

MAT02036 - Amostragem 2

Aula 14 - Amostragem por Conglomerados - Correlação Intraclasses

Markus Stein

Departamento de Estatística, IME/UFRGS

2022/2

Housekeeping

- Aproveitem o momento presencial para tirar dúvidas
- Se estivéssemos no ensino remoto ou à distância
 - vocês poderiam estar somente ouvindo, sem interação
 - ou assistindo vídeos e material em outro momento
- Depois das aulas, rever material da aula passada
 - fazer exercícios
 - se preparar para a próxima aula

Aula passada

Amostragem por Conglomerados em 1 Estágio Simples

Exercício

Um certo país possui $M = 10$ companhias aéreas com N_i aviões cada. As milhas (*em milhares*) de cada avião (y_{ij}) num determinado período de tempo foram registradas.

Cia (i)	No. aviões (N_i)	T_i	\bar{Y}_i
1	10	40	4
2	15	75	5
3	15	75	5
4	15	60	4
5	10	60	6
6	15	90	6
7	15	75	5
8	10	70	7
9	10	40	4
10	15	90	6

Aula passada

Amostragem por Conglomerados em 1 Estágio Simples

Exercício

- a. Calcule os parâmetros total (T) e média, individual (\bar{Y}) e por conglomerados (\bar{Y}_C) e a variância entre totais dos conglomerados S_{ec}^2 .
- b*. Calcule o viés dos estimadores **HT** e de **razão** para o total, \hat{T} (ou média, \bar{y}), são não viesados para os respectivos parâmetros que se destinam a estimar, T e \bar{Y} . (Obs. mostrar analiticamente ou com os dados do exercício)
- b. Assumindo o plano **AC1S** com **AASs** de conglomerados, para amostras de tamanho $m = 4$, calcule a variância do estimador natural (HT) do total, $Var_{AC1S}(\hat{T}^{HT})$, e a variância do estimador da média, $Var_{AC1S}(\bar{y}^{HT})$.
- c. Repetir o item (b) para estimador o estimador de razão.
- d. Escolha um estimador para o total (ou para a média), selecione uma amostra e estime o parâmetro com base na amostra observada.

Efeito do plano amostral

Efeito do plano amostral

- O **Efeito do Plano Amostral - EPA** é uma medida para comparar a **eficiência** de duas estratégias, E_1 e E_2 , formadas pelas combinações de **plano amostral** e **estimador**, para um **mesmo tamanho de amostra**.

$$EPA(E_1; E_2) = V_{E_1}(\hat{\theta}_1)/V_{E_2}(\hat{\theta}_2)$$

- O termo original em inglês é **Design Effect - deff** e foi sugerido por @Kish1965.
- Outra medida que dá uma indicação **semelhante** ao EPA é o **Fator do Plano Amostral - FPA**, que vem do inglês *Design Factor*, definido como:

$$FPA(E_1; E_2) = \sqrt{EPA(E_1; E_2)} = DP_{E_1}(\hat{\theta}_1)/DP_{E_2}(\hat{\theta}_2)$$

- O *FPA* compara diretamente o *desvio padrão* dos estimadores sob duas estratégias diferentes de amostragem.
 - É mais comum o uso do EPA que do FPA, sendo o FPA mais diretamente relacionado com a margem de erro das estimativas, enquanto o uso do EPA é mais conveniente quando se trata de planejar e dimensionar amostras.

Efeito do plano amostral

Exemplo - Efeito do plano amostral ao estimar a média populacional por unidade elementar, através do estimador HT com amostragem conglomerada simples em um estágio, em relação ao uso de uma AAS de igual tamanho.

Neste caso, as duas estratégias cuja eficiência se quer comparar são:

- **Estratégia 1:** Amostragem conglomerada em um estágio simples - AC1S, com o estimador natural $\bar{y}_{AC1S/HT}$.
- **Estratégia 2:** Amostragem aleatória simples - AAS de mesmo tamanho total (n), com o estimador usual de média $\bar{y} = \frac{1}{n} \sum_{i \in s} y_i$.

