

Estatística Numérica Computacional
Trabalho nº2
Grupo V

Marta Paz nº49861
Rafael Almeida nº49788
Rafael Gameiro nº50677
Ricardo Pinto nº49811

October, 2018

Exercício 1

Considere as duas amostras de duas v. aleatórias X e Y

Amostras de X	92	90	85	96	92	88	96	88
Amostras de Y	89	90	88	93	90	85	95	90

- (a) Use $\hat{\rho} = \frac{\sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^N (X_i - \bar{X})^2 \sum_{i=1}^N (Y_i - \bar{Y})^2}}$ para estimar o coeficiente de correlação entre as duas variáveis.

Resolução

A resolução desta alínea foi feita somente em R.

```
1  #####
2
3
4  #EXERCICIO 1
5
6  x <- c(92,90,85,96,92,88,96,88)
7  y <- c(89,90,88,93,90,85,95,90)
8
9  alpha <- 0.025
10 n <- 999
11
12
13 #EXERCICIO 1 - ALINEA A
14
15 # calculo de Pearson Coefficient
16 func <- function(x,y){
17   sum((x-mean(x))*(y-mean(y)))/sqrt((sum((x-mean(x))^2))*(sum((y-mean(y))^2)))
18 }
19
20 pCoefficient = func(x,y)
21
```

R: $\hat{\rho} \approx 0.798$

- (b) Calcule a variância de jackknife de $\hat{\rho}$.

Resolução

Para calcular a variância de jackknife temos a seguinte fórmula:

$$Var_{jack} = \frac{1}{n(n-1)} \left(\sum_{j=1}^N l_{jack}^2 - nb_{jack}^2 \right)$$

Para tal, também precisamos de calcular o enviesamento.

Definição

Suponhamos que $\rho = t(F)$ e que F é uma função de distribuição cumulativa e $t(\cdot)$ algo funcional. Então, a função de influência de t em F é dada por:

$$L_t(y, F) = \lim_{\varepsilon \rightarrow 0} \frac{t[(1 - \varepsilon)F + \varepsilon H_y] - t(F)}{\varepsilon}$$

e

$$H_y(u) = \begin{cases} 0, & u < y \\ 1, & u \geq y \end{cases}$$

Suponhamos que x_1, x_2, \dots, x_n é uma amostra e \hat{F} é uma função empírica da amostra. A função de influência empírica é definida como:

$$l(y) = L_t(y, \hat{F})$$

Os valores da função de influência empírica nos pontos de dados são chamadas de valores de influência empírica.

$$l_j = l(x_j) = L_t(x_j, \hat{F}), \quad j = 1, 2, \dots, n$$

Uma extensão do Teorema de Taylor diz que para funções t e medidas G e F ,

$$t(G) \approx t(F) + \int L_t(y, F) dG(y)$$

Isto é um resultado exato se t é uma estatística linear. Aplicando essa fórmula com F à função de distribuição cumulativa e à função $G = \hat{F}$ obtemos

$$\begin{aligned} t(\hat{F}) &\approx t(F) + \int L_t(y, F) d\hat{F}(y) \quad , \text{que é} \\ t(\hat{F}) &\approx t(F) + \frac{1}{n} \sum_{j=1}^N L_t(x_j, \hat{F}) = t(F) + \frac{1}{n} \sum_{j=1}^N l_j \\ &\Leftrightarrow \rho - \hat{\rho} = -\frac{1}{n} \sum_{j=1}^N l_j \end{aligned}$$

O Jackknife fornece uma maneira de aproximar valores de influência empíricos através da reamostragem dos dados.

$$L_t(y, F) = \lim_{\varepsilon \rightarrow 0} \frac{t \left[(1 - \varepsilon) \hat{F} + \varepsilon H_y \right] - t(\hat{F})}{\varepsilon}$$

$$\approx \frac{t \left[(1 - \varepsilon) \hat{F} + \varepsilon H_y \right] - t(\hat{F})}{\varepsilon}$$

