

Amostragem aleatória simples sem reposição (parte 2)

Prof. Caio Azevedo

Estimação da proporção populacional

- População: observações univariadas - y_1, \dots, y_N (variáveis não aleatórias), em que y_i é a observação relativa ao indivíduo i (podemos também considerar observações multivariadas).
- Temos que $y_i = 1$ se o indivíduo i possui a característica de interesse e 0 caso contrário.

Estimação da proporção

- Exemplos: presença de alguma doença, procedência (1 se é oriundo de determinado lugar, 0, caso contrário), inadimplência (1 se inadimplente, 0 caso contrário).
- Os procedimentos definidos anteriormente, em princípio, se mantêm (slides AAS_s , parte 1, [link](#)). A principal diferença, de forma geral, reside na estrutura da variável de interesse.
- Parâmetro de interesse: $p = \frac{1}{N} \sum_{i=1}^N y_i$.
- Lembremos que $y_i = y_i^k, \forall k \in \mathbb{R}^+$.

Estimação da proporção

- Estimador “natural”:

$$\begin{aligned}\hat{p} &= \frac{1}{n} \sum_{i=1}^n Y_i = \frac{1}{n} \sum_{i \in \mathcal{S}} y_i \\ &= \frac{1}{n} \sum_{i=1}^N F_i y_i\end{aligned}$$

- Note que, neste caso, a variância populacional $\sigma^2 = \frac{1}{N} \sum_{i=1}^N (y_i - p)^2$, toma a seguinte forma:

$$\begin{aligned}\sigma^2 &= \frac{1}{N} \sum_{i=1}^N (y_i^2 - 2y_i p + p^2) = \frac{1}{N} (Np - Np^2) \\ &= p(1 - p) = pq, q = 1 - p\end{aligned}$$

- Consequentemente, $s^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - p)^2 = \frac{N}{N-1} \sigma^2 = \frac{N}{N-1} pq$

Propriedades do estimador

- Note que, essencialmente, \hat{p} é uma média amostral (de variáveis binárias), semelhante à $\hat{\mu}$ ([aqui](#)).
- Portanto, as propriedades de \hat{p} são semelhantes as de $\hat{\mu}$ (lembrando que $f = \frac{n}{N}$), sob AAS_s , por exemplo:
 - $\mathcal{E}_{A_2}(\hat{p}) = \mathcal{E}_{A_2}(\hat{\mu}) = \frac{1}{n} \sum_{i=1}^N y_i \mathcal{E}(F_i) = \frac{1}{n} \sum_{i=1}^N y_i \frac{n}{N} = \frac{1}{N} \sum_{i=1}^N y_i = p$.
 - $\mathcal{V}_{A_2}(\hat{p}) = \mathcal{V}_{A_2}(\hat{\mu}) = (1-f) \frac{s^2}{n} = (1-f) \frac{N}{N-1} \frac{pq}{n} = \frac{N-n}{N-1} \frac{pq}{n}$.
- Estimativa: $\tilde{p} = \frac{1}{n} \sum_{i \in s} y_i = \frac{1}{n} \sum_{i=1}^N f_i y_i$

Propriedades do estimador

- Vimos também que um estimador não viciado para a variância populacional ($s^2 = \frac{N}{N-1}pq$) é dado por

$$\begin{aligned}\hat{s}^2 &= \frac{1}{n-1} \sum_{i=1}^n (Y_i - \hat{\mu})^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \hat{p})^2 \\ &= \frac{1}{n-1} \sum_{i=1}^N F_i (y_i - \hat{p})^2\end{aligned}$$

- Note, no entanto, que neste caso

$$\begin{aligned}\hat{s}^2 &= \frac{1}{n-1} \sum_{i=1}^n (Y_i^2 - 2Y_i\hat{p} + \hat{p}^2) = \frac{1}{n-1} (n\hat{p} - n\hat{p}^2) \\ &= \frac{n}{n-1} \hat{p}\hat{q}, \hat{q} = 1 - \hat{p}\end{aligned}$$

Propriedades do estimador

- Consequentemente, um estimador não viciado para a variância do estimador é dado por:

$$(1 - f) \frac{\widehat{S}^2}{n} = (1 - f) \frac{n\widehat{p}\widehat{q}}{n(n-1)} = (1 - f) \frac{\widehat{p}\widehat{q}}{n-1}$$

- Analogamente ao caso da média, temos que

$$\frac{\widehat{p} - p}{\sqrt{\frac{N-n}{N-1} \frac{pq}{n}}} \xrightarrow[n \rightarrow \infty, N-n \rightarrow \infty]{D} N(0, 1); \quad \frac{\widehat{p} - p}{\sqrt{(1-f) \frac{\widehat{p}\widehat{q}}{(n-1)}}} \xrightarrow[n \rightarrow \infty, N-n \rightarrow \infty]{D} N(0, 1)$$

Comparação dos estimadores sob os planos A_1 e A_2

- O estimador para a proporção sob AAS_c ou AAS_s , é o mesmo.
- Temos que $\mathcal{E}_{A_i}(\hat{p}) = p$, $i = 1, 2$.
- Portanto, o efeito do planejamento (EPA), do estimador sob o plano A_2 em relação ao plano A_1 , é dado por:

$$EPA = \frac{\mathcal{V}_{A_2}(p)}{\mathcal{V}_{A_1}(p)} = \frac{(1-f) \frac{N}{N-1} \frac{pq}{n}}{\frac{pq}{n}} = (1-f) \frac{N}{N-1}$$

- Assim, quando $N \rightarrow \infty$, $\mathcal{V}_{A_2}(\hat{p}) \rightarrow \mathcal{V}_{A_1}(\hat{p})$ (semelhante ao que ocorre com $\hat{\mu}$).

Comparação dos estimadores sob os planos A_1 e A_2

- Consequentemente, temos que o plano AAS_s é melhor do que AAS_c , tendendo ambos a serem equivalentes, à medida que o tamanho da população tende a infinito.
- Exercício: Calcular o EPA para $\hat{\tau}$.
- Defina um estimador para o EPA, da seguinte forma: $\widehat{EPA} = \frac{\widehat{V}_{A_2}(\hat{\theta})}{\widehat{V}_{A_1}(\hat{\theta})}$, em que $\hat{\theta}$ é um estimador escolhido, para o parâmetro de interesse (θ) e $\widehat{V}_A(\hat{\theta})$ é o estimador da variância do estimador de $\hat{\theta}$. Exercício: calcular o \widehat{EPA} para os estimadores $\hat{\mu}$, $\hat{\tau}$ e \hat{p} .
- OBS: Se for necessário, é possível ser mais específico com relação à notação do EPA (\widehat{EPA}), ou seja, escrevendo $EPA(\theta)(\widehat{EPA}(\theta))$.

Intervalo de Confiança

- Erro-padrão do estimador: $EP_{A_2}(\hat{p}) = \sqrt{\mathcal{V}_{A_2}(\hat{p})}$.
- Um estimador da variância do estimador: $\hat{\mathcal{V}}_{A_2}(\hat{p}) = (1 - f) \frac{\hat{p}\hat{q}}{n-1}$.
- Um estimador do erro-padrão do estimador: $\widehat{EP}_{A_2}(\hat{p}) = \sqrt{\hat{\mathcal{V}}_{A_2}(\hat{p})}$.

