

# 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. 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  $y_i = y_i^k, \forall k \in \mathbb{R}^+$ .

- Estimador “natural”:

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

- 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  $\hat{\mu}$  em

[http://www.ime.unicamp.br/~cnaber/aula\\_AAS%20sem%20reposicao%20parte%201%20Amost%20S%202018.pdf](http://www.ime.unicamp.br/~cnaber/aula_AAS%20sem%20reposicao%20parte%201%20Amost%20S%202018.pdf).

- 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

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

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

## 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  $\mathcal{E}_{A_i}(\hat{p}) = p$ ,  $i = 1, 2$ .
- Além disso,  $\hat{\mathcal{V}}_{A_1}(\hat{p}) = \frac{\hat{p}\hat{q}}{n-1}$  e  $\hat{\mathcal{V}}_{A_2}(\hat{p}) = (1-f)\frac{\hat{p}\hat{q}}{n-1}$ , em que  $f \in (0, 1)$ .
- Portanto, o efeito do planejamento (EPA), do estimador sob o plano  $A_2$  em relação ao plano  $A_1$ , é dado por:

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

- Ademais, quando  $N \rightarrow \infty$ ,  $\mathcal{V}_{A_2}(\hat{p}) \rightarrow \mathcal{V}_{A_1}(\hat{p})$ .
- 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.



# Intervalo de Confiança

- Assim, dois intervalos de confiança (assintóticos) com coeficiente de confiança de aproximadamente  $\gamma$ , 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).

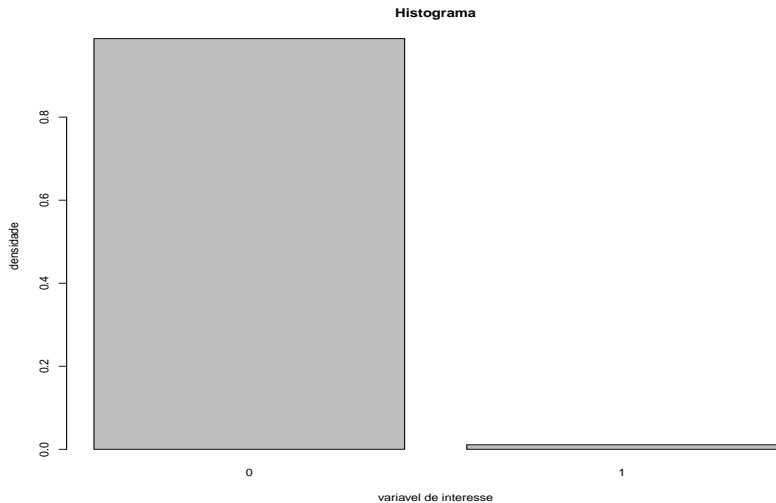
# Estudos de simulação

- Distribuição assintótica do estimador para a proporção. Tamanho da população  $N = 100000$ .
- 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

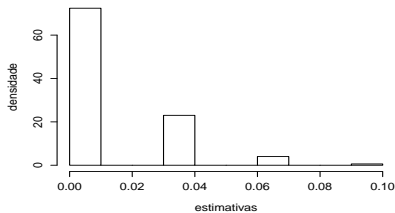
- Quatro tamanhos amostrais (30, 50, 100, 1000), em termos percentuais, com relação ao tamanho da população (0,03%,0,05%,0,1%,1%).
- Estudar a distribuição amostral (empírica) com base em  $R = 1000$  réplicas (amostras selecionadas da população de interesse).

$p=0,01$

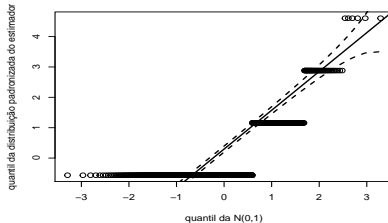


$p=0,01$

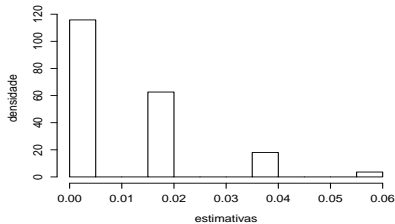
$n = 30$ ,  $p$ -valor (teste-SW) = 0



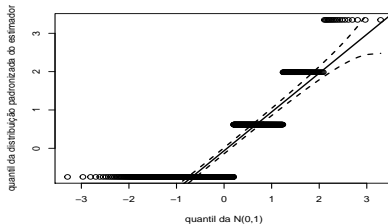
$n = 30$ ,  $p$ -valor (teste-SW) = 0



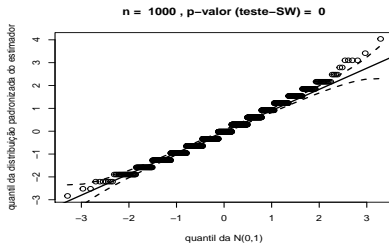
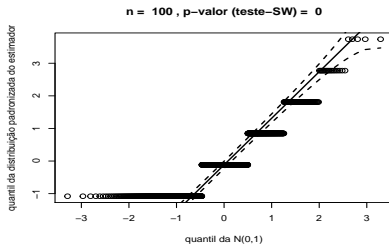
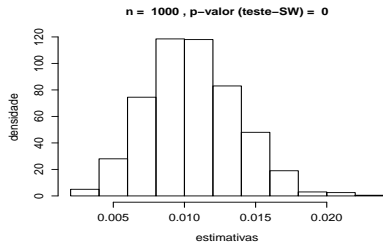
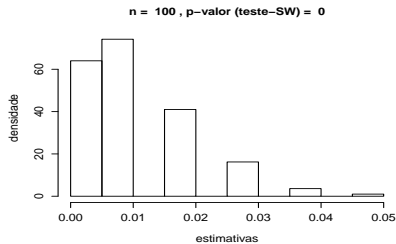
$n = 50$ ,  $p$ -valor (teste-SW) = 0



