = \mu = p$.

Como já vimos, em (2.62) e (2.44), $\sigma^2(\bar{Y}) = \frac{\sigma^2}{n}$ e $\sigma^2 = p(1-p)$, pelo que a variância de p^* é agora

$$(2.63) \quad \sigma^2(p^*) = \frac{p(1-p)}{n}.$$

2.5. PRÉ - AMOSTRAGEM

Como já vimos, a amostragem com reposição, apesar de ser de mais fácil aplicação, conduz a estimadores com maior variância.

Quando se pretende obter resultados com uma dada precisão, realiza-se uma pré-amostragem com reposição para determinar a dimensão a dar à amostra definitiva.

Sejam Y_1, \dots, Y_n as observações da pré-amostra, pode admitir-se que

$$(2.64) \quad D = \sum_{i=1}^n Y_i^2 - \frac{\left(\sum_{i=1}^n Y_i\right)^2}{n}$$

se distribui, pelo menos aproximadamente, como um chi-quadrado com $n-1$ graus de liberdade e que

$$(2.65) \quad t = \frac{\bar{Y} - \mu}{\sqrt{\frac{D}{n(n-1)}}}$$

tem densidade t de Student, com $n-1$ graus de liberdade.

Assim os intervalos de confiança (aproximados) para μ são

$$(2.66) \quad \bar{Y} - t_{n-1,q} \sqrt{\frac{D}{n(n-1)}} \leq \mu \leq \bar{Y} + t_{n-1,q} \sqrt{\frac{D}{n(n-1)}}$$

com a semi-amplitude

$$(2.67) \quad A = t_{n-1,q} \sqrt{\frac{D}{n(n-1)}}.$$

Suponhamos que se pretende uma dada semi-amplitude \bar{A} . Terá então que se resolver a equação em ordem a n , de modo a determinar a dimensão de amostra conveniente para obter a precisão pré-definida.

Seja

$$(2.68) \quad \bar{A} = t_{n-1,q} \sqrt{\frac{D}{n(n-1)}}$$

fazendo

$$(2.69) \quad k = \left(\frac{\bar{A}}{t_{n-1,q}} \right)^2$$

ou seja

$$(2.70) \quad k = \frac{D}{n(n-1)}$$

obtém-se

$$(2.71) \quad kn^2 - kn - D = 0$$

vindo

$$(2.72) \quad n = \frac{k + \sqrt{k^2 + 4kD}}{2k}$$

(a outra solução não serve por ser negativa).

Na prática o valor de n obtido por esta fórmula é arredondado para cima. Quando $n > 120$ a densidade t coincide com a normal reduzida, podendo-se utilizar as tabelas desta última. Quando $n > 30$ pode-se, em geral, realizar testes t para hipóteses definidas a partir de μ .

No caso de a amostragem ser por atributos, quando $30 \leq n \leq 120$ pode admitir-se que

$$(2.73) \quad t = \frac{p^* - p}{\sqrt{\frac{p^*(1-p^*)}{n}}}$$

tem a densidade t com $n-1$ graus de liberdade e que, para $n > 120$, é normal reduzida.

Obtêm-se então os intervalos de confiança para p dados por

$$(2.74) \quad p^* - t_{n-1,q} \sqrt{\frac{p^*(1-p^*)}{n}} \leq p \leq p^* + t_{n-1,q} \sqrt{\frac{p^*(1-p^*)}{n}}$$

podendo igualmente realizar-se testes t para hipóteses definidas a partir de p .

A semi - amplitude do intervalo anterior é

$$(2.75) \quad A = t_{n-1,q} \sqrt{\frac{p^*(1-p^*)}{n}}$$

Se pretendermos obter uma determinada semi - amplitude , digamos \bar{A} , terá de se resolver a equação seguinte em ordem a n

$$(2.76) \quad \bar{A} = t_{n-1,q} \sqrt{\frac{p^*(1-p^*)}{n}}$$

Fazendo

$$(2.77) \quad k = \frac{p^*(1-p^*)}{n}$$

vem

$$(2.78) \quad n = \frac{p^*(1-p^*)}{k}$$

valor este que, tal como no caso referido anteriormente, se arredonda para cima.

Consideremos a função $g(p) = p(1-p)$ cujo gráfico está representado na figura 1 .

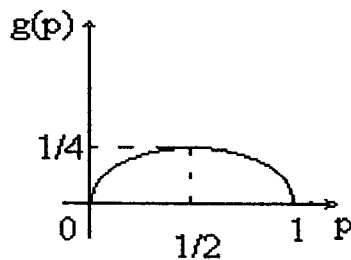


Figura 1

verificando-se

$$(2.79) \quad \begin{cases} g\left(\frac{1}{2}-q\right) = g\left(\frac{1}{2}+q\right) \\ g(0) = g(1) = 0 \\ g\left(\frac{1}{2}\right) = \frac{1}{4} \end{cases}$$

Concluimos então que, quanto mais próximo p^* estiver de $\frac{1}{2}$, mais observações é preciso tomar, para se obter a precisão pretendida.

2.6. UTILIZAÇÃO DE TABELAS DE NÚMEROS ALEATÓRIOS

A utilização de tabelas de números aleatórios é um método simples e satisfatório de obter uma amostra aleatória simples, desde que se consiga indexar os elementos da população por qualquer ordem conveniente.

As tabelas de números aleatórios tanto podem ser lidas na horizontal como na vertical, podendo ser tomados algarismos isolados como combinações de 2, 3 ou mais dígitos.

Partindo de determinado ponto (escolhido aleatoriamente), percorrem-se as linhas (colunas) da tabela num qualquer sentido, escolhendo-se os números que aparecem ao fim de um intervalo previamente estabelecido. Esses números indicarão quais os elementos da população que irão compor a amostra.

Apesar da facilidade de utilização destas tabelas há que ter em consideração que, em utilizações repetidas, se deve variar o ponto de partida, o sentido e o passo com que os números são seleccionados.

Há ainda que evitar utilizar a mesma tabela um número exagerado de vezes. Este procedimento assegura que todas as possíveis combinações têm igual probabilidade de serem seleccionadas.

3. AMOSTRAGEM SIMPLES COM PROBABILIDADES VARIÁVEIS

A amostragem simples, com probabilidades variáveis, pressupõe um conhecimento, à priori, da população e o uso desse conhecimento para atribuir, às unidades estatísticas, diferentes probabilidades de pertencerem à amostra.

Nos esquemas de selecção amostral discutidos até agora, todas as unidades da população tinham a mesma hipótese de serem seleccionadas para a amostra. Não é essencial que assim seja, a única exigência, que se deve fazer, é que as probabilidades de inclusão sejam conhecidas e não nulas. Pode, de facto, ser demonstrado que é possível atingir uma precisão mais elevada com probabilidades variáveis.

3.1. SEM REPOSIÇÃO

Como já vimos anteriormente, na amostragem sem reposição, os elementos já escolhidos deixam de estar disponíveis para as tiragens seguintes.

Consideremos

$$(3.1) \quad P_i^j = P(Y_j = y_i) ; j = 1, \dots, n \quad , \quad i = 1, \dots, N$$

a probabilidade do i -ésimo elemento da população ser escolhido na j -ésima tiragem, e

$$(3.2) \quad \Pi_i = \sum_{j=1}^n P_i^j ; i = 1, \dots, N .$$

Seja

$$(3.3) \quad A_i = \bigcup_{j=1}^n (Y_j = y_i),$$

o acontecimento que se verifica quando o i -ésimo elemento da população é escolhido. Como os acontecimentos que figuram nesta reunião são incompatíveis dois a dois, tem-se:

$$(3.4) \quad P(A_i) = \sum_{j=1}^n P_i^j = \Pi_i ; i = 1, \dots, N .$$

Seja X_i , $i=1,\dots,N$, uma variável aleatória que toma os valores 1 ou 0 consoante o i -ésimo elemento da população é ou não escolhido. Esta variável tem o esquema:

$$(3.5) \quad X_i = \begin{cases} 0 & 1 \\ 1-\Pi_i & \Pi_i \end{cases}, \quad i=1,\dots,N$$

logo

$$(3.6) \quad \mu(X_i) = \Pi_i \quad ; \quad i=1,\dots,N .$$

Por outro lado tem-se:

$$(3.7) \quad \sum_{i=1}^N X_i = n$$

visto haver sempre n elementos escolhidos. Assim:

$$(3.8) \quad \sum_{i=1}^n \Pi_i = \sum_{i=1}^N \mu(X_i) = \mu\left(\sum_{i=1}^N X_i\right) = n .$$

Como já vimos, existem

$$(3.9) \quad L = \frac{N!}{(N-n)!}$$

amostras possíveis, às quais podemos atribuir índices r , $r=1,\dots,L$.

Seja C_i , $i=1,\dots,N$ o conjunto dos índices das amostras que contêm o i -ésimo elemento. Dada uma função $h(Y)$, dos valores da característica numérica para a qual estamos a amostrar, sendo q_r a probabilidade de se colher a amostra com índice r , teremos:

$$(3.10) \quad \mu\left(\sum_{j=1}^n h(Y_j)\right) = \sum_{r=1}^L q_r \left(\sum_{j=1}^n h(Y_j)\right)_{(r)}$$

sendo $\left(\sum_{j=1}^n h(Y_j)\right)_{(r)}$ a soma dos valores da função relativamente à amostra com índice r . Logo,

agrupando, no segundo membro, os termos em que intervém o i -ésimo elemento da população vem:

$$(3.11) \quad \mu\left(\sum_{j=1}^n h(Y_j)\right) = \sum_{i=1}^N \left(\sum_{r \in C_i} q_r\right) h(y_i) = \sum_{i=1}^N \Pi_i h(y_i)$$

visto que Π_i , a probabilidade de o i -ésimo elemento ser escolhido, é a soma das probabilidades das amostras que contêm esse elemento.

Utilizando estas igualdades, podemos agora obter um estimador linear centrado de μ (o valor médio da característica em estudo, $\mu = \frac{1}{N} \sum_{i=1}^N y_i$).

Seja

$$(3.12) \quad h(Y_j) = \frac{Y_j}{\Pi_j} ; j = 1, \dots, n ,$$

então a expressão do estimador pretendido será:

$$(3.13) \quad \mu^* = \frac{1}{N} \sum_{j=1}^n \frac{Y_j}{\Pi_j}$$

pois como se pode ver

$$(3.14) \quad \mu(\mu^*) = \frac{1}{N} \mu \left(\sum_{j=1}^n \frac{Y_j}{\Pi_j} \right) = \frac{1}{N} \sum_{i=1}^N \Pi_i \frac{y_i}{\Pi_i} = \frac{1}{N} \sum_{i=1}^N y_i = \mu .$$

Por outro lado, sendo $\Pi_{i,i'}$, com $i \neq i'$, a probabilidade de os elementos com índices i e i' serem ambos escolhidos para a amostra, tem-se:

$$(3.15) \quad \Pi_{i,i'} = \sum_{r \in C_i \cap C_{i'}} q_r ; i, i' = 1, \dots, N ; i \neq i' ,$$

ou seja, $\Pi_{i,i'}$ é a soma das probabilidades correspondentes às amostras que contêm ambos os elementos indicados e essas amostras têm os índices em $C_i \cap C_{i'}$.

Dada agora uma função $g(Y, Y')$ de pares de valores da característica, temos

$$(3.16) \quad \mu \left(\sum_{j=1}^n \sum_{j' \neq j}^n g(Y_j; Y_{j'}) \right) = \sum_{r=1}^L q_r \left(\sum_{j=1}^n \sum_{j' \neq j}^n g(Y_j; Y_{j'}) \right)_{(r)}$$

Agrupando os termos do segundo membro em que intervém o mesmo par de elementos da população vem:

$$(3.17) \quad \mu \left(\sum_{j=1}^n \sum_{j' \neq j}^n g(Y_j; Y_{j'}) \right) = \sum_{i=1}^N \sum_{i' \neq i}^N \left(\sum_{r \in C_i \cap C_{i'}} q_r \right) g(y_i; y_{i'}) = \sum_{i=1}^N \sum_{i' \neq i}^N \Pi_{i,i'} g(y_i; y_{i'})$$

Consideremos uma variável aleatória

$$(3.18) \quad U = \sum_{j=1}^n h(Y_j)$$

para a qual se calcula a variância. Então:

$$(3.19) \quad U^2 = \sum_{j=1}^n h^2(Y_j) + \sum_{\substack{j=1 \\ j \neq j'}}^n \sum_{\substack{j'=1 \\ j' \neq j}}^n h(Y_j)h(Y_{j'}).$$

De (3.11) vem

$$(3.20) \quad \mu(U) = \sum_{i=1}^N \Pi_i h(y_i)$$

e considerando também (3.17) vem

$$(3.21) \quad \mu(U^2) = \sum_{i=1}^N \Pi_i h^2(y_i) + \sum_{\substack{i=1 \\ i \neq i'}}^N \sum_{\substack{i'=1 \\ i' \neq i}}^N \Pi_{i,i'} h(y_i)h(y_{i'})$$

e

$$(3.22) \quad \mu^2(U) = \sum_{i=1}^N \Pi_i^2 h^2(y_i) + \sum_{\substack{i=1 \\ i \neq i'}}^N \sum_{\substack{i'=1 \\ i' \neq i}}^N \Pi_i \Pi_{i'} h(y_i)h(y_{i'}) .$$

Temos então

$$(3.23) \quad \begin{aligned} \sigma^2(U) &= \mu(U^2) - \mu^2(U) = \\ &= \sum_{i=1}^N \Pi_i h^2(y_i) + \sum_{\substack{i=1 \\ i \neq i'}}^N \sum_{\substack{i'=1 \\ i' \neq i}}^N \Pi_{i,i'} h(y_i)h(y_{i'}) - \sum_{i=1}^N \Pi_i^2 h^2(y_i) - \sum_{\substack{i=1 \\ i \neq i'}}^N \sum_{\substack{i'=1 \\ i' \neq i}}^N \Pi_i \Pi_{i'} h(y_i)h(y_{i'}) = \\ &= \sum_{i=1}^N \Pi_i (1 - \Pi_i) h^2(y_i) + \sum_{\substack{i=1 \\ i \neq i'}}^N \sum_{\substack{i'=1 \\ i' \neq i}}^N (\Pi_{i,i'} - \Pi_i \Pi_{i'}) h(y_i)h(y_{i'}) \end{aligned}$$

Considerando (3.11) e (3.17), tem-se $\mu\left(\sum_{j=1}^n (1 - \Pi_j) h^2(Y_j)\right) = \sum_{i=1}^N \Pi_i (1 - \Pi_i) h^2(y_i)$, e

$$\mu\left(\sum_{\substack{j=1 \\ j \neq j'}}^n \sum_{\substack{j'=1 \\ j' \neq j}}^n \frac{\Pi_{j,j'} - \Pi_j \Pi_{j'}}{\Pi_{j,j'}} h(Y_j)h(Y_{j'})\right) = \sum_{\substack{i=1 \\ i \neq i'}}^N \sum_{\substack{i'=1 \\ i' \neq i}}^N (\Pi_{i,i'} - \Pi_i \Pi_{i'}) h(y_i)h(y_{i'}). \text{ Então, dado o valor}$$

médio de uma soma ser a soma dos valores médios, vê-se que,

$$(3.24) \quad \sigma^{2*}(U) = \sum_{j=1}^n (1 - \Pi_j) h^2(y_j) + \sum_{j=1}^n \sum_{\substack{j'=1 \\ j' \neq j}}^n \frac{\Pi_{j,j'} - \Pi_j \Pi_{j'}}{\Pi_{j,j'}} h(y_j)h(y_{j'})$$

é estimador centrado de $\sigma^2(U)$.

Usando (3.23) e (3.24) é possível calcular $\sigma^2(\mu^*)$ e obter um estimador centrado para esta variância.