Se tomarmos $\varepsilon = \frac{1}{n-1}$, então

$$(1 - \varepsilon) \hat{F} + \varepsilon H_{x_j} = \frac{n}{n(n-1)} \hat{F} - \frac{1}{n-1} H_{x_j} = \hat{F}_{-j}$$

é uma distribuição sem peso no ponto x_j e peso $\frac{1}{n-1}$ no resto da amostra. Isto é equivalente a ter apenas a amostra de tamanho $n-1$ encontrada omitindo x_j da amostra original. Então, a aproximação Jackknife ao valor da influência empírico l_j é

$$l_{jack:j} = (n-1)[t(\hat{F}) - t(\hat{F}_{-j})] = (n-1)(\hat{\rho} - \rho_{-j})$$

As estimativas imparciais do enviesamento e da variância, que usam os valores de influência empíricos do Jackknife a que chamamos de viés de Jackknife e variância de Jackknife, são:

$$b_{jack} = -\frac{1}{n} \sum_{j=1}^N l_{jack:j}$$

$$Var_{jack} = \frac{1}{n(n-1)} \left(\sum_{j=1}^N l_{jack:j}^2 - n b_{jack}^2 \right)$$

```

25 #*****
26
27
28 #EXERCICIO 1 - ALINEA B
29
30 # Criacao das matrizes xAux e yAux
31 # Ciclo for que preenche as matrizes xAux e yAux com os valores "teta_chapeu exceto o elemento na posicao j"
32 xAux <- matrix(nrow=8, ncol = 7)
33 yAux <- matrix(nrow=8, ncol = 7)
34
35 for(l in 1:8){
36   tempX=x[-l] # Remocao do elemento l do vetor x, com posterior criacao do vetor tempX ja sem esse elemento
37   tempY=y[-l] # Remocao do elemento l do vetor y, com posterior criacao do vetor tempY ja sem esse elemento
38
39   for(r in 1:7){
40     xAux[row=l,col=r]=tempX[r]
41     yAux[row=l,col=r]=tempY[r]
42   }
43 }
44
45 # Calculo do ljack_j: (n-1) * (teta_chapeu - teta_chapeu_excetoj)
46 # Seja n a dimensao das amostras, teta_chapeu o valor de pcoefficient calculado na alinea a)
47 # Consequente calculo do somatorio do ljack de todas as posicoes
48 ljack = 0
49 somaBjack=0
50
51 for(i in 1:8){
52   ljack = 7 * (pcoefficient - func(xAux[i, ], yAux[i, ]))
53   somaBjack = somaBjack + ljack
54 }
55
56 # calculo do bjack: -(1/n) * somatorio(ljack_j)
57 bjack = -(1/8) * somaBjack
58
59 # calculo do varjack: (1/(n(n-1))) * (somatorio(ljack^2 - n*bjack^2))
60 somaVarJack=0
61
62 for (i in 1:8){
63   ljack = 7 * (pcoefficient - func(xAux[i, ], yAux[i, ]))
64   somaVarJack = somaVarJack + (ljack^2)
65 }
66
67 varJack = (1/(8*(8-1))) * (somaVarJack-(8*(bjack^2)))
68
69
70
71

```

R: $Var_{jack} \approx 0.014$

(c) Construa um intervalo de confiança bootstrap básico para ρ .

Resolução

Para calcular um intervalo de confiança de Bootstrap tivemos de gerar várias amostras bootstrap, tanto para X como para Y . Para tal, aplicámos a cada amostra a fórmula do coeficiente de correlação, e de seguida o enviesamento a cada resultado obtido. Ordenámos os valores e utilizando a estatística pivot,