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

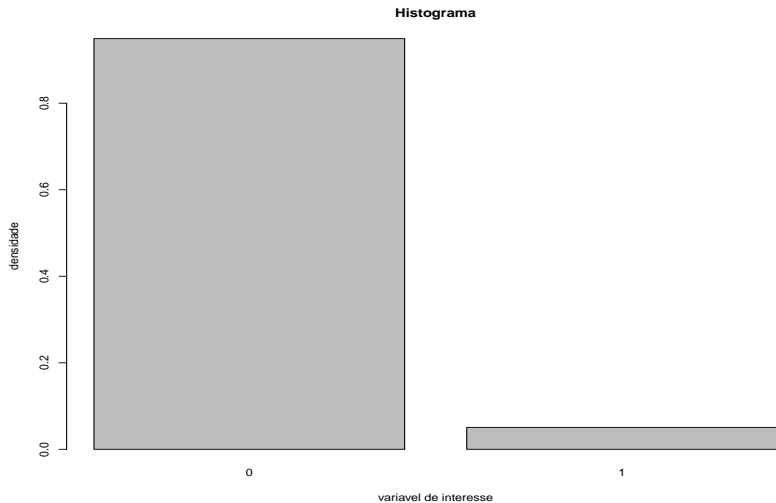


$p=0,01$



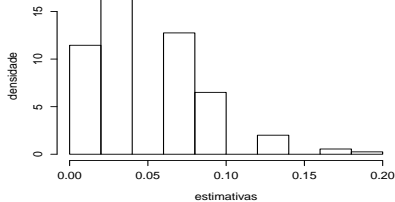


$p=0,05$

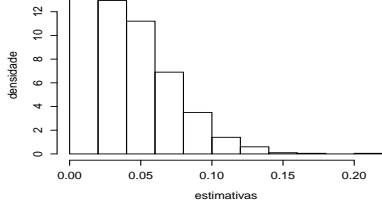


$p=0,05$

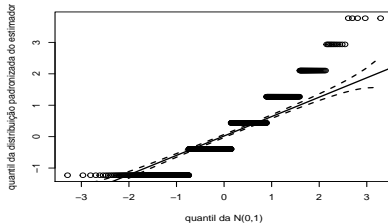
$n = 30$  ,  $p$ -valor (teste-SW) = 0



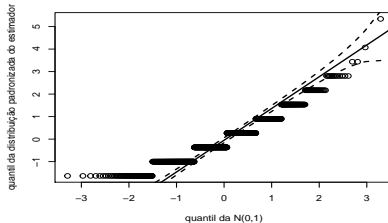
$n = 50$  ,  $p$ -valor (teste-SW) = 0



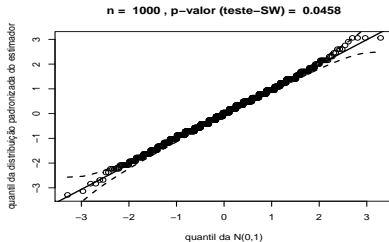
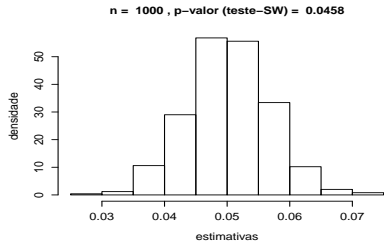
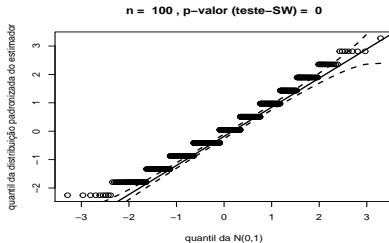
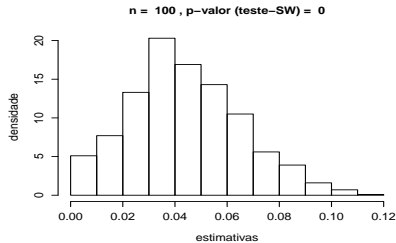
$n = 30$  ,  $p$ -valor (teste-SW) = 0



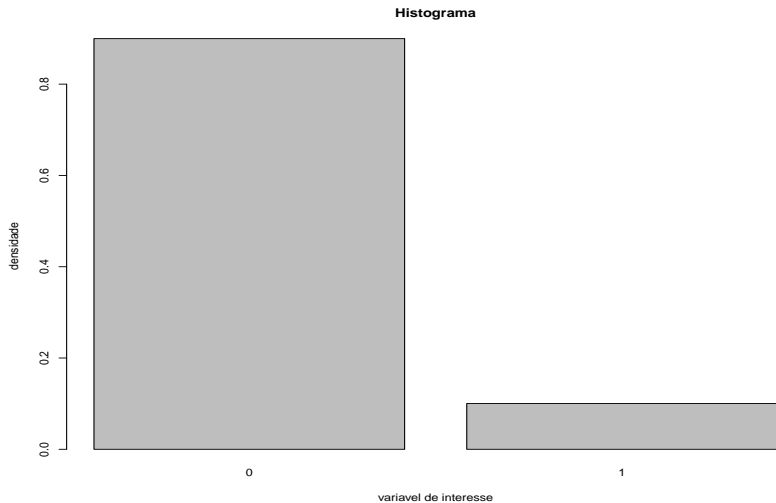
$n = 50$  ,  $p$ -valor (teste-SW) = 0



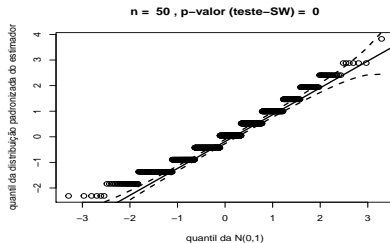
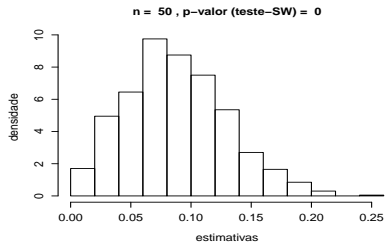
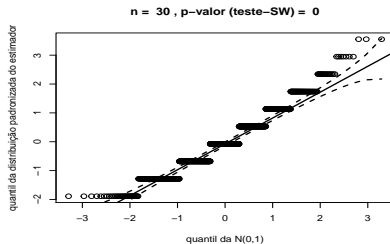
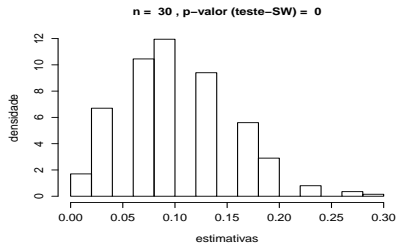
$p=0,05$