Seja

$$(3.25) \quad h(Y_j) = \frac{Y_j}{\Pi_j} \quad ; \quad j=1, \dots, n$$

atendendo a (3.23) obtém-se

$$(3.26) \quad \sigma^2(\mu^*) = \frac{1}{N^2} \sigma^2 \left(\sum_{j=1}^n \frac{Y_j}{\Pi_j} \right) = \frac{1}{N^2} \left[\sum_{i=1}^N (1 - \Pi_i) \frac{y_i^2}{\Pi_i} + \sum_{i=1}^N \sum_{i'=1}^N (\Pi_{i,i'} - \Pi_i \Pi_{i'}) \frac{y_i}{\Pi_i} \frac{y_{i'}}{\Pi_{i'}} \right]$$

e atendendo a (3.24) obtém-se

$$(3.27) \quad \sigma^{2*}(\mu^*) = \frac{1}{N^2} \sigma^{2*} \left(\sum_{j=1}^n \frac{Y_j}{\Pi_j} \right) = \frac{1}{N^2} \left[\sum_{j=1}^n (1 - \Pi_j) \frac{Y_j^2}{\Pi_j^2} + \sum_{j=1}^n \sum_{j'=1}^n \frac{\Pi_{j,j'} - \Pi_j \Pi_{j'}}{\Pi_{j,j'}} \frac{Y_j}{\Pi_j} \frac{Y_{j'}}{\Pi_{j'}} \right].$$

$j \neq j'$

Para escolhermos as probabilidades que levam a minimizar $\sigma^2(\mu^*)$, temos que obter uma expressão alternativa para (3.23). Para tal comecemos por observar que a variável aleatória $X_{i,i'} = X_i X_{i'}$; $i, i' = 1, \dots, N$ e $i \neq i'$, toma o valor 1 quando ambos os elementos da população com índices i e i' são escolhidos e o valor 0 quando pelo menos um deles não é escolhido. Assim esta variável terá o esquema

$$(3.28) \quad X_{i,i'} \begin{cases} 0 & 1 \\ 1 - \Pi_{i,i'} & \Pi_{i,i'} \end{cases} ; \quad i, i' = 1, \dots, N, \quad i \neq i'$$

vindo

$$(3.29) \quad \mu(X_{i,i'}) = \Pi_{i,i'} ; \quad i, i' = 1, \dots, N, \quad i \neq i'.$$

Observe-se agora que

$$(3.30) \quad \sum_{\substack{i'=1 \\ i' \neq i}}^N X_{i,i'} = \sum_{\substack{i'=1 \\ i' \neq i}}^N X_i X_{i'} = X_i \sum_{\substack{i'=1 \\ i' \neq i}}^N X_{i'} = X_i (n - X_i) = nX_i - X_i^2 = \\ = nX_i - X_i = (n-1)X_i \quad ; \quad i = 1, \dots, N$$

visto que $\sum_{i=1}^N X_i = n$ e que, como X_i toma apenas os valores 0 e 1, $X_i^2 = X_i$. Logo, atendendo à expressão anterior:

$$(3.31) \quad \sum_{\substack{i'=1 \\ i' \neq i}}^N \Pi_{i,i'} = \sum_{\substack{i'=1 \\ i' \neq i}}^N \mu(X_{i,i'}) = \mu \left(\sum_{\substack{i'=1 \\ i' \neq i}}^N X_{i,i'} \right) = \mu[(n-1)X_i] = \\ = (n-1)\mu(X_i) = (n-1)\Pi_i \quad ; \quad i = 1, \dots, N$$

visto que $\mu(X_i) = \Pi_i$.

Como $\sum_{i=1}^N \Pi_i = n$ vem

$$(3.32) \quad \sum_{\substack{i=1 \\ i' \neq i}}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) = \Pi_i \sum_{\substack{i'=1 \\ i' \neq i}}^N \Pi_{i'} - \sum_{\substack{i'=1 \\ i' \neq i}}^N \Pi_{i,i'} = \Pi_i (n - \Pi_i) - (n - 1)\Pi_i = \Pi_i (1 - \Pi_i).$$

Atendendo à expressão anterior, (3.23) passa a ter a forma:

$$(3.33) \quad \sigma^2(U) = \sum_{i=1}^N \left[\sum_{\substack{i'=1 \\ i' \neq i}}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \right] h^2(y_i) - \sum_{i=1}^N \sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) h(y_i) h(y_{i'})$$

e como

$$(3.34) \quad \sum_{\substack{i=1 \\ i' \neq i}}^N \left[\sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \right] h^2(y_i) = \sum_{\substack{i=1 \\ i' \neq i}}^N \left[\sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \right] h^2(y_{i'})$$

vem

$$(3.35) \quad \begin{aligned} \sigma^2(U) &= \frac{1}{2} \sum_{i=1}^N \left[\sum_{\substack{i'=1 \\ i' \neq i}}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \right] h^2(y_i) - \sum_{i=1}^N \sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) h(y_i) h(y_{i'}) + \\ &\quad + \frac{1}{2} \sum_{\substack{i=1 \\ i' \neq i}}^N \left[\sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \right] h^2(y_{i'}) = \\ &= \frac{1}{2} \sum_{\substack{i=1 \\ i' \neq i}}^N \left[\sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \right] \left[h^2(y_i) - 2h(y_i)h(y_{i'}) + h^2(y_{i'}) \right] = \\ &= \frac{1}{2} \sum_{\substack{i=1 \\ i' \neq i}}^N \left[\sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \right] \left[h(y_i) - h(y_{i'}) \right]^2 \end{aligned}$$

Assim, como para μ^* se tem $h(Y_j) = \frac{Y_j}{\Pi_j}$, vem

$$(3.36) \quad \sigma^2(\mu^*) = \frac{1}{2N^2} \sum_{i=1}^N \sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \left[\frac{y_i}{\Pi_i} - \frac{y_{i'}}{\Pi_{i'}} \right]^2,$$

que mostra que, se as probabilidades Π_i forem proporcionais aos y_i , $\sigma^2(\mu^*)$ é nula.

Atendendo a (3.24) e (3.32), obtemos

$$(3.37) \quad \sigma^{2*}(U) = \frac{1}{2} \sum_{j=1}^n \sum_{j'=1}^n \frac{\Pi_j \Pi_{j'} - \Pi_{j,j'}}{\Pi_{j,j'}} \left[h(Y_j) - h(Y_{j'}) \right]^2$$

então

$$(3.38) \quad \sigma^{2*}(\mu^*) = \frac{1}{2N^2} \sum_{j=1}^n \sum_{j'=1}^n \frac{\Pi_j \Pi_{j'} - \Pi_{j,j'}}{\Pi_{j,j'}} \left[\frac{Y_j}{\Pi_j} - \frac{Y_{j'}}{\Pi_{j'}} \right]^2 .$$

3.2. COM REPOSIÇÃO

Como já vimos anteriormente, na amostragem com reposição, os elementos que já foram escolhidos continuam disponíveis para as tiragens seguintes. O número de amostras possíveis é portanto $L = N^n$, pois, para cada uma das tiragens, existem N hipóteses de escolha.

Consideremos novamente as expressões (3.1) e (3.2).

Como, neste caso, há reposição, os acontecimentos $(Y_j = y_i)$, $j = 1, \dots, n$, não são, para um dado i , incompatíveis dois a dois, logo não se podem interpretar os Π_i como probabilidades.

Veremos a seguir que, os Π_i são agora os valores médios das variáveis que dão o número de vezes que cada elemento da população é escolhido.

Seja X_i^j a variável aleatória que toma o valor 1 ou 0 conforme o i -ésimo elemento da população é ou não escolhido, na j -ésima tiragem.

A variável aleatória X_i , que dá o número de vezes que o i -ésimo elemento é escolhido, é então da forma :

$$(3.39) \quad X_i = \sum_{j=1}^n X_i^j .$$

Visto que, em cada tiragem, apenas um elemento é escolhido, temos:

$$(3.40) \quad \sum_{i=1}^N X_i^j = 1 ; j = 1, \dots, n .$$

À variável aleatória X_i^j corresponde então o esquema:

$$(3.41) \quad X_i^j \begin{cases} 0 & 1 \\ 1 - P_i^j & P_i^j \end{cases} ; j = 1, \dots, n , i = 1, \dots, N .$$

Logo

$$(3.42) \quad \mu(X_i^j) = P_i^j ; j = 1, \dots, n , i = 1, \dots, N ,$$

pelo que

$$(3.43) \quad \Pi_i = \sum_{j=1}^n P_i^j = \mu \left(\sum_{j=1}^n X_i^j \right) = \mu(X_i) ; \quad i = 1, \dots, N .$$

Embora, neste caso, a variável aleatória não tome só os valores 0 ou 1 , mas possa tomar valores de 1 até n, continuamos a ter

$$(3.44) \quad \sum_{i=1}^N X_i = n$$

visto que, em cada tiragem, é escolhido apenas um único elemento. Então:

$$(3.45) \quad \sum_{i=1}^N \Pi_i = \sum_{i=1}^N \mu(X_i) = \mu \left(\sum_{i=1}^N X_i \right) = n .$$

Seja $C_{i,s}$, $s = 0, \dots, n$, o conjunto dos índices das amostras em que o i -ésimo elemento da população figura s vezes. A soma das probabilidades de se obterem essas amostras dá-nos a probabilidade do i -ésimo elemento da população ser escolhido s vezes, logo

$$(3.46) \quad P(X_i = s) = \sum_{r \in C_{i,s}} q_r ; \quad s = 0, 1, \dots, n , \quad i = 1, \dots, N .$$

Consideremos ainda uma função $h(Y)$, dos valores da característica numérica para a qual se está a amostrar. Continua a ter-se

$$(3.47) \quad \mu \left(\sum_{j=1}^n h(Y_j) \right) = \sum_{r=1}^L q_r \left(\sum_{j=1}^n h(Y_j) \right)_{(r)} ,$$

mas, se quisermos reagrupar os termos do segundo membro, para cada elemento da população, teremos de atender ao número de vezes que o elemento é escolhido. Como no segundo membro $h(Y_j)$ figura s vezes, nas amostras pertencentes a $C_{i,s}$, vem então:

$$(3.48) \quad \begin{aligned} \mu \left(\sum_{j=1}^n h(Y_j) \right) &= \sum_{i=1}^N \left(\sum_{s=1}^n s \cdot \left(\sum_{r \in C_{i,s}} q_r \right) \right) h(y_i) = \sum_{i=1}^N \left[\sum_{s=1}^n P(X_i = s) \cdot s \right] h(y_i) = \\ &= \sum_{i=1}^N \mu(X_i) h(y_i) = \sum_{i=1}^N \Pi_i h(y_i) \end{aligned}$$

o que mostra que a expressão (3.11) se mantém, só que agora os Π_i são valores médios e não probabilidades.

Raciocinando como para a amostragem sem reposição, verifica-se que (3.13) é, também neste caso, a expressão do estimador centrado para o valor médio μ .

Vamos agora mostrar que a expressão (3.17) também se aplica neste tipo de amostragem. Como estamos a considerar a reposição dos elementos, com $j \neq j'$, os acontecimentos $(Y_j = y_i)$ e $(Y_{j'} = y_{i'})$ são independentes, pelo que as variáveis aleatórias X_i^j e $X_{i'}^{j'}$ também são independentes.

Recordemos que X_i^j toma os valores 1 e 0 consoante o i -ésimo elemento da população é ou não escolhido na j -ésima tiragem, continua pois a verificar-se a expressão (3.42).

Note-se que, como X_i^j e $X_{i'}^{j'}$ são, com $j \neq j'$, independentes se tem

$$(3.49) \quad \mu(X_i^j X_{i'}^{j'}) = \mu(X_i^j) \mu(X_{i'}^{j'}) ; j, j' = 1, \dots, N, \quad i, i' = 1, \dots, N ; j \neq j'$$

enquanto que, como, em cada tiragem, só é escolhido um elemento,

$$(3.50) \quad X_i^j \cdot X_i^{j'} = 0 ; j = 1, \dots, n, \quad i \neq i' .$$

Finalmente, como X_i^j apenas toma os valores 0 e 1, temos

$$(3.51) \quad X_i^{j^2} = X_i^j ; j = 1, \dots, n, \quad i = 1, \dots, N .$$

Consideremos agora

$$(3.52) \quad \Pi_{i, i'} = \sum_{j=1}^n \sum_{j'=1}^n P_i^j P_{i'}^{j'} ; i, i' = 1, \dots, N ,$$

quando $i \neq i'$ tem-se, com $j \neq j'$

$$(3.53) \quad \begin{aligned} \Pi_{i, i'} &= \sum_{j=1}^n \sum_{j'=1}^n \mu(X_i^j) \mu(X_{i'}^{j'}) = \sum_{j=1}^n \sum_{j'=1}^n \mu(X_i^j X_{i'}^{j'}) = \\ &= \mu \left(\sum_{j=1}^n \sum_{j'=1}^n X_i^j X_{i'}^{j'} \right) = \mu \left(\sum_{j=1}^n \sum_{j'=1}^n X_i^j X_{i'}^{j'} \right) = \\ &= \mu \left(\sum_{j=1}^n \sum_{j'=1}^n X_i^j X_{i'}^{j'} \right) = \mu \left[\left(\sum_{j=1}^n X_i^j \right) \left(\sum_{j'=1}^n X_{i'}^{j'} \right) \right] = \mu(X_i X_{i'}) \end{aligned}$$

enquanto que, se $i = i'$, vem, com $j \neq j'$

$$(3.54) \quad \begin{aligned} \Pi_{i, i'} &= \sum_{j=1}^n \sum_{j'=1}^n \mu(X_i^j) \mu(X_{i'}^{j'}) = \sum_{j=1}^n \sum_{j'=1}^n \mu(X_i^j X_{i'}^{j'}) = \\ &= \mu \left(\sum_{j=1}^n \sum_{j'=1}^n X_i^j X_{i'}^{j'} \right) = \mu \left(\sum_{j=1}^n \sum_{j'=1}^n X_i^j X_{i'}^{j'} \right) = \\ &= \mu \left(\sum_{j=1}^n \sum_{j'=1}^n X_i^j X_{i'}^{j'} - \sum_{j=1}^n X_i^{j^2} \right) = \mu \left[\left(\sum_{j=1}^n X_i^j \right) \left(\sum_{j'=1}^n X_{i'}^{j'} \right) - \sum_{j=1}^n X_i^j \right] = \\ &= \mu(X_i^2 - X_{i'}) = \mu[X_i(X_i - 1)] \end{aligned}$$

Dada agora uma função $g(Y, Y')$ de pares de valores da característica numérica para a qual se está a amostrar, temos



$$(3.55) \quad \mu \left(\sum_{j=1}^n \sum_{j'=1}^n g(Y_j, Y_{j'}) \right) = \sum_{r=1}^{L'} q_r \left(\sum_{j=1}^n \sum_{j'=1}^n g(y_j, y_{j'}) \right)_{(r)}.$$

Ao agruparmos, os termos do segundo membro, segundo os pares de elementos da população, temos de considerar, os pares de elementos distintos em separado dos pares de elementos repetidos. Para uma amostra com $r \in C_{i,s} \cap C_{i',s'}$, em que o i -ésimo elemento aparece s vezes e o i' -ésimo elemento aparece s' vezes, o par $(y_i, y_{i'})$ ocorre $s \cdot s'$ vezes ($s + s' \leq n$). Para uma amostra com $r \in C_{i,s}$, o par $(y_i, y_{i'})$ aparece $s(s-1)$ vezes, logo

$$(3.56) \quad \begin{aligned} \mu \left(\sum_{\substack{j=1 \\ j \neq j'}}^n \sum_{j'=1}^n g(Y_j, Y_{j'}) \right) &= \sum_{i=1}^N \sum_{i'=1}^N \left[\sum_{s=1}^{n-1} \sum_{s'=1}^{n-s} \left(\sum_{r \in C_{i,s} \cap C_{i',s'}} q_r \right) s \cdot s' \right] g(y_i, y_{i'}) + \sum_{i=1}^N \left[\sum_{s=1}^N \left(\sum_{r \in C_{i,s}} q_r \right) s(s-1) \right] g(y_i, y_{i'}) = \\ &= \sum_{\substack{i=1 \\ i \neq i'}}^N \sum_{i'=1}^N \mu(X_i X_{i'}) g(y_i, y_{i'}) + \sum_{i=1}^N \mu[X_i(X_i - 1)] g(y_i, y_i) = \\ &= \sum_{\substack{i=1 \\ i \neq i'}}^N \sum_{i'=1}^N \Pi_{i,i'} g(y_i, y_{i'}) + \sum_{i=1}^N \Pi_{i,i} g(y_i, y_i) = \\ &= \sum_{i=1}^N \sum_{i'=1}^N \Pi_{i,i'} g(y_i, y_{i'}) \end{aligned}$$

pelo que a expressão (3.17) continua a aplicar-se, só que agora, por ser com reposição, se pode ter $i = i'$.

Consideremos novamente as expressões (3.18), (3.19), (3.20), (3.21) e (3.22), vê-se que:

$$(3.57) \quad \begin{cases} \mu \left(\sum_{j=1}^n h^2(Y_j) \right) = \sum_{i=1}^N \Pi_i h^2(y_i) \\ \mu \left(\sum_{j=1}^n \sum_{j'=1}^n \frac{\Pi_{j,j'} - \Pi_j \Pi_{j'}}{\Pi_{j,j'}} h(Y_j) h(Y_{j'}) \right) = \sum_{i=1}^N \sum_{i'=1}^N (\Pi_{i,i'} - \Pi_i \Pi_{i'}) h(y_i) h(y_{i'}) \end{cases}$$

o que mostra que

$$(3.58) \quad \sigma^{2*}(U) = \sum_{j=1}^n h^2(Y_j) + \sum_{j=1}^n \sum_{j'=1}^n \frac{\Pi_{j,j'} - \Pi_j \Pi_{j'}}{\Pi_{j,j'}} h(Y_j) h(Y_{j'})$$

é um estimador centrado de $\sigma^2(U)$.