Intervalo de Confiança

- Assim, dois intervalos de confiança (assintóticos) com coeficiente de confiança de aproximadamente γ , são dados por

$$IC(\mu, \gamma) \approx \left[\hat{p} - z_\gamma \sqrt{(1-f) \frac{\hat{p}\hat{q}}{n-1}}; \hat{p} + z_\gamma \sqrt{(1-f) \frac{\hat{p}\hat{q}}{n-1}} \right] \quad (1)$$

$$IC(\mu, \gamma) \approx \left[\hat{p} - z_\gamma \sqrt{\frac{1-f}{4(n-1)}}; \hat{p} + z_\gamma \sqrt{\frac{1-f}{4(n-1)}} \right] \quad (2)$$

em que $P(Z \leq z_\gamma) = \frac{1+\gamma}{2}$ e $Z \sim N(0, 1)$.

- Erro da estimativa: $z_\gamma \sqrt{(1-f) \frac{\hat{p}\hat{q}}{n-1}}$ ou $z_\gamma \sqrt{\frac{1-f}{4(n-1)}}$.
- O comprimento do intervalo (2) sempre será maior (ou igual) ao comprimento do intervalo (1).

Testes de Hipótese

- Hipóteses usuais (p_0 conhecido, $q_0 = 1 - p_0$)

1 $H_0 : p = p_0$ vs $H_1 : p < p_0$.

2 $H_0 : p = p_0$ vs $H_0 : p > p_0$.

3 $H_0 : p = p_0$ vs $H_0 : p \neq p_0$.

- Estatística do teste $Z_t = \frac{\hat{p} - p_0}{\sqrt{[(N-n)/(N-1)]p_0q_0/n}}$.

- Sob H_0 , vimos que $Z_t \approx N(0, 1)$, para n e $N-n$ suficientemente grandes.

- Defina $z_t = \frac{\tilde{p} - p_0}{\sqrt{[(N-n)/(N-1)]p_0q_0/n}}$ o valor calculado da estatística do teste e z_c o(s) valor(es) crítico(s).

- Defina ainda $Z \sim N(0, 1)$. Os procedimentos são análogos ao caso da média, com as devidas adaptações.

Determinação do tamanho amostral: erro da estimativa

Analogamente ao caso da média populacional, temos que

$$\delta = z_\gamma \sqrt{(1-f) \frac{pq}{n}} \rightarrow n = \frac{1}{\frac{\delta^2}{z_\gamma^2 pq} + \frac{1}{N}} \quad (3)$$

Podemos usar estimativas de p obtidas em pesquisas anteriores, sob uma amostra piloto ou, considerar o pior caso, em termos da variabilidade dos dados. Neste último caso, temos que:

$$n = \frac{1}{\frac{4\delta^2}{z_\gamma^2} + \frac{1}{N}} \quad (4)$$

Isto vale para qualquer um dos dois critérios: erro da estimativa e precisão. Note que o tamanho da amostra fornecido por (4) será maior ou igual àquele fornecido por (3).

Estudos de simulação

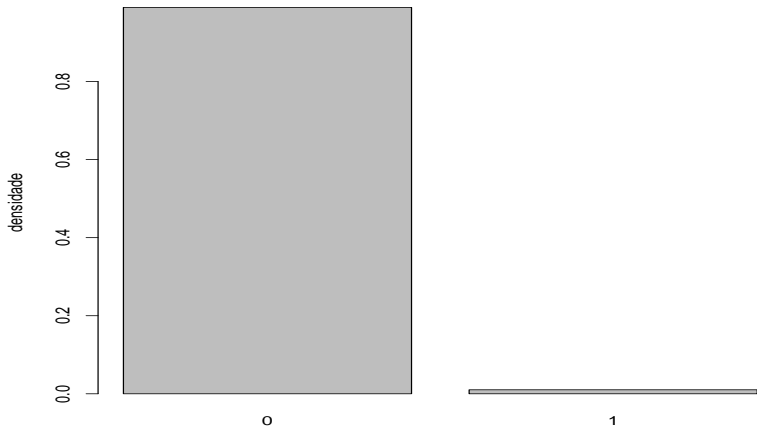
- Distribuição assintótica do estimador para a proporção. Tamanho da população $N = 100.000$.
- Vários cenários, variando em função do valor verdadeiro da proporção populacional p .
- $p =$
 $(0,01; 0,05; 0,10; 0,25; 0,35; 0,50; 0,65; 0,75; 0,9; 0,95; 0,99)^T$.
- A distribuição, (em princípio), da variável de interesse é Bernoulli(p).

Estudos de simulação

- Quatro tamanhos amostrais (30, 50, 100, 1000), em termos percentuais, com relação ao tamanho da população (0,03%,0,05%,0,1%,1%).
- Estudar a distribuição amostral (empírica) com base em $R = 1.000$ réplicas (amostras selecionadas da população de interesse).

$p=0,01$

Gráfico de Colunas

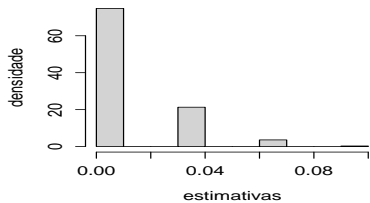


variavel de interesse

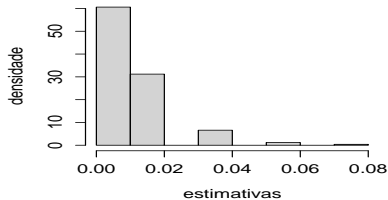


$p=0,01$

$n = 30$, p -valor (teste-SW) = 0

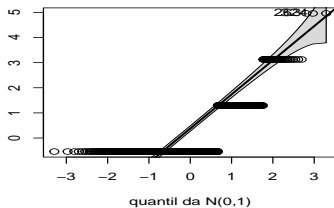


$n = 50$, p -valor (teste-SW) = 0



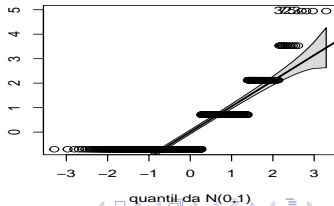
quantil da distribuição padronizada do estimador

$n = 30$, p -valor (teste-SW) = 0



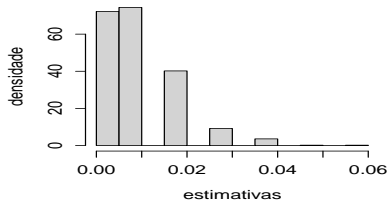
quantil da distribuição padronizada do estimador