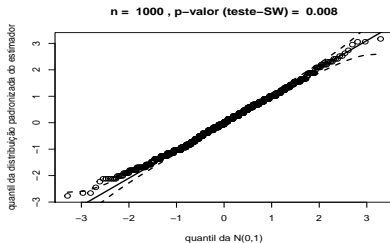
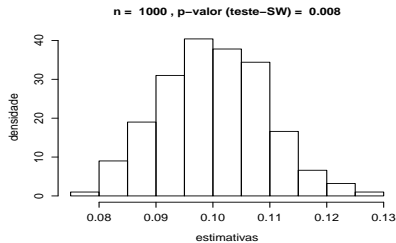
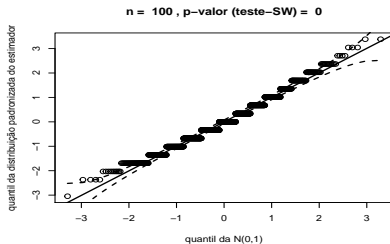
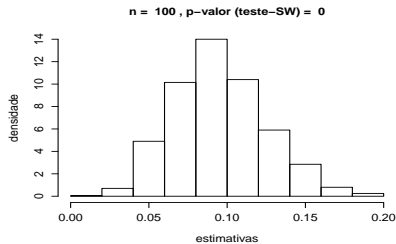
$$p=0,10$$



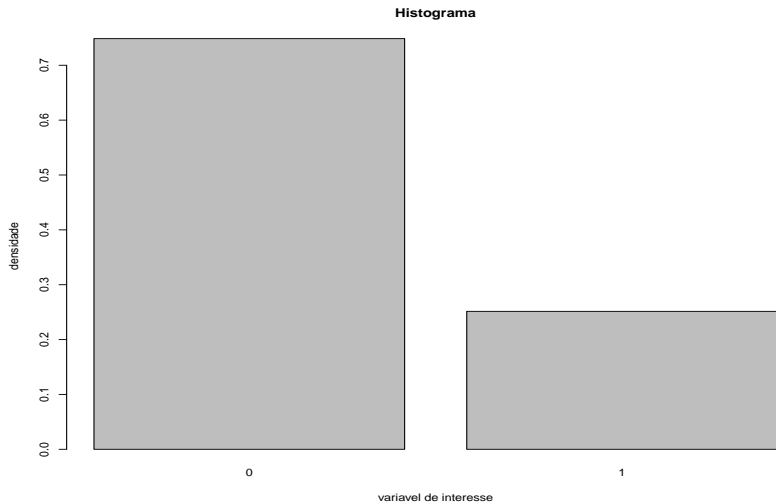
$p=0,10$



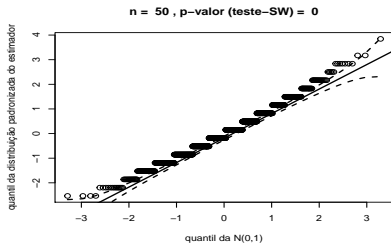
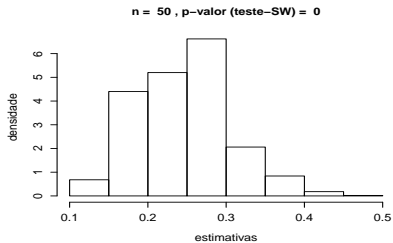
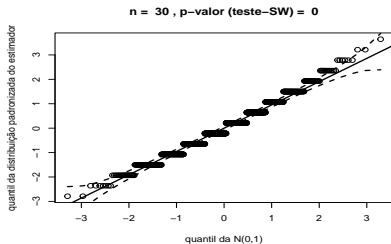
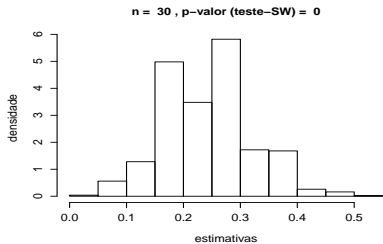
$p=0,10$



$p=0,25$



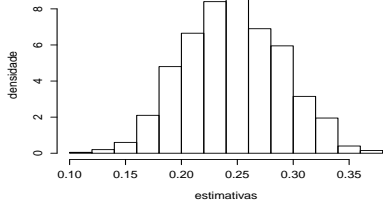
$p=0,25$



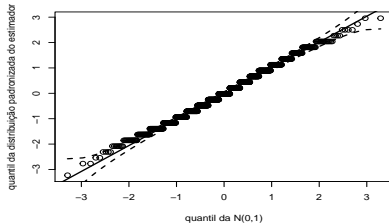


$p=0,25$

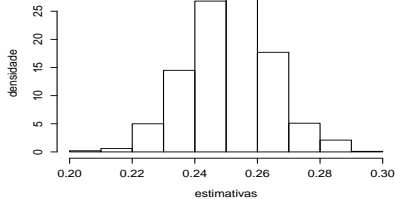
$n = 100$  ,  $p$ -valor (teste-SW) =  $3e-04$



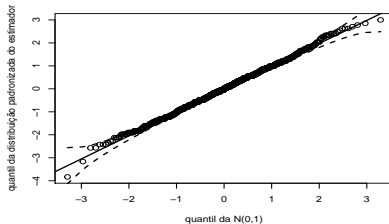
$n = 100$  ,  $p$ -valor (teste-SW) =  $3e-04$