Em particular, para μ^* , com $h(Y_j) = \frac{Y_j}{\Pi_j}$; $j = 1, \dots, n$, as expressões (3.57) e (3.58) dão

$$(3.59) \quad \begin{cases} \sigma^2(\mu^*) = \frac{1}{N^2} \left[\sum_{i=1}^N \frac{y_i^2}{\Pi_i} + \sum_{i=1}^N \sum_{i'=1}^N (\Pi_{i,i'} - \Pi_i \Pi_{i'}) h(y_i) h(y_{i'}) \right] \\ \sigma^{2*}(\mu^*) = \frac{1}{N^2} \left[\sum_{j=1}^n \frac{Y_j^2}{\Pi_j^2} + \sum_{j=1}^n \sum_{j'=1}^n \frac{\Pi_{j,j'} - \Pi_j \Pi_{j'}}{\Pi_{j,j'}} \frac{Y_j}{\Pi_j} \frac{Y_{j'}}{\Pi_{j'}} \right] \end{cases}$$

Tal como foi feito para a amostragem com reposição, procuremos uma expressão alternativa para $\sigma^2(U)$.

Assim:

(3.60)

$$\begin{aligned} \sum_{i'=1}^N \Pi_{i,i'} &= \Pi_{i,i} + \sum_{\substack{i'=1 \\ i' \neq i}}^N \Pi_{i,i'} = \Pi_{i,i} + \sum_{\substack{i'=1 \\ i' \neq i}}^N \mu(X_i X_{i'}) = \\ &= \Pi_{i,i} + \mu \left(\sum_{\substack{i'=1 \\ i' \neq i}}^N X_i X_{i'} \right) = \Pi_{i,i} + \mu \left(X_i \sum_{\substack{i'=1 \\ i' \neq i}}^N X_{i'} \right) = \Pi_{i,i} + \mu [X_i (n - X_i)] = \\ &= \Pi_{i,i} + n\mu(X_i) - \mu(X_i^2) = \Pi_{i,i} + n\Pi_i - (\Pi_{i,i} + \Pi_i) = (n-1)\Pi_i \end{aligned}$$

$$i = 1, \dots, N$$

logo

$$(3.61) \quad \sum_{i=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) = \Pi_i \sum_{i'=1}^N \Pi_{i'} - \sum_{i'=1}^N \Pi_{i,i'} = n\Pi_i - (n-1)\Pi_i = \Pi_i ; \quad i = 1, \dots, N$$

pele que

(3.62)

$$\begin{aligned} \sigma^2(U) &= \frac{1}{2} \sum_{i=1}^N \left[\sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \right] h^2(y_i) - \sum_{i=1}^N \sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) h(y_i) h(y_{i'}) + \\ &+ \frac{1}{2} \sum_{i'=1}^N \left[\sum_{i=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \right] h^2(y_{i'}) = \\ &= \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) [h^2(y_i) - 2h(y_i)h(y_{i'}) + h^2(y_{i'})] = \\ &= \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) [h(y_i) - h(y_{i'})]^2 \end{aligned}$$

assim, para μ^* ter-se-à

$$(3.63) \quad \sigma^2(\mu^*) = \frac{1}{2N^2} \sum_{i=1}^N \sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \left[\frac{y_i}{\Pi_i} - \frac{y_{i'}}{\Pi_{i'}} \right]^2$$

pelo que, se os valores médios Π_i forem proporcionais aos y_i , $\sigma^2(\mu^*)$ seria nula.

Atendendo a (3.58) e a (3.62),obtemos

$$(3.64) \quad \sigma^{2*}(U) = \frac{1}{2} \sum_{j=1}^n \sum_{\substack{j'=1 \\ j \neq j'}}^n \frac{\Pi_j \Pi_{j'} - \Pi_{j,j'}}{\Pi_{j,j'}} \left[h(Y_j) - h(Y_{j'}) \right]^2$$

vindo,

$$(3.65) \quad \sigma^{2*}(\mu^*) = \frac{1}{2N^2} \sum_{j=1}^n \sum_{\substack{j'=1 \\ j \neq j'}}^n \frac{\Pi_j \Pi_{j'} - \Pi_{j,j'}}{\Pi_{j,j'}} \left[\frac{Y_j}{\Pi_j} - \frac{Y_{j'}}{\Pi_{j'}} \right]^2.$$

Considerando agora um caso particular, em que as probabilidades de escolha dos vários elementos da população não variam de tiragem para tiragem (embora possam variar de elemento para elemento), ou seja, $P_i^j = p_i$; $j = 1, \dots, n$, $i = 1, \dots, N$. Temos então:

$$(3.66) \quad \begin{cases} \Pi_i = np_i & ; \quad i = 1, \dots, N \\ \Pi_{i,i'} = n(n-1)p_i p_{i'} & ; \quad i, i' = 1, \dots, N \end{cases}$$

logo

$$(3.67) \quad \Pi_i \Pi_{i'} - \Pi_{i,i'} = np_i p_{i'} & ; \quad i, i' = 1, \dots, N ,$$

tendo em conta (3.62) e (3.64) obtemos

$$(3.68) \quad \sigma^2(U) = \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N np_i p_{i'} \left[h(y_i) - h(y_{i'}) \right]^2$$

e

$$(3.69) \quad \sigma^{2*}(U) = \frac{1}{2} \sum_{\substack{j=1 \\ j \neq j'}}^n \sum_{j'=1}^n \frac{\left[h(Y_j) - h(Y_{j'}) \right]^2}{n-1}.$$

Neste caso tem-se

$$(3.70) \quad \mu^* = \frac{1}{N} \sum_{j=1}^n \frac{Y_j}{np_j} ,$$

$$(3.71) \quad \sigma^2(\mu^*) = \frac{1}{2N^2 n} \sum_{i=1}^N \sum_{i'=1}^N p_i p_{i'} \left(\frac{y_i}{p_i} - \frac{y_{i'}}{p_{i'}} \right)^2$$

e

$$(3.72) \quad \sigma^{2*}(\mu^*) = \frac{1}{2N^2 n^2 (n-1)} \sum_{\substack{j=1 \\ j \neq j'}}^n \sum_{j'=1}^n \left(\frac{Y_j}{p_j} - \frac{Y_{j'}}{p_{j'}} \right)^2$$

Se as probabilidades p_i forem proporcionais nos y_i , $\sigma^2(\mu^*)$ é nula.

4. AMOSTRAGEM ESTRATIFICADA

Já foi visto que, na amostragem aleatória simples, a variância do estimador, da média da população, \bar{Y} , depende da variabilidade da característica Y na população, independentemente da dimensão da amostra.

Se a população é bastante heterogénea, e a característica em estudo difere, consideravelmente, entre certos subconjuntos da população inicial, é aconselhável dividir a população. Essa divisão, chamada estratificação, deve ser feita de forma a que cada um desses subconjuntos (os estratos) seja constituído por indivíduos de alguma forma semelhantes. Os estratos podem ser definidos considerando diversas características. Em populações humanas podemos referir, por exemplo, a idade, o sexo, a classe social ou o grau de instrução.

A estratificação pode produzir um ganho de precisão nos estimadores das características da população inteira. Ao dividir uma população heterogénea em estratos, internamente homogéneos, nos quais existem apenas pequenas variações de indivíduo para indivíduo, é possível obter estimadores precisos, em cada estrato, usando uma pequena amostra de elementos de cada um desses estratos.

Obtidos os estimadores para cada estrato, estes podem então ser combinados num estimador preciso para a população inteira.

O principal objectivo, da teoria da amostragem estratificada, é portanto saber qual a melhor escolha para as dimensões das amostras, retiradas de cada estrato, por forma a obter a melhor precisão.

Consideremos uma divisão de uma população, de dimensão N, em k estratos, com dimensões N_1, \dots, N_k , então:

$$(4.1) \quad N = \sum_{t=1}^k N_t \quad .$$

Se forem colhidas, dos vários estratos, sub-amostras independentes com dimensões n_1, \dots, n_k , a amostra total terá dimensão

$$(4.2) \quad n = \sum_{t=1}^k n_t \quad .$$

Admitindo que as sub-amostras são colhidas sem reposição e com probabilidades iguais, a probabilidade de um elemento, do estrato t, ser escolhido será dada por:

$$(4.3) \quad \Pi_t = \frac{n_t}{N_t} \quad ; \quad t = 1, \dots, N_t$$

visto que, em cada uma das n_t tiragens, realizadas nesse estrato, cada um dos seus elementos tem probabilidade $\frac{1}{N_t}$ de ser escolhido.

Sejam $Y_{t,j}$; $j=1, \dots, n_t$, as variáveis aleatórias que dão os resultados obtidos para o estrato t , atendendo à expressão (3.13), obtemos

$$(4.4) \quad \mu^* = \frac{1}{N} \sum_{t=1}^k \sum_{j=1}^{n_t} \frac{Y_{t,j}}{\Pi_t} = \frac{1}{N} \sum_{t=1}^k \frac{N_t}{n_t} \sum_{j=1}^{n_t} Y_{t,j} = \frac{1}{N} \sum_{t=1}^k N_t \bar{Y}_t$$

onde

$$(4.5) \quad \bar{Y}_t = \frac{1}{n_t} \sum_{j=1}^{n_t} Y_{t,j} ; t=1, \dots, k$$

é a média da sub-amostra colhida do estrato t .

Devido à grande homogeneidade de cada estrato, podemos considerá-lo como uma população. Assim, tendo em conta as devidas adaptações, podemos escrever a variância do estrato t utilizando (2.40). Temos então:

$$(4.6) \quad \sigma^2(\bar{Y}_t) = \frac{\sigma_t^2}{n_t} \frac{N_t - n_t}{N_t - 1} ; t=1, \dots, k$$

Visto que as sub-amostras são colhidas independentemente, em cada estrato, as médias, $\bar{Y}_1, \dots, \bar{Y}_k$, são independentes, pelo que

$$(4.7) \quad \sigma^2(\mu^*) = \sigma^2\left(\frac{1}{N} \sum_{t=1}^k N_t \bar{Y}_t\right) = \frac{1}{N^2} \sum_{t=1}^k N_t^2 \sigma^2(\bar{Y}_t) = \frac{1}{N^2} \sum_{t=1}^k N_t^2 \frac{\sigma_t^2}{n_t} \frac{N_t - n_t}{N_t - 1}$$

Para N_t grande, $\frac{1}{N_t - 1} \cong \frac{1}{N_t}$, pelo que podemos considerar a aproximação, a $\sigma^2(\mu^*)$, dada por:

$$(4.8) \quad \begin{aligned} \sigma^2(\mu^*) &= \frac{1}{N^2} \sum_{t=1}^k N_t^2 \frac{\sigma_t^2}{n_t} \left(1 - \frac{n_t}{N_t}\right) = \frac{1}{N^2} \sum_{t=1}^k \left(N_t^2 \frac{\sigma_t^2}{n_t} - N_t \sigma_t^2\right) = \\ &= -\frac{1}{N^2} \sum_{t=1}^k N_t \sigma_t^2 + \sum_{t=1}^k \frac{N_t^2 \sigma_t^2}{N^2} \frac{1}{n_t} \end{aligned}$$

Fazendo

$$(4.9) \quad \begin{cases} V_0 = -\frac{1}{N^2} \sum_{t=1}^k N_t \sigma_t^2 \\ V_t = \frac{N_t^2}{N^2} \sigma_t^2 ; t=1, \dots, k \\ W_t = n_t ; t=1, \dots, k \end{cases}$$

vem

$$(4.10) \quad \sigma^2(\mu^*) = V_0 + \sum_{t=1}^k \frac{V_t}{W_t} .$$

Note-se que V_0, V_1, \dots, V_k apenas dependem dos parâmetros dos estratos (N_t e σ_t^2), enquanto que os W_1, \dots, W_k coincidem com as dimensões das sub-amostras.

Por outro lado, quando se realiza uma amostragem estratificada, para além de um custo fixo C_0 , podemos ter custos C_1, \dots, C_k para cada estrato. Assim o custo total será

$$(4.11) \quad C = C_0 + \sum_{t=1}^k C_t n_t = \sum_{t=0}^k C_t n_t \quad , \quad \text{com } n_0 = 1$$

Pela desigualdade de Cauchy, podemos escrever

$$(4.12) \quad (\sigma^2(\mu^*) - V_0)(C - C_0) = \left(\sum_{t=1}^k \frac{V_t}{W_t} \right) \left(\sum_{t=1}^k C_t W_t \right) \geq \left(\sqrt{\sum_{t=1}^k C_t V_t} \right)^2 \geq \sum_{t=1}^k C_t V_t$$

Visto que, a igualdade em (4.12) é obtida apenas quando se verifica a condição

$$(4.13) \quad \frac{\left(\frac{V_t}{W_t} \right)}{W_t C_t} = k^2 \quad (\text{constante para todos os valores de } t),$$

e que segundo membro de (4.12), é independente de W_t , então a escolha de W_t para satisfazer esta condição, que rescrevemos

$$(4.14) \quad W_t^2 = \frac{V_t}{k^2 C_t} \quad , \quad \forall t \quad ,$$

minimizará $\sigma^2(\mu^*) \cdot C$.

Considerando (4.9) e (4.14), vem:

$$(4.15) \quad n_t = W_t = \frac{1}{k} \sqrt{\frac{V_t}{C_t}} = \frac{N_t \sigma_t}{kN \sqrt{C_t}} \quad ; \quad t = 1, \dots, k$$

e portanto

$$(4.16) \quad n = \sum_{t=1}^k n_t = \frac{1}{kN} \sum_{t=1}^k \frac{N_t \sigma_t}{\sqrt{C_t}}$$

o que dá

$$(4.17) \quad \frac{n_t}{n} = \frac{\frac{N_t \sigma_t}{\sqrt{C_t}}}{\sum_{t'=1}^k \frac{N_{t'} \sigma_{t'}}{\sqrt{C_{t'}}}} ; t=1, \dots, k$$

ou seja

$$(4.18) \quad n_t = n \frac{\frac{N_t \sigma_t}{\sqrt{C_t}}}{\sum_{t'=1}^k \frac{N_{t'} \sigma_{t'}}{\sqrt{C_{t'}}}} ; t=1, \dots, k .$$

Note-se que $\sigma^2(\mu^*)$ e C foram expressas como somas de duas componentes, uma fixa V_0 e C_0 e outra dada pelas funções $\sum_{t=1}^k \frac{V_t}{n_t}$ e $\sum_{t=1}^k C_t n_t$.

Fixando o custo total C, para minimizar $(\sigma^2(\mu^*) - V_0)(C - C_0)$ basta minimizar $\sigma^2(\mu^*)$. Se fixarmos $\sigma^2(\mu^*)$ basta minimizar C.

Se os custos C_1, \dots, C_k , correspondentes aos diferentes estratos, forem desprezados ou considerando iguais, $C_t = C$, obtêm-se expressões simplificadas para as dimensões dos estratos, n_t , da forma:

$$(4.19) \quad n_t = n \frac{N_t \sigma_t}{\sum_{t'=1}^k N_{t'} \sigma_{t'}} ; t=1, \dots, k .$$

As expressões obtidas aqui podem também aplicar-se à amostragem estratificada por atributos.

Consideremos p_t a probabilidade de os elementos de um qualquer estrato t possuírem o atributo que se está a estudar, tem-se então:

$$(4.20) \quad \sigma_t = \sqrt{p_t(1-p_t)}$$

vindo, quando são considerados os custos C_t correspondentes aos diferentes estratos,

$$(4.21) \quad n_t = n \frac{\frac{N_t \sqrt{p_t(1-p_t)}}{\sqrt{C_t}}}{\sum_{t'=1}^k \frac{N_{t'} \sqrt{p_{t'}(1-p_{t'})}}{\sqrt{C_{t'}}}} ; t=1, \dots, k$$

ou então, ignorando esses custos,

$$(4.22) \quad n_t = n \frac{N_t \sqrt{p_t(1-p_t)}}{\sum_{t'=1}^k N_{t'} \sqrt{p_{t'}(1-p_{t'})}} ; t = 1, \dots, k .$$

Comparemos agora as variâncias, dos estimadores de μ , obtidas usando a amostragem estratificada $\sigma_E^2(\mu^*)$ e a amostragem, sem reposição, com probabilidades iguais $\sigma_s^2(\mu^*)$, utilizando em ambos os casos amostras da mesma dimensão, n .