$n = 50$, p -valor (teste-SW) = 0



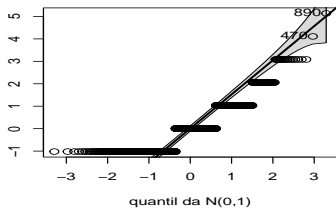
$p=0,01$

$n = 100$, p -valor (teste-SW) = 0

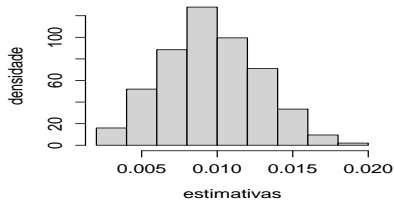


$n = 100$, p -valor (teste-SW) = 0

quantil da distribuição padronizada do estimador

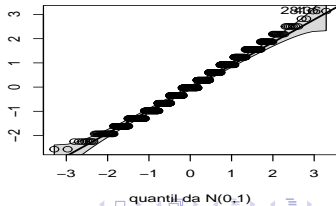


$n = 1000$, p -valor (teste-SW) = 0



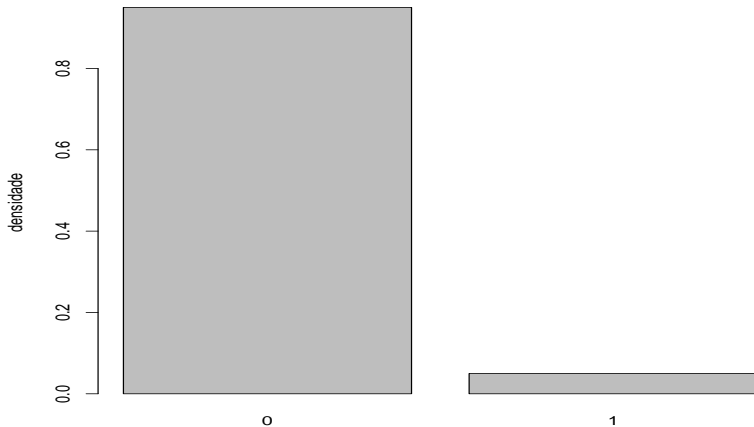
$n = 1000$, p -valor (teste-SW) = 0

quantil da distribuição padronizada do estimador



$p=0,05$

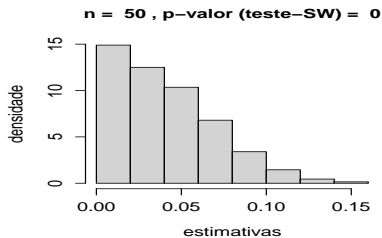
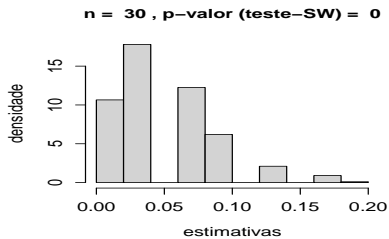
Gráfico de Colunas



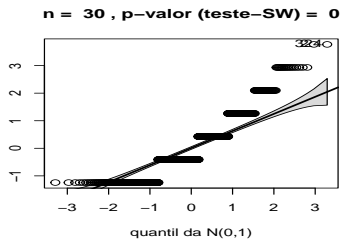
variável de interesse



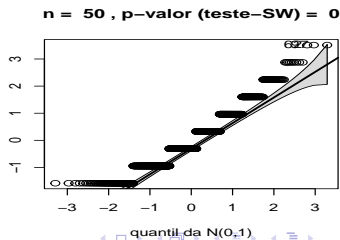
$p=0,05$



quantil da distribuição padronizada do estimador

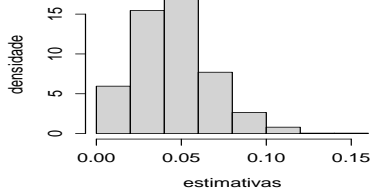


quantil da distribuição padronizada do estimador



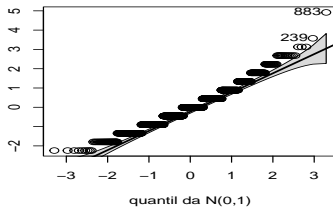
$p=0,05$

$n = 100$, p -valor (teste-SW) = 0

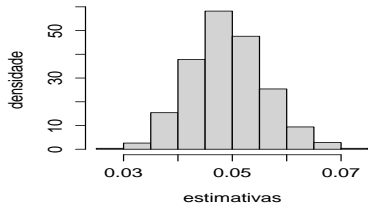


quantil da distribuição padronizada do estimador

$n = 100$, p -valor (teste-SW) = 0

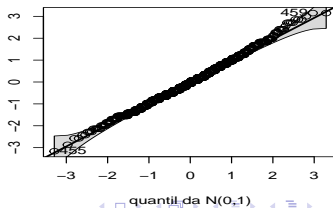


$n = 1000$, p -valor (teste-SW) = 0.0025



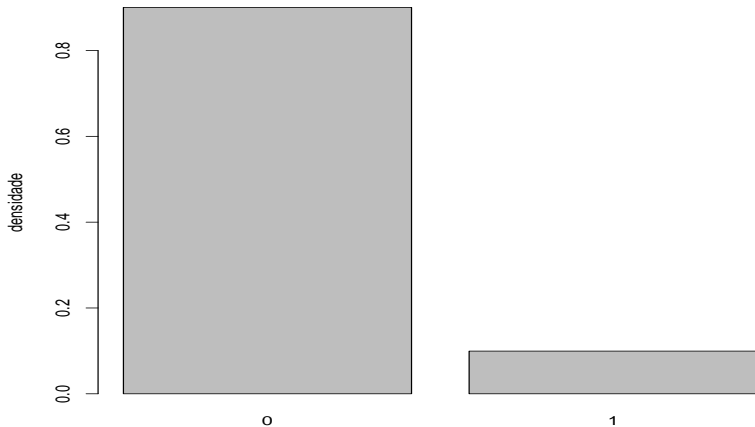
quantil da distribuição padronizada do estimador

$n = 1000$, p -valor (teste-SW) = 0.0025



$p=0,10$

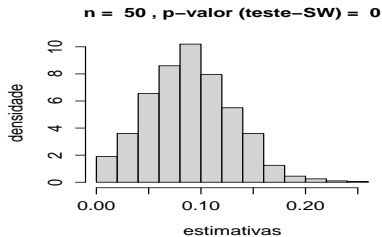
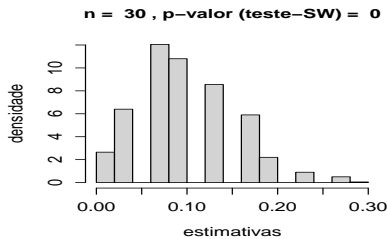
Gráfico de Colunas