O efeito do plano amostral (neste caso, conglomeração) ao estimar a média populacional por unidade elementar é:

$$EPA(AC1S/HT; AAS) = \frac{V_{AC1S}(\bar{y}_{AC1S/HT})}{V_{AAS}(\bar{y})}$$

Efeito do plano amostral

O EPA mede o quanto a variância do estimador é maior (ou menor) por usar, neste caso, AC1S em lugar de AAS.

- $EPA < 1 \Rightarrow$ *ganho de precisão*, devido ao uso de amostragem conglomerada.
- $EPA = 1 \Rightarrow$ *mesma precisão*, não há diferença de precisão, pode-se optar pelo plano operacionalmente mais vantajoso.
- $EPA > 1 \Rightarrow$ *perda de precisão*, devido ao uso de amostragem conglomerada.

Um valor de $EPA = 5$, por exemplo, indicaria que a variância sob amostragem conglomerada seria cinco vezes maior que a variância de uma AAS de igual tamanho total.

Efeito do plano amostral

Exemplo 2 - Efeito do plano amostral ao estimar a média populacional por unidade elementar, através do estimador tipo razão com AC1S em relação ao uso da AAS.

- **Estratégia 1:** Amostragem conglomerada em um estágio simples - AC1S, com estimador tipo razão $\bar{y}_{AC1S}^R = \frac{1}{n} \sum_{i \in a} Y_i$ para a média.
- **Estratégia 2:** Amostragem aleatória simples - AAS de mesmo tamanho total (n), com o estimador usual de média $\bar{y} = \frac{1}{n} \sum_{i \in s} y_i$.

O efeito do plano amostral (neste caso, conglomeração) ao estimar a média populacional por unidade elementar é:

$$EPA(AC1S^R; AAS) = \frac{V_{AC1S}(\bar{y}_{AC1S}^R)}{V_{AAS}(\bar{y})}$$

Nota: Os estimadores pontuais são idênticos; somente os planos amostrais (e as variâncias) são diferentes.

Efeito do plano e Correlação Intraclass

Coeficiente de Correlação Intraclass

Vimos que a eficiência na **AC1S** depende do grau de similaridade dos seus elementos.

- Uma medida para indicar esse grau de similaridade é o coeficiente de correlação intraclass (**CCI**), *intraclass correlation coefficient* ou *intracluster correlation coefficient* (**ICC**).
- Considere a população dividida em M conglomerados:
 - Dentro do i -ésimo conglomerado existem $N_i(N_i - 1)$ pares de valores distintos da variável Y .

Elemento	$(i, 1)$	$(i, 1)$...	(i, j)	...	(i, N_i)
$(i, 1)$	-	(Y_{i1}, Y_{i2})	...	(Y_{i1}, Y_{ij})	...	(Y_{i1}, Y_{iN_i})
$(i, 2)$	(Y_{i2}, Y_{i1})	-	...	(Y_{i2}, Y_{ij})	...	(Y_{i2}, Y_{iN_i})
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
(i, j)	(Y_{ij}, Y_{i1})	(Y_{ij}, Y_{i2})	...	-	...	(Y_{ij}, Y_{iN_i})
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
(i, N_i)	(Y_{iN_i}, Y_{i1})	(Y_{iN_i}, Y_{i2})	...	(Y_{iN_i}, Y_{ij})	...	-

Coeficiente de Correlação Intraclass

- O **CCI** é coeficiente de correlação de Pearson para todos os $\sum_{i \in C} N_i(N_i - 1)$ pares do tipo (Y'_1, Y'_2) ,

$$\rho = \frac{Cov(Y'_1, Y'_2)}{\sqrt{Var(Y'_1)Var(Y'_2)}}$$

- Y'_1 significa possíveis valores da primeira posição do par
- Y'_2 significa possíveis valores da segunda posição do par.