Sejam $y_{t,i}$ os valores da característica numérica, para a qual se está a amostrar, nos elementos do estrato t , os valores médios dos estratos serão:

$$(4.23) \quad \mu_t = \frac{1}{N_t} \sum_{i=1}^{N_t} y_{t,i} ; t = 1, \dots, k$$

e o valor médio da população será dado por:

$$(4.24) \quad \mu = \frac{1}{N} \sum_{t=1}^k \sum_{i=1}^{N_t} y_{t,i} .$$

As variâncias dos estratos e da população podem agora ser redefinidas como

$$(4.25) \quad \begin{cases} \sigma_t^2 = \frac{1}{N_t - 1} \sum_{i=1}^{N_t} (y_{t,i} - \mu_t)^2 ; t = 1, \dots, k \\ \sigma^2 = \frac{1}{N - 1} \sum_{t=1}^k \sum_{i=1}^{N_t} (y_{t,i} - \mu)^2 \end{cases}$$

Obtém-se assim, considerando (2.40) e

$$(4.26) \quad \begin{aligned} (N - 1)\sigma^2 &= \sum_{t=1}^k \sum_{i=1}^{N_t} (y_{t,i} - \mu)^2 = \sum_{t=1}^k \sum_{i=1}^{N_t} [(y_{t,i} - \mu_t) + (\mu_t - \mu)]^2 = \\ &= \sum_{t=1}^k \sum_{i=1}^{N_t} [(y_{t,i} - \mu_t)^2 + 2(\mu_t - \mu)(y_{t,i} - \mu_t) + (\mu_t - \mu)^2] = \\ &= \sum_{t=1}^k \sum_{i=1}^{N_t} (y_{t,i} - \mu_t)^2 + 2 \sum_{t=1}^k \left[(\mu_t - \mu) \sum_{i=1}^{N_t} (y_{t,i} - \mu_t) \right] + \sum_{t=1}^k N_t (\mu_t - \mu)^2 = \\ &= \sum_{t=1}^k (N_t - 1)\sigma_t^2 + \sum_{t=1}^k N_t (\mu_t - \mu)^2 \end{aligned}$$

com $\sum_{i=1}^{N_t} (y_{t,i} - \mu_t) = 0 ; t = 1, \dots, k$ (visto que a soma dos desvios para a média é nula),

$$(4.27) \quad \sigma_s^2(\mu^*) = \frac{\sigma^2}{n} \frac{N-n}{N} = \frac{N-n}{nN(N-1)} \left[\sum_{t=1}^k (N_t - 1) \sigma_t^2 + \sum_{t=1}^k N_t (\mu_t - \mu)^2 \right] .$$

Consideremos agora que a amostragem estratificada é realizada com critério proporcional, isto é, que

$$(4.28) \quad \frac{n_t}{N_t} = \frac{n}{N} ; \quad t = 1, \dots, k .$$

Atendendo a (4.8) vem então:

$$(4.29) \quad \sigma_E^2(\mu^*) = \frac{1}{N^2} \sum_{t=1}^k N_t^2 \frac{\sigma_t^2}{n_t} \left(1 - \frac{n_t}{N_t} \right) = \frac{1}{N^2} \frac{N}{n} \left(1 - \frac{n}{N} \right) \sum_{t=1}^k N_t \sigma_t^2 = \frac{N-n}{nN^2} \sum_{t=1}^k N_t \sigma_t^2 ,$$

logo

$$(4.30) \quad \begin{aligned} \sigma_s^2(\mu^*) - \sigma_E^2(\mu^*) &= \frac{N-n}{nN} \left\{ \frac{1}{N-1} \left[\sum_{t=1}^k (N_t - 1) \sigma_t^2 + \sum_{t=1}^k N_t (\mu_t - \mu)^2 \right] - \frac{1}{N} \sum_{t=1}^k N_t \sigma_t^2 \right\} = \\ &= \frac{N-n}{nN} \left\{ \left[\frac{1}{N-1} \sum_{t=1}^k (N_t - 1) \sigma_t^2 - \frac{1}{N} \sum_{t=1}^k N_t \sigma_t^2 \right] + \frac{1}{N-1} \sum_{t=1}^k N_t (\mu_t - \mu)^2 \right\} \end{aligned}$$

onde

$$(4.31) \quad \frac{1}{N-1} \sum_{t=1}^k (N_t - 1) \sigma_t^2 - \frac{1}{N} \sum_{t=1}^k N_t \sigma_t^2 = \frac{1}{N(N-1)} \sum_{t=1}^k (N_t - N) \sigma_t^2 < 0 .$$

Considerando, o caso particular em que $\mu_1 = \dots = \mu_k$, isto é, $\mu_t = \mu$, tem-se

$$(4.32) \quad \frac{1}{N-1} \sum_{t=1}^k N_t (\mu_t - \mu)^2 = 0 ,$$

pelo que,

$$(4.33) \quad \sigma_s^2(\mu^*) < \sigma_E^2(\mu^*) .$$

Vê-se assim que, para que a amostragem estratificada conduza a bons resultados, os valores médios, dos estratos, devem diferir bastante entre si, isto é, a variância entre os estratos, $\frac{1}{N-1} \sum_{t=1}^k N_t (\mu_t - \mu)^2$, deve ser grande, enquanto que cada estrato deve ser homogéneo, isto é, a variância, dentro de cada estrato, deve ser pequena.

5. AMOSTRAGEM PARA AGREGADOS

De modo a não alongar, desnecessariamente, a discussão deste tipo de amostragem, consideramos apenas o seu caso mais elementar, que, apesar da sua simplicidade, ilustra, perfeitamente, a teoria da amostragem para agregados.

Consideremos agora uma divisão da população em agregados disjuntos dois a dois e escrevamos $i\theta i'$ quando “ i ” e “ i' ” são índices de elementos que pertencem ao mesmo agregado.

Sejam Π_i a probabilidade de o elemento com índice i ser escolhido e $\Pi_{i,i'}$ a probabilidade de o par de elementos i e i' ser escolhido.

Se $\Pi_i = \Pi_{i'} = \Pi_{i,i'}$, com $i\theta i'$, vê-se que a probabilidade de o elemento i' ser escolhido, quando o elemento i é escolhido, é dada por

$$(5.1) \quad \frac{\Pi_{i,i'}}{\Pi_i} = 1$$

logo quando um elemento de um agregado é escolhido os restantes elementos do agregado são escolhidos (com probabilidade 1). Assim, para se ter $\Pi_i = \Pi_{i'} = \Pi_{i,i'}$, com $i\theta i'$, a amostra tem de ser constituída por agregados completos, isto é, que abranjam toda a população.

Inversamente, se a amostra for constituída por agregados completos, dado $i\theta i'$, as probabilidades de os elementos i e i' serem escolhidos, isoladamente ou em conjunto, são iguais à probabilidade de ser escolhido o agregado a que pertencem.

Assim, com $i\theta i'$, tem-se $\Pi_i = \Pi_{i'} = \Pi_{i,i'}$ se e só se a amostra for constituída por agregados completos.

Admitamos que se constituem N_1 agregados com N_2 elementos cada e que a amostra é formada por n_1 agregados, tem-se então

$$(5.2) \quad \begin{cases} N = N_1 N_2 \\ n = n_1 N_2 \end{cases}$$

logo

$$(5.3) \quad \frac{n_1}{N_1} = \frac{n}{N}$$

A escolha da amostra reduz-se à escolha dos agregados que a vão compor. Trata-se pois de escolher uma amostra de n_1 agregados dentro da população dos agregados.

Suponhamos que esta amostragem é realizada com probabilidades iguais e sem reposição. A probabilidade de se escolher qualquer agregado será $\frac{n_1}{N_1}$. Como a probabilidade de se escolher um elemento é igual à probabilidade de se escolher o agregado a que pertence, vem

$$(5.4) \quad \Pi_i = \frac{n_1}{N_1} = \frac{n}{N} \quad ; \quad i = 1, \dots, N$$

vindo igualmente

$$(5.5) \quad i\theta_{i'} \quad ; \quad \Pi_{i'} = \frac{n_1}{N_1} = \frac{n}{N}$$

Se $i\theta_{i'}$, isto é, se os elementos i e i' pertencem a agregados diferentes, a probabilidade de serem ambos escolhidos é igual à probabilidade de os dois agregados a que pertencem serem escolhidos. Admitindo que os agregados são escolhidos pela técnica da amostragem aleatória simples, com probabilidades iguais e sem reposição, havendo N_1 agregados e destes sendo escolhidos n_1 , haverá, adaptando (3.9),

$$(5.6) \quad L = \frac{N_1!}{(N_1 - n_1)!}$$

amostras possíveis formadas por n_1 agregados, todos com probabilidades iguais.

Dados dois agregados, podemos agrupar as amostras, que os contêm, segundo as posições dos dois agregados nessa amostra. Haverá $n(n-1)$ destas classes de amostras cada uma com

$$(5.7) \quad L' = \frac{(N_1 - 2)!}{(N_1 - n_1)!}$$

amostras, visto que as amostras de cada classe correspondem às amostras de dimensão n_1-2 que se podem extrair da população dos agregados, uma vez excluídos os dois agregados dados. Assim haverá $n_1(n_1-1)L'$ amostras contendo os dois agregados. Como todas as amostras têm probabilidades iguais teremos

$$(5.8) \quad i\theta_{i'} \quad ; \quad \Pi_{i'} = \frac{n_1(n_1-1) \frac{(N_1-2)!}{(N_1-n_1)!}}{\frac{N_1!}{(N_1-n_1)!}} = \frac{n_1(n_1-1)}{N_1(N_1-1)}$$

O estimador de μ é então dado por

$$(5.9) \quad \mu^* = \frac{1}{N} \sum_{j=1}^n \frac{Y_j}{\Pi_j} = \frac{1}{N} \sum_{j=1}^n \frac{N}{n} Y_j = \bar{Y}$$

onde \bar{Y} é a média da amostra.

Por outro lado, atendendo a (3.36), com $\Pi_i = \Pi_{i'} = \frac{n}{N}$, temos

$$(5.10) \quad \sigma^2(\mu^*) = \frac{1}{2N^2} \sum_{i=1}^N \sum_{i'=1}^N (\Pi_i \Pi_{i'} - \Pi_{i,i'}) \frac{N^2}{n^2} (y_i - y_{i'})^2.$$

Como

$$(5.11) \quad \begin{cases} i\theta i' ; \Pi_i \Pi_{i'} - \Pi_{i,i'} = \frac{n^2}{N^2} - \frac{n}{N} = \frac{n}{N} \left(\frac{n}{N} - 1 \right) = \frac{n_1}{N_1} \left(\frac{n_1}{N_1} - 1 \right) \\ i\theta i ; \Pi_i \Pi_{i'} - \Pi_{i,i'} = \frac{n^2}{N^2} - \frac{n_1}{N_1} \left(\frac{n_1 - 1}{N_1 - 1} \right) = \frac{n_1^2}{N_1^2} - \frac{n_1}{N_1} = \frac{n_1}{N_1} \left(\frac{N_1 - n_1}{N_1(N_1 - 1)} \right) \end{cases}$$

se pusermos

$$(5.12) \quad \begin{cases} s_1 = \sum \sum_{i\theta i'} (y_i - y_{i'})^2 \\ s_2 = \sum \sum_{i\theta i} (y_i - y_{i'})^2 \end{cases}$$

onde os dois somatórios são feitos para os pares (i, i') que satisfazem ou não a relação θ , vemos que

$$(5.13) \quad s = \sum_{\substack{i=1 \\ i \neq i'}}^N \sum_{i'=1}^N (y_i - y_{i'})^2 = s_1 + s_2$$

obtendo-se ainda, devido a (5.11):

$$(5.14) \quad \begin{aligned} \sigma^2(\mu^*) &= \frac{1}{2N^2} \frac{N^2}{n^2} \left[\frac{n_1}{N_1} \left(\frac{n_1}{N_1} - 1 \right) s_1 + \frac{n_1}{N_1} \frac{N_1 - n_1}{N_1(N_1 - 1)} s_2 \right] = \\ &= \frac{1}{2n^2} \left[\frac{n_1}{N_1} \left(\frac{n_1}{N_1} - 1 \right) s_1 - \frac{n_1}{N_1} \frac{N_1 - n_1}{N_1(N_1 - 1)} s_2 \right] = \\ &= \frac{1}{2n^2} \frac{n}{N} \left[\left(\frac{n_1}{N_1} - 1 \right) s_1 - \frac{N_1 - n_1}{N_1(N_1 - 1)} (s - s_1) \right] = \\ &= \frac{1}{2nN} \left[\frac{N_1 - n_1}{N_1(N_1 - 1)} s - \frac{1}{N_1} \left(\frac{N_1 - n_1}{N_1 - 1} + N_1 - n_1 \right) s_1 \right] = \\ &= \frac{1}{2nN} \left[\frac{N_1 - n_1}{N_1(N_1 - 1)} s - \frac{(N_1 - n_1)N_1^2}{N_1^2(N_1 - 1)} s_1 \right] = \\ &= \frac{1}{2nN} \frac{N_1 - n_1}{N_1(N_1 - 1)} (s - N_1 s_1) \end{aligned}$$

o que mostra que $\sigma^2(\mu^*)$ é tanto menor quanto maior for s_1 . Observe-se que s depende apenas da população, enquanto que o valor de s_1 resulta da maneira como os agregados são constituídos. Vê-se assim que os agregados devem ser tão heterogêneos quanto possível, por forma a maximizar s_1 e, conseqüentemente, a minimizar $\sigma^2(\mu^*)$.

Voltando à expressão (5.10), como estamos a realizar uma amostragem sem reposição, vemos que quando $n_1 > 1$

$$(5.15) \quad \sigma^{2*}(\mu^*) = \frac{1}{2n^2} \sum_{j=1}^n \sum_{j'=1}^n \frac{\Pi_j \Pi_{j'} - \Pi_{j,j'}}{\Pi_{j,j'}} (Y_j - Y_{j'})^2$$

é o estimador centrado de $\sigma^2(\mu^*)$. Quando $n=1$ vê-se que, devido a (5.11), este estimador toma valores negativos pelo que não se deve utilizar. Em certos casos dispõe-se de um modelo matemático para a população, o que nos permite estimar $\sigma^2(\mu^*)$ mesmo tomando $n_1=1$.

Comparando a amostragem estratificada com a amostragem para agregados, vemos que a principal diferença entre elas é que, na primeira todos os estratos são amostrados, enquanto que na segunda os próprios agregados são sujeitos a um procedimento de selecção. É este facto que conduz os princípios, para os dois métodos, em direcções opostas:

- Na amostragem estratificada, a variabilidade amostral está confinada ao interior dos estratos.
- Na amostragem para agregados., existe apenas variabilidade entre os agregados, uma vez que, todos eles, são inteiramente amostrados.

Um caso particular da amostragem para agregados é a amostragem sistemática. Este tipo de amostragem é recomendado e adoptado, na prática, devido à sua atractiva simplicidade e utilidade, particularmente no estudo de populações instáveis, onde a amostragem tradicional não é aplicável.

Consideremos uma população, a qual é disposta (fisicamente ou por meio de uma lista) numa sequência simples de $N=N_1 \cdot N_2$ indivíduos. O agregado com índice h , $h=1, \dots, N_1$, é constituído pelos elementos da população com índices $i = h + (m-1)N_1$; $m = 1, \dots, N_2$ (Quadro 1).

AGREGADO	COMPOSIÇÃO DO AGREGADO
1	1 , N_1+1 , $2N_1+1$, ... , $(N_2-1)N_1+1$
.	
.	
.	
h	h , $h+N_1$, $h+2N_1$, ... , $h+(N_2-1)N_1$
.	
.	
.	
N_1	N_1 , $2N_1$, $3N_1$, ... , N_2N_1

Quadro 1

Neste tipo de amostragem, toma-se, em geral, um único agregado para constituir a amostra. A probabilidade de seleccionar um agregado é $\frac{1}{N_1}$, ou seja, é a probabilidade com que cada membro do agregado é seleccionado para a amostra.

Para se escolher o agregado que irá compor a amostra, de entre os primeiros N_1 elementos da população é seleccionado um, aleatoriamente, com probabilidades de selecção iguais. (Nesta escolha podem-se utilizar tabelas de números aleatórios).

Consideremos que foi o p -ésimo elemento o seleccionado, tomando-o como ponto de partida, os restantes elementos, que vão compor amostra, são obtidos por progressão aritmética de razão N_1 . Assim, a amostra consiste no agregado de índice p , ou seja, nos elementos nas posições $p, p + N_1, p + 2 N_1, \dots, p + (N_2-1) N_1$.