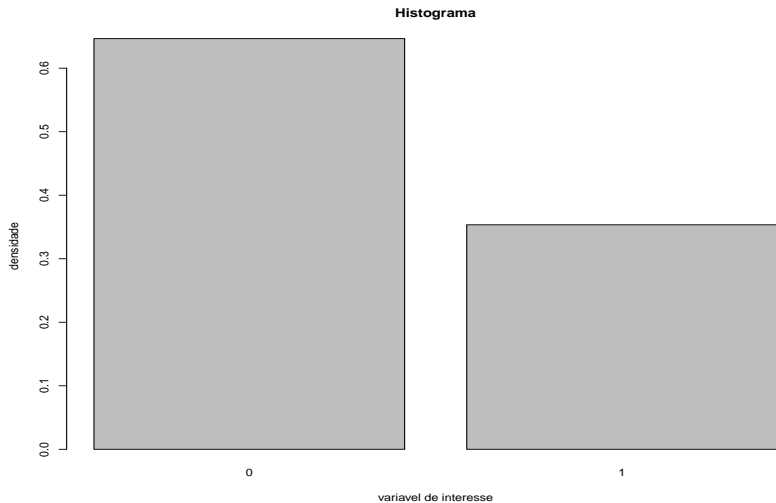
$n = 1000$  ,  $p$ -valor (teste-SW) = 0.2888



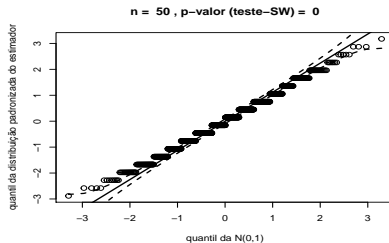
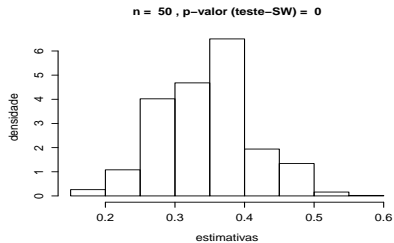
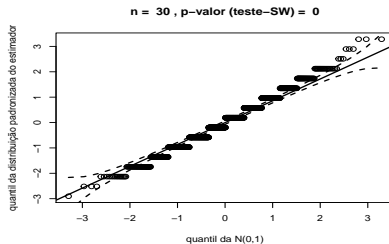
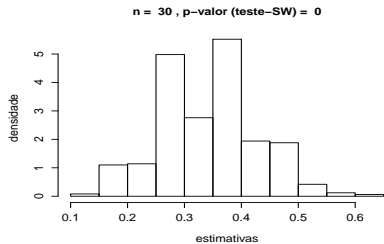
$n = 1000$  ,  $p$ -valor (teste-SW) = 0.2888



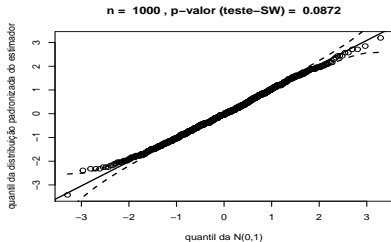
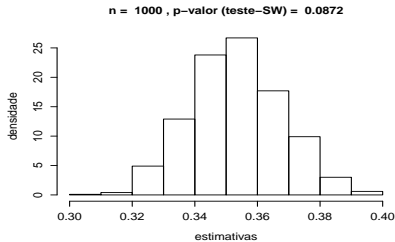
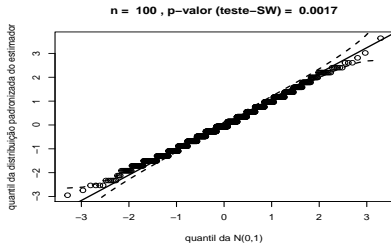
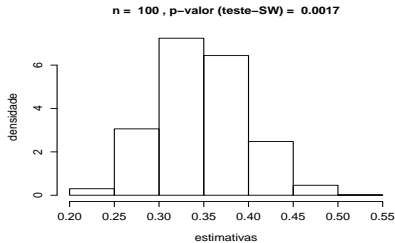
$$p=0,35$$