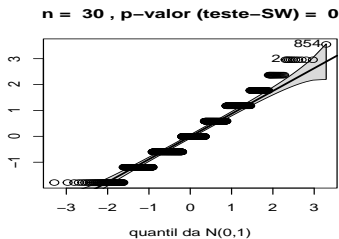
variável de interesse



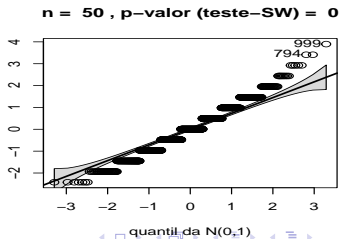
$p=0,10$



quantil da distribuição padronizada do estimador

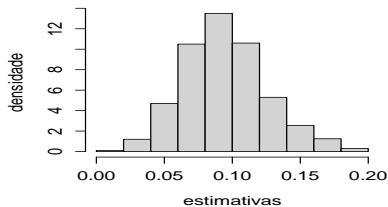


quantil da distribuição padronizada do estimador



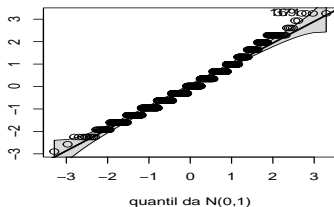
$p=0,10$

$n = 100$, p -valor (teste-SW) = 0

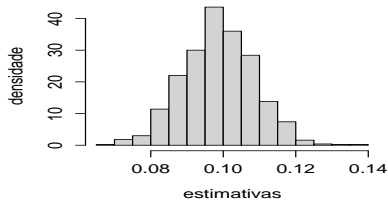


$n = 100$, p -valor (teste-SW) = 0

quantil da distribuição padronizada do estimador

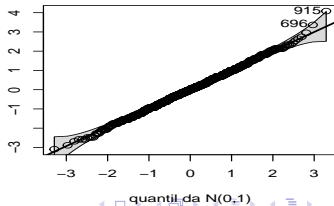


$n = 1000$, p -valor (teste-SW) = 0.2383



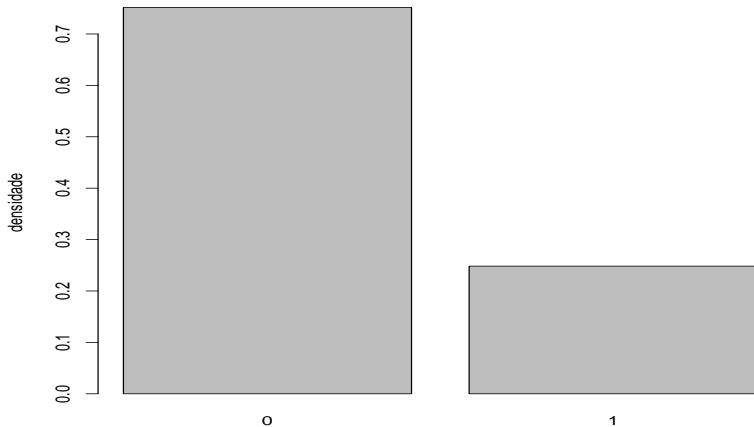
$n = 1000$, p -valor (teste-SW) = 0.2383

quantil da distribuição padronizada do estimador



$p=0,25$

Gráfico de Colunas

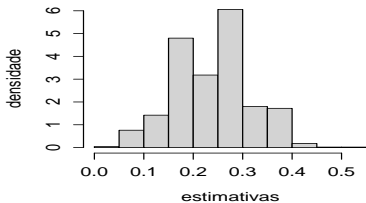


variavel de interesse



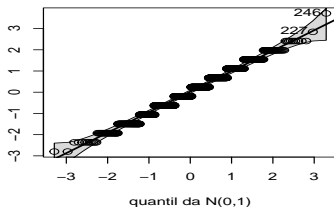
$p=0,25$

$n = 30$, p -valor (teste-SW) = 0

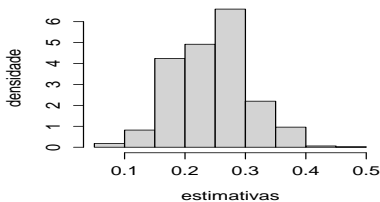


$n = 30$, p -valor (teste-SW) = 0

quantil da distribuição padronizada do estimador

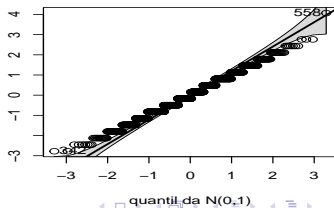


$n = 50$, p -valor (teste-SW) = 0



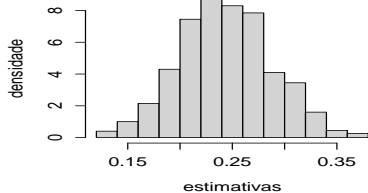
$n = 50$, p -valor (teste-SW) = 0

quantil da distribuição padronizada do estimador

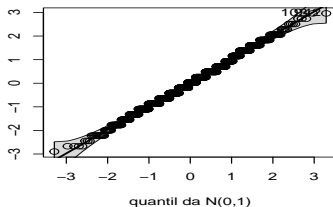


$p=0,25$

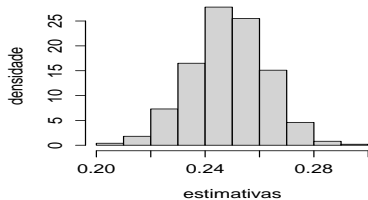
$n = 100$, p -valor (teste-SW) = $6e-04$



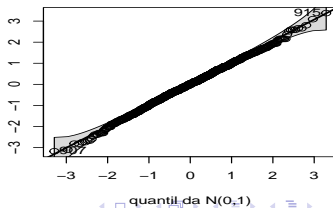
$n = 100$, p -valor (teste-SW) = $6e-04$



$n = 1000$, p -valor (teste-SW) = 0.4711

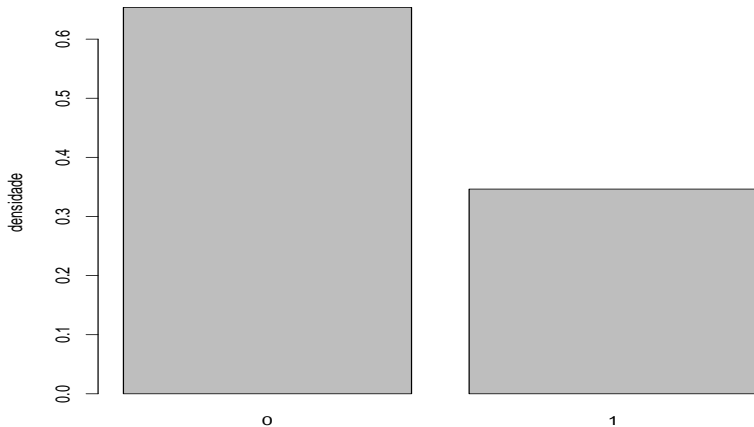


$n = 1000$, p -valor (teste-SW) = 0.4711



$p=0,35$

Gráfico de Colunas

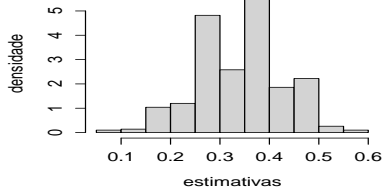


variavel de interesse



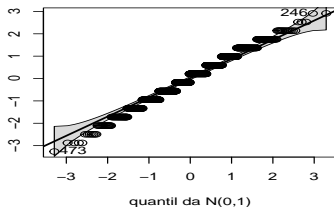
$p=0,35$

$n = 30$, p -valor (teste-SW) = 0

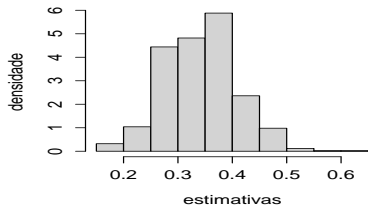


quantil da distribuição padronizada do estimador

$n = 30$, p -valor (teste-SW) = 0

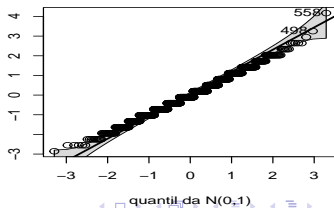


$n = 50$, p -valor (teste-SW) = 0



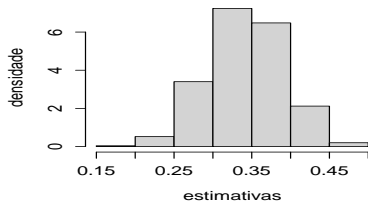
quantil da distribuição padronizada do estimador