A conveniência deste tipo de amostragem recai no facto de a selecção do primeiro membro da amostra determinar automaticamente a amostra inteira.

Prova-se que se a população for autocorrelacionada, isto é, se os elementos próximos tiverem tendência a ser semelhantes, a amostragem sistemática conduz a bons resultados.

Nas diversas técnicas de amostragem estudadas até aqui, considerámos sempre que, os elementos, que vão compor a amostra, devem ser retirados directamente da população. Factores económicos e outros, baseados na dificuldade ou impossibilidade de o estudo abranger, de igual forma, toda a população, desencorajam a aplicação destas técnicas. Estes factores conduzem à necessidade de seleccionar grupos de indivíduos em vez de indivíduos directamente da população. mas isso apenas é válido em alguns tipos de pesquisa.

Suponhamos que cada grupo da população, aos quais chamaremos unidades primárias, pode ser dividido num determinado número de sub-unidades. Se ao seleccionar uma amostra de n unidades primárias, as sub-unidades, dentro de cada unidades seleccionada, dão resultados semelhantes, parece ser pouco económico medi-las todas.

Uma prática comum é seleccionar e medir uma amostra de sub-unidades de uma qualquer unidades escolhida, ou seja, o conjunto de elementos a ser estudado é obtido por sub-amostragem.

A sub-amostragem tem uma grande variedade de aplicações, que vão bem para além do objectivo imediato dos estudos por amostragem. Sempre que qualquer processo envolva testes químicos, físicos ou biológicos, ele pode ser executado numa pequena quantidade de material, é portanto indiferente que os resultados obtidos sejam provenientes de uma sub-amostra de uma grande quantidade ou de uma amostra.

Quanto à escolha das unidades que vão ser observadas, tanto pode ser feita com probabilidades iguais como com probabilidades proporcionais à sua dimensão.

Várias regras se podem considerar para determinar as fracções de amostragem e sub-amostragem, e vários métodos de estimação estão disponíveis. As vantagens dos diferentes métodos dependem da natureza da população, dos custos e dos dados suplementares que estão à nossa disposição.

Apesar de o tipo de sub-amostragem mais usual ser para três níveis, este processo pode se estendido sucessivamente a mais um nível tanto quanto for necessário, se as sub-unidades anteriormente obtidas em vez de serem enumeradas completamente, forem elas também amostradas.

Tendo em conta que a teoria de sub-amostragem com probabilidades proporcionais existente, praticamente só se encontra desenvolvida para um número reduzido de níveis de fraccionamento, pretende-se obter uma teoria geral deste tipo de sub-amostragem, válida para qualquer número de níveis, partindo da generalização da sub-amostragem com probabilidades iguais.

6.1. RESULTADOS PRÉVIOS

Seja X uma variável aleatória cuja distribuição depende de uma condição C . Esta variável aleatória terá o valor médio condicional $\mu(X|C)$ e a variância condicional

$$(6.1) \quad \sigma^2(X|C) = \mu(X^2|C) - \mu^2(X|C).$$

O valor médio, não condicional de X , obtém-se descondicionando $\mu(X|C)$, vindo

$$(6.2) \quad \mu(X) = \mu_c(\mu(X|C)),$$

obtendo-se, analogamente, $\mu(X^2) = \mu_c(\mu(X^2|C))$, pelo que

$$(6.3) \quad \begin{aligned} \sigma^2(X) &= \mu(X^2) - \mu^2(X) = \\ &= \mu_c(\mu(X^2|C)) - \mu_c^2(\mu(X|C)) = \\ &= \mu_c(\mu(X^2|C)) - \mu_c(\mu^2(X|C)) + \mu_c(\mu^2(X|C)) - \mu_c^2(\mu(X|C)) = \\ &= \mu_c(\mu(X^2|C) - \mu^2(X|C)) + \mu_c(\mu^2(X|C)) - \mu_c^2(\mu(X|C)) = \\ &= \mu_c(\sigma^2(X|C)) + \mu_c(\mu^2(X|C)) - \mu_c^2(\mu(X|C)) \end{aligned}$$

fazendo $Z = \mu(X|C)$ temos

$$\begin{aligned} \mu_c(\sigma^2(X|C)) + \mu_c(\mu^2(X|C)) - \mu_c^2(\mu(X|C)) &= \mu_c(\sigma^2(X|C)) + \mu_c(Z^2) - \mu_c^2(Z) = \\ &= \mu_c(\sigma^2(X|C)) + \sigma_c^2(Z) \end{aligned}$$

e portanto

$$(6.4) \quad \sigma^2(X) = \mu_c(\sigma^2(X|C)) + \sigma_c^2(\mu(X|C))$$

6.2. NOTAÇÕES

De modo a ser possível a generalização da teoria da sub-amostragem para qualquer número de níveis, é necessário criar uma terminologia de índices que permita referenciar, sem ambiguidades, a posição e génese de uma dada sub-população, qualquer que ela seja, isto é, qual o nível de estrutura a que essa sub-população pertence, dentro desse nível qual a posição que ocupa e quais as sub-populações, que por decomposições sucessivas, lhe deram origem.

Vamos considerar neste estudo que a sub-amostragem é efectuada em k níveis de fraccionamento.

Por forma a clarificar, à partida, como a população inicial é decomposta, sucessivamente, até ao nível k , recorreu-se à elaboração de um esquema que está representado na figura 2.

A identificação das sub-populações, dos vários níveis, a partir do primeiro é feita através da atribuição de vectores. Consideremos as sub-populações de nível h , com $h=1, \dots, k$. Podemos identificá-las por vectores \vec{i}^h com h componentes, ou seja:

$$(6.5) \quad \vec{i}^h = (i_1, \dots, i_h) ; \quad h = 1, \dots, k ,$$

indicando as componentes de \vec{i}^h as sucessivas opções que levam à escolha da população correspondente.

Consideremos ainda o índice p , $1 \leq p \leq h$, e representemos por \vec{i}_p^h o vector obtido tomando as p primeiras componentes do vector \vec{i}^h , temos então:

$$(6.6) \quad \vec{i}_p^h = (\vec{i}_p^h | i_{p+1}, \dots, i_h) ; \quad 1 \leq p \leq h .$$

No caso em que temos $p = h-1$ vem

$$(6.7) \quad \vec{i}_p^h = (\vec{i}^{h-1} | i_h) ; \quad 1 \leq p \leq h .$$

(Na situação particular de $p=0$ toma-se $\vec{i}_0^h = \vec{i}^0$).

Atendendo às expressões (6.5) e (6.7) as sub-populações de nível 1 serão identificadas pelo vector $\vec{i}^1 = (i_1)$, as sub-populações de nível 2 pelo vector $\vec{i}^2 = (i_1, i_2) = (\vec{i}^1 | i_2)$, as de nível 3 pelo vector $\vec{i}^3 = (i_1, i_2, i_3) = (\vec{i}^2 | i_3)$ e assim sucessivamente até às populações de nível k que serão identificadas pelo vector $\vec{i}^k = (\vec{i}^{k-1} | i_k)$.

Deste modo, $P(\vec{i}^1)$ representará qualquer sub-população de nível 1, $P(\vec{i}^2)$ representará qualquer sub-população de nível 2, e assim sucessivamente até $P(\vec{i}^{k-1})$, que representará por sua vez qualquer sub-população de nível $k-1$, que finalmente, por decomposição, origina, directamente, os elementos. Por extensão de linguagem, estes elementos serão considerados como uma sub-população, de nível k , $P(\vec{i}^k)$, assim como a própria população inicial é considerada como uma sub-população de nível zero, isto é, $P(\vec{i}^0)$.

Facilmente se verifica que, dada uma sub-população de nível h , ela está contida noutras sub-populações existentes em níveis anteriores a h , ou seja, com $p < h$, $P(\vec{i}^h)$ é sub-população de $P(\vec{i}_p^h)$.

Seja $X (X_h, h = 0, 1, \dots, k)$ o conjunto dos vectores representativos das sub-populações (de nível $h, h=0, 1, \dots, k$). Tem-se $X_0 = \{\vec{i}^0\}$ e o seu único vector corresponde à população.

Dado $\vec{i}^h (h = 0, 1, \dots, k-1)$, $\Omega(\vec{i}^h)$ representa o conjunto dos vectores \vec{i}^{h+1} cujas primeiras h componentes são as componentes de \vec{i}^h . Tais vectores correspondem às sub-populações em que se decompõe uma determinada população genérica $P(\vec{i}^h)$. Relembrando que

\vec{i}_p^h representa os vectores obtidos tomando as p primeiras componentes de \vec{i}^h , no caso de $p=h-1$ tem-se $\vec{i}^h \in \Omega(\vec{i}_{h-1}^h)$.

Consideremos agora

$$(6.8) \quad \begin{cases} N(\vec{i}^h) = \# [P(\vec{i}^h)] ; h = 0, 1, \dots, k \\ N'(\vec{i}^h) = \# [\Omega(\vec{i}^h)] ; h = 0, 1, \dots, k-1 \end{cases}$$

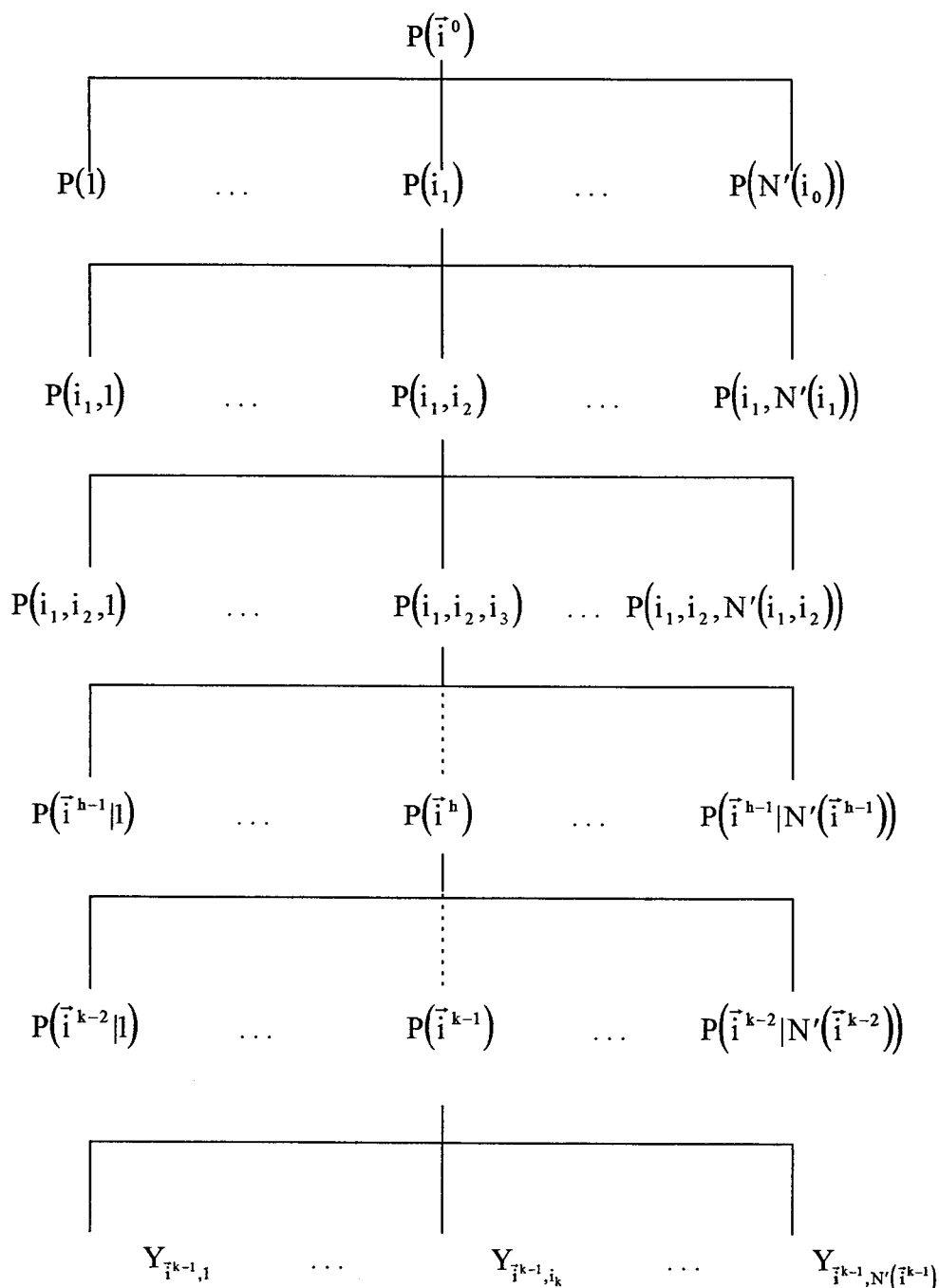


Figura 2 - Caso geral da decomposição de uma população em k níveis de fraccionamento.

ou seja, $N(\tilde{i}^h)$ representa a dimensão de $P(\tilde{i}^h)$, estando definida para $h=0,1,\dots,k$ e $N'(\tilde{i}^h)$ dá o número de sub-populações em que $P(\tilde{i}^h)$ se decompõe. $N(\tilde{i}^h)$ está definida apenas para $h=0,1,\dots,k-1$, visto a passagem do penúltimo para o último nível de estrutura, $k-1$ e k respectivamente, resultar numa decomposição directa de qualquer das sub-populações $P(\tilde{i}^{k-1})$ nos respectivos elementos, $P(\tilde{i}^k)$. Então:

$$(6.9) \quad \begin{cases} N(\tilde{i}^k) = 1 \\ N(\tilde{i}^{k-1}) = N'(\tilde{i}^{k-1}) \end{cases}$$

e

$$(6.10) \quad N(\tilde{i}^h) = \sum_{\tilde{i}^{h+1} \in \Omega(\tilde{i}^h)} N(\tilde{i}^{h+1}),$$

isto é, a dimensão de uma sub-população é igual à soma das dimensões das várias sub-populações resultantes da sua decomposição.

Atendendo a (6.10), obtemos a expressão

$$(6.11) \quad N(\tilde{i}^0) = \sum_{\tilde{i}^1 \in \Omega(\tilde{i}^0)} N(\tilde{i}^1),$$

que mostra que a dimensão da população, $N(\tilde{i}^0)$, pode ser dada em função das dimensões, $N(\tilde{i}^1)$, das sub-populações de nível 1.

Se quisermos reportar a dimensão da população ao nível k (último nível, sendo este constituído por elementos), é válida a seguinte expressão:

$$(6.12) \quad N(\tilde{i}^0) = \sum_{\tilde{i}^1 \in X_1} \sum_{\tilde{i}^2 \in X_2} \dots \sum_{\tilde{i}^{k-2} \in X_{k-2}} \sum_{\tilde{i}^{k-1} \in \Omega(\tilde{i}^{k-2})} N(\tilde{i}^{k-1}).$$

As expressões seguintes são apresentadas com o objectivo de melhor clarificar o significado de X_h e de $\Omega(\tilde{i}^h)$.

$$(6.13) \quad \left\{ \begin{array}{l} \mathbf{X}_0 = \{\tilde{i}^0\} \\ \mathbf{X}_1 = \{\tilde{i}^1\} = \{(i_1)\} = \Omega(\tilde{i}^0) = \{1, \dots, N'(\tilde{i}^0)\} \\ \mathbf{X}_2 = \{\tilde{i}^2\} = \{(i_1, i_2)\} = \bigcup_{\tilde{i}^1 \in \mathbf{X}_1} \Omega(\tilde{i}^1) = \bigcup_{\tilde{i}^1 \in \mathbf{X}_1} \{(i^1 | 1), \dots, (i^1 | N'(\tilde{i}^1))\} \\ \mathbf{X}_3 = \{\tilde{i}^3\} = \{(i_1, i_2, i_3)\} = \bigcup_{\tilde{i}^2 \in \mathbf{X}_2} \Omega(\tilde{i}^2) = \bigcup_{\tilde{i}^2 \in \mathbf{X}_2} \{(i^2 | 1), \dots, (i^2 | N'(\tilde{i}^2))\} \\ \vdots \\ \mathbf{X}_{h+1} = \{\tilde{i}^{h+1}\} = \{(i_1, \dots, i_{h+1})\} = \bigcup_{\tilde{i}^h \in \mathbf{X}_h} \Omega(\tilde{i}^h) = \bigcup_{\tilde{i}^h \in \mathbf{X}_h} \{(i^h | 1), \dots, (i^h | N'(\tilde{i}^h))\} \\ \vdots \\ \mathbf{X}_k = \{\tilde{i}^k\} = \{(i_1, \dots, i_k)\} = \bigcup_{\tilde{i}^{k-1} \in \mathbf{X}_{k-1}} \Omega(\tilde{i}^{k-1}) = \bigcup_{\tilde{i}^{k-1} \in \mathbf{X}_{k-1}} \{(i^{k-1} | 1), \dots, (i^{k-1} | N'(\tilde{i}^{k-1}))\} \end{array} \right.$$

À semelhança do que foi feito anteriormente para as sub-populações, iremos agora criar uma notação para as amostras.