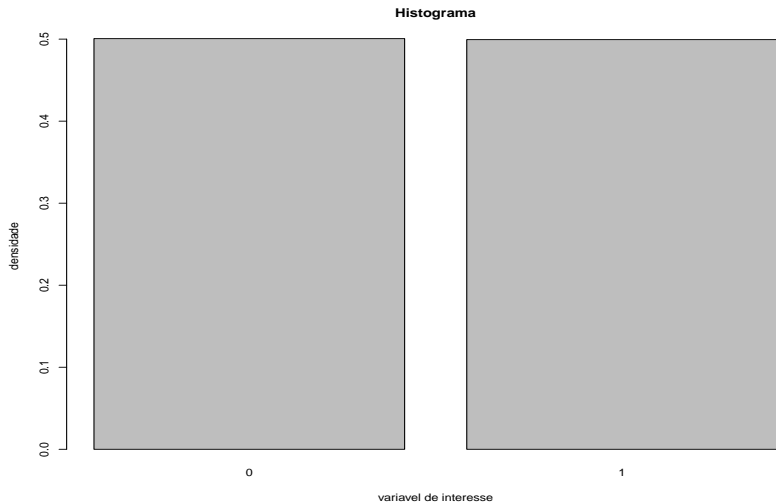
$p=0,35$



$p=0,35$

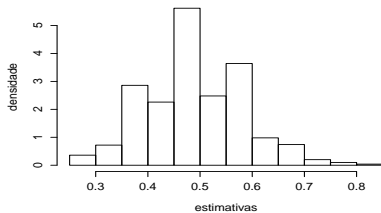


$p=0,50$

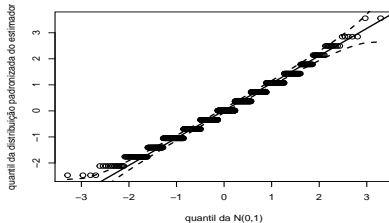


$p=0,50$

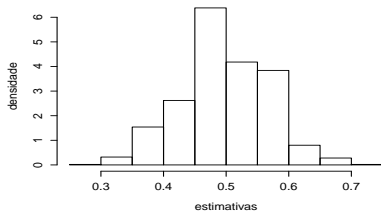
$n = 30$ ,  $p$ -valor (teste-SW) = 0



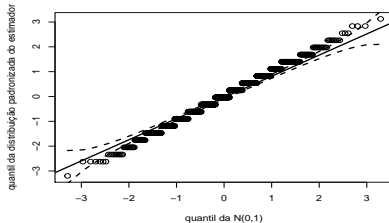
$n = 30$ ,  $p$ -valor (teste-SW) = 0



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

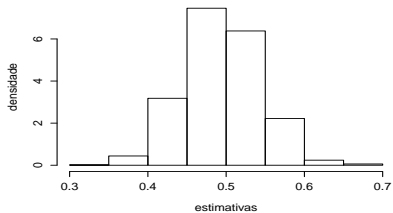


$n = 50$ ,  $p$ -valor (teste-SW) = 0

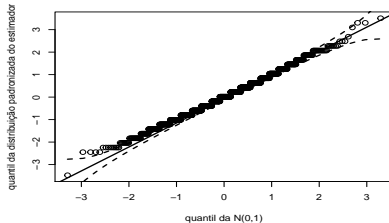


$p=0,50$

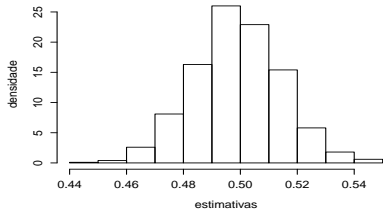
$n = 100$  ,  $p$ -valor (teste-SW) =  $8e-04$



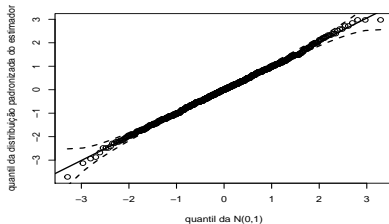
$n = 100$  ,  $p$ -valor (teste-SW) =  $8e-04$



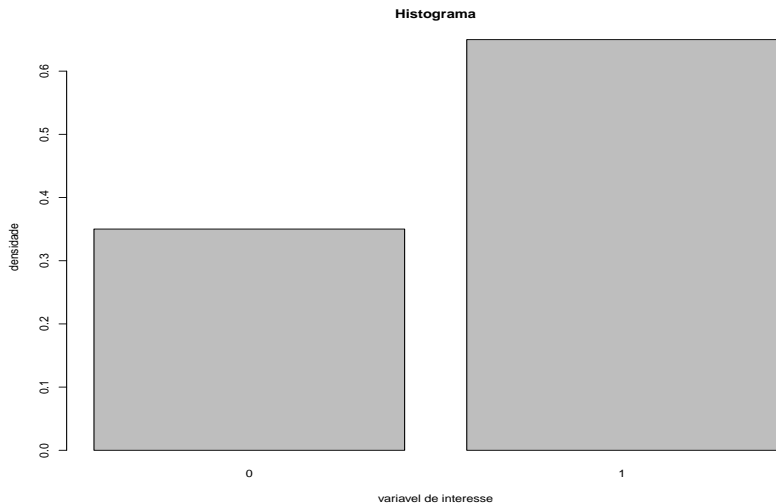
$n = 1000$  ,  $p$ -valor (teste-SW) =  $0.6002$



$n = 1000$  ,  $p$ -valor (teste-SW) =  $0.6002$



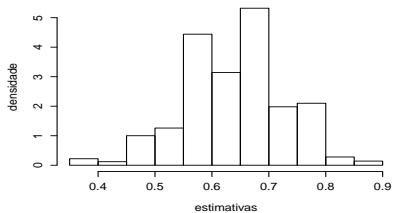
$p=0,65$



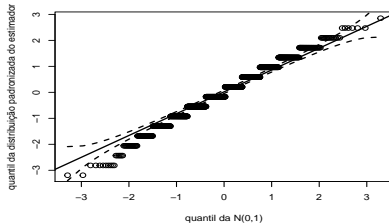


$p=0,65$

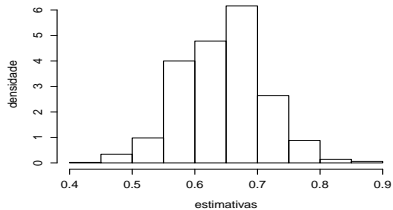
$n = 30$ ,  $p$ -valor (teste-SW) = 0



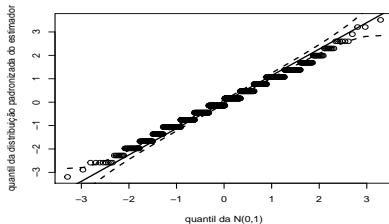
$n = 30$ ,  $p$ -valor (teste-SW) = 0



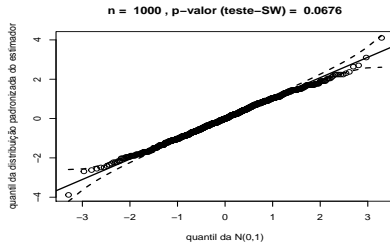
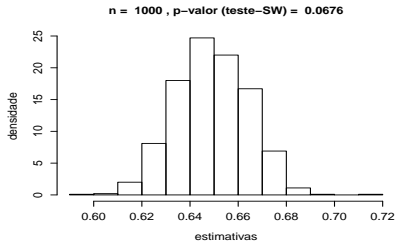
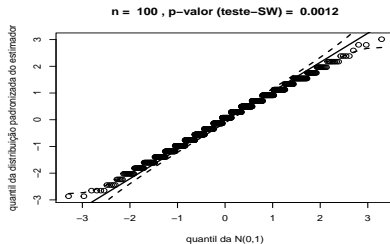
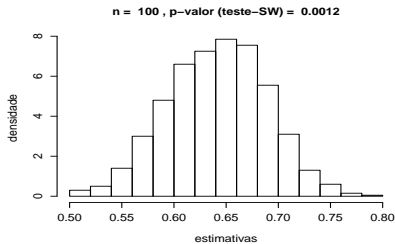
$n = 50$ ,  $p$ -valor (teste-SW) = 0