Exemplo: (cont. Exemplo slides Aula 12)

Considere a população de tamanho $N = 6$ agrupada em $M = 3$ conglomerados, de três maneiras diferentes:

$$\mathbf{Y}_A = ((7, 8); (9, 10); (12, 14))$$

$$\mathbf{Y}_B = ((7, 10); (12, 8); (9, 14))$$

$$\mathbf{Y}_C = ((7, 14); (12, 8); (9, 10))$$

Coeficiente de Correlação Intraclass

Tamanhos dos conglomerados *iguais*

- Se $N_i = \bar{N}$, $\forall i = 1, \dots, M$, sob **AASc** dentro dos conglomerados,

$$Cov(Y_1', Y_2') = \frac{\sum_{i \in C} \sum_{j \in C_i} \sum_{k \neq j \in C_i} (y_{ij} - \bar{Y}) (y_{ik} - \bar{Y})}{M \bar{N} (\bar{N} - 1)}$$

e

$$Var(Y_1') = Var(Y_2') = Var_y,$$

então

$$\rho = \frac{Var_{ec} - \frac{Var_{dc}}{\bar{N} - 1}}{Var_y}$$

Mostrar???

Coeficiente de Correlação Intraclass

Tamanhos dos conglomerados *iguais*

- Sob **AASs** dentro dos conglomerados, de acordo com @Cochran1977, página 242, tem-se:

$$\rho = \frac{\sum_{i \in C} \sum_{j \in C_i} \sum_{k \neq j \in C_i} (y_{ij} - \bar{Y}) (y_{ik} - \bar{Y})}{(\bar{N} - 1) (M\bar{N} - 1) S_y^2} \doteq 1 - \frac{S_d^2}{S_y^2}$$

é o **coeficiente de correlação intraconglomerado** ou **intraclass**; e S_d^2 é a medida da variância *dentro* dos conglomerados, dada por:

$$S_{dc}^2 = \frac{1}{M} \sum_{i \in C} \frac{1}{\bar{N} - 1} \sum_{j \in C_i} (y_{ij} - \bar{Y}_i)^2 = \frac{1}{M} \sum_{i \in C} S_i^2$$

com:

$$S_i^2 = \frac{1}{\bar{N} - 1} \sum_{j \in C_i} (y_{ij} - \bar{Y}_i)^2$$

Coeficiente de Correlação Intraclass

Tamanhos dos conglomerados *iguais*

- Sob **AASs** de conglomerados, lembrando, variância **total** é dada por:

$$S_y^2 = \frac{1}{N-1} \sum_{i \in C} \sum_{j \in C_i} (y_{ij} - \bar{Y})^2$$

- Vimos que a variância **total** também pode ser expressa em função das variâncias *entre* conglomerados, S_{ec}^2 , e *dentro* dos conglomerados, S_{dc}^2 , através da expressão:

$$S_y^2 = \frac{(\bar{N} - 1)MS_{dc}^2 + \bar{N}(M - 1)\bar{S}_{ec}^2}{M\bar{N} - 1}$$

onde, $\bar{S}_{ec}^2 = \frac{S_{ec}^2}{\bar{N}}$.

Coeficiente de Correlação Intraclass

Tamanhos dos conglomerados *iguais*

- Sob AASs de conglomerados

$$EPA(AC1S^R; AAS) \doteq 1 + (\bar{N} - 1)\rho$$

A expressão para o $EPA(AC1S^R; AAS)$ resulta do uso das expressões de acordo com @Cochran1977, página 241:

$$Var_{AC1S}(\bar{y}_{AC1S}^R) \doteq \left(\frac{1}{m\bar{N}} - \frac{1}{M\bar{N}} \right) S_y^2 [1 + (\bar{N} - 1)\rho]$$

$$Var_{AAS}(\bar{y}) = \left(\frac{1}{m\bar{N}} - \frac{1}{M\bar{N}} \right) S_y^2$$

Coeficiente de Correlação Intraclass

Algumas considerações relacionadas com a variação do *EPA* para AC1S:

1. Se os conglomerados tiverem variância dentro grande, isto é, se $S_{dc}^2 \doteq S_y^2$, então $\rho \doteq 0$ e portanto, $EPA(AC1S^R; AAS) \doteq 1 + (\bar{N} - 1) \times 0 = 1$.

