

---

Estatística: Aplicação ao Sensoriamento Remoto

SER 204 - ANO 2024

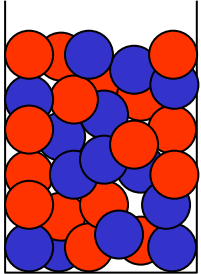
Estimação Pontual

Camilo Daleles Rennó

camilo.renno@inpe.br

<http://www.dpi.inpe.br/~camilo/estatistica/>

# Inferência Estatística



Considere o experimento: retiram-se 3 bolas de uma urna (com reposição). Define-se uma v.a.  $X$  cujo valor representa o número total de bolas vermelhas dentre as 3 escolhidas. Qual a média e variância de  $X$ ?

Quais os valores possíveis de  $X$ ?

Valores inteiros (número de tentativas bem-sucedidas)

Mínimo 0 (nenhuma bola vermelha)

Máximo 3 (todas 3 são bolas vermelhas)

$X: \{0, 1, 2, 3\}$

Qual a distribuição de probabilidade de  $X$ ?

$X$  é discreto

A probabilidade de sucesso  $p$  é igual para todas tentativas (sorteio com reposição)

O número de sorteios é pré-definido ( $n = 3$ ) e o número de sucessos é a v.a.  $X$

Distribuição: Binomial

Quais os parâmetros que definem esta Binomial?

$n$  e  $p$

$n = 3$

$p = ?$  (precisaria conhecer toda a população)

**DISTRIBUIÇÃO CONHECIDA**  
**PARÂMETRO(S) DESCONHECIDO(S)**

# Inferência Estatística



Numa imagem, um *pixel* é selecionado ao acaso. Define-se uma v.a.  $X$  cujo valor representa seu valor digital. Qual a probabilidade deste *pixel* possuir valor entre 100 e 150?

Quais os valores possíveis de  $X$ ?

Considerando uma imagem 8 *bits*...

Mínimo 0 (região escura)

Máximo 255 (região clara)

$X: \{0, 1, \dots, 255\}$

Qual a distribuição de probabilidade de  $X$ ?

$X$  é discreto

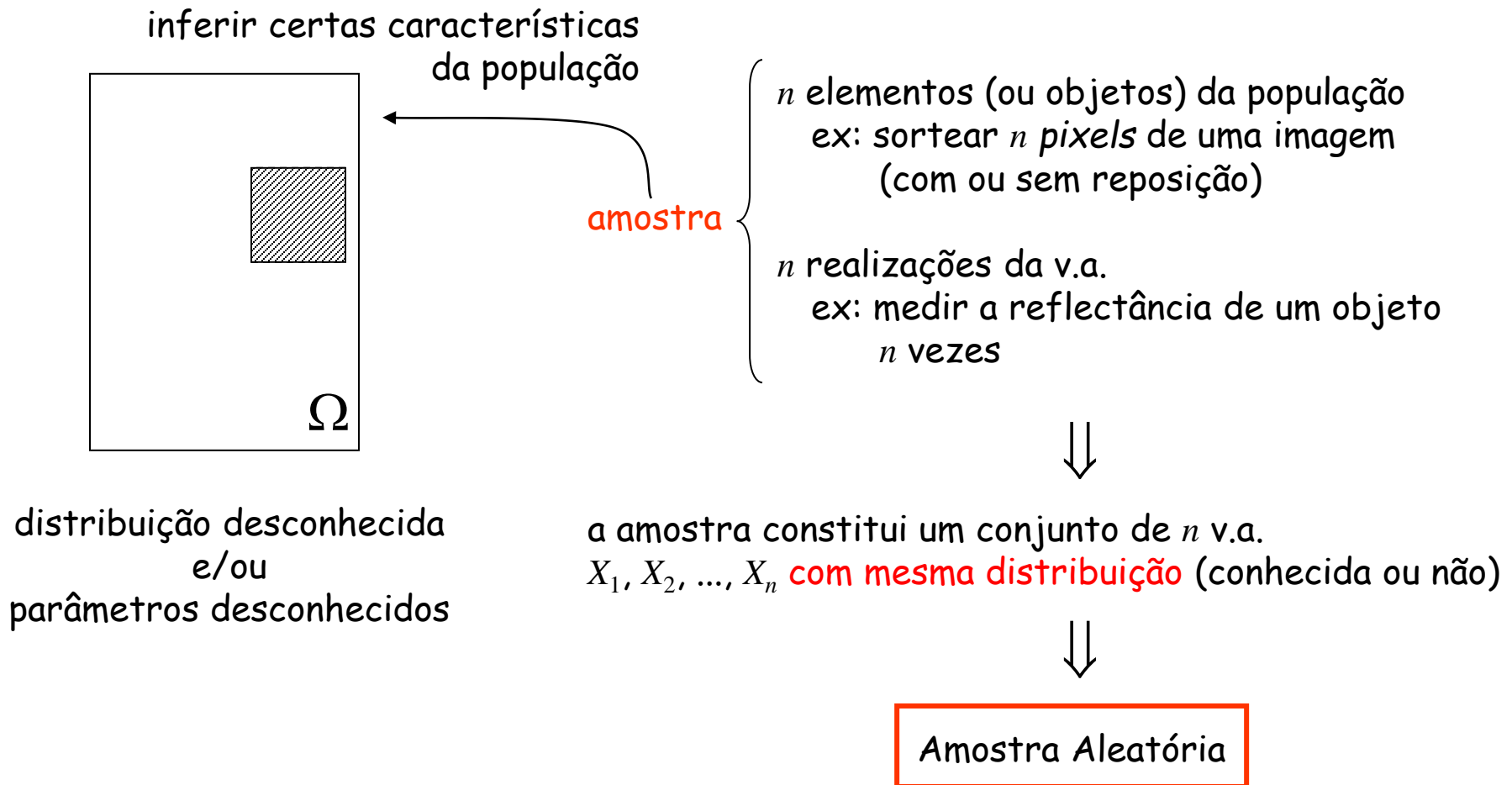
Distribuição: Desconhecida (discreta)

Que parâmetros são necessários para definir esta distribuição?

??????

**DISTRIBUIÇÃO DESCONHECIDA**

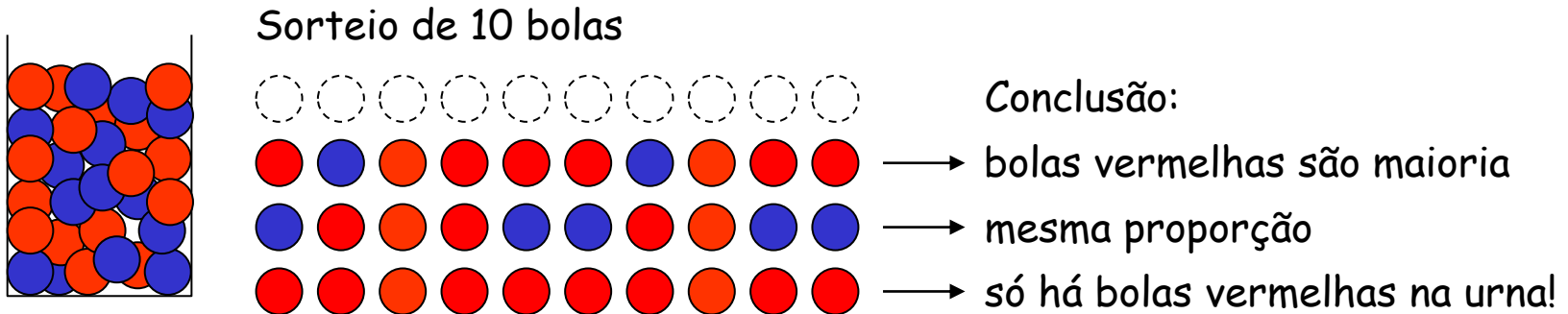
# Inferência Estatística



# Amostra Aleatória

Como uma amostra aleatória é um conjunto de  $n$  v.a.:

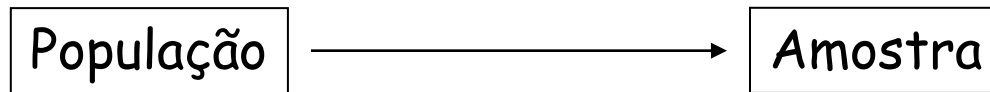
cada amostragem resulta num conjunto distinto de valores e portanto pode levar a uma conclusão distinta



Grandes questões:

- Quão representativa é a amostra disponível para a análise?  
tamanho de amostra e métodos de obtenção das amostras
- Que características devem ser observadas para representar a população?  
estimação de parâmetros
- Quão confiável é conclusão obtida pela pesquisa?  
erros

# Estimação de Parâmetros



Distribuição de Probabilidade (ou FDP)

Distribuição Amostral (Frequências)

Parâmetros  
(valor fixo) ← *estimar*

Estatísticas  
(variável aleatória)

Estimação {  
  pontual (*estatísticas*)  
  por intervalo (*intervalos de confiança*)

OBS: *estatística*: é a v.a. que estima (pontualmente) um parâmetro (populacional)  
as vezes é chamada simplesmente de *estimador*  
*estimativa*: é o valor do estimador obtido para uma amostra específica

# Distribuição Amostral

---

É muito comum utilizar um conjunto de valores observados (amostra) para tentar “enxergar” a verdadeira distribuição da população.

Esta capacidade, é claro, depende do tamanho e representatividade da amostra.

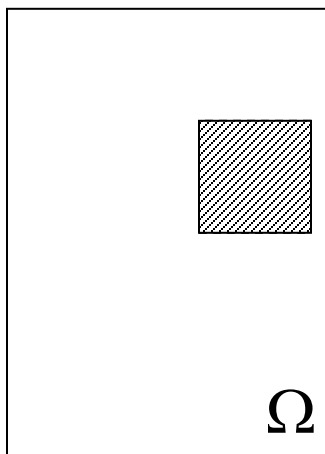
**Obs1:** Tipos de amostragem e tamanho ideal de uma amostra serão discutidos em “Teoria de Amostragem”.