$n = 50$, p -valor (teste-SW) = 0

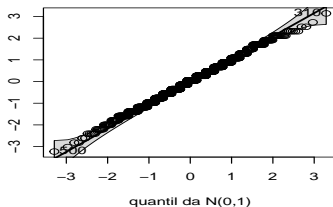


$p=0,35$

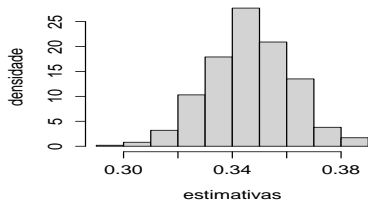
$n = 100$, p -valor (teste-SW) = 0.0056



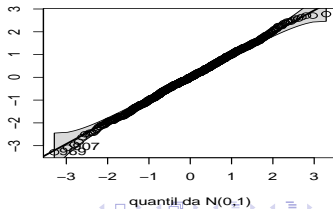
$n = 100$, p -valor (teste-SW) = 0.0056



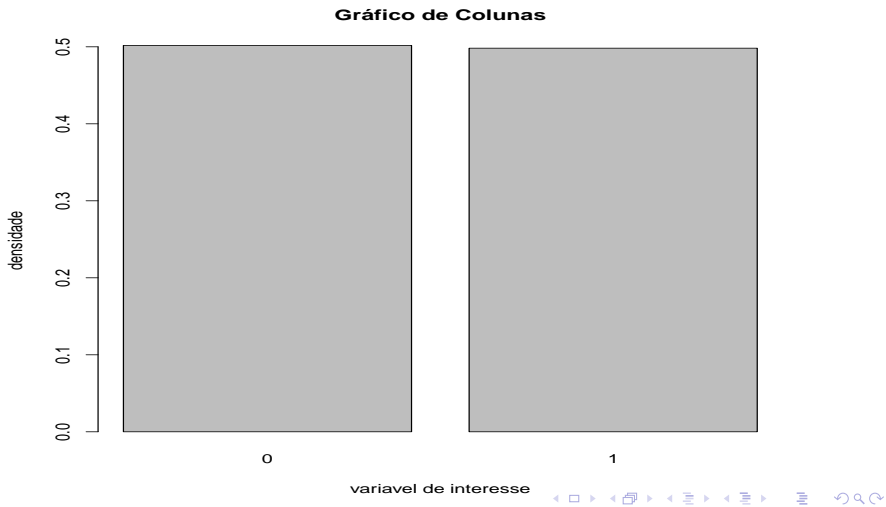
$n = 1000$, p -valor (teste-SW) = 0.5187



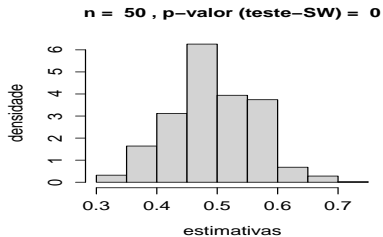
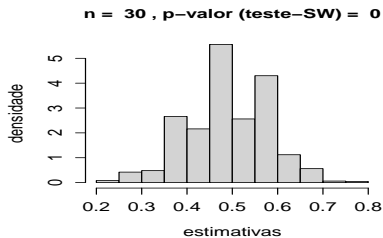
$n = 1000$, p -valor (teste-SW) = 0.5187