$$\hat{a}_\alpha = \hat{\rho}_{(R+1).\alpha}^* - \hat{\rho}$$

$$\hat{a}_{1-\alpha} = \hat{\rho}_{(R+1).(1-\alpha)}^* - \hat{\rho}$$

obtivemos o intervalo de confiança na forma,

$$]\hat{\rho} - \hat{a}_{1-\alpha}, \hat{\rho} - \hat{a}_\alpha[$$

```

73 #*****
74
75
76 #EXERCICIO 1 - ALINEA C
77
78 # Geracao de amostras Bootstrap
79 amBootsX <- matrix(NA, nrow=n, ncol = 8)
80 amBootsY <- matrix(NA, nrow=n, ncol = 8)
81
82 for(r in 1:n) {
83   amBootsX[r,] <- sample(x, replace = TRUE)
84   amBootsY[r,] <- sample(y, replace = TRUE)
85 }
86
87
88 # Geracao de amostras auxiliares
89 amostra_Vies <- 1:n
90 amostra_Ro <- 1:n
91
92 for(r in 1:n) {
93   amostra_Ro[r] <- func(amBootsX[r,], amBootsY[r,])
94   amostra_Vies[r] <- amostra_Ro[r] - pcoefficient
95 }
96
97
98 # Calculo do Intervalo de Confianca
99 # Limite Inferior: teta_chapeu - a_1menosAlpha
100 # Limite Superior: teta_chapeu + a_alpha
101 # Seja o alpha = 0.025
102
103 # Passo 1: Gerar amostra usada para calcular o Intervalo de Confianca
104 # Seja amIC a amostra do Intervalo de Confianca
105 amIC <- 1:n+1
106 for(r in 1:r) {
107   amIC[r] <- amostra_Vies[r]
108 }
109
110 amIC[n+1] <- 0
111 amIC <- sort(amIC)
112
113 # Passo 2: Calculo do a_alpha e do a_1menosAlpha
114 a_alpha <- amIC[(n+1)*alpha]
115 a_1menosAlpha <- amIC[(n+1)*(1-alpha)]
116
117 # Passo 3: Estabelecer limites inferior e superior
118 limite_inferior_Boots <- pcoefficient-a_1menosAlpha
119 limite_superior_Boots <- pcoefficient+a_alpha
120

```

R:]0.86, 2.33[

- (d) Construa um intervalo de confiança t-bootstrap para ρ . Use a variância bootstrap para estimar a variância de $\hat{\rho}_r^*$.

Resolução

Para o cálculo do intervalo de confiança bootstrap studentized utilizámos uma estatística pivot diferente da usada no cálculo do bootstrap básico:

$$\hat{z} = \frac{\hat{\rho}^* - \hat{\rho}}{\sqrt{Var_{Boots}(\hat{\rho}^*)}}$$

$$Var_{Boots}(\hat{\rho}^*) = \frac{1}{n-1} \sum_{j=1}^N (\hat{\rho}_j^* - \bar{\hat{\rho}}^*)^2$$

Primeiramente, usámos as amostras bootstrap geradas da alínea c) para o cálculo dos $\hat{\rho}_r^*$. Para obtermos a variância bootstrap de cada $\hat{\rho}_r^*$, gerámos um novo conjunto de amostras bootstrap com base em cada amostra previamente usada na alínea anterior. O resto do procedimento, foi a reordenação dos diferentes valores de z , e a determinação do intervalo de confiança, usando a fórmula