**Obs2:** Testes estatísticos formais que visam comprovar se uma população segue ou não uma distribuição específica serão discutidos durante o curso.

Existem pelos menos 3 representações gráficas que podem ser utilizadas para avaliar a distribuição amostral:

- histograma (gráfico de frequências)
- frequência acumulada
- boxplot

# Estimação Pontual de um Parâmetro



amostra composta por  $n$  valores

parâmetro desconhecido  $\theta$

De que maneira os valores da amostra podem ser combinados a fim de se produzir uma "boa" estimativa desse parâmetro  $\theta$ ?

$\left\{ \begin{array}{l} \text{método dos momentos} \\ \text{método da máxima verossimilhança} \end{array} \right. \Rightarrow \text{é preciso conhecer a distribuição!}$



# Estimação Pontual de um Parâmetro

Considere que seja possível produzir  $m$  diferentes estimadores para  $\theta$ , sendo que  $\hat{\theta}_i$  representa o  $i$ -ésimo estimador de  $\theta$  ( $i = 1, \dots, m$ )

Como escolher qual estimador é melhor?

## Importante:

- lembre-se que todo estimador é uma v.a. e portanto seu valor (estimativa) varia de amostra para amostra
- dificilmente (ou é improvável) que uma amostra forneça uma estimativa igual ao parâmetro que se deseja estimar

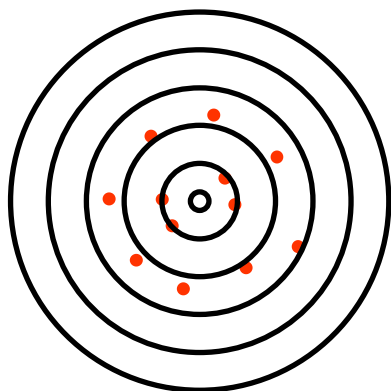
# Estimação Pontual de um Parâmetro

Considere que seja possível produzir  $m$  diferentes estimadores para  $\theta$ , sendo que  $\hat{\theta}_i$  representa o  $i$ -ésimo estimador de  $\theta$  ( $i = 1, \dots, m$ )

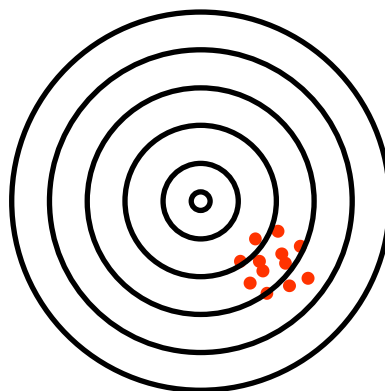
Para que  $\hat{\theta}_k$  seja o melhor, então esse estimador deveria

- ser não tendencioso (**exatidão**)  $\Rightarrow$  A média das estimativas de todas as amostras de tamanho  $n$  possíveis de serem retiradas da população é igual ao verdadeiro valor do parâmetro  
$$E(\hat{\theta}_k) = \theta$$
- ter variância mínima (**precisão**)  $\Rightarrow$  O melhor estimador irá produzir estimativas mais próximas entre si (idealmente próximas ao verdadeiro valor do parâmetro)  
$$Var(\hat{\theta}_k) < Var(\hat{\theta}_j) \quad \forall k \neq j$$

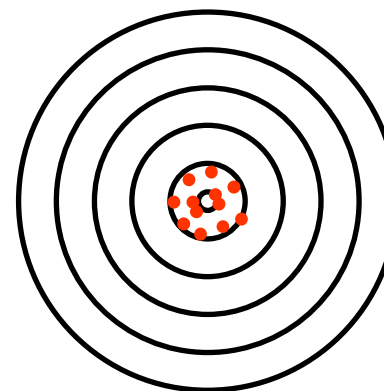
Tiro ao alvo



**Exato**  
**Impreciso**



**Inexato**  
**Preciso**



**Exato**  
**Preciso**

# Estimação Pontual de $\mu$

Seja  $X$  uma v.a. com distribuição qualquer com média ( $\mu$ ) e variância ( $\sigma^2$ ) também desconhecidas. Retira-se uma amostra de tamanho  $n$  com a finalidade de se estimar  $\mu$ .

- média populacional  $\mu$

Ex. amostra com  $n = 5$       $\{3,4; 4,5; 2,6; 3,8; 6,0\}$       $\bar{X} = \frac{3,4 + 4,5 + 2,6 + 3,8 + 6,0}{5} = 4,06$

De que maneira os valores da amostra podem ser combinados a fim de se produzir uma "boa" estimativa de  $\mu$ ?

Como não há nenhuma razão para acreditar que um valor da amostra é mais importante do que o outro:

$$\hat{\mu} = \frac{\sum_{i=1}^n x_i}{n} = \bar{X} = \underbrace{\sum_{j=1}^N x_j FR(X = x_j)}_{\text{dados agrupados (v.a. discreta)}} \quad \text{média amostral}$$

# Estimação Pontual de $\mu$

Seja  $X$  uma v.a. com distribuição qualquer com média ( $\mu$ ) e variância ( $\sigma^2$ ) também desconhecidas. Retira-se uma amostra de tamanho  $n$  com a finalidade de se estimar  $\mu$ .

- média populacional  $\mu$

Verificando a tendenciosidade de  $\bar{X}$

$$\begin{aligned} E(\bar{X}) &= E\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right) = \frac{E(X_1 + X_2 + \dots + X_n)}{n} \\ &= \frac{E(X_1) + E(X_2) + \dots + E(X_n)}{n} = \frac{\cancel{n}\mu}{\cancel{n}} = \mu \end{aligned}$$

estimador  
não tendencioso

Interpretação (teórica): se calculássemos a média dos  $\bar{X}$  de todas amostras (de tamanho  $n$ ) possíveis de serem obtidas, o resultado seria  $\mu$

# Estimação Pontual de $\mu$

Seja  $X$  uma v.a. com distribuição qualquer com média ( $\mu$ ) e variância ( $\sigma^2$ ) também desconhecidas. Retira-se uma amostra de tamanho  $n$  com a finalidade de se estimar  $\mu$ .

- média populacional  $\mu$

Calculando a variância de  $\bar{X}$  (avaliação de precisão)

$$\begin{aligned} \text{Var}(\bar{X}) &= \text{Var}\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right) = \frac{\text{Var}(X_1 + X_2 + \dots + X_n)}{n^2} \\ &= \frac{\text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_n)}{n^2} = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n} \end{aligned}$$

$\sigma^2 \quad + \quad \sigma^2 \quad + \dots + \quad \sigma^2$

Se as amostras forem independentes, ou seja, se elas não guardarem nenhuma relação entre si.

# Estimação Pontual de $\mu$

Seja  $X$  uma v.a. com distribuição qualquer com média ( $\mu$ ) e variância ( $\sigma^2$ ) também desconhecidas. Retira-se uma amostra de tamanho  $n$  com a finalidade de se estimar  $\mu$ .

- média populacional  $\mu$

$$\bar{X} = \frac{\sum_{i=1}^n x_i}{n}$$

$$E(\bar{X}) = \mu$$

$$Var(\bar{X}) = \frac{\sigma^2}{n}$$

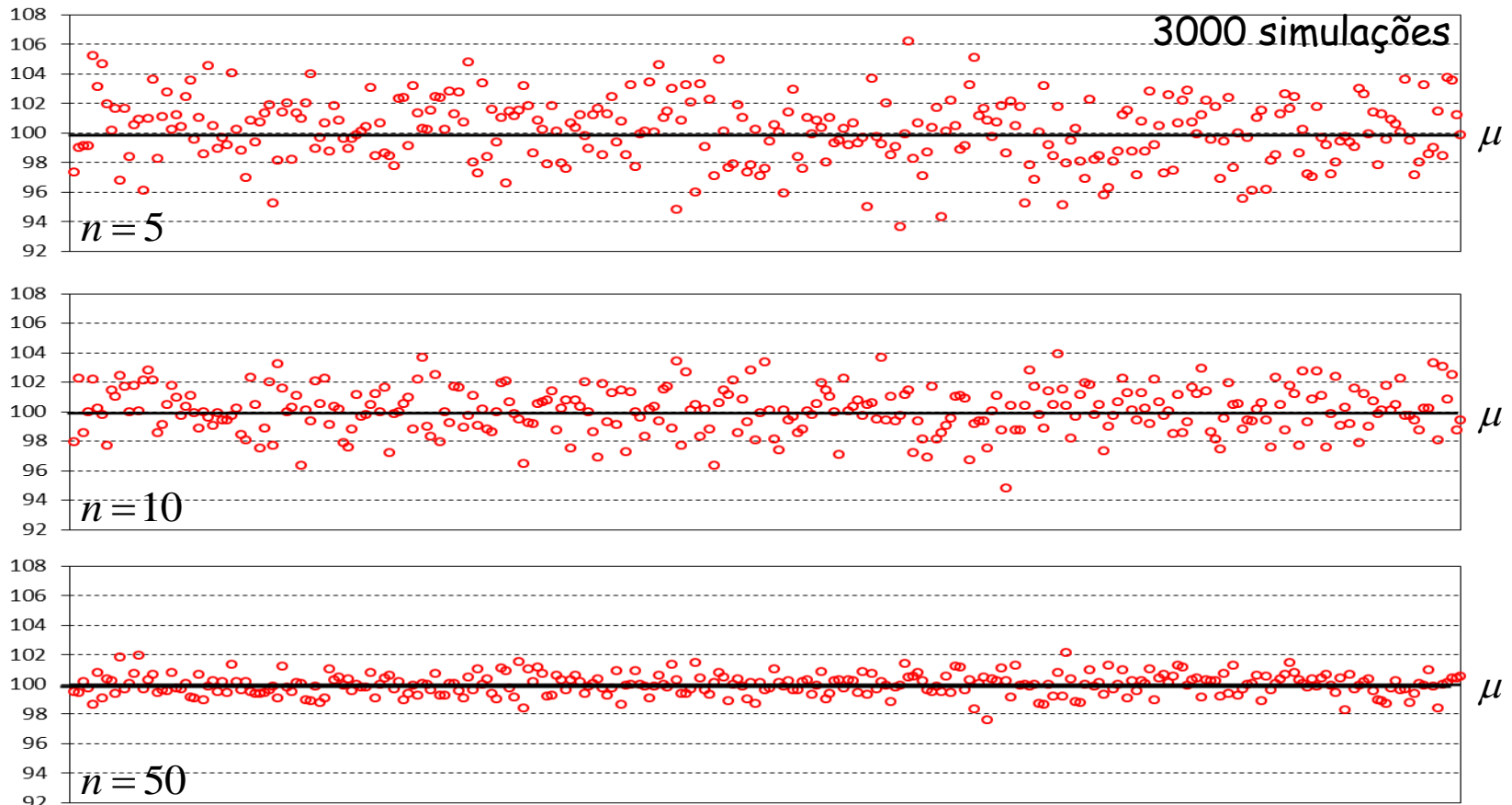
A precisão da média amostral depende da variação original dos dados ( $\sigma^2$ ) e do tamanho da amostra ( $n$ )