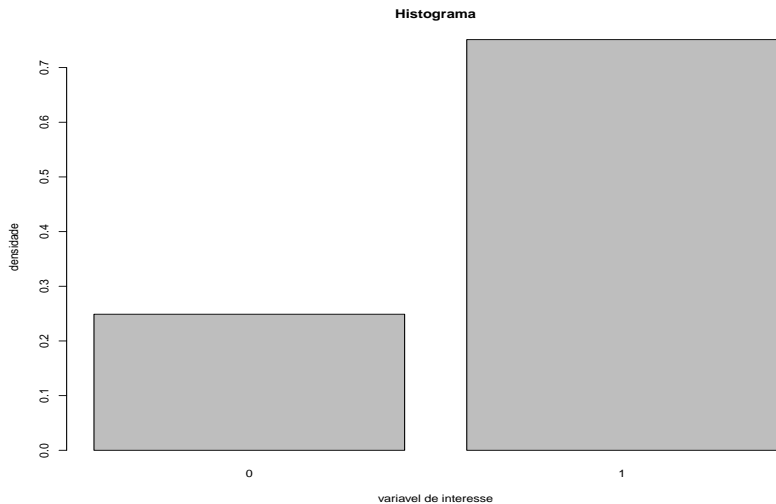
$n = 50$ ,  $p$ -valor (teste-SW) = 0



$p=0,65$

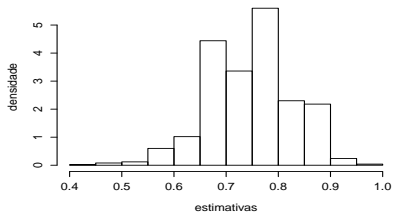


$$p=0,75$$

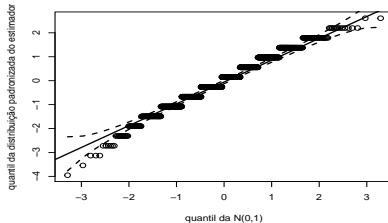


$p=0,75$

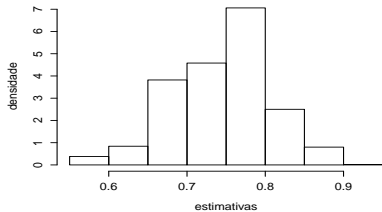
$n = 30$ ,  $p$ -valor (teste-SW) = 0



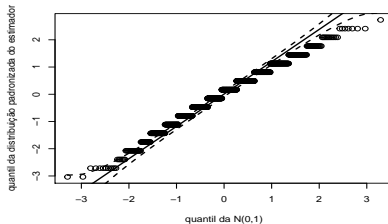
$n = 30$ ,  $p$ -valor (teste-SW) = 0



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

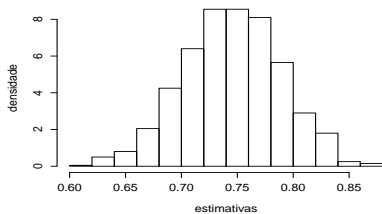


$n = 50$ ,  $p$ -valor (teste-SW) = 0

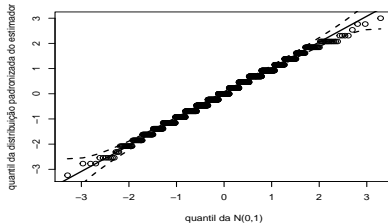


$p=0,75$

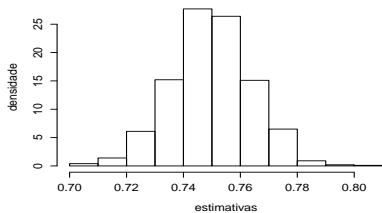
$n = 100$  ,  $p$ -valor (teste-SW) =  $7e-04$



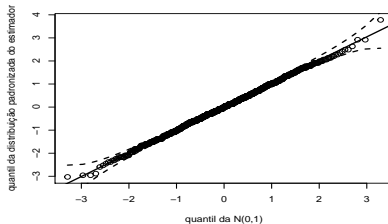
$n = 100$  ,  $p$ -valor (teste-SW) =  $7e-04$



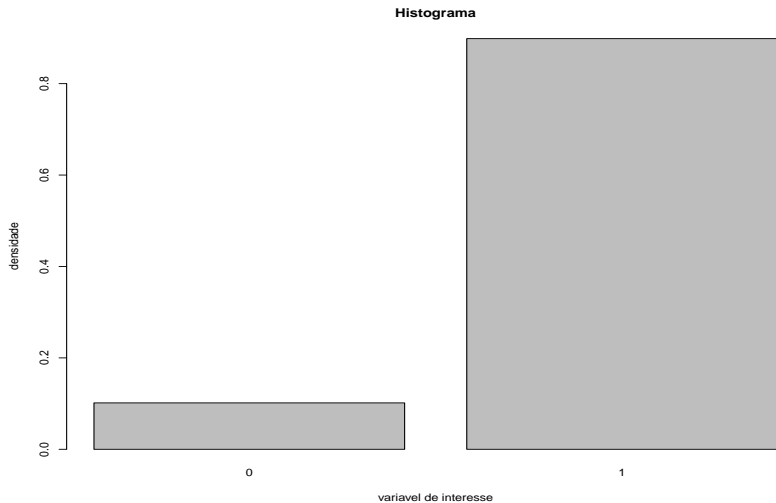
$n = 1000$  ,  $p$ -valor (teste-SW) = 0.7532