Dado existirem k níveis de fraccionamento, a colheita da amostra é feita em k etapas. Em cada uma destas etapas colhem-se sub-populações, do nível correspondente, de entre as contidas nas sub-populações escolhidas na etapa precedente.

Representaremos pois, as sub-populações escolhidas por $P(\tilde{j}^h)$ sendo $n(\tilde{j}^h)$, $T(\tilde{j}^h)$, $\bar{Y}(\tilde{j}^h)$ respectivamente o número, o total e a média das observações colhidas em $P(\tilde{j}^h)$. Observe-se que $T(\tilde{j}^k)$ e $\bar{Y}(\tilde{j}^k)$ coincidem com a única observação $Y(\tilde{j}^k)$ colhida em $P(\tilde{j}^k)$. Os vectores \tilde{j}^h são aleatórios, podendo ocorrer mais de uma vez, se se estiver a amostrar com reposição.

Representaremos ainda por \mathbf{X}_h^* o conjunto dos vectores \tilde{j}^h e por $A(\tilde{i}^h)$ o acontecimento que se verifica quando $P(\tilde{i}^h)$ é escolhida. Observe-se que, com $p < h$, $A(\tilde{i}^h)$ implica $A(\tilde{i}^p)$.

Importa ainda referir que, as várias amostragens feitas na mesma etapa são independentes.

6.3. SUB-AMOSTRAGEM COM PROBABILIDADES IGUAIS SEM REPOSIÇÃO

Admitamos que a população se pode decompor em $N'(\tilde{i}^0)$ sub-populações, de nível 1, das quais a i_1 -ésima sub-população contem $N'(\tilde{i}^1)$ sub-populações, de nível 2, e assim sucessivamente até se chegar à sub-população de nível $k-1$ que é composta por $N'(\tilde{i}^{k-2})$ elementos.

Consideremos o acontecimento $A(\tilde{i}^h)$, com $h < k$, que se verifica quando a população $P(\tilde{i}^h)$ é escolhida.

Na fase, da sub-amostragem, em que se verifica $A(\tilde{i}^h)$, são escolhidas, com probabilidades iguais e sem reposição, $n'(\tilde{i}^h)$ sub-populações, de nível $h+1$, sub-populações estas que estão contidas em $P(\tilde{i}^h)$.

Como é obvio, a obtenção das sub-populações de nível $h+1$ depende, directamente, da sub-população, de nível h , que lhe deu origem, assim, temos as probabilidades condicionais:

$$(6.14) \quad \Pi(\tilde{i}^{h+1} | \tilde{i}_h^{h+1}) = P[A(\tilde{i}^{h+1}) | A(\tilde{i}_h^{h+1})] = \frac{n'(\tilde{i}_h^{h+1})}{N'(\tilde{i}_h^{h+1})}$$

Tendo em atenção o que se referiu para o nível $h+1$, facilmente se vê que o mesmo se verifica para qualquer nível. Logo tem-se:

$$A(\tilde{i}^{h+1}) \Rightarrow A(\tilde{i}_h^{h+1}) \Rightarrow \dots \Rightarrow A(\tilde{i}_p^{h+1}), \text{ com } p = 0, \dots, h,$$

isto é, a ocorrência de um acontecimento, num dado nível, resulta da ocorrência de um acontecimento do nível anterior, sucessivamente até ao nível zero.

Obtemos então

$$(6.15) \quad \Pi(\tilde{i}^{h+1}) = \Pi(\tilde{i}^{h+1} | \tilde{i}_h^{h+1}) \dots \Pi(\tilde{i}^{h+1} | \tilde{i}_p^{h+1}) \Pi(\tilde{i}_p^{h+1}) = \frac{n'(\tilde{i}_h^{h+1})}{N'(\tilde{i}_h^{h+1})} \dots \frac{n'(\tilde{i}_p^{h+1})}{N'(\tilde{i}_p^{h+1})} \Pi(\tilde{i}_p^{h+1}),$$

representando $\Pi(\tilde{i}^h)$ a probabilidade de $A(\tilde{i}^h)$. Como $\Pi(\tilde{i}^0) = 1$, podemos escrever

$$(6.16) \quad \Pi(\tilde{i}^{h+1}) = \prod_{p=0}^h \frac{n'(\tilde{i}_p^{h+1})}{N'(\tilde{i}_p^{h+1})}$$

tendo-se ainda

$$(6.17) \quad \Pi(\tilde{i}^{h+1}) = \Pi(\tilde{i}_h^{h+1}) \frac{n'(\tilde{i}_h^{h+1})}{N'(\tilde{i}_h^{h+1})}$$

6.3.1. CONSTRUÇÃO DO ESTIMADOR

Tendo em conta que, cada sub-população, de qualquer nível, apenas pode ser escolhida uma vez, dado que estamos a amostrar sem reposição, obtém-se o seguinte estimador para valor médio:

$$(6.18) \quad \mu^* = \frac{1}{N(\tilde{j}^0)} \sum_{j_1=1}^{n'(\tilde{j}^0)} \sum_{j_2=1}^{n'(\tilde{j}^1)} \dots \sum_{j_k=1}^{n'(\tilde{j}^{k-1})} \frac{Y(j_1, j_2, \dots, j_k)}{\Pi(j_1, j_2, \dots, j_k)} = \frac{1}{N(\tilde{j}^0)} \sum_{j^k \in X_k} \frac{Y(\tilde{j}^k)}{\Pi(\tilde{j}^k)}$$

onde $N(\tilde{j}^0)$ representa a dimensão da população.

De modo a simplificar a notação, representaremos, a partir daqui, a dimensão da população, $N(\tilde{j}^0)$, apenas por N .

Os vectores que representam as sub-populações de nível $h+1$, escolhidas entre as contidas em $P(\tilde{j}^h)$, pertencem a $\Omega(\tilde{j}^h) \cap X_{h+1}^*$. Usando (6.17), tem-se pois:

$$(6.19) \quad \mu^* = \frac{1}{N} \sum_{\tilde{j}^{k-1} \in X_{k-1}^*} \frac{1}{\Pi(\tilde{j}^{k-1})} \frac{N'(\tilde{j}^{k-1})}{n'(\tilde{j}^{k-1})} \sum_{j^k \in \Omega(\tilde{j}^{k-1}) \cap X_k^*} Y(\tilde{j}^k) = \frac{1}{N} \sum_{\tilde{j}^{k-1} \in X_{k-1}^*} \frac{N'(\tilde{j}^{k-1})}{\Pi(\tilde{j}^{k-1})} \bar{Y}(\tilde{j}^{k-1}).$$

Considerando que as sub-populações, pertencentes a qualquer nível h ($h=0,1,\dots,k-1$), se subdividem todas em igual número de sub-populações do nível seguinte, isto é,

$$(6.20) \quad \begin{cases} N'(\tilde{i}^h) = N'_h & ; \quad h = 0, 1, \dots, k-1 \\ n'(\tilde{i}^h) = n'_h & ; \quad h = 0, 1, \dots, k-1 \end{cases},$$

é então possível obter uma expressão simplificada do estimador μ^* .

Assim, sendo n a dimensão da amostra,

$$(6.21) \quad \begin{cases} N = \prod_{h=0}^{k-1} N'_h \\ n = \prod_{h=0}^{k-1} n'_h \end{cases},$$

logo, de (6.16), (6.20) e (6.21) vem $\Pi(\tilde{i}^k) = \frac{n}{N}$, e portanto o estimador para o valor médio vem:

$$(6.22) \quad \mu^* = \frac{1}{N} \sum_{\tilde{j}^k \in X_k^*} \frac{N}{n} Y(\tilde{j}^k) = \bar{Y},$$

onde \bar{Y} é a média global da amostra.

6.3.2. VARIÂNCIA DO ESTIMADOR

Como foi visto, em (6.4), a variância de X é a soma do valor médio da variância condicional com a variância do valor médio condicional.

Quando se chega ao nível k (k -ésima etapa da amostragem) a condição C corresponde aos resultados dos $k-1$ níveis anteriores pelo que a representamos por C_{k-1} , então:

$$(6.23) \quad \sigma^2(\mu^*) = \mu_{C_{k-1}} \left[\sigma^2(\mu^* | C_{k-1}) \right] + \sigma_{C_{k-1}}^2 \left[\mu(\mu^* | C_{k-1}) \right].$$

Utilizando os símbolos μ_h e σ_h^2 para indicar o valor médio e a variância, relativos à etapa h , sendo

$$(6.24) \quad \begin{cases} \sigma^2(\mu^* | C_{k-1}) = \sigma_k^2(\mu^*) \\ \mu(\mu^* | C_{k-1}) = \mu_k(\mu^*) \end{cases},$$

como podemos considerar $\mu_{(k)}(\mu^*)$ uma variável aleatória, cuja distribuição depende de C_{k-2} , vem

$$(6.25) \quad \sigma_{C_{k-1}}^2 \left[\mu_k(\mu^*) \right] = \mu_{C_{k-2}} \left[\sigma^2(\mu_k(\mu^*) | C_{k-2}) \right] + \sigma_{C_{k-2}}^2 \left[\mu(\mu_k(\mu^*) | C_{k-2}) \right].$$

Analogamente ao que fizemos anteriormente vamos pôr

$$(6.26) \quad \begin{cases} \sigma^2 \left[\mu_k(\mu^*) | C_{k-2} \right] = \sigma_{k-1}^2 \left[\mu_k(\mu^*) \right] \\ \mu \left[\mu_k(\mu^*) | C_{k-2} \right] = \mu_{k-1} \left[\mu_k(\mu^*) \right] \end{cases}$$

obtendo-se assim

$$(6.27) \quad \sigma^2(\mu^*) = \mu_{C_{k-1}} \left[\sigma_k^2(\mu^*) + \mu_{C_{k-2}} \left(\sigma_{k-1}^2(\mu_k) \right) \right] + \sigma_{C_{k-2}}^2 \left[\mu_{k-1} \mu_k(\mu^*) \right]$$

Pode ainda concluir-se que:

$$(6.28) \quad \mu_{C_h}(w) = \mu_1(\dots \mu_h(w)) ; \quad h = 1, \dots, k.$$

Representando por ∇_h ; $h = 1, \dots, k$, o cálculo dos valores médios, para as etapas com índices $h' \neq h$ e da variância para a etapa com índice h , obtém-se:

$$(6.29) \quad \sigma^2(\mu^*) = \sum_{h=1}^k \nabla_h(\mu^*),$$

ou seja, a variância do estimador $\sigma^2(\mu^*)$ é composta por tantas componentes quantas as etapas da sub-amostragem.

Note-se que, a validade deste resultado se verifica para qualquer técnica de sub-amostragem.

Começando por calcular a componente $\nabla_k(\mu^*)$, usando (6.19) e (2.40) temos então:

$$(6.30) \quad \begin{aligned} \sigma_k^2(\mu^*) &= \frac{1}{N^2} \sum_{\vec{j}^{k-1} \in X_{k-1}^*} \frac{N^{i^2}(\vec{j}^{k-1})}{\Pi^2(\vec{j}^{k-1})} \sigma^2(\bar{Y}(\vec{j}^{k-1})) = \\ &= \frac{1}{N^2} \sum_{\vec{j}^{k-1} \in X_{k-1}^*} \frac{N^{i^2}(\vec{j}^{k-1})}{\Pi^2(\vec{j}^{k-1})} \frac{\sigma^2(\vec{j}^{k-1})}{n^i(\vec{j}^{k-1})} \frac{N(\vec{j}^{k-1}) - n^i(\vec{j}^{k-1})}{N(\vec{j}^{k-1}) - 1} \end{aligned}$$

onde

$$(6.31) \quad \sigma^2(\bar{j}^{k-1}) = \frac{1}{N(\bar{j}^{k-1})} \left[\sum_{\bar{i}^k \in \Omega(\bar{j}^{k-1})} Y^2(\bar{i}^k) - \frac{T^2(\bar{j}^{k-1})}{N(\bar{j}^{k-1})} \right]$$

é a variância da sub-população $P(\bar{j}^{k-1})$, visto que os vectores, que correspondem aos seus elementos, são vectores de $\Omega(\bar{j}^{k-1})$, onde a característica em estudo toma os valores $Y(\bar{i}^k)$ com $\bar{i}^k \in \Omega(\bar{j}^{k-1})$.

Como os vectores representativos das sub-populações de nível $k-1$, escolhidos entre as contidas em $P(\bar{j}^{k-2})$, são os vectores que pertencem a $\Omega(\bar{j}^{k-2}) \cap X_{k-1}^*$ e como, devido a (6.17), se tem:

$$(6.32) \quad \Pi(\bar{i}^{h+1}) = \Pi(\bar{i}^h) \frac{n'(\bar{i}^h)}{N'(\bar{i}^h)} ; \text{ com } \bar{i}^{h+1} \in \Omega(\bar{i}^h) ,$$

podemos escrever (6.30) na forma

$$(6.33) \quad \sigma_k^2(\mu^*) = \frac{1}{N^2} \sum_{\bar{j}^{k-2} \in X_{k-2}^*} \frac{1}{\Pi^2(\bar{j}^{k-2})} \frac{N'(\bar{j}^{k-2})}{n'(\bar{j}^{k-2})} \cdot \frac{\sum_{\bar{j}^{k-1} \in \Omega(\bar{j}^{k-2}) \cap X_{k-1}} \frac{N'^2(\bar{j}^{k-1}) \sigma^2(\bar{j}^{k-1}) [N(\bar{j}^{k-1}) - n'(\bar{j}^{k-1})]}{n'(\bar{j}^{k-1}) [N(\bar{j}^{k-1}) - 1]}}{\frac{n'(\bar{j}^{k-2})}{N'(\bar{j}^{k-2})}} .$$

Para calcularmos os sucessivos valores médios, que acabarão por dar $\nabla_k(\mu^*)$, procederemos sempre da mesma maneira, utilizando (3.11) e (6.32). Assim temos:

(6.34)

$$\begin{aligned} \mu_{k-1}(\sigma_k^2(\mu^*)) &= \\ &= \frac{1}{N^2} \sum_{\bar{j}^{k-2} \in X_{k-2}^*} \frac{1}{P^2(\bar{j}^{k-2})} \frac{N'(\bar{j}^{k-2})}{n'(\bar{j}^{k-2})} \cdot \sum_{\bar{i}^{k-1} \in \Omega(\bar{j}^{k-2})} N^2(\bar{i}^{k-1}) \frac{s^2(\bar{i}^{k-1})}{n'(\bar{i}^{k-1})} \frac{N(\bar{i}^{k-1}) - n'(\bar{i}^{k-1})}{N(\bar{i}^{k-1}) - 1} = \end{aligned}$$

$$= \frac{1}{N^2} \sum_{\vec{j}^{k-3} \in X_{k-3}^*} \frac{1}{P^2(\vec{j}^{k-3})} \frac{N'(\vec{j}^{k-3})}{n'(\vec{j}^{k-3})} \sum_{\vec{j}^{k-2} \in X_{k-2} \cap \Omega(\vec{j}^{k-3})} \frac{N'(\vec{j}^{k-2})}{n'(\vec{j}^{k-2})} \cdot \sum_{\vec{i}^{k-1} \in \Omega(\vec{j}^{k-2})} \frac{N^2(\vec{i}^{k-1}) \frac{s^2(\vec{i}^{k-1}) N(\vec{i}^{k-1}) - n'(\vec{i}^{k-1})}{n'(\vec{i}^{k-1})} \frac{N(\vec{i}^{k-1}) - 1}{N(\vec{i}^{k-1}) - 1}}{\frac{n'(\vec{j}^{k-3})}{N'(\vec{j}^{k-3})}}$$

Analogamente tem-se

$$(6.35) \quad \mu_{k-2} \left[\mu_{k-1} \left(\sigma_k^2(\mu^*) \right) \right] = \frac{1}{N^2} \sum_{\vec{j}^{k-3} \in X_{k-3}^*} \frac{1}{\Pi^2(\vec{j}^{k-3})} \frac{N'(\vec{j}^{k-3})}{n'(\vec{j}^{k-3})} \sum_{\vec{i}^{k-2} \in \Omega(\vec{j}^{k-3})} \frac{N'(\vec{i}^{k-2})}{n'(\vec{i}^{k-2})} \cdot \sum_{\vec{i}^{k-1} \in \Omega(\vec{i}^{k-2})} N^2(\vec{i}^{k-1}) \frac{\sigma^2(\vec{i}^{k-1}) N(\vec{i}^{k-1}) - n'(\vec{i}^{k-1})}{n'(\vec{i}^{k-1})} \frac{N(\vec{i}^{k-1}) - 1}{N(\vec{i}^{k-1}) - 1}$$

e portanto

$$(6.36) \quad \nabla_k(\mu^*) = \mu_1 \left[\dots \mu_{k-1} \left(\sigma_k^2(\mu^*) \right) \right] = \frac{N'(\vec{i}^0)}{N^2 n'(\vec{i}^0)} \sum_{\vec{i}^1 \in X_1} \frac{N'(\vec{i}^1)}{n'(\vec{i}^1)} \sum_{\vec{i}^2 \in \Omega(\vec{i}^1)} \frac{N'(\vec{i}^2)}{n'(\vec{i}^2)} \dots \dots \sum_{\vec{i}^{k-2} \in \Omega(\vec{i}^{k-3})} \frac{N'(\vec{i}^{k-2})}{n'(\vec{i}^{k-2})} \cdot \sum_{\vec{i}^{k-1} \in \Omega(\vec{i}^{k-2})} N^2(\vec{i}^{k-1}) \frac{\sigma^2(\vec{i}^{k-1}) N(\vec{i}^{k-1}) - n'(\vec{i}^{k-1})}{n'(\vec{i}^{k-1})} \frac{N(\vec{i}^{k-1}) - 1}{N(\vec{i}^{k-1}) - 1}$$

Calculamos agora uma das $\nabla_h(\mu^*)$ com $1 < h < k$.