Quanto maior o tamanho da amostra ( $n$ ), mais precisa será a estimativa de  $\mu$

# Estimação Pontual de $\mu$

- média populacional  $\mu$      $E(\bar{X}) = \mu$      $Var(\bar{X}) = \frac{\sigma^2}{n}$

Simulando-se  $\bar{X}$  a partir de amostras de uma v.a.  $X \sim N(\mu = 100, \sigma^2 = 25)$



# Estimação Pontual de $\sigma^2$

Seja  $X$  uma v.a. com distribuição qualquer com média ( $\mu$ ) e variância ( $\sigma^2$ ) também desconhecidas. Retira-se uma amostra de tamanho  $n$  com a finalidade de se estimar  $\sigma^2$ .

- **variância populacional  $\sigma^2$**

De que maneira os valores da amostra podem ser combinados a fim de se produzir uma "boa" estimativa de  $\sigma^2$ ?

Como não há nenhuma razão para acreditar que um valor da amostra é mais importante do que o outro e  $\mu$  é desconhecido:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n}$$

Mas será um estimador tendencioso?



# Estimação Pontual de $\sigma^2$

Seja  $X$  uma v.a. com distribuição qualquer com média ( $\mu$ ) e variância ( $\sigma^2$ ) também desconhecidas. Retira-se uma amostra de tamanho  $n$  com a finalidade de se estimar  $\sigma^2$ .

- variância populacional  $\sigma^2$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n}$$
$$\sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n (X_i^2 - 2\bar{X}X_i + \bar{X}^2)$$
$$= \sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2$$
$$= \sum_{i=1}^n X_i^2 - 2n\bar{X}^2 + n\bar{X}^2$$
$$= \sum_{i=1}^n X_i^2 - n\bar{X}^2$$
$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n} \Rightarrow \sum_{i=1}^n X_i = n\bar{X}$$

# Estimação Pontual de $\sigma^2$

Seja  $X$  uma v.a. com distribuição qualquer com média ( $\mu$ ) e variância ( $\sigma^2$ ) também desconhecidas. Retira-se uma amostra de tamanho  $n$  com a finalidade de se estimar  $\sigma^2$ .

- variância populacional  $\sigma^2$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n}$$
$$E(\hat{\sigma}^2) = E\left(\frac{\sum_{i=1}^n X_i^2 - n\bar{X}^2}{n}\right) = \frac{1}{n} E\left(\sum_{i=1}^n X_i^2\right) - E(\bar{X}^2)$$
$$= \frac{1}{n} \sum_{i=1}^n E(X_i^2) - E(\bar{X}^2)$$

$$Var(X_i) = \sigma^2 = E(X_i^2) - (E(X_i))^2 = E(X_i^2) - \mu^2 \Rightarrow E(X_i^2) = \sigma^2 + \mu^2$$

# Estimação Pontual de $\sigma^2$

Seja  $X$  uma v.a. com distribuição qualquer com média ( $\mu$ ) e variância ( $\sigma^2$ ) também desconhecidas. Retira-se uma amostra de tamanho  $n$  com a finalidade de se estimar  $\sigma^2$ .

- variância populacional  $\sigma^2$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n}$$
$$E(\hat{\sigma}^2) = E\left(\frac{\sum_{i=1}^n X_i^2 - n\bar{X}^2}{n}\right) = \frac{1}{n} E\left(\sum_{i=1}^n X_i^2\right) - E(\bar{X}^2)$$
$$= \frac{1}{n} \sum_{i=1}^n E(X_i^2) - E(\bar{X}^2)$$

$$Var(\bar{X}) = \frac{\sigma^2}{n} = E(\bar{X}^2) - (E(\bar{X}))^2 = E(\bar{X}^2) - \mu^2 \Rightarrow E(\bar{X}^2) = \frac{\sigma^2}{n} + \mu^2$$

# Estimação Pontual de $\sigma^2$

Seja  $X$  uma v.a. com distribuição qualquer com média ( $\mu$ ) e variância ( $\sigma^2$ ) também desconhecidas. Retira-se uma amostra de tamanho  $n$  com a finalidade de se estimar  $\sigma^2$ .

- variância populacional  $\sigma^2$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n}$$

$$E(\hat{\sigma}^2) = E\left(\frac{\sum_{i=1}^n X_i^2 - n\bar{X}^2}{n}\right) = \frac{1}{n} E\left(\sum_{i=1}^n X_i^2\right) - E(\bar{X}^2)$$

$$E(X_i^2) = \sigma^2 + \mu^2$$

$$E(\bar{X}^2) = \frac{\sigma^2}{n} + \mu^2$$

$$= \frac{1}{n} \sum_{i=1}^n (E(X_i^2) - E(\bar{X}^2))$$

$$= \sigma^2 + \cancel{\mu^2} - \frac{\sigma^2}{n} - \cancel{\mu^2} = \frac{n\sigma^2 - \sigma^2}{n} = \frac{n-1}{n} \sigma^2$$

estimador  
tendencioso!

# Estimação Pontual de $\sigma^2$

Seja  $X$  uma v.a. com distribuição qualquer com média ( $\mu$ ) e variância ( $\sigma^2$ ) também desconhecidas. Retira-se uma amostra de tamanho  $n$  com a finalidade de se estimar  $\sigma^2$ .

- **variância populacional  $\sigma^2$**

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{\cancel{n} \quad \cancel{n} \quad n-1}$$

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n-1}$$

**variância amostral**

$$E(s^2) = \sigma^2$$

**estimador  
não tendencioso**

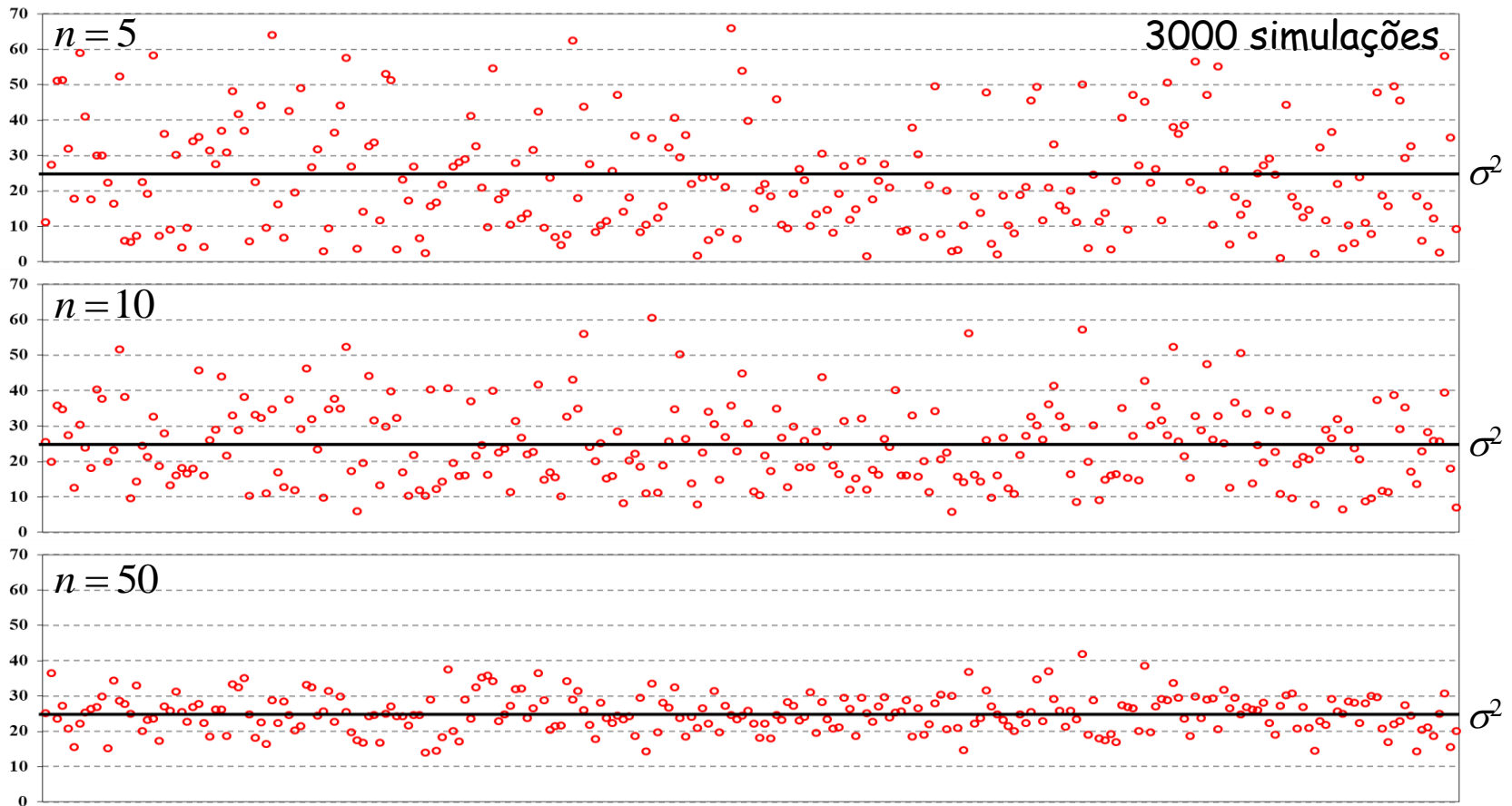
Interpretação (teórica): se calculássemos a média dos  $s^2$  de todas amostras (de tamanho  $n$ ) possíveis de serem obtidas, o resultado seria  $\sigma^2$

Curiosidade:  $Var(s^2) = \frac{2\sigma^4}{n-1}$  (precisão aumenta com o tamanho da amostra!)