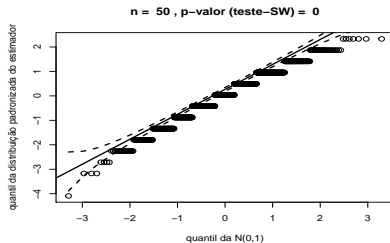
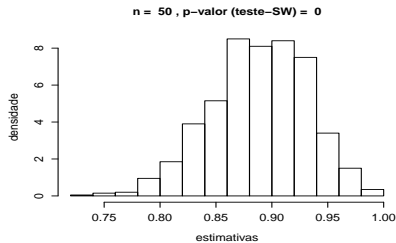
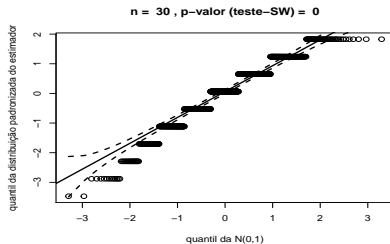
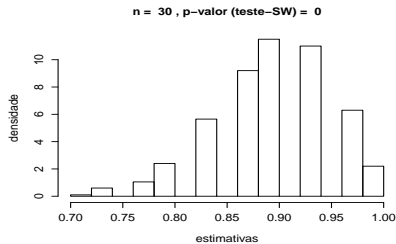
$n = 1000$  ,  $p$ -valor (teste-SW) = 0.7532



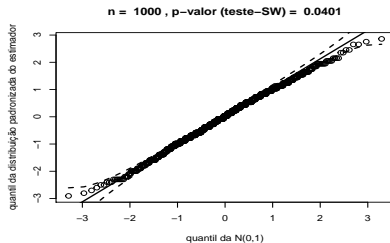
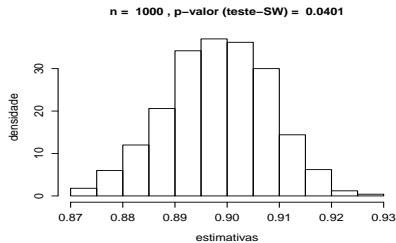
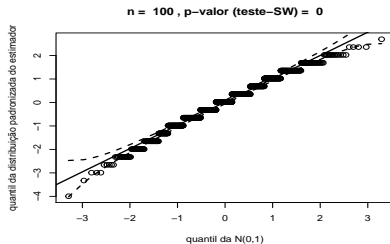
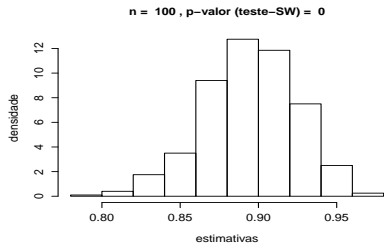
$$p=0,90$$



$p=0,90$

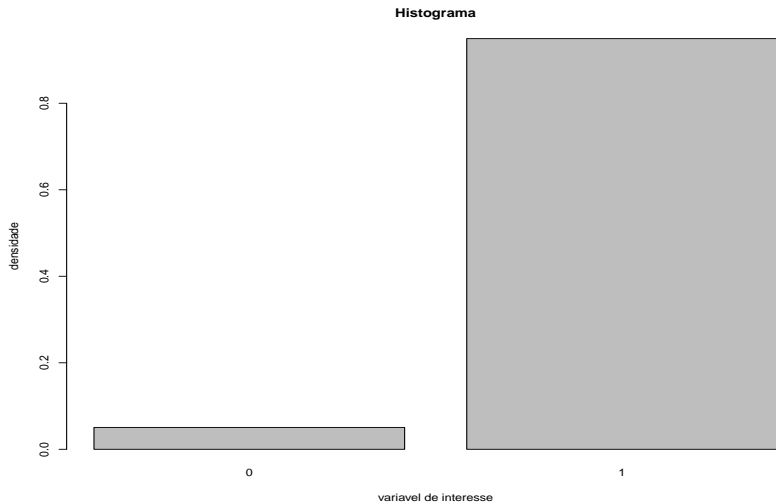


$p=0,90$

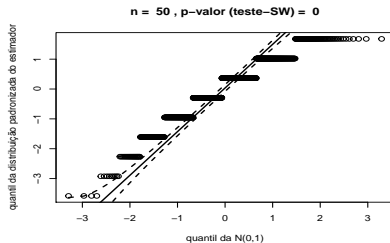
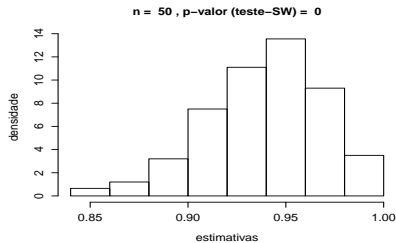
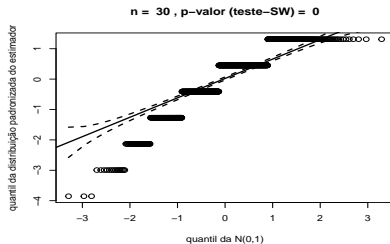
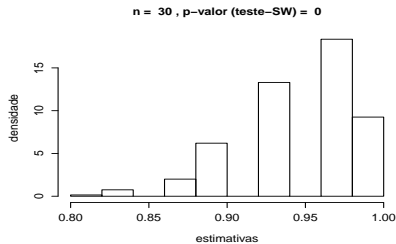