- Nesse caso, não ocorreria perda de precisão devido ao uso de amostragem conglomerada.

2. Pode ser demonstrado (?) que $\left[- \left(\frac{M-1}{N-M} \right) \right] \leq \rho \leq 1$, ou ainda, de forma aproximada, que $\frac{-1}{N-1} \leq \rho \leq 1$.

- Em muitas aplicações práticas, $\rho > 0$, porque os conglomerados tendem a ser mais homogêneos internamente do que a população em geral.
- Consequência: $EPA(AC1S^R; AAS) > 1$ na maioria das vezes.

Coeficiente de Correlação Intraclass

3. Raramente $\rho < 0$, caso em que AC1S seria mais eficiente que AAS.

4. Num caso extremo, $\rho = 1$ e portanto $EPA(AC1S^R; AAS) = \overline{N}$ e $Var_{AC1S}(\overline{y}_{AC1S}^R) = EPA(AC1S^R; AAS) Var_{AAS}(\overline{y}) \doteq \overline{N} \frac{S_y^2}{m\overline{N}} = \frac{S_y^2}{m}$

- Nesse caso, a precisão da amostra conglomerada de tamanho total igual a $m\overline{N}$ é equivalente apenas àquela obtida com uma amostra aleatória simples de tamanho m !!!

A tabela abaixo apresenta EPAs sob **AC1S** de acordo com valores selecionados dos tamanhos dos conglomerados e do coeficiente de correlação intraclass.

Tamanho	$\rho = 0,01$	$\rho = 0,05$	$\rho = 0,1$	$\rho = 0,2$	$\rho = 0,3$	$\rho = 0,5$
2	1	1	1	1	1	2
5	1	1	1	2	2	3
11	1	2	2	3	4	6
21	1	2	3	5	7	11
31	1	3	4	7	10	16
51	2	4	6	11	16	26

Coeficiente de Correlação Intraclass

Estimação

- Numa **amostra** retirada **com reposição**, o **coeficiente de correlação intraclass** pode ser **estimado** por:

$$r = \frac{s_{ec}^2 \frac{s_{dc}^2}{N}}{s_{ec}^2 + s_{dc}^2}$$

em que

$$s_{ec}^2 = \frac{1}{m-1} \sum_{i \in a} \left(\bar{Y}_i - \bar{y}^{HT} \right)^2 \quad \text{e} \quad s_{dc}^2 = \frac{1}{m} \sum_{i \in a} \frac{N_i}{N} \text{Var}_i$$

Coeficiente de Correlação Intraclass

Exemplo - continuação: Considerando a divisão A:

Para conglomerados de tamanhos diferentes.

Exemplo -

Considere a população de tamanho $N = 6$ agrupada em $M = 3$ conglomerados da seguinte forma: $Y = ((12); (7, 9, 14); (8, 10))$

a) Veja neste exemplo que $\frac{\gamma^2}{\sigma^2}$.

b) Calcule o CCI pela definição e pelo método proposto do Bolfarine e Bussab.

Para casa

- Fazer a lista 2 de exercícios.
- Ler o capítulo 2 da apostila da Profa. Vanessa.
- Rever os slides.

Próxima aula

- Acompanhar o material no moodle.

Amostragem por Conglomerados

- Tamanho de amostra e Intervalos de confiança

Muito obrigado!



Fonte: imagem do livro *Combined Survey Sampling Inference: Weighing of Basu's Elephants*.

Referências

- Amostragem: Teoria e Prática Usando o R
- **Elementos de Amostragem**, Bolfarine e Bussab.
- Cochran(1977)

Resumo da notação