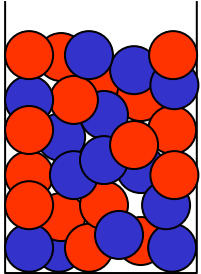
# Estimação Pontual de $\sigma^2$

- variância populacional  $\sigma^2$   $E(s^2) = \sigma^2$

Simulando-se  $s^2$  a partir de amostras de uma v.a.  $X \sim N(\mu = 100, \sigma^2 = 25)$



# Estimação Pontual de $p$



Numa urna, há  $N$  bolas, sendo  $K$  vermelhas e  $N - K$  azuis. Assim, pode-se dizer que  $K/N$  representa a proporção  $p$  de bolas vermelhas na urna (que por sua vez, representa a probabilidade de se selecionar uma bola vermelha desta urna).

Mas se  $N$  e  $K$  são desconhecidos, como estimar  $p$ ?

Considere que  $n$  bolas são escolhidas ao acaso (com reposição), definindo-se  $Y$  como o número de bolas vermelhas entre as  $n$  selecionadas, qual a distribuição de  $Y$ ?

$$Y \sim \text{Binomial}(n, p) \quad Y = \sum_{i=1}^n X_i \quad X_i \sim \text{Bernoulli} \quad p = P(X_i = 1) \Leftrightarrow P(\text{sucesso})$$

Qual é o melhor estimador pontual de  $p$ ?

$$\frac{Y}{n} = \hat{p} \quad \text{Proporção Amostral}$$

$$E(\hat{p}) = E\left(\frac{Y}{n}\right) = \frac{E(Y)}{n} = \frac{np}{n} = p \quad \begin{array}{l} \text{estimador} \\ \text{não tendencioso} \end{array}$$

$$\text{Var}(\hat{p}) = \text{Var}\left(\frac{Y}{n}\right) = \frac{\text{Var}(Y)}{n^2} = \frac{npq}{n^2} = \frac{pq}{n}$$

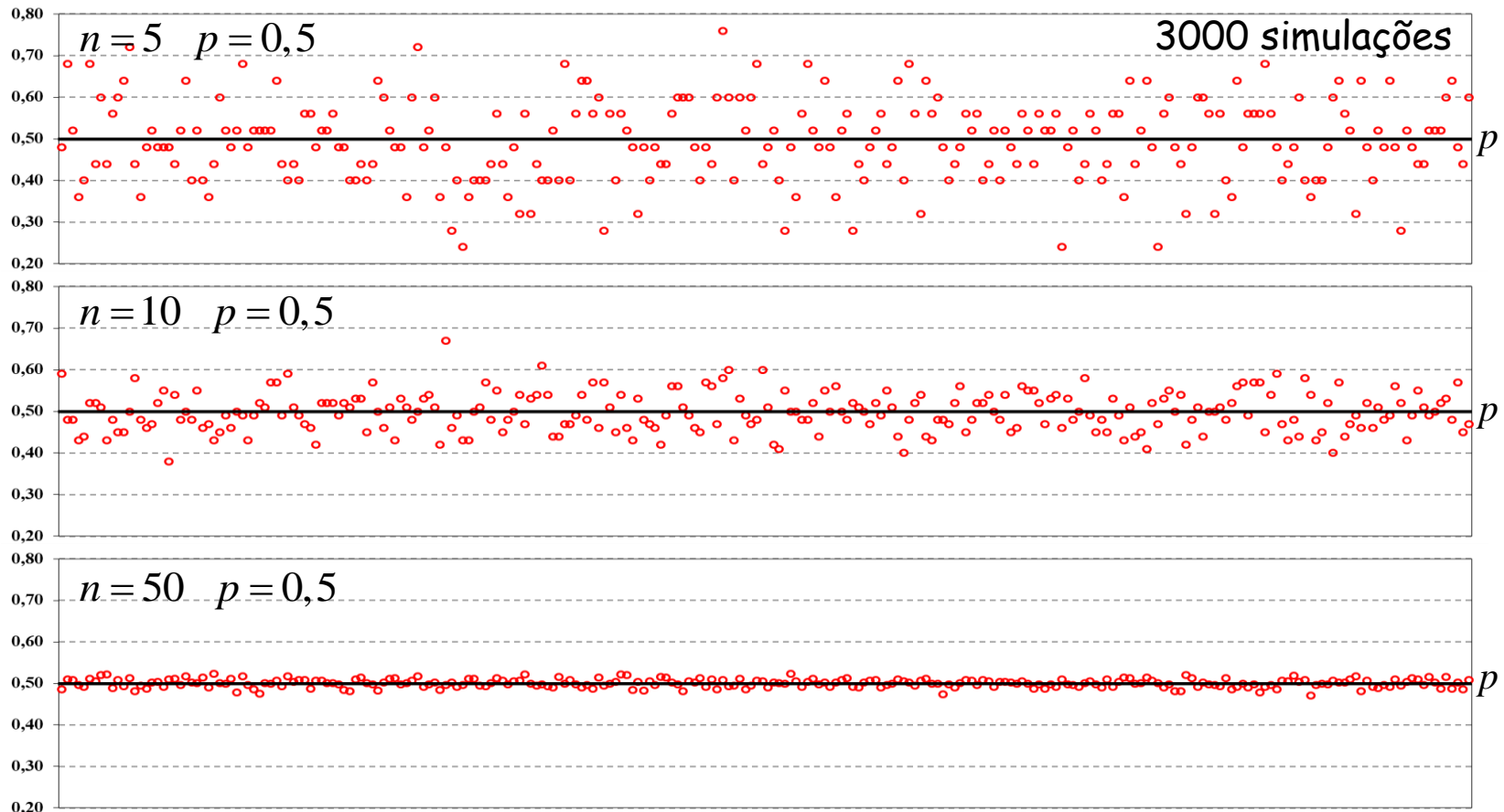
Quanto maior o tamanho da amostra ( $n$ ), mais precisa será a estimativa de  $p$

Quanto mais  $p$  se aproxima de 0,5 (50%), menos precisa será sua estimativa

# Estimação Pontual de $p$

- proporção populacional  $p$        $E(\hat{p}) = p$        $Var(\hat{p}) = \frac{pq}{n}$

Simulando-se  $\hat{p}$  a partir de amostras de uma v.a.  $X \sim \text{Binomial}(n, p = 0,5)$

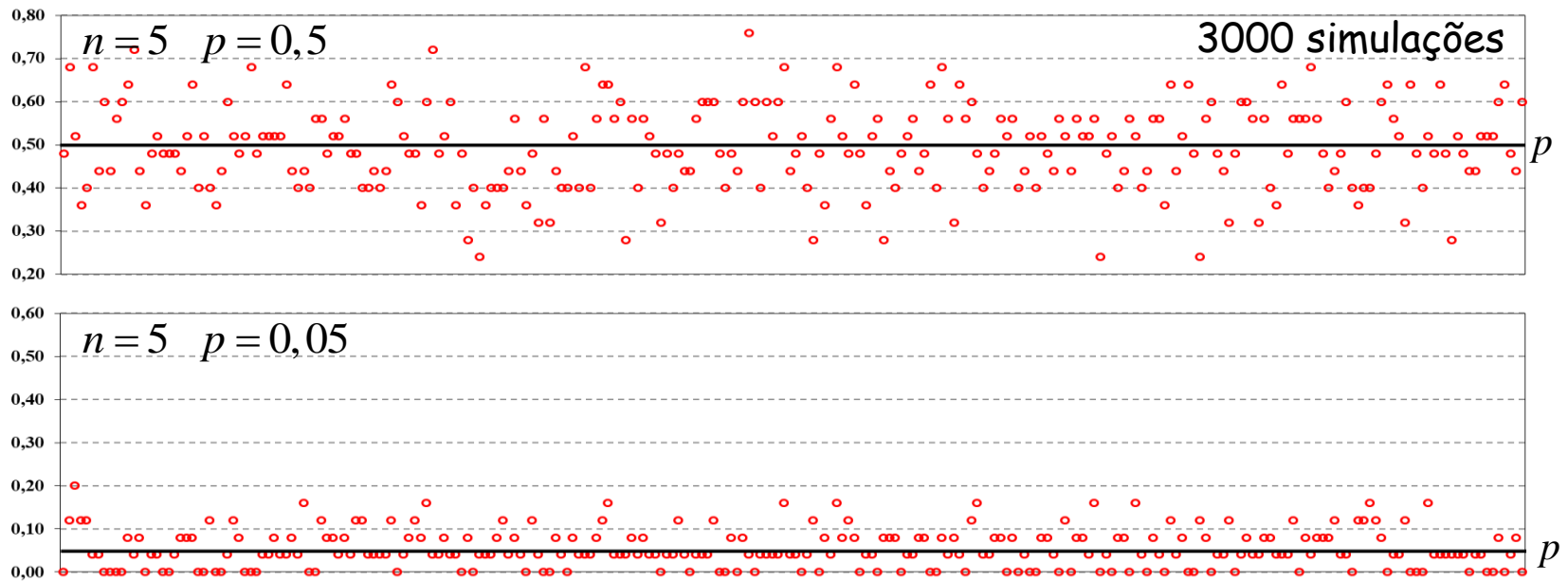




# Estimação Pontual de $p$

- proporção populacional  $p$        $E(\hat{p}) = p$        $Var(\hat{p}) = \frac{pq}{n}$

Simulando-se  $\hat{p}$  a partir de amostras de uma v.a.  $X \sim \text{Binomial}(n = 5, p)$



# Estimação Pontual de $\mu$ e $\sigma^2$

Exemplo: uma amostra ( $n = 12$ ) é retirada de uma população e os seguintes valores são observados: 0, 2, 3, 5, 2, 1, 2, 1, 3, 3, 4, 2. Calcule a média e variância amostrais.