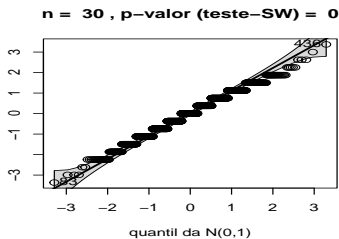
$p=0,50$



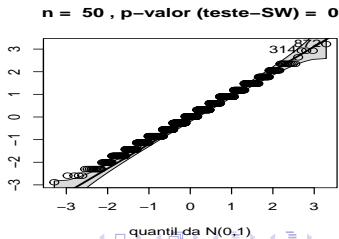
$p=0,50$



quantil da distribuição padronizada do estimador

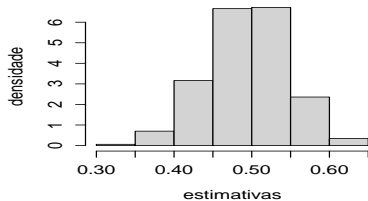


quantil da distribuição padronizada do estimador

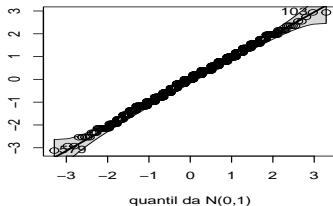


$p=0,50$

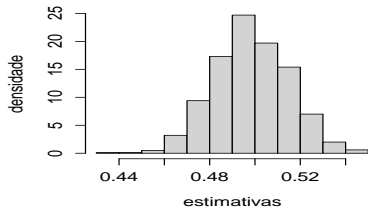
$n = 100$, p -valor (teste-SW) = 0.0051



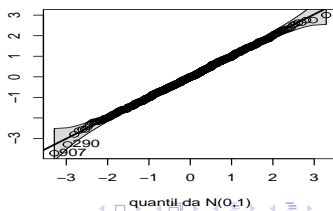
$n = 100$, p -valor (teste-SW) = 0.0051



$n = 1000$, p -valor (teste-SW) = 0.7653

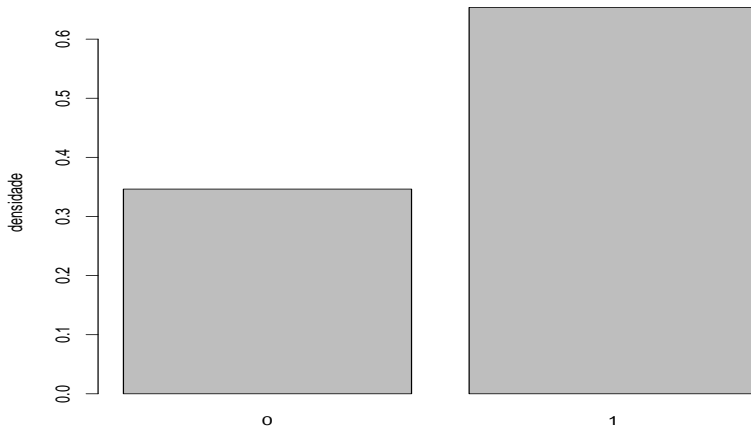


$n = 1000$, p -valor (teste-SW) = 0.7653



$p=0,65$

Gráfico de Colunas

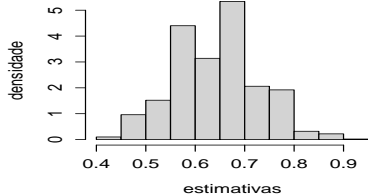


variavel de interesse



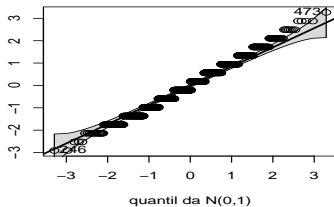
$p=0,65$

$n = 30$, p -valor (teste-SW) = 0

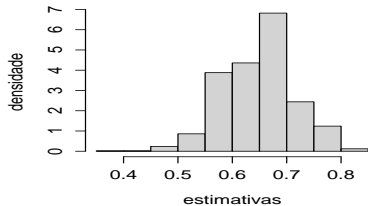


quantil da distribuição padronizada do estimador

$n = 30$, p -valor (teste-SW) = 0

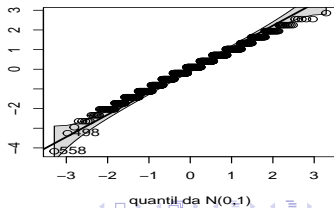


$n = 50$, p -valor (teste-SW) = 0



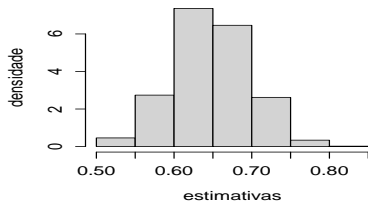
quantil da distribuição padronizada do estimador

$n = 50$, p -valor (teste-SW) = 0



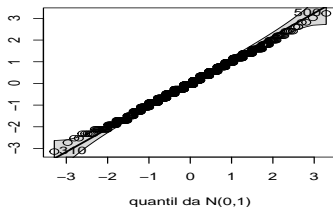
$p=0,65$

$n = 100$, p -valor (teste-SW) = 0.0056

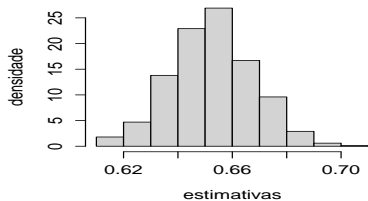


$n = 100$, p -valor (teste-SW) = 0.0056

quantil da distribuição padronizada do estimador

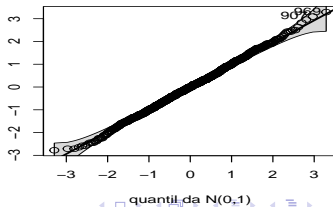


$n = 1000$, p -valor (teste-SW) = 0.5187



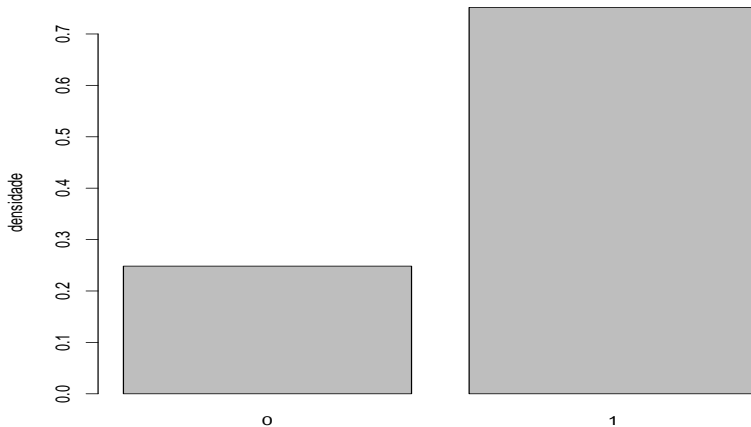
$n = 1000$, p -valor (teste-SW) = 0.5187

quantil da distribuição padronizada do estimador



$p=0,75$

Gráfico de Colunas

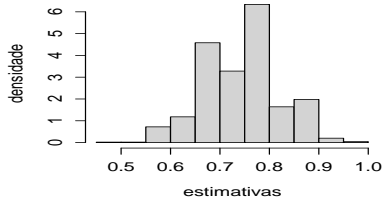


variavel de interesse



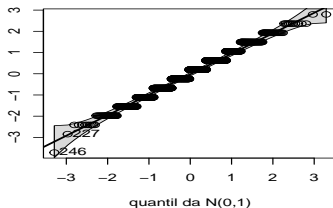
$p=0,75$

$n = 30$, p -valor (teste-SW) = 0

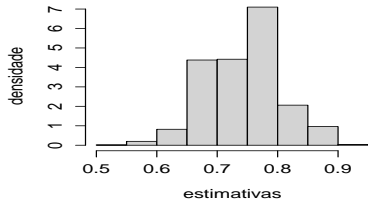


quantil da distribuição padronizada do estimador

$n = 30$, p -valor (teste-SW) = 0

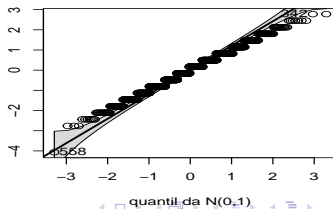


$n = 50$, p -valor (teste-SW) = 0



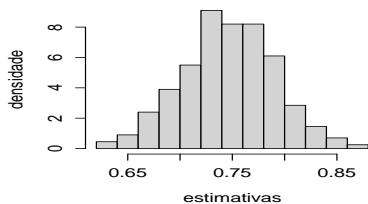
quantil da distribuição padronizada do estimador

$n = 50$, p -valor (teste-SW) = 0

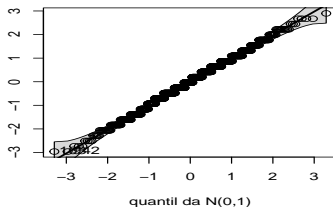


$p=0,75$

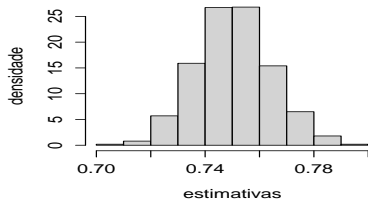
$n = 100$, p -valor (teste-SW) = $6e-04$



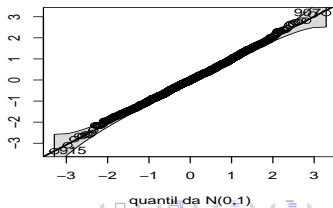
$n = 100$, p -valor (teste-SW) = $6e-04$



$n = 1000$, p -valor (teste-SW) = 0.4711

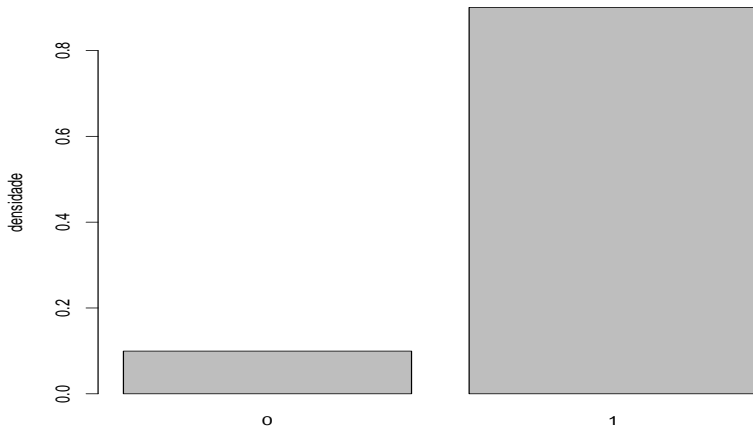


$n = 1000$, p -valor (teste-SW) = 0.4711



$p=0,90$

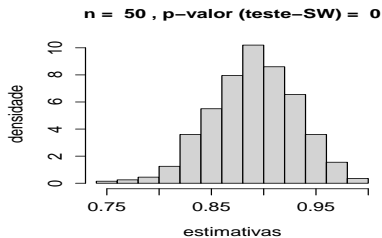
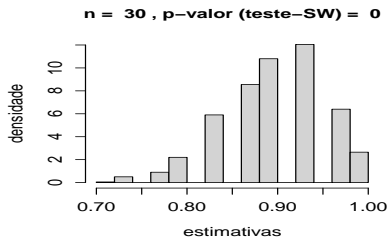
Gráfico de Colunas