De (6.19) vem

$$(6.37) \quad \mu_k(\mu^*) = \frac{1}{N} \sum_{\vec{j}^{k-1} \in X_{k-1}^*} \frac{N'(\vec{j}^{k-1})}{\Pi(\vec{j}^{k-1})} \mu \left[\bar{Y}(\vec{j}^{k-1}) \right] = \frac{1}{N} \sum_{\vec{j}^{k-1} \in X_{k-1}^*} \frac{T(\vec{j}^{k-1})}{\Pi(\vec{j}^{k-1})}$$

visto que

$$(6.38) \quad N'(\vec{j}^{k-1}) \mu(\bar{Y}(\vec{j}^{k-1})) = N(\vec{j}^{k-1}) \mu(\bar{Y}(\vec{j}^{k-1})) = T(\vec{j}^{k-1}).$$

Portanto

$$(6.39) \quad \mu_{k-1}(\mu_k(\mu^*)) = \mu_{k-1} \left(\frac{1}{N} \sum_{\vec{j}^{k-1} \in X_{k-1}^*} \frac{T(\vec{j}^{k-1})}{\Pi(\vec{j}^{k-1})} \right) =$$

$$\begin{aligned}
&= \frac{1}{N} \mu_{k-1} \left(\sum_{\tilde{j}^{k-2} \in X_{k-2}^*} \frac{1}{\Pi(\tilde{j}^{k-2})} \sum_{\tilde{j}^{k-1} \in X_{k-1}^* \cap \Omega(\tilde{j}^{k-2})} \frac{T(\tilde{j}^{k-1})}{n'(\tilde{j}^{k-2})} \right) = \\
&= \frac{1}{N} \sum_{\tilde{j}^{k-2} \in X_{k-2}^*} \frac{1}{\Pi(\tilde{j}^{k-2})} \mu_{k-1} \left(\sum_{\tilde{j}^{k-1} \in X_{k-1}^* \cap \Omega(\tilde{j}^{k-2})} \frac{T(\tilde{j}^{k-1})}{\Pi(\tilde{j}^{k-2})} \right) = \\
&= \frac{1}{N} \sum_{\tilde{j}^{k-2} \in X_{k-2}^*} \frac{1}{\Pi(\tilde{j}^{k-2})} \sum_{\tilde{j}^{k-1} \in \Omega(\tilde{j}^{k-2})} T(\tilde{j}^{k-1}) = \\
&= \frac{1}{N} \sum_{\tilde{j}^{k-2} \in X_{k-2}^*} \frac{T(\tilde{j}^{k-2})}{\Pi(\tilde{j}^{k-2})}
\end{aligned}$$

e

$$(6.40) \quad \mu_{h+1}(\dots \mu_k(\mu^*)) = \frac{1}{N} \sum_{\tilde{j}^h \in X_h^*} \frac{T(\tilde{j}^h)}{\Pi(\tilde{j}^h)} = \frac{1}{N} \sum_{\tilde{j}^{h-1} \in X_{h-1}^*} \frac{1}{\Pi(\tilde{j}^{h-1})} \frac{N'(\tilde{j}^{h-1})}{n'(\tilde{j}^{h-1})} \sum_{\tilde{j}^h \in X_h^* \cap \Omega(\tilde{j}^{h-1})} T(\tilde{j}^h)$$

Como

$$(6.41) \quad \sigma^2 \left[\frac{1}{n'(\tilde{j}^{h-1})} \sum_{\tilde{j}^h \in X_h^* \cap \Omega(\tilde{j}^{h-1})} T(\tilde{j}^h) \right] = \frac{\sigma^2(\tilde{j}^{h-1})}{n'(\tilde{j}^{h-1})} \frac{N'(\tilde{j}^{h-1}) - n'(\tilde{j}^{h-1})}{N'(\tilde{j}^{h-1}) - 1},$$

visto tratar-se da média de uma amostra de totais de sub-populações, obtida sem reposição e com probabilidades iguais, tem-se ainda:

$$(6.42) \quad \sigma^2(\tilde{j}^{h-1}) = \frac{1}{N'(\tilde{j}^{h-1})} \left[\sum_{\tilde{i}^h \in \Omega(\tilde{j}^{h-1})} T^2(\tilde{i}^h) - \frac{T^2(\tilde{j}^{h-1})}{N'(\tilde{j}^{h-1})} \right]$$

logo

$$(6.43) \quad \sigma_h^2 \left[\mu_{h+1}(\dots \mu_k(\mu^*)) \right] = \frac{1}{N^2} \sum_{\tilde{j}^{h-1} \in X_{h-1}^*} \frac{1}{\Pi^2(\tilde{j}^{h-1})} \frac{N'^2(\tilde{j}^{h-1}) \sigma^2(\tilde{j}^{h-1}) [N'(\tilde{j}^{h-1}) - n'(\tilde{j}^{h-1})]}{n'(\tilde{j}^{h-1}) [N'(\tilde{j}^{h-1}) - 1]},$$

Raciocinando como atrás, obtemos:

(6.44)

$$\begin{aligned} \nabla_h(\mu^*) &= \frac{N'(\bar{i}^0)}{N^2 n'(\bar{i}^0)} \sum_{\bar{i}^1 \in X_1} \frac{N'(\bar{i}^1)}{n'(\bar{i}^1)} \sum_{\bar{i}^2 \in \Omega(\bar{i}^1)} \frac{N'(\bar{i}^2)}{n'(\bar{i}^2)} \dots \\ &\cdot \sum_{\bar{i}^{h-1} \in \Omega(\bar{i}^{h-2})} N^2(\bar{i}^{h-1}) \frac{\sigma^2(\bar{i}^{h-1})}{n'(\bar{i}^{h-1})} \frac{N'(\bar{i}^{h-1}) - n'(\bar{i}^{h-1})}{N'(\bar{i}^{h-1}) - 1} \end{aligned}$$

com $h = 2, \dots, k$.

Finalmente para obtermos $\nabla_1(\mu^*)$ vemos que se tem, atendendo à expressão (6.40),

$$(6.45) \quad \mu_2(\dots \mu_k(\mu^*)) = \frac{1}{N} \sum_{\bar{j}^1 \in X_1} \frac{1}{\Pi(\bar{j}^0)} \frac{N'(\bar{i}^0)}{n'(\bar{i}^0)} \sum_{\bar{j}^1 \in X_1} T(\bar{j}^1),$$

pelo que, raciocinando como para obter (6.44) e atendendo a que $\bar{j}^0 = \bar{i}^0$ e que $\Pi(\bar{i}^0) = 1$, vem

$$(6.46) \quad \nabla_1(\mu^*) = \frac{1}{N^2} N'^2(\bar{i}^0) \frac{\sigma^2(\bar{i}^0)}{n'(\bar{i}^0)} \frac{N'(\bar{i}^0) - n'(\bar{i}^0)}{N'(\bar{i}^0) - 1}$$

$$\text{com } \sigma^2(\bar{i}^0) = \frac{1}{N'(\bar{i}^0)} \left[\sum_{\bar{i}^1 \in X_1} T^2(\bar{i}^0) - \frac{T^2(\bar{i}^0)}{N'(\bar{i}^0)} \right], \text{ onde } \Omega(\bar{i}^0) = X_1.$$

6.4. SUB-AMOSTRAGEM COM PROBABILIDADES PROPORCIONAIS

Nesta técnica as sub-amostras da etapa h , $h=1, \dots, k$, têm dimensões n'_h , $h=1, \dots, k$, sendo obtidas com reposição e com probabilidades proporcionais, salvo na última etapa em que a selecção é feita sem reposição e com probabilidades iguais.

Tendo-se escolhido $P(\bar{j}^h)$, com $h=0, \dots, k-2$, as probabilidades das sub-populações $P(\bar{i}^{h+1})$, com $\bar{i}^{h+1} \in \Omega(\bar{i}^h)$, serem escolhidas, em cada tiragem, são $\frac{N(\bar{i}^{h+1})}{N(\bar{i}^h)}$. Facto que determina o nome adoptado para esta técnica.

Seja agora $\Pi(\bar{i}^h)$ o valor médio da variável aleatória $X(\bar{i}^h)$, que nos dá o número de vezes que $P(\bar{i}^h)$ é escolhida. Observe-se que dado estarmos a trabalhar com reposição, salvo na última etapa, se $X(\bar{i}^h) > 1$ e se $h < k-1$ então colhem-se $X(\bar{i}^h)$ sub-amostras em $P(\bar{i}^h)$.

Se $\bar{i}^{h+1} \in \Omega(\bar{i}^h)$, o número de vezes que $P(\bar{i}^{h+1})$ é escolhida, é a soma dos números de vezes que correspondem às várias sub-amostras colhidas em $P(\bar{i}^h)$. Sendo $X(\bar{i}^{h+1}, m)$,

com $m=1, \dots, X(\tilde{i}^h)$, as variáveis aleatórias que dão o número de vezes que $P(\tilde{i}^{h+1})$ é escolhida, teremos

$$(6.47) \quad X(\tilde{i}^{h+1}) = \sum_{m=1}^{X(\tilde{i}^h)} X(\tilde{i}^{h+1}, m) ; \quad h=1, \dots, k-1.$$

Sabendo que

$$(6.48) \quad \mu[X(\tilde{i}^{h+1}, m)] = n'_{h+1} \frac{N(\tilde{i}^{h+1})}{N(\tilde{i}^h)} ; \quad h=1, \dots, k-1,$$

e que, dado o valor tomado por $X(\tilde{i}^h)$ poder ser interpretado como uma condição de que depende a distribuição de $X(\tilde{i}^{h+1})$, com $\tilde{i}^{h+1} \in \Omega(\tilde{i}^h)$, teremos portanto

$$(6.49) \quad \begin{aligned} \Pi(\tilde{i}^{h+1}) &= \mu[X(\tilde{i}^{h+1})] = \sum_t P[X(\tilde{i}^h) = t] \cdot \mu[X(\tilde{i}^{h+1}) | X(\tilde{i}^h) = t] = \\ &= \sum_t P[X(\tilde{i}^h) = t] \cdot \mu\left[\sum_{m=1}^t X(\tilde{i}^{h+1}, m)\right] = \sum_t P[X(\tilde{i}^h) = t] \cdot t \cdot n'_{h+1} \frac{N(\tilde{i}^{h+1})}{N(\tilde{i}^h)} \\ &= n'_{h+1} \frac{N(\tilde{i}^{h+1})}{N(\tilde{i}^h)} \mu[X(\tilde{i}^h)] = n'_{h+1} \frac{N(\tilde{i}^{h+1})}{N(\tilde{i}^h)} \Pi(\tilde{i}^h), \quad \text{com } h=1, \dots, k-1. \end{aligned}$$

Por outro lado, utilizando os resultados apresentados anteriormente, obtém-se

$$(6.50) \quad \Pi(\tilde{i}^1) = n'_1 \frac{N(\tilde{i}^1)}{N(\tilde{i}^0)}.$$

De (6.49) e (6.50), vem

$$(6.51) \quad \Pi(\tilde{i}^{h+1}) = n'_1 \dots n'_{h+1} \frac{N(\tilde{i}^{h+1})}{N(\tilde{i}^0)} ; \quad h=1, \dots, k-1.$$

Observemos que a dimensão da amostra é dada por

$$(6.52) \quad n = n'_1 \dots n'_k$$

logo

$$(6.53) \quad \Pi(\tilde{i}^k) = \frac{n}{N(\tilde{i}^0)}.$$

6.4.1. CONSTRUÇÃO DO ESTIMADOR

Os resultados anteriores permitem construir a seguinte expressão para o estimador do valor médio:

$$(6.54) \quad \mu^* = \frac{1}{N(\vec{i}^0)} \sum_{\vec{j}^k \in X_k^*} \frac{Y(\vec{j}^k)}{\Pi(\vec{j}^k)} = \frac{1}{n} \sum_{\vec{j}^k \in X_k^*} Y(\vec{j}^k) = \bar{Y},$$

que indica que, o cálculo do estimador se reduz à determinação da média aritmética da amostra.

Relembremos, contudo, que esta generalização, para qualquer número (k) de níveis, só é possível desde que se satisfaçam as seguintes condições:

- Em todas as etapas, excepto a última, a amostragem é com probabilidades proporcionais às dimensões das sub-populações;
- Na última etapa, a amostragem é com probabilidades iguais.

6.4.2. VARIÂNCIA DO ESTIMADOR

Tal como sucede no caso da sub-amostragem com probabilidades iguais, a variância do estimador é dada por uma decomposição, em tantas componentes quantas as etapas da sub-amostragem, isto é,

$$(6.55) \quad \sigma^2(\bar{Y}) = \sigma^2(\mu^*) = \sum_{h=1}^k \nabla_{(h)}(\mu^*)$$

onde $\nabla_{(h)}(\mu^*)$; $h = 1, \dots, k$ é a componente da etapa h ($h = 1, \dots, k$) , para $\sigma^2(\mu^*)$.

Observe-se agora que, com $n_{k-1} = n'_1 \dots n'_{k-1}$, se tem

$$(6.56) \quad \mu^* = \frac{1}{n} \sum_{\vec{j}^k \in X_k^*} Y(\vec{j}^k) = \frac{1}{n_{k-1}} \sum_{\vec{j}^{k-1} \in X_{k-1}^*} \frac{1}{n'_k} \sum_{\vec{j}^k \in Q(\vec{j}^{k-1}) \cap X_k^*} Y(\vec{j}^k) = \frac{1}{n_{k-1}} \sum_{\vec{j}^{k-1} \in X_{k-1}^*} \bar{Y}(\vec{j}^{k-1})$$

onde $\bar{Y}(\vec{j}^{k-1})$ representa a média das observações, tomadas em $P(\vec{j}^{k-1})$. Como as sub-amostras da última etapa são colhidas independentemente uma das outras sem reposição, as $\bar{Y}(\vec{j}^{k-1})$ são independentes.