• média amostral  $\bar{X}$

$$\bar{X} = \frac{\sum_{i=1}^n x_i}{n} \quad (\text{dados brutos})$$

$$\bar{X} = \frac{0+2+3+\dots+2}{12} = \frac{7}{3}$$

$$\bar{X} = \frac{\sum_{j=1}^N x_j FA(X = x_j)}{n} = \sum_{j=1}^N x_j FR(X = x_j) \quad (\text{dados agrupados})$$

$$\bar{X} = \frac{0*1+1*2+2*4+3*3+4*1+5*1}{12} = \frac{7}{3} \quad (\text{usando FA})$$

$$\bar{X} = 0*\frac{1}{12} + 1*\frac{1}{6} + 2*\frac{1}{3} + 3*\frac{1}{4} + 4*\frac{1}{12} + 5*\frac{1}{12} = \frac{7}{3} \quad (\text{usando FR})$$

distribuição amostral

Valor	Freq. Absoluta	Freq. Relativa
0	1	1/12
1	2	1/6
2	4	1/3
3	3	1/4
4	1	1/12
5	1	1/12
Total	12	1

# Estimação Pontual de $\mu$ e $\sigma^2$

Exemplo: uma amostra ( $n = 12$ ) é retirada de uma população e os seguintes valores são observados: 0, 2, 3, 5, 2, 1, 2, 1, 3, 3, 4, 2. Calcule a média e variância amostrais.

• variância amostral  $s^2$   $\bar{X} = \frac{7}{3}$

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n-1} = \frac{\sum_{i=1}^n x_i^2 - n\bar{X}^2}{n-1} \quad (\text{dados brutos})$$

$$s^2 = \frac{(0 - \frac{7}{3})^2 + (2 - \frac{7}{3})^2 + \dots + (2 - \frac{7}{3})^2}{11} = \frac{(0^2 + 2^2 + \dots + 2^2) - 12 * (\frac{7}{3})^2}{11} = 1,88$$

$$s^2 = \frac{\sum_{j=1}^N (x_j - \bar{X})^2 FA(X = x_j)}{n-1} = \frac{\sum_{j=1}^N x_j^2 FA(X = x_j) - n\bar{X}^2}{n-1} \quad (\text{dados agrupados})$$

$$s^2 = \frac{(0 - \frac{7}{3})^2 * 1 + (1 - \frac{7}{3})^2 * 2 + \dots + (5 - \frac{7}{3})^2 * 1}{11} = \frac{(0^2 * 1 + 1^2 * 2 + \dots + 5^2 * 1) - 12 * (\frac{7}{3})^2}{11} = 1,88$$

distribuição amostral

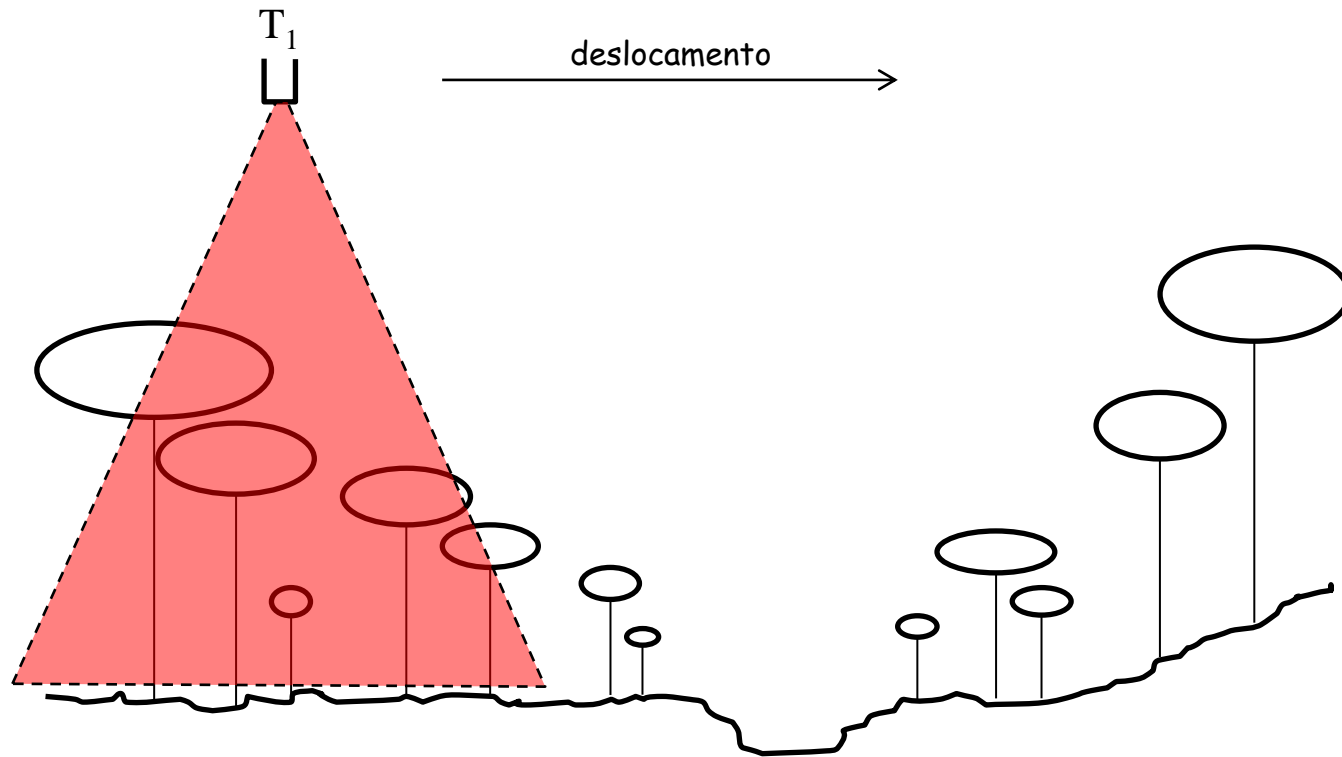
Valor	Freq. Absoluta	Freq. Relativa
0	1	1/12
1	2	1/6
2	4	1/3
3	3	1/4
4	1	1/12
5	1	1/12
Total	12	1

# Estimação Pontual de $\mu$ , $\sigma^2$ e $p$

Observações:

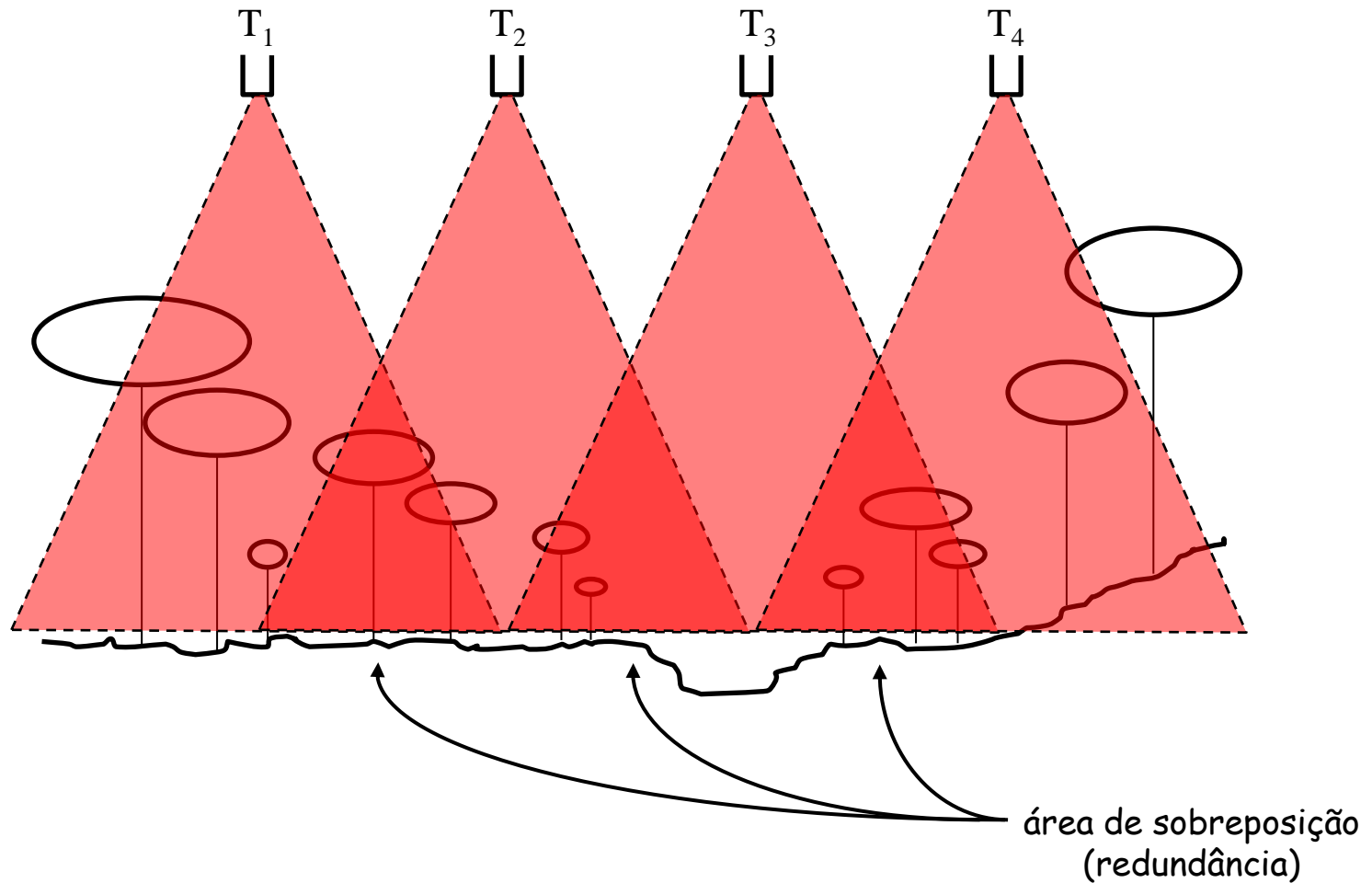
- $\mu$ ,  $\sigma^2$  e  $p$  são parâmetros que representam a população e portanto são valores fixos sendo, em geral, desconhecidos
- $\bar{X}$ ,  $s^2$  e  $\hat{p}$  são estatísticas calculadas a partir da amostra e representam variáveis aleatórias (cada conjunto de amostras pode apresentar um valor diferente)
- Não confunda **variância amostral** ( $s^2$ ) com **variância da média amostral** ( $\text{Var}(\bar{X})$ )
- De modo geral, as amostras devem ser obtidas de modo independente uma das outras, ou seja, o valor de uma amostra não deve ter relação com o(s) valor(es) das outras amostras (exceção em estudos de séries temporais ou dados espaciais, onde estuda-se exatamente esta relação)