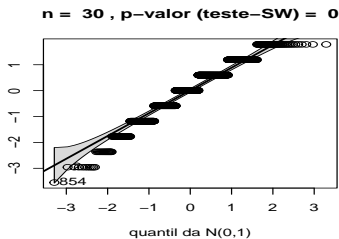
variavel de interesse



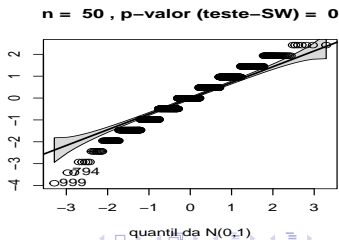
$p=0,90$



quantil da distribuição padronizada do estimador

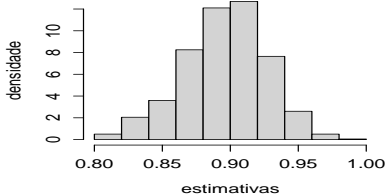


quantil da distribuição padronizada do estimador



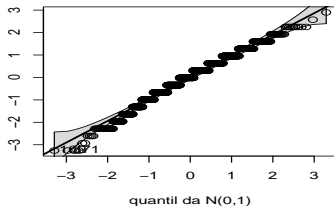
$p=0,90$

$n = 100$, p -valor (teste-SW) = 0

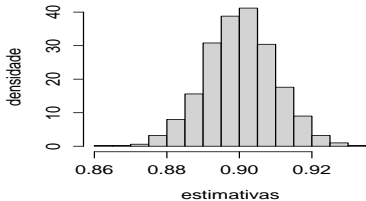


$n = 100$, p -valor (teste-SW) = 0

quantil da distribuição padronizada do estimador

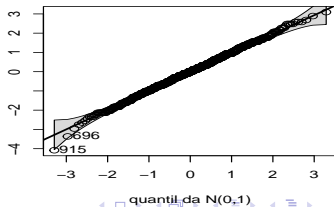


$n = 1000$, p -valor (teste-SW) = 0.2383



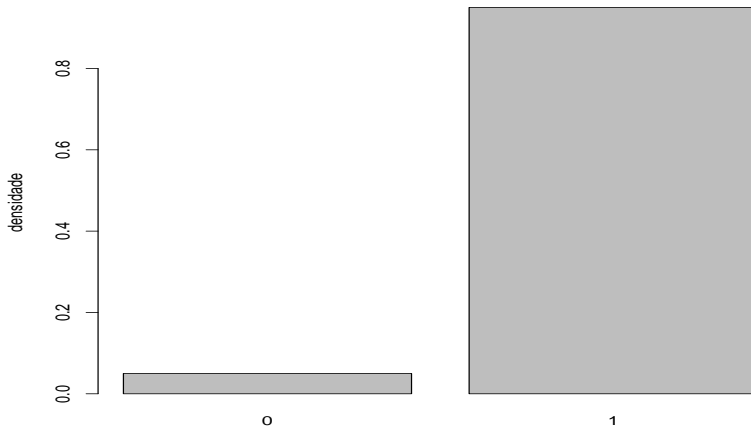
$n = 1000$, p -valor (teste-SW) = 0.2383

quantil da distribuição padronizada do estimador



$p=0,95$

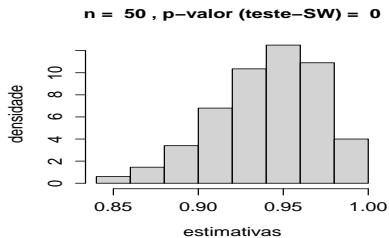
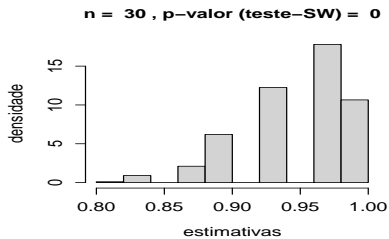
Gráfico de Colunas



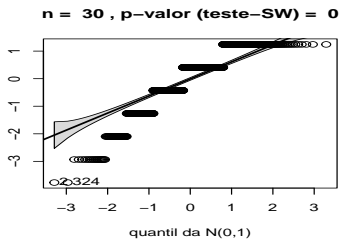
variavel de interesse



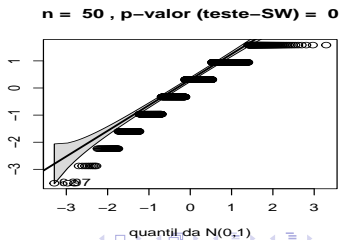
$p=0,95$



quantil da distribuição padronizada do estimador

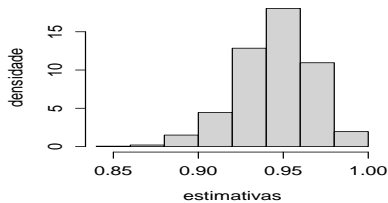


quantil da distribuição padronizada do estimador



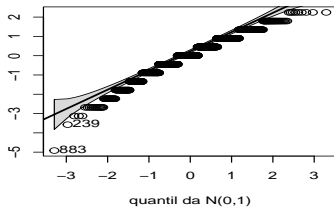
$p=0,95$

$n = 100$, p -valor (teste-SW) = 0

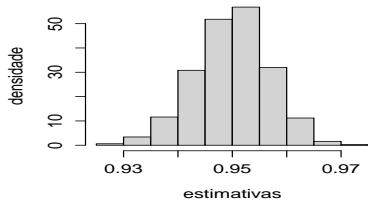


quantil da distribuição padronizada do estimador

$n = 100$, p -valor (teste-SW) = 0

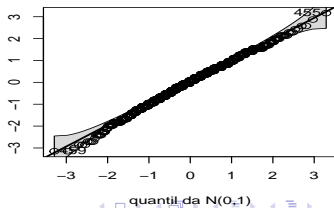


$n = 1000$, p -valor (teste-SW) = 0.0025



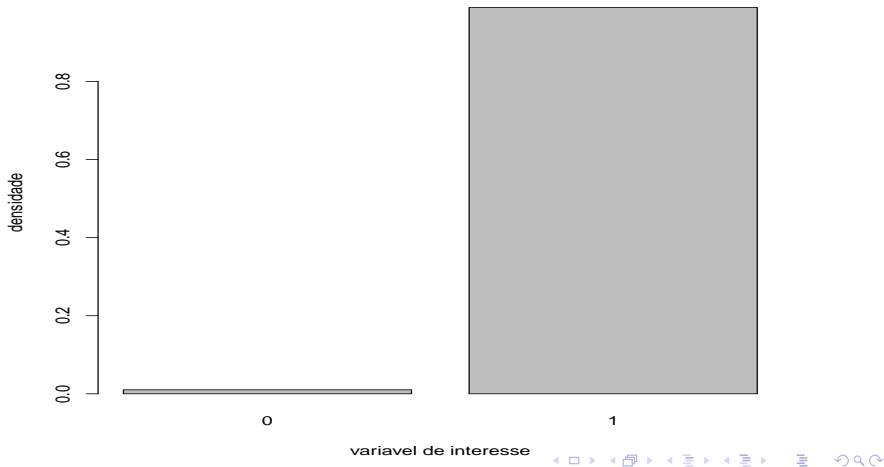
quantil da distribuição padronizada do estimador