$$]\hat{\rho} - \hat{z}_{(R+1)(1-\alpha)}^* \sqrt{Var_{Boots}(\hat{\rho})}, \hat{\rho} - \hat{z}_{(R+1)(\alpha)}^* \sqrt{Var_{Boots}(\hat{\rho})}[$$

```

126 #EXERCICIO 1 - ALINEA D
127
128 # Cálculo do theta da amostra
129 amRho <- 1:n
130
131
132 for(r in 1:n){
133   amRho[r] <- func(amBootsX[r,], amBootsV[r,])
134 }
135
136 meanAmRho <- mean(amRho, na.rm = TRUE)
137
138
139 # Cálculo do varBoots: (1/(r-1)) * (somatorio( (teta_chapeu_estrela_R - media_teta_chapeu_estrela)^2) )
140 # Seja r a dimensão n da amostra bootstrap
141 varboots <- (1/(n-1))*sum((amRho - meanAmRho)^2, na.rm = TRUE)
142
143
144 # Criação de uma matriz com os elementos das amostras das variâncias
145 # Cálculo de todos os z's e suas variâncias
146 ryr <- 1:8
147 rxr <- 1:8
148 z <- 1:(n+1)
149 amRhoStar <- 1:n
150
151 for(r in 1:n){
152   for(i in 1:n){
153     rxr <- sample(amBootsX[r, ], replace = TRUE)
154     ryr <- sample(amBootsV[r, ], replace = TRUE)
155
156     amRhoStar[i] <- func(rxr, ryr)
157   }
158
159   meanAmRhoStar <- mean(amRhoStar, na.rm = TRUE)
160
161   # Cálculo da variância para o z atual
162   varbootsStar <- (1/(n-1))*sum((amRhoStar - meanAmRhoStar)^2, na.rm = TRUE)
163   z[r] <- (amRho[r] - pCoefficient)/sqrt(varbootsStar)
164 }
165
166
167 z[n+1] <- 0
168 amIC <- 1:(n+1)
169
170
171 # Cálculo do Intervalo de Confiança
172 # Limite Inferior: teta_chapeu - z_estrela_((R+1)(1-alpha)) * raiz(variancia)
173 # Limite Superior: teta_chapeu - z_estrela_((R+1)*alpha) * raiz(variancia)
174 amIC <- sort(z)
175
176 limite_inferior_tBoots <- pCoefficient - (amIC[(n+1)*(1-alpha)]*sqrt(varboots))
177 limite_superior_tBoots <- pCoefficient - (amIC[(n+1)*alpha]*sqrt(varboots))

```

R:]0.89, 2.27[

Discussão

Comparativamente com o método Bootstrap, o método Jackknife é mais rápido a executar e permite calcular estatísticas não-paramétricas. Por outro lado, uma vez que o intervalo de confiança baseia-se em aproximações, o melhor método a ser usado é o de Bootstrap.

O intervalo de confiança bootstrap studentized tem menor amplitude que o intervalo de confiança bootstrap básico. Uma vez que quanto menor for a amplitude de um intervalo, maior será a sua precisão, então podemos concluir que o intervalo de confiança bootstrap studentized é mais preciso que o de bootstrap básico. O intervalo de confiança bootstrap básico não se encontra como seria esperado, pois o valor ρ calculado anteriormente não pertence a esse intervalo. O intervalo de confiança bootstrap studentized também não está como era esperado, pois o valor ρ calculado anteriormente não pertence a esse intervalo.

Exercício 2

Considere a amostra seguinte:

1	3	3	1	1	3	2	2	3	0	2	4	2	6	4	2	4	3	2	4
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Considere a hipótese nula, H_0 : A v.a. subjacente a esta amostra tem distribuição binomial de parâmetros $(8, 0.3)$. Use a função de distribuição empírica da estatística de teste para estimar o p -value da estatística de teste qui-quadrado, T que permite testar a hipótese nula.

$$T = \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i}$$

$i \in \{1, \dots, K\}$ - suporte da v.a.

$\{p_i, i = 1, \dots, K\}$ - função de probabilidade da v.a.

n - dimensão da amostra.

N_i - número de observações na amostra que tomam o valor i .

Resolução

Neste exercício, foi-nos dada uma amostra e a seguinte hipótese nula:

H_0 : A variável aleatória subjacente a esta amostra tem distribuição binomial de parâmetros $(8, 0.3)$.

De seguida, era pedido que usássemos a função de distribuição empírica da estatística de teste para estimar o p -value da estatística de teste qui-quadrado, T , que permite testar a hipótese nula fornecida. Desta forma, para o teste de hipótese que realizámos, considerámos as seguintes hipóteses:

H_0 : A variável aleatória subjacente a esta amostra tem distribuição binomial de parâmetros $(8, 0.3)$.