# Utilização de amostras não independentes

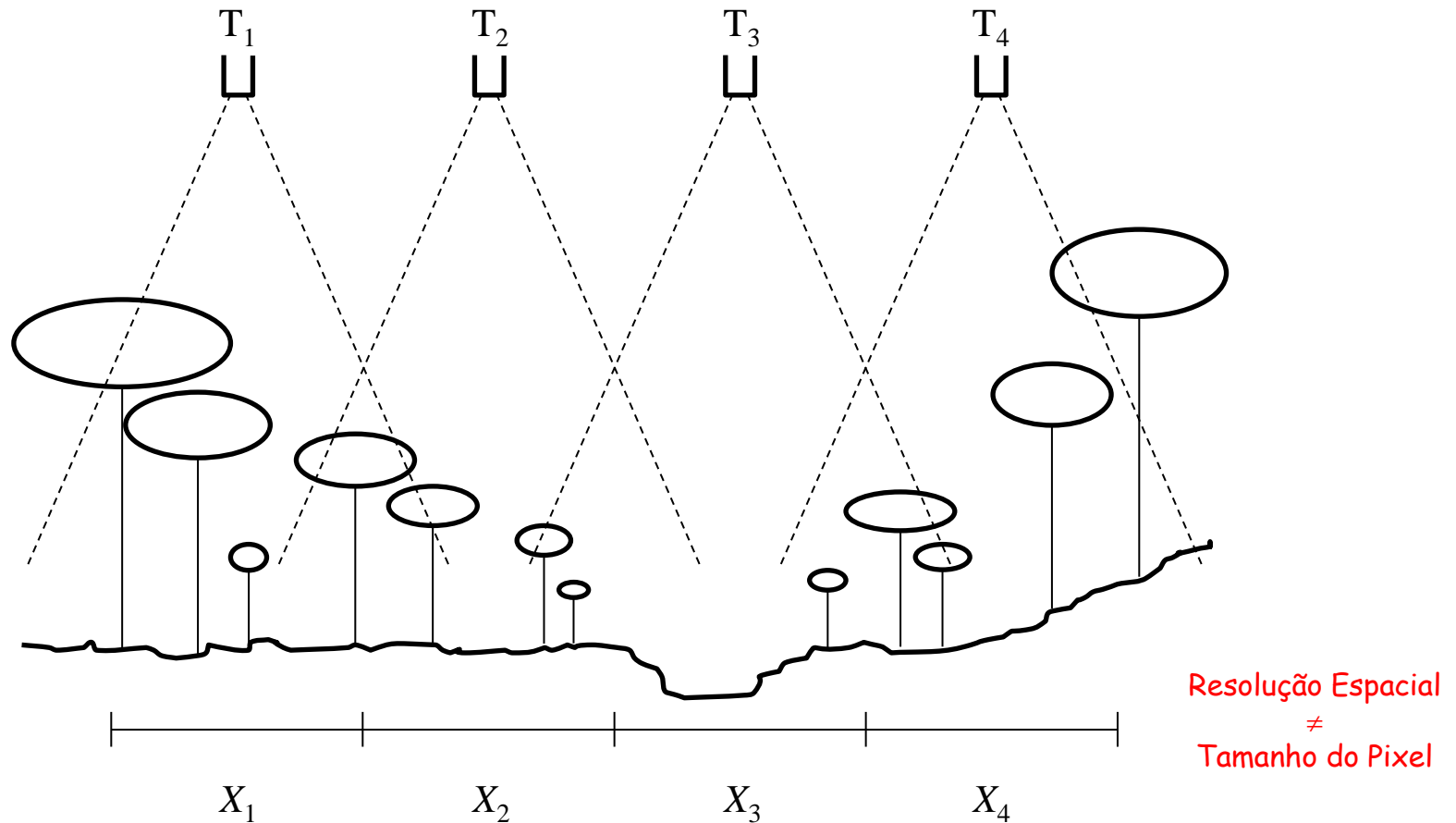


A resposta do sensor representa a integração das respostas de todos objetos que estão no campo de visada

# Utilização de amostras não independentes



# Utilização de amostras não independentes



$X_i$ : valor que representa a resposta do sensor no tempo  $T_i$  ( $\equiv$  elemento de resolução)

Note que estes valores **não são independentes** (devido a sobreposição)

# Utilização de amostras não independentes

Suponha que  $X$  representa um conjunto de amostras **independentes** de uma v.a. qualquer obtidas numa determinada sequência (série temporal por exemplo)

$X$	$X'$	$X''$	$X'''$
9			
6			
2			
3	3,2		
6			
2			
6			
10			
6			
5			
7			
1			
7			
8			
5			

$X'$ ,  $X''$  e  $X'''$  resultam do cálculo de médias móveis (tamanho 3) aplicado sobre  $X$

$$X'_i = 0,1X_{i-1} + 0,8X_i + 0,1X_{i+1}$$



# Utilização de amostras não independentes

Suponha que  $X$  representa um conjunto de amostras **independentes** de uma v.a. qualquer obtidas numa determinada sequência (série temporal por exemplo)

$X$	$X'$	$X''$	$X'''$
9			
6	5,9		
2	2,5		
3	3,2		
6	5,3		
2	2,8		
6	6		
10	9,2		
6	6,3		
5	5,3		
7	6,2		
1	2,2		
7	6,5		
8	7,6		
5			

$X'$ ,  $X''$  e  $X'''$  resultam do cálculo de médias móveis (tamanho 3) aplicado sobre  $X$

$$X'_i = 0,1X_{i-1} + 0,8X_i + 0,1X_{i+1}$$

$$X''_i = 0,2X_{i-1} + 0,6X_i + 0,2X_{i+1}$$

$$X'''_i = (X_{i-1} + X_i + X_{i+1}) / 3$$

# Utilização de amostras não independentes

Suponha que  $X$  representa um conjunto de amostras **independentes** de uma v.a. qualquer obtidas numa determinada sequência (série temporal por exemplo)

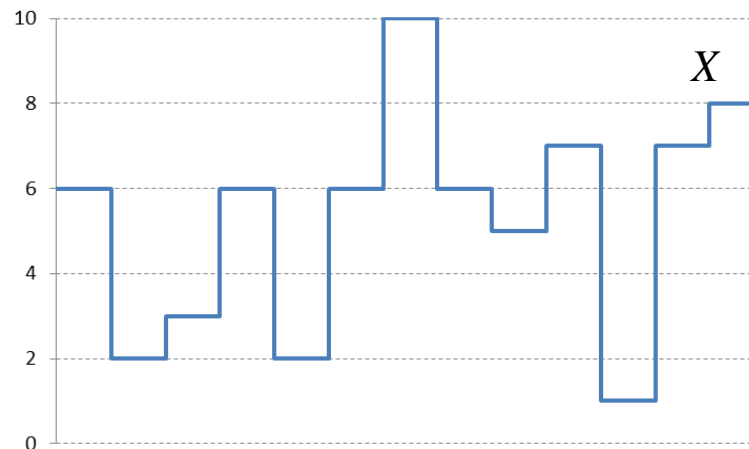
$X$	$X'$	$X''$	$X'''$
<del>9</del>			
6	5,9	5,8	5,67
2	2,5	3	3,67
3	3,2	3,4	3,67
6	5,3	4,6	3,67
2	2,8	3,6	4,67
6	6	6	6
10	9,2	8,4	7,33
6	6,3	6,6	7
5	5,3	5,6	6
7	6,2	5,4	4,33
1	2,2	3,4	5
7	6,5	6	5,33
8	7,6	7,2	6,67
<del>5</del>			

$X'$ ,  $X''$  e  $X'''$  resultam do cálculo de médias móveis (tamanho 3) aplicado sobre  $X$

$$X'_i = 0,1X_{i-1} + 0,8X_i + 0,1X_{i+1}$$

$$X''_i = 0,2X_{i-1} + 0,6X_i + 0,2X_{i+1}$$

$$X'''_i = (X_{i-1} + X_i + X_{i+1}) / 3$$



# Utilização de amostras não independentes

Suponha que  $X$  representa um conjunto de amostras **independentes** de uma v.a. qualquer obtidas numa determinada sequência (série temporal por exemplo)