$n = 1000$, p -valor (teste-SW) = 0.0025

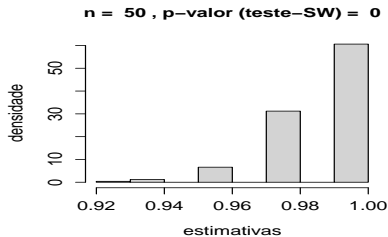
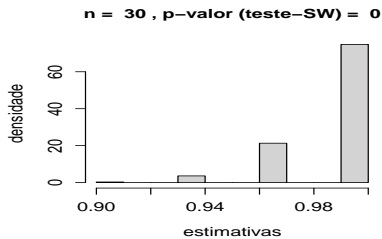


$p=0,99$

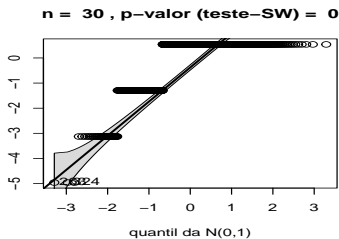
Gráfico de Colunas



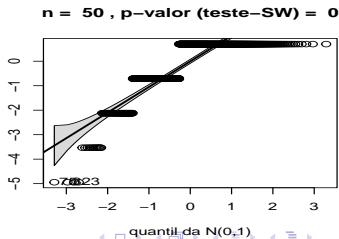
$p=0,99$



quantil da distribuição padronizada do estimador

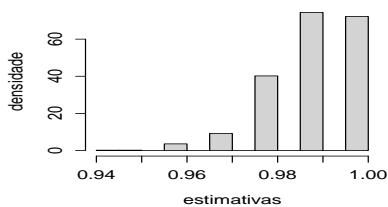


quantil da distribuição padronizada do estimador



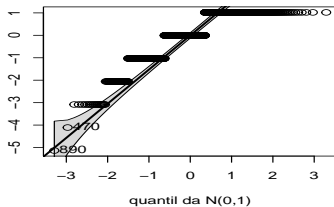
$p=0,99$

$n = 100$, p -valor (teste-SW) = 0

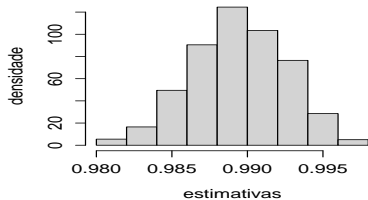


quantil da distribuição padronizada do estimador

$n = 100$, p -valor (teste-SW) = 0

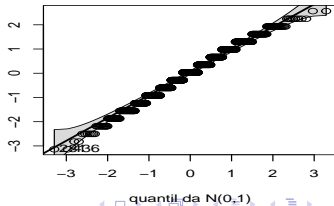


$n = 1000$, p -valor (teste-SW) = 0



quantil da distribuição padronizada do estimador

$n = 1000$, p -valor (teste-SW) = 0



Otimidade dos estimadores

- Vamos nos concentrar na média amostral e na classe de estimadores não viciados que sejam combinações lineares das variáveis aleatórias (Y_1, \dots, Y_n) .
- Os resultados para os outros parâmetros são análogos.
- A forma geral do estimador em questão é dada por

$$\hat{\mu}_{sc} = \sum_{i=1}^n c_i Y_i$$

- Note que, sob $AAS_s(A_2)$ temos que os Y_i 's não são mais independentes (embora ainda sejam “identicamente distribuídas” (pelo menos em relação a alguns momentos)) Exercício 2.9, livro: “Elementos de Amostragem”). Temos que $\mathcal{E}(Y_i) = \mu$, $\mathcal{V}(Y_i) = \frac{N-1}{N} s^2$, $Cov(Y_i, Y_j) = -\frac{s^2}{N}$, $\forall i \neq j$.

Otimidade dos estimadores

- Note que $\mathcal{E}(\hat{\mu}_{sc}) = \sum_{i=1}^n c_i \mathcal{E}(Y_i) = \mu \sum_{i=1}^n c_i$.
- Exercício: provar que $\hat{\mu}_{sc}$ é um estimador não viciado se e somente se

$$\sum_{i=1}^n c_i = 1 \quad (5)$$

- Defina $\bar{c} = \frac{1}{n} \sum_{i=1}^n c_i$. Pode-se provar (usando (5)) que $\sum_{i,j=1,2,\dots,N}^{i \neq j} c_i c_j = 1 - \sum_{i=1}^n c_i^2$.

Otimidade dos estimadores

- Assim, temos que:

$$\begin{aligned}\mathcal{V}(\hat{\mu}_{sc}) &= \sum_{i=1}^n c_i^2 \mathcal{V}(Y_i) + \sum_{\substack{i \neq j \\ i, j=1, 2, \dots, N}} \text{Cov}(c_i Y_i, c_j Y_j) \\ &= \frac{N-1}{N} s^2 \sum_{i=1}^n c_i^2 - \frac{s^2}{N} \sum_{\substack{i \neq j \\ i, j=1, 2, \dots, N}} c_i c_j \\ &= s^2 \left(\sum_{i=1}^n c_i^2 - \frac{1}{N} \right)\end{aligned}$$

- Devemos então minimizar $g(\mathbf{c}) = \sum_{i=1}^n c_i^2 - \frac{1}{N} + \lambda (\sum_{i=1}^n c_i - 1)$, em $c_i = 1, 2, \dots, n$, $\mathbf{c} = (c_1, \dots, c_N)^T$.

Otimidade dos estimadores

- Derivando $g(\mathbf{c})$ com relação à λ e $c_i, i = 1, 2, \dots, n$, e igualando cada uma das $n+1$ equações a 0, vem que

$$\sum_{i=1}^n c_i - 1 = 0 \quad (6)$$

$$2c_i + \lambda = 0 \quad (7)$$

- De (6) vem que $\sum_{i=1}^n c_i = 1(*)$. Assim, utilizando (*) em (7), $i=1,2,\dots,n$, tem-se que $\lambda = -\frac{2}{n}(**)$.
- Logo utilizando (**) em (7), tem-se que $c_i = \frac{1}{n}, i = 1, 2, \dots, n$. O que implica que o estimador ótimo é dado por $\hat{\mu}_{sc} = \frac{1}{n} \sum_{i=1}^n Y_i$.
- Utilizando desenvolvimentos análogos, obtemos resultados semelhantes para a estimação de τ e ρ .