H_1 : A variável aleatória subjacente a esta amostra não tem distribuição binomial de parâmetros $(8, 0.3)$

Para calcular o valor do p -value, utilizámos a seguinte fórmula:

$$p - value = \hat{P}(T > t_{obs} | H_0) = 1 - \hat{F}_{H_0}(t_{obs}) = \frac{\#\{t_j : t_j \geq t_{obs}\} + 1}{m + 1}$$

Para podermos utilizar a fórmula acima, tivemos de calcular os seguintes valores:

1. t_{obs} , que corresponde ao valor da estatística de teste dos valores observados (valores estes que foram dados no enunciado). Recorremos à seguinte fórmula para o cálculo da estatística de teste (que foi dada no enunciado do problema):

$$T = \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i}$$

E realizámos os seguintes passos no R:

- (a) Percorremos a amostra dada no enunciado e contamos o número de ocorrências de cada elemento (0 a 6, pois a amostra só tem elementos compreendidos entre esses valores). Para guardar o número de ocorrências usámos um vetor, `nObs`, em que cada posição correspondia a um valor do intervalo entre 0 e 6.
- (b) Fizemos um ciclo onde é calculado a probabilidade binomial, com os parâmetros fornecidos, de cada ocorrência.
- (c) Por último, aplicámos a fórmula da estatística de teste acima ao conjunto das probabilidades calculadas no ponto anterior e daí obtivemos que $t_{obs} = 2.205978$.

2. t_j , que corresponde ao valor da estatística de teste assumindo que a hipótese nula é verdade. Para tal, criamos uma matriz, `amAux`, onde cada linha vai corresponder a uma amostra gerada com o comando `rbinom` do R usando os parâmetros da distribuição binomial que nos deram no enunciado.

De seguida criamos uma matriz, `nObsAux`, onde guardamos o número de ocorrências decada, em cada amostra gerada. Depois, calculamos a probabilidade de cada posição de `nObsAux` acontecer e, com os valores obtidos, pudemos calcular os t_j 's. Para cada t_j , comparamos o seu valor com t_{obs} e contamos quantas vezes t_j era igual ou superior que t_{obs} . Esse valor foi guardado na nossa variável `soma`.

3. Por último, aplicamos a fórmula do $p - value$ aos valores que obtivemos, ou seja:

$$p - value = \frac{soma + 1}{(999 + 1) + 1} \approx 0.86$$

Para verificar se aceitamos a hipótese nula, consideramos três valores diferentes para o nível de significância: $\alpha = 0.01$, $\alpha = 0.025$ e $\alpha = 0.5$.

Para os três valores de alfa considerados, verificamos que $p - value \geq \alpha$, logo não rejeitamos a hipótese nula nos três níveis de significância. Isto é, existem evidências estatísticas para afirmar que a variável aleatória subjacente à amostra fornecida tem distribuição binomial de parâmetros $(8, 0.3)$.

```

184 #EXERCICIO 2
185
186 amostra <- c(1,3,3,1,1,3,2,2,3,0,2,4,2,6,4,2,4,3,2,4)
187 nobs <- c(1,3,6,5,4,1) # Numero de vezes que os valores 0, 1, 2, 3, 4 e 6 se repetem, respetivamente
188 k <- 6
189
190
191 # Funcao Auxiliar: Funcao Fatorial
192 fact <- function(x) {
193   if(x == 0) {
194     1
195   } else {
196     x * fact(x-1)
197   }
198 }
199
200
201 # calculo da Probabilidade da Distribuicao Binomial: probBinomial = (n!/(k!(n-k)!)) * (p^k) * ((1-p)^(n-k))
202 # Seja B(n,p) = B(8,0.3), temos n=8 e p=0.3.
203 # Seja k cada valor de nobs
204 probBinomial <- 1:6
205 for(i in 1:k) {
206   probBinomial[i] <- (fact(8)/(fact(nobs[i])*fact(8-nobs[i]))) * (0.3^nobs[i]) * ((0.7)^(8-nobs[i]))
207 }
208
209
210 # calculo do T, usando a funcao de probabilidade dada: estaTest = somatorio( ((N_i-n*p_i)^2) / (n*p_i) )
211 # Seja N_i cada valor de nobs, n a dimensao da amostra e p_i cada valor da funcao de probabilidade
212 estaTest <- 0
213 for(i in 1:k) {
214   estaTest <- sum( ((nobs[i]-20*probBinomial[i])^2) / (20*probBinomial[i]) )
215 }
216
217
218
219 # Calculo do p-value: (#{tj : tj>=tobs} + 1) / (m+1)
220 # Seja tj cada valor da amostra inicial, tobs a estatistica de teste e m a dimensao de nobs
221 soma=0
222 probaux <- 1:p
223 amaux <- matrix(nrow=(n+1), ncol=p)
224
225 #geracao de 1000 amostra com probabilidade binomial (8,0.3)
226 for(i in 1:(n+1)) {
227   probaux <- rbinom(8,8,0.3)
228   for(r in 1:p){
229     amaux[i,r] <- probaux[r]
230   }
231 }
232
233
234 #calcular o numero de ocorrencias de cada valor na amostra gerada
235 nobsAux <- matrix(nrow=(n+1), ncol=9)
236 for(r in 1:(n+1)){
237   for(c in 1:9){
238     nobsAux[r,c] <- 0
239   }
240 }
241
242