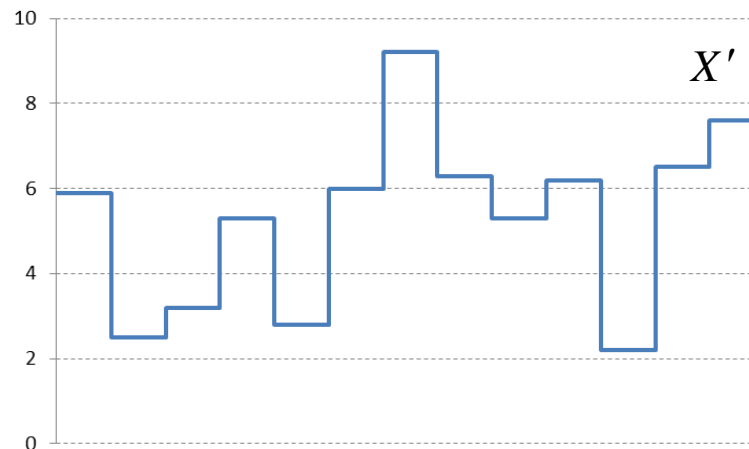
$X$	$X'$	$X''$	$X'''$
<del>9</del>			
6	5,9	5,8	5,67
2	2,5	3	3,67
3	3,2	3,4	3,67
6	5,3	4,6	3,67
2	2,8	3,6	4,67
6	6	6	6
10	9,2	8,4	7,33
6	6,3	6,6	7
5	5,3	5,6	6
7	6,2	5,4	4,33
1	2,2	3,4	5
7	6,5	6	5,33
8	7,6	7,2	6,67
<del>5</del>			

$X'$ ,  $X''$  e  $X'''$  resultam do cálculo de médias móveis (tamanho 3) aplicado sobre  $X$

$$X'_i = 0,1X_{i-1} + 0,8X_i + 0,1X_{i+1}$$

$$X''_i = 0,2X_{i-1} + 0,6X_i + 0,2X_{i+1}$$

$$X'''_i = (X_{i-1} + X_i + X_{i+1}) / 3$$



# Utilização de amostras não independentes

Suponha que  $X$  representa um conjunto de amostras **independentes** de uma v.a. qualquer obtidas numa determinada sequência (série temporal por exemplo)

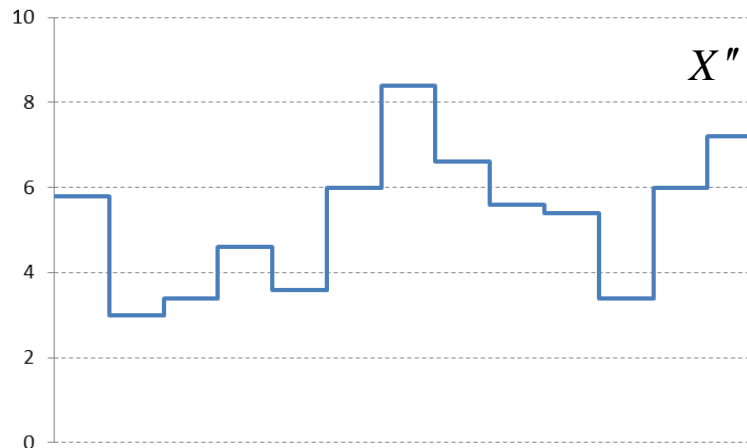
$X$	$X'$	$X''$	$X'''$
<del>9</del>			
6	5,9	5,8	5,67
2	2,5	3	3,67
3	3,2	3,4	3,67
6	5,3	4,6	3,67
2	2,8	3,6	4,67
6	6	6	6
10	9,2	8,4	7,33
6	6,3	6,6	7
5	5,3	5,6	6
7	6,2	5,4	4,33
1	2,2	3,4	5
7	6,5	6	5,33
8	7,6	7,2	6,67
<del>5</del>			

$X'$ ,  $X''$  e  $X'''$  resultam do cálculo de médias móveis (tamanho 3) aplicado sobre  $X$

$$X'_i = 0,1X_{i-1} + 0,8X_i + 0,1X_{i+1}$$

$$X''_i = 0,2X_{i-1} + 0,6X_i + 0,2X_{i+1}$$

$$X'''_i = (X_{i-1} + X_i + X_{i+1}) / 3$$



# Utilização de amostras não independentes

Suponha que  $X$  representa um conjunto de amostras **independentes** de uma v.a. qualquer obtidas numa determinada sequência (série temporal por exemplo)

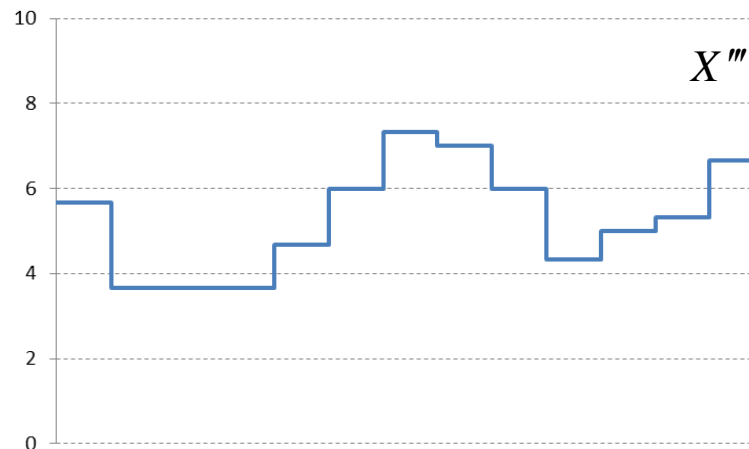
$X$	$X'$	$X''$	$X'''$
<del>9</del>			
6	5,9	5,8	5,67
2	2,5	3	3,67
3	3,2	3,4	3,67
6	5,3	4,6	3,67
2	2,8	3,6	4,67
6	6	6	6
10	9,2	8,4	7,33
6	6,3	6,6	7
5	5,3	5,6	6
7	6,2	5,4	4,33
1	2,2	3,4	5
7	6,5	6	5,33
8	7,6	7,2	6,67
<del>5</del>			

$X'$ ,  $X''$  e  $X'''$  resultam do cálculo de médias móveis (tamanho 3) aplicado sobre  $X$

$$X'_i = 0,1X_{i-1} + 0,8X_i + 0,1X_{i+1}$$

$$X''_i = 0,2X_{i-1} + 0,6X_i + 0,2X_{i+1}$$

$$X'''_i = (X_{i-1} + X_i + X_{i+1}) / 3$$



# Utilização de amostras não independentes

Suponha que  $X$  representa um conjunto de amostras **independentes** de uma v.a. qualquer obtidas numa determinada sequência (série temporal por exemplo)

$X$	$X'$	$X''$	$X'''$
<del>9</del>			
6	5,9	5,8	5,67
2	2,5	3	3,67
3	3,2	3,4	3,67
6	5,3	4,6	3,67
2	2,8	3,6	4,67
6	6	6	6
10	9,2	8,4	7,33
6	6,3	6,6	7
5	5,3	5,6	6
7	6,2	5,4	4,33
1	2,2	3,4	5
7	6,5	6	5,33
8	7,6	7,2	6,67
<del>5</del>			

$X'$ ,  $X''$  e  $X'''$  resultam do cálculo de médias móveis (tamanho 3) aplicado sobre  $X$

$$X'_i = 0,1X_{i-1} + 0,8X_i + 0,1X_{i+1}$$

$$X''_i = 0,2X_{i-1} + 0,6X_i + 0,2X_{i+1}$$

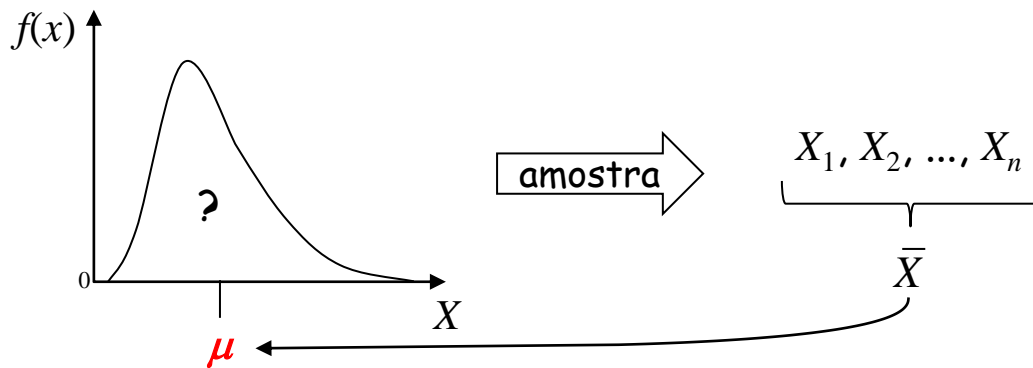
$$X'''_i = (X_{i-1} + X_i + X_{i+1}) / 3$$

	$\bar{X}$	$s^2$
$X$	5,31	6,90
$X'$	5,31	4,39
$X''$	5,31	2,68
$X'''$	5,31	1,62

Conclusão: a utilização de amostras não independentes (autocorrelacionadas) afetam mais a estimação da variância do que a estimação da média

# Distribuições amostrais

Um parâmetro pode ser estimado através de um único valor (estimador pontual)



Se amostras de tamanho  $n$  fossem obtidas e para cada uma fosse calculada  $\bar{X}$ , poderíamos esperar que todas tivessem o mesmo valor? **Muito pouco provável!**

Como o estimador é uma v.a., então há, pelo menos teoricamente, uma distribuição associada a esse estimador.

Conhecer essas distribuições é fundamental para se entender o quão próximas ou distintas poderão ser as estimativas obtidas para as diferentes amostras, ou seja, entender qual a relação existente entre o estimador e o parâmetro que se deseja estimar.

# Distribuição amostral relacionada com $\bar{X}$

$\{X_1, X_2, \dots, X_n\}$  amostra aleatória

$X_i \sim ?(\mu, \sigma^2)$  distribuição desconhecida,  $\mu$  desconhecido, mas  $\sigma^2$  conhecido

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

Se  $X_i \sim N(\mu, \sigma^2)$ :  $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$       $E(\bar{X}) = \mu$       $Var(\bar{X}) = \frac{\sigma^2}{n}$

Se  $n$  for grande (ou seja, adotando-se o TLC):

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

$$\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0,1)$$

é necessário  
conhecer  $\sigma^2$

se  $X$  tiver distribuição Normal ou  
se  $n$  for grande (TLC válida)

Conclusão: sempre que precisarmos entender a relação entre  $\bar{X}$  e  $\mu$ , iremos usar a distribuição Normal Padrão, desde que  $\sigma^2$  seja conhecida.