Adaptando a este caso a expressão (2.40), a variância da média das observações é dada por

$$(6.57) \quad \sigma^2[\bar{Y}(\vec{j}^{k-1})] = \frac{\sigma^2(\vec{j}^{k-1})}{n'_k} \frac{N(\vec{j}^{k-1}) - n'_k}{N(\vec{j}^{k-1}) - 1},$$

onde $\sigma^2(\vec{j}^{k-1})$ continua a representar a variância da sub-população $P(\vec{j}^{k-1})$, logo

$$(6.58) \quad \sigma_k^2(\mu^*) = \frac{1}{n_{k-1}^2} \sum_{\vec{j}^{k-1} \in X_{k-1}^*} \sigma^2[\bar{Y}(\vec{j}^{k-1})] = \frac{1}{n_{k-1}^2} \sum_{\vec{j}^{k-1} \in X_{k-1}^*} \frac{\sigma^2(\vec{j}^{k-1}) N(\vec{j}^{k-1}) - n'_k}{n'_k} \frac{N(\vec{j}^{k-1}) - n'_k}{N(\vec{j}^{k-1}) - 1},$$

então, recorrendo à expressão (3.11) utilizada na amostragem, com probabilidades variáveis, com reposição, vem:

(6.59)

$$\begin{aligned} \mu_{k-1}[\sigma_k^2(\mu^*)] &= \mu_{k-1} \left(\frac{1}{n_{k-1}^2} \sum_{\vec{j}^{k-2} \in X_{k-2}^*} \sum_{\vec{j}^{k-1} \in X_{k-1} \cap \Omega(\vec{j}^{k-2})} \frac{\sigma^2(\vec{j}^{k-1}) N(\vec{j}^{k-1}) - n'_k}{n'_k} \frac{N(\vec{j}^{k-1}) - n'_k}{N(\vec{j}^{k-1}) - 1} \right) = \\ &= \frac{1}{n_{k-1}^2} \sum_{\vec{j}^{k-2} \in X_{k-2}^*} \mu_{k-1} \left[\sum_{\vec{j}^{k-1} \in X_{k-1} \cap \Omega(\vec{j}^{k-2})} \frac{\sigma^2(\vec{j}^{k-1}) N(\vec{j}^{k-1}) - n'_k}{n'_k} \frac{N(\vec{j}^{k-1}) - n'_k}{N(\vec{j}^{k-1}) - 1} \right] = \\ &= \frac{1}{n_{k-1}^2} \sum_{\vec{j}^{k-2} \in X_{k-2}^*} \sum_{\vec{i}^{k-1} \in \Omega(\vec{j}^{k-2})} \frac{n'_{k-1} N(\vec{i}^{k-1})}{N(\vec{j}^{k-2})} \frac{\sigma^2(\vec{i}^{k-1}) N(\vec{i}^{k-1}) - n'_k}{n'_k} \frac{N(\vec{i}^{k-1}) - n'_k}{N(\vec{i}^{k-1}) - 1}. \end{aligned}$$

Tendo em atenção a expressão(6.48), o valor médio da variável aleatória que dá o número de vezes que $P(\vec{j}^{h+1})$ é escolhida, cada vez que se amostra para $P(\vec{j}^{k-2})$ com $\vec{j}^{k-2} \in \Omega(\vec{j}^{k-1})$, é

$$(6.60) \quad \mu[X(\vec{j}^{k-1}, m)] = n'_{k-1} \frac{N(\vec{j}^{k-1})}{N(\vec{j}^{k-2})} ; \quad m = 1, \dots, X(\vec{j}^{k-2}).$$

Continuando a raciocinar desta forma obtemos

(6.61)

$$\begin{aligned} \nabla_k(\mu^*) &= \mu_1[\dots \sigma_k^2(\mu^*)] = \frac{1}{n_{k-1}^2} \sum_{\vec{i}^1 \in X_1} \frac{n'_1 N(\vec{i}^1)}{N(\vec{i}^0)} \sum_{\vec{i}^2 \in \Omega(\vec{i}^1)} \frac{n'_2 N(\vec{i}^2)}{N(\vec{i}^1)} \dots \\ &\dots \sum_{\vec{i}^{k-2} \in \Omega(\vec{i}^{k-3})} \frac{n'_{k-2} N(\vec{i}^{k-2})}{N(\vec{i}^{k-3})} \sum_{\vec{i}^{k-1} \in \Omega(\vec{i}^{k-2})} \frac{n'_{k-1} N(\vec{i}^{k-1})}{N(\vec{i}^{k-2})} \frac{\sigma^2(\vec{i}^{k-1}) N(\vec{i}^{k-1}) - n'_k}{n'_k} \frac{N(\vec{i}^{k-1}) - n'_k}{N(\vec{i}^{k-1}) - 1} = \\ &= \frac{1}{nN(\vec{i}^0)} \sum_{\vec{i}^{k-1} \in X_{k-1}} N(\vec{i}^{k-1}) \sigma^2(\vec{i}^{k-1}) \frac{N(\vec{i}^{k-1}) - n'_k}{N(\vec{i}^{k-1}) - 1}. \end{aligned}$$

Procuramos agora obter a expressão das $\nabla_h(\mu^*)$, $h = 2, \dots, k-1$.

De (6.38) vem que $\mu[\bar{Y}(\bar{j}^{k-1})] = \frac{T(\bar{j}^{k-1})}{N(\bar{j}^{k-1})}$ é o valor médio de $P(\bar{j}^{h+1})$, obtém-se então:

$$(6.62) \quad \mu_k(\mu^*) = \frac{1}{n_{k-1}} \sum_{\bar{j}^{k-1} \in X_{k-1}^*} \frac{T(\bar{j}^{k-1})}{N(\bar{j}^{k-1})}.$$

Recorrendo novamente à expressão (3.11) vem

$$(6.63) \quad \begin{aligned} \mu_{k-1}[\mu_k(\mu^*)] &= \mu_{k-1} \left[\frac{1}{n'_1 \dots n'_{k-1}} \sum_{\bar{j}^{k-2} \in X_{k-2}^*} \sum_{\bar{j}^{k-1} \in X_{k-1}^* \cap \Omega(\bar{j}^{k-2})} \frac{T(\bar{j}^{k-1})}{N(\bar{j}^{k-1})} \right] = \\ &= \frac{1}{n'_1 \dots n'_{k-1}} \sum_{\bar{j}^{k-2} \in X_{k-2}^*} \mu_{k-1} \left[\sum_{\bar{j}^{k-1} \in X_{k-1}^* \cap \Omega(\bar{j}^{k-2})} \frac{T(\bar{j}^{k-1})}{N(\bar{j}^{k-1})} \right] = \\ &= \frac{1}{n'_1 \dots n'_{k-1}} \sum_{\bar{j}^{k-2} \in X_{k-2}^*} \sum_{\bar{i}^{k-1} \in \Omega(\bar{j}^{k-2})} n'_{k-1} \frac{N(\bar{i}^{k-1})}{N(\bar{j}^{k-2})} \frac{T(\bar{i}^{k-1})}{N(\bar{i}^{k-1})} = \\ &= \frac{1}{n'_1 \dots n'_{k-2}} \sum_{\bar{j}^{k-2} \in X_{k-2}^*} \frac{T(\bar{j}^{k-2})}{N(\bar{j}^{k-2})} \end{aligned}$$

com $\sum_{\bar{i}^{k-1} \in \Omega(\bar{j}^{k-2})} T(\bar{i}^{k-1}) = T(\bar{j}^{k-2})$, e também

$$(6.64) \quad \mu_{h+1}[\dots \mu_k(\mu^*)] = \frac{1}{n'_1 \dots n'_h} \sum_{\bar{j}^h \in X_h^*} \frac{T(\bar{j}^h)}{N(\bar{j}^h)} = \frac{1}{n'_1 \dots n'_h} \sum_{\bar{j}^h \in X_h^*} \mu(\bar{j}^h)$$

onde $\mu(\bar{j}^h)$ é o valor médio de $P(\bar{j}^h)$. Atendendo à independência das sub-amostragens realizadas na mesma etapa, vem

$$(6.65) \quad \begin{aligned} \sigma_h^2[\mu_{h+1}[\dots \mu_k(\mu^*)]] &= \frac{1}{n_1'^2 \dots n_{h-1}'^2} \sum_{\bar{j}^{h-1} \in X_{h-1}^*} \frac{1}{n'_h} \sum_{\bar{j}^h \in X_h^* \cap \Omega(\bar{j}^{h-1})} \mu(\bar{j}^h) = \\ &= \frac{1}{n_1'^2 \dots n_{h-1}'^2} \sum_{\bar{j}^{h-1} \in X_{h-1}^*} \sigma^2[W(\bar{j}^{h-1})] \end{aligned}$$

onde

$$(6.66) \quad W(\bar{j}^{h-1}) = \frac{1}{n'_h} \sum_{\bar{j}^h \in X_h^* \cap \Omega(\bar{j}^{h-1})} \mu(\bar{j}^h)$$

é a média de uma amostra com dimensão n'_h , de valores médios de sub-populações $P(\tilde{i}^h)$, com $\tilde{i}^h \in \Omega(\tilde{i}^{h-1})$. Como esta amostra é colhida com reposição e probabilidades iguais tem-se

$$(6.67) \quad \sigma^2[W(\tilde{j}^{h-1})] = \frac{1}{n'_h} \sum_{\tilde{i}^h \in \Omega(\tilde{j}^{h-1})} \frac{N(\tilde{i}^h)}{N(\tilde{j}^{h-1})} [\mu(\tilde{i}^h) - \mu(\tilde{j}^{h-1})]^2$$

vindo

$$(6.68) \quad \sigma_h^2[\mu_{h+1}[\dots\mu_k(\mu^*)]] = \frac{1}{n_1'^2 \dots n_{h-1}'^2 \dots n_h'} \sum_{\tilde{j}^{h-1} \in X_{h-1}^*} \sum_{\tilde{i}^h \in \Omega(\tilde{j}^{h-1})} \frac{N(\tilde{i}^h)}{N(\tilde{j}^{h-1})} [\mu(\tilde{i}^h) - \mu(\tilde{j}^{h-1})]^2$$

Utilizando mais uma vez a expressão (3.11), temos

$$(6.69) \quad \begin{aligned} \mu_{h-1} \left[\sigma_h^2[\mu_{h+1}(\dots\mu_k(\mu^*))] \right] &= \\ &= \frac{1}{n_1'^2 \dots n_{h-1}'^2} \sum_{\tilde{j}^{h-2} \in X_{h-2}^*} \mu_{h-1} \left[\sum_{\tilde{j}^{h-1} \in X_{h-1}^* \cap \Omega(\tilde{j}^{h-2})} \frac{1}{n'_h} \sum_{\tilde{i}^h \in \Omega(\tilde{j}^{h-1})} \frac{N(\tilde{i}^h)}{N(\tilde{j}^{h-1})} [\mu(\tilde{i}^h) - \mu(\tilde{j}^{h-1})]^2 \right] = \\ &= \frac{1}{n_1'^2 \dots n_{h-1}'^2} \sum_{\tilde{j}^{h-2} \in X_{h-2}^*} \sum_{\tilde{i}^{h-1} \in \Omega(\tilde{j}^{h-2})} \frac{n'_{h-1} N(\tilde{i}^{h-1})}{N(\tilde{j}^{h-2})} \frac{1}{n'_h} \sum_{\tilde{i}^h \in \Omega(\tilde{i}^{h-1})} \frac{N(\tilde{i}^h)}{N(\tilde{i}^{h-1})} [\mu(\tilde{i}^h) - \mu(\tilde{i}^{h-1})]^2 = \\ &= \frac{1}{n_1'^2 \dots n_{h-2}'^2 n'_{h-1} n'_h} \sum_{\tilde{j}^{h-2} \in X_{h-2}^*} \frac{1}{N(\tilde{j}^{h-2})} \sum_{\tilde{i}^{h-1} \in \Omega(\tilde{j}^{h-2})} \sum_{\tilde{i}^h \in \Omega(\tilde{i}^{h-1})} N(\tilde{i}^h) [\mu(\tilde{i}^h) - \mu(\tilde{i}^{h-1})]^2 \end{aligned}$$

e, em seguida, com $n_h = n'_1 \dots n'_h$, $h = 1, \dots, k$ tem-se

$$(6.70) \quad \begin{aligned} \nabla_h(\mu^*) &= \mu_1 \left[\dots \mu_{h-1} \left[\sigma_h^2(\mu_{h+1} \dots \mu_k(\mu^*)) \right] \right] = \\ &= \frac{1}{n'_h N(\tilde{i}^0)} \sum_{\tilde{i}^h \in X_h} N(\tilde{i}^h) [\mu(\tilde{i}^h) - \mu(\tilde{i}^{h-1})]^2, \quad h = 2, \dots, k \end{aligned}$$

Para obtermos $\nabla_1(\mu^*)$, observamos que, devido a (6.59) e (6.61), se tem

$$(6.71) \quad \mu_2[\dots\mu_k(\mu^*)] = \frac{1}{n'_1} \sum_{\tilde{j}^1 \in X_1} \mu(\tilde{j}^1) = W(\tilde{j}^0)$$

deduzindo-se então

$$(6.72) \quad \nabla_1(\mu^*) = \sigma_1^2 [\mu_2, \dots, \mu_k(\mu^*)] = \frac{1}{n'_1} \sum_{\tilde{i}^1 \in X_1} \frac{N(\tilde{i}^1)}{N(\tilde{i}^0)} \left[\mu(\tilde{i}^1) - \mu(\tilde{i}^0) \right]^2.$$

Se considerarmos

$$(6.73) \quad \begin{cases} A_h = \frac{1}{N(\tilde{i}^0)} \sum_{\tilde{i}^h \in X_h} N(\tilde{i}^h) \left[\mu(\tilde{i}^h) - \mu(\tilde{i}^{h-1}) \right]^2 ; h = 1, \dots, k-1 \\ A_k = \frac{1}{N(\tilde{i}^0)} \sum_{\tilde{i}^{k-1} \in X_k} N(\tilde{i}^{k-1}) \sigma^2(\tilde{i}^{k-1}) \frac{N(\tilde{i}^{k-1}) - n'_k}{N(\tilde{i}^{k-1}) - 1} \end{cases}$$

obtemos, atendendo a (6.50), (6.56), (6.65) e (6.67),

$$(6.74) \quad \sigma^2(\mu^*) = \sum_{h=1}^k \frac{A_h}{n_h} = \sum_{h=1}^k \frac{A_h}{\prod_{m=1}^h n'_m}$$

6.5. CUSTO TOTAL

Nesta secção procura-se resolver o problema da escolha de n'_1, n'_2, \dots, n'_k , para um custo total dado, adoptando para o mesmo uma expressão da forma

$$(6.75) \quad C = C_0 + \sum_{h=1}^k C_h n_h = C_0 + \sum_{h=1}^k C_h \prod_{m=1}^h n'_m.$$

Dado C_0 ser um custo fixo, a obtenção do valor óptimo para n'_m , $m=1, \dots, k$, não depende do seu valor. Pode pois considerar-se $C_0=0$.

Assim, de (6.74) e (6.75), usando a desigualdade de Cauchy, temos

$$(6.76) \quad \left(\sum_{h=1}^k \frac{A_h}{\prod_{m=1}^h n'_m} \right) \left(\sum_{h=1}^k C_h \prod_{m=1}^h n'_m \right) \geq \left(\sqrt{\sum_{h=1}^k A_h C_h} \right)^2 \geq \sum_{h=1}^k A_h C_h,$$

atingindo-se o mínimo quando

$$(6.77) \quad \frac{\frac{A_h}{\prod_{m=1}^h n'_m}}{C_h \prod_{m=1}^h n'_m} = k \text{ (constante),}$$

ou seja,

$$(6.78) \quad \left(\prod_{m=1}^h n'_m \right)^2 = \frac{A_h}{C_h k}.$$

Como

$$(6.79) \quad \prod_{m=1}^h n'_m = n_h$$

temos

$$(6.80) \quad \frac{n_{h+1}}{n_h} = n'_{h+1}.$$

Substituindo (6.79) em (6.78) obtém-se

$$(6.81) \quad n_h^2 = \frac{A_h}{C_h k}$$

e também

$$(6.82) \quad n_{h+1}^2 = \frac{A_{h+1}}{C_{h+1} k}.$$

Atendendo a (6.80), usando (6.81) e (6.82), temos

$$(6.83) \quad n_{h+1}'^2 = \frac{\frac{A_{h+1}}{k C_{h+1}}}{\frac{A_h}{k C_h}} = \frac{A_{h+1}}{C_{h+1}} \frac{C_h}{A_h}.$$

Então o custo mínimo é atingido com

$$(6.84) \quad n'_{h+1} = \sqrt{\frac{A_{h+1} C_h}{C_{h+1} A_h}}, \quad h = 1, \dots, k-1.$$

- (1) Hansen, M. H. & Hurwitz, W. N. - The problem of non-response in sample survey, J.A.S.A., n.91 p.517 , 1948
- (2) Hardy, G.H. , Littlewood , J.E. & Polya - Inequalities - Cambridge University Press , Cambridge , 1934
- (3) Kendall, M. ; Stuart, A. & Ord, J.K. - The Advanced Theory of Statistics - Vol III , Charles Griffins & Co , London , 1983
- (4) Kish , L. & Hess , I. - On variance of ratios and their differences in multistage samples - J.A.S.A., n.54 p.416 , 1954
- (5) Kish , L. & Anderson,D.W. - Multivariate and multipurpose stratification , J.A.S.A. , n.73 p.24 - 1978
- (6) Kokan, A.R. & Khan,S. - Optimum allocation in multivariate surveys: na analytical solution, J.R. Stat. Soc. , B29 - 1967

Das referências bibliográficas apresentadas, destaca-se a indicada em (3) pois serviu de base a este estudo.

Da bibliografia existente sobre amostragem, terão especial interesse os artigos indicados em (1) , (4) , (5) e (6) .

É também referido, em (2) o livro onde vem desenvolvido o estudo da desigualdade de Cauchy