```

```

243 ~ for(j in 1:(n+1)){
244 ~   for(i in 1:p){
245 ~     if(amaux[j,i] == 0){
246 ~       nobsaux[j, 1] = nobsaux[j,1]+1
247 ~     }
248 ~     else if(amaux[j,i] == 1){
249 ~       nobsaux[j, 2] = nobsaux[j,2]+1
250 ~     }
251 ~     else if(amaux[j,i] == 2){
252 ~       nobsaux[j, 3] = nobsaux[j,3]+1
253 ~     }
254 ~     else if(amaux[j,i] == 3){
255 ~       nobsaux[j, 4] = nobsaux[j,4]+1
256 ~     }
257 ~     else if(amaux[j,i] == 4){
258 ~       nobsaux[j, 5] = nobsaux[j,5]+1
259 ~     }
260 ~     else if(amaux[j,i] == 5){
261 ~       nobsaux[j, 6] = nobsaux[j,6]+1
262 ~     }
263 ~     else if(amaux[j,i] == 6){
264 ~       nobsaux[j, 7] = nobsaux[j,7]+1
265 ~     }
266 ~     else if(amaux[j,i] == 7){
267 ~       nobsaux[j, 8] = nobsaux[j,8]+1
268 ~     }
269 ~     else if(amaux[j,i] == 8){
270 ~       nobsaux[j, 9] = nobsaux[j,9]+1
271 ~     }
272 ~   }
273 ~ }
274 ~
275 ~

```

```

276 ~ estaTestAux <- 0
277 ~ probBinomialAux <- 1:9
278 ~
279 ~ #faz o calculo de soma = #{tj : tj>=tobs}
280 ~ for(i in 1:(n+1)){
281 ~   #calcula a probabilidade da ocorrencia de cada elemento gerado, assumindo que H0 e verdade
282 ~   for(j in 1:9){
283 ~     if(nobsaux[i,j] != 0){
284 ~       probBinomialAux[j] <- (fact(8)/(fact(nobsaux[i,j])*fact(8-nobsaux[i,j])))^(0.3^nobsaux[i,j])*((0.7)^(8-nobsaux[i,j]))
285 ~     }else{
286 ~       probBinomialAux[j] <- 0
287 ~     }
288 ~   }
289 ~ }
290 ~
291 ~ for(j in 1:9){
292 ~   if(probBinomialAux[j] == 0){
293 ~     estaTestAux <- estaTestAux + 0
294 ~   }else{
295 ~     estaTestAux <- sum(((nobsaux[i,j]-20*probBinomialAux[j])^2)/(20*probBinomialAux[j]))
296 ~   }
297 ~ }
298 ~
299 ~ if(estaTestAux >= estaTestObs){
300 ~   soma = soma+1
301 ~ }
302 ~
303 ~ estaTestAux <- 0
304 ~ }
305 ~
306 ~
307 ~ # calcula o pvalue = soma/m+1, com m = numero de elementos gerados
308 ~ pvalue <- (soma+1)/(n+2)
309 ~
310 ~

```