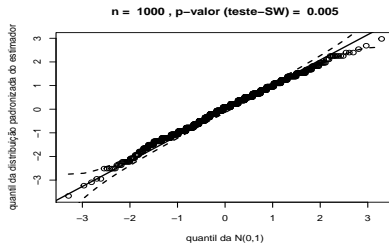
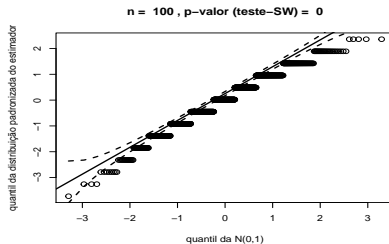
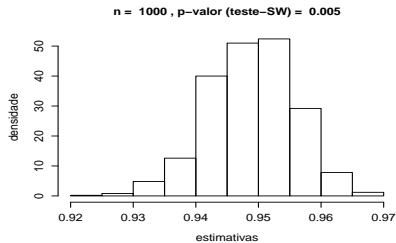
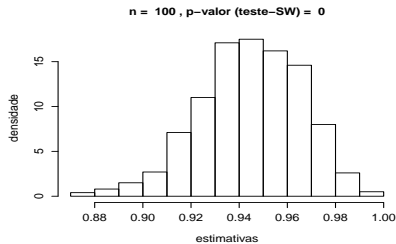
$p=0,95$



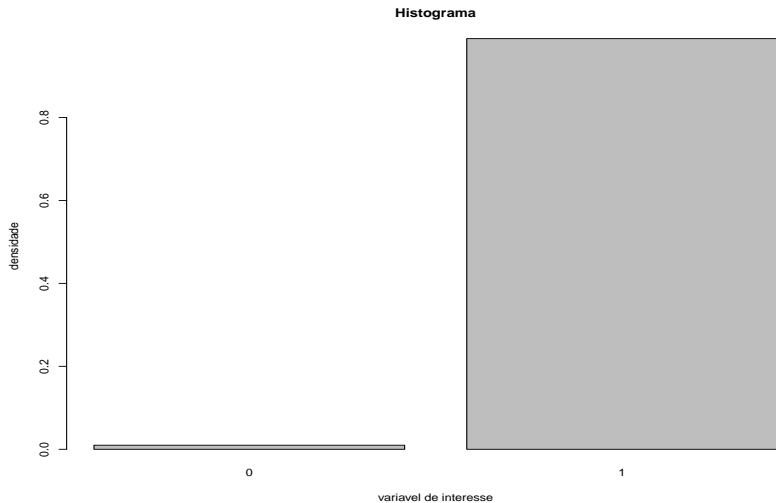
$p=0,95$



$p=0,95$

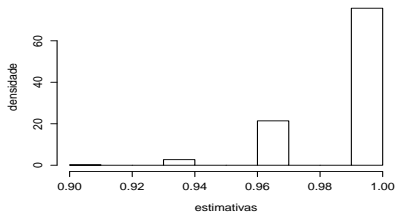


$p=0,99$

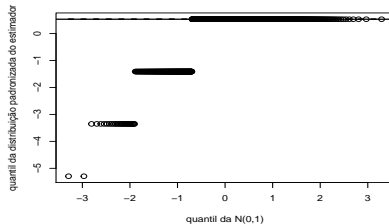


$p=0,99$

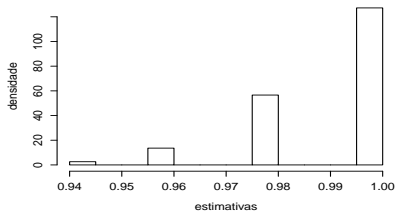
$n = 30$ ,  $p$ -valor (teste-SW) = 0



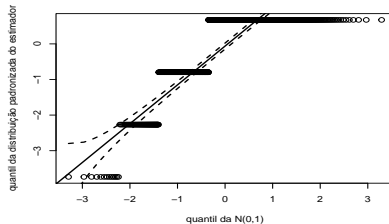
$n = 30$ ,  $p$ -valor (teste-SW) = 0



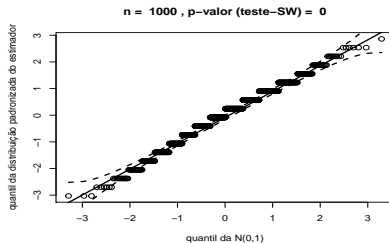
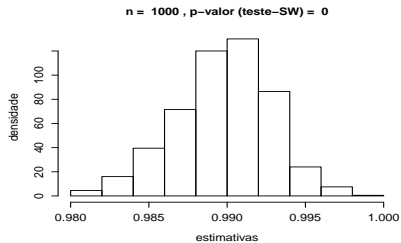
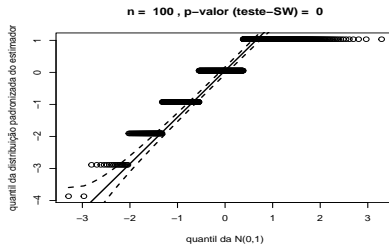
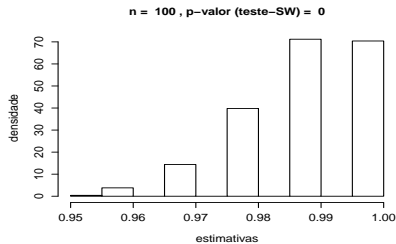
$n = 50$ ,  $p$ -valor (teste-SW) = 0



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



$p=0,99$



## 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  $Y_i$ 's não são mais independentes (embora ainda sejam identicamente distribuídas Exercício 2.9, livro-texto). Temos que  $\mathcal{E}(Y_i) = \mu$ ,  $\mathcal{V}(Y_i) = \frac{N-1}{N} s^2$ ,  
 $Cov(Y_i, Y_j) = -\frac{s^2}{N}, \forall i \neq j$ .

# Otimalidade 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 que (usando (5) e  $\sum_{i,j=1,2,\dots,N}^{i \neq j} c_i c_j = 1 - \sum_{i=1}^n c_i^2$ ).

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

- Devemos então buscar minimizar  $\sum_{i=1}^n c_i^2 - \frac{1}{N} + \lambda (\sum_{i=1}^n c_i - 1)$ , em  $c_i = 1, 2, \dots, n$ .

- Derivando com relação à  $\lambda$  e  $c_i, i = 1, 2, \dots, n$ , e igualando as  $n+1$  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  $\tau$  e  $p$ .