# Distribuição amostral relacionada com $s^2$

$X = \{X_1, X_2, \dots, X_n\}$  amostra aleatória

Se  $X_i \sim N(\mu, \sigma^2)$  distribuição **normal** com  $\mu$  e  $\sigma^2$  desconhecidos

$$\frac{X_i - \mu}{\sigma} \sim N(0,1) \quad \frac{(X_i - \mu)^2}{\sigma^2} \sim \chi_1^2$$

$$\frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2} \sim \chi_n^2 \quad \Rightarrow \text{precisaria conhecer } \mu!$$

Substituindo-se  $\mu$  por  $\bar{X}$  tem-se que

$$\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} \sim \chi_{n-1}^2 \quad (\text{perde-se 1 grau de liberdade})$$

$$\text{mas } s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1} \Rightarrow \sum_{i=1}^n (X_i - \bar{X})^2 = (n-1)s^2 \Rightarrow \boxed{\frac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2}$$

Conclusão: sempre que precisarmos entender a relação entre  $s^2$  e  $\sigma^2$ , iremos usar a distribuição Qui-quadrado (neste caso, é necessário que  **$X \sim \text{Normal}$** ).

# Graus de liberdade

De modo geral, pode-se entender “grau de liberdade” como o número de valores que, no final de um cálculo de uma estatística, estão “livres para variarem”, ou seja, que têm o comportamento de variáveis aleatórias.

Por exemplo: deseja-se avaliar os desvios em torno da média a partir de uma amostra de 3 valores retirados de uma população normalmente distribuída.

$$X \sim N(\mu, \sigma^2) \quad \text{Amostra: } \{X_1 - \mu, X_2 - \mu, X_3 - \mu\} \quad (\mu \text{ é conhecida})$$

Quais são os valores possíveis de serem obtidos nesta amostra?

R: Neste caso, posso escolher “livremente” quaisquer 3 valores entre  $-\infty$  e  $+\infty$

$$\underline{-3,5} \quad \underline{0,5} \quad \underline{100,9}$$

# Graus de liberdade

Agora, se  $\mu$  é desconhecido e o substituímos por  $\bar{X}$  ...

$$\text{Amostra: } \{X_1 - \bar{X}, X_2 - \bar{X}, X_3 - \bar{X}\}$$

$$\text{Como } \bar{X} = (X_1 + X_2 + X_3) / 3 \text{ então } \sum_{i=1}^3 (X_i - \bar{X}) = 0$$

Assim, ao se escolher os dois primeiros valores, o terceiro é necessariamente conhecido. Neste caso, perde-se 1 grau de liberdade

$$\begin{array}{ccc} -1,5 & 0,1 & \underline{1,4} \\ \hline \end{array} \quad \begin{array}{l} \nearrow \\ -1,5 + 0,1 + (X_3 - \bar{X}) = 0 \end{array}$$

As perdas de graus de liberdade acontecem sempre que um parâmetro é substituído por seu estimador

# Distribuição amostral relacionada com $\bar{X}$

$X = \{X_1, X_2, \dots, X_n\}$  amostra aleatória

Se  $X_i \sim N(\mu, \sigma^2)$  distribuição **normal** com  $\mu$  e  $\sigma^2$  **desconhecidos**

$$\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0,1) \quad \text{Lembrete: } \bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2 \quad \frac{\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}}{\sqrt{\frac{(n-1)s^2}{(n-1)\sigma^2}}} = \frac{\frac{\bar{X} - \mu}{\cancel{\sigma}/\sqrt{n}}}{\cancel{\sigma}} = \boxed{\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \sim t_{n-1}} \quad \text{Se } n \text{ é grande } (n > 100):$$
$$\frac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \sim N(0,1)$$

Conclusão: sempre que precisarmos entender a relação entre  $\bar{X}$  e  $\mu$  mas  $\sigma^2$  for desconhecida, iremos usar a distribuição *t-Student* (neste caso, é necessário que  $X \sim \text{Normal}$ ). Se  $n$  for grande, pode-se usar a Normal Padrão.

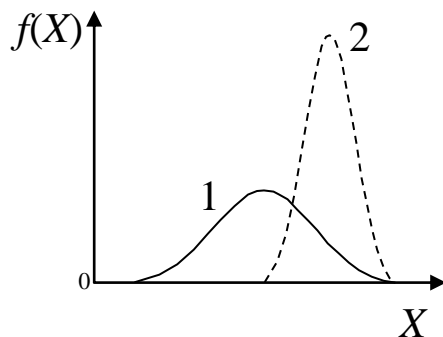
# Distribuição amostral relacionada com $\bar{X}_1$ e $\bar{X}_2$

$$X_1 = \{X_{1,1}, X_{1,2}, \dots, X_{1,n_1}\}$$

2 amostras aleatórias independentes

$$X_2 = \{X_{2,1}, X_{2,2}, \dots, X_{2,n_2}\}$$

$$X_{1,i} \sim N(\mu_1, \sigma_1^2) \quad X_{2,i} \sim N(\mu_2, \sigma_2^2) \quad \mu_j \text{ desconhecidas, mas } \sigma_j^2 \text{ conhecidas}$$



$$\bar{X}_1 \sim N\left(\mu_1, \frac{\sigma_1^2}{n_1}\right)$$

Quão próximos estão  $\mu_1$  e  $\mu_2$  ?

$$\bar{X}_2 \sim N\left(\mu_2, \frac{\sigma_2^2}{n_2}\right)$$

$$\bar{X}_1 - \bar{X}_2 \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$$

$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0,1)$$

Normal Padrão

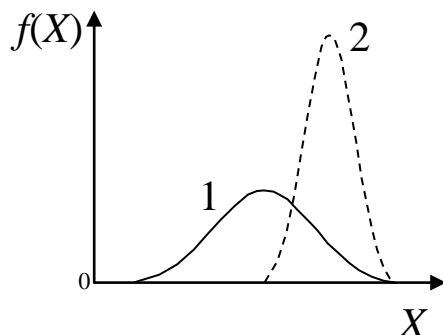
# Distribuição amostral relacionada com $\bar{X}_1$ e $\bar{X}_2$

$$X_1 = \{X_{1,1}, X_{1,2}, \dots, X_{1,n_1}\}$$

2 amostras aleatórias independentes

$$X_2 = \{X_{2,1}, X_{2,2}, \dots, X_{2,n_2}\}$$

$$X_{1,i} \sim N(\mu_1, \sigma_1^2) \quad X_{2,i} \sim N(\mu_2, \sigma_2^2) \quad \mu_j \text{ e } \sigma_j^2 \text{ desconhecidas}$$



$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0,1)$$

$$\frac{(n_1 - 1)s_1^2}{\sigma_1^2} + \frac{(n_2 - 1)s_2^2}{\sigma_2^2} \sim \chi_{n_1+n_2-2}^2$$

$$\frac{\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}}{\sqrt{\frac{\frac{(n_1 - 1)s_1^2}{\sigma_1^2} + \frac{(n_2 - 1)s_2^2}{\sigma_2^2}}{n_1 + n_2 - 2}}} \sim t_{n_1+n_2-2}$$

a princípio sem solução pois  $\sigma_1^2$   
e  $\sigma_2^2$  são desconhecidos!

# Distribuição amostral relacionada com $\bar{X}_1$ e $\bar{X}_2$

$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim t_{n_1+n_2-2}$$

$$\sqrt{\frac{(n_1-1)s_1^2}{\sigma_1^2} + \frac{(n_2-1)s_2^2}{\sigma_2^2}} \sim t_{n_1+n_2-2}$$

(fazendo  $\sigma_1^2 = \sigma_2^2 \Rightarrow \sigma^2$ )  
abordagem homocedástica

$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\cancel{\sigma} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1+n_2-2}$$

$$\cancel{\sigma} \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}} \sim t_{n_1+n_2-2}$$

rearranjando os termos...

$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \sim t_{n_1+n_2-2}$$

# Distribuição amostral relacionada com $\bar{X}_1$ e $\bar{X}_2$

$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0,1)$$

(considerando  $\sigma_1^2 \neq \sigma_2^2$ )  
abordagem heterocedástica

$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \sim t_g$$

$$g \approx \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{\left(\frac{s_1^2}{n_1}\right)^2}{n_1 - 1} + \frac{\left(\frac{s_2^2}{n_2}\right)^2}{n_2 - 1}}$$

**Importante:** a seleção de qual abordagem (homo ou heterocedástica) deve ser adotada é feita verificando-se previamente se as variâncias populacionais podem ou não ser consideradas iguais (teste de hipóteses)



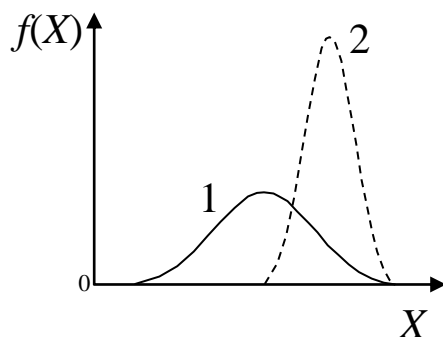
# Distribuição amostral relacionada com $s_1^2$ e $s_2^2$

$$X_1 = \{X_{1,1}, X_{1,2}, \dots, X_{1,n_1}\}$$

2 amostras aleatórias independentes

$$X_2 = \{X_{2,1}, X_{2,2}, \dots, X_{2,n_2}\}$$

$$X_{1,i} \sim N(\mu_1, \sigma_1^2) \quad X_{2,i} \sim N(\mu_2, \sigma_2^2) \quad \mu_j \text{ e } \sigma_j^2 \text{ desconhecidas}$$



Quão semelhantes são  $\sigma_1^2$  e  $\sigma_2^2$ ?

$$\frac{(n_1 - 1)s_1^2}{\sigma_1^2} \sim \chi_{n_1-1}^2$$

$$\frac{(n_2 - 1)s_2^2}{\sigma_2^2} \sim \chi_{n_2-1}^2$$

$$\frac{\frac{(n_1 - 1)s_1^2}{\sigma_1^2}}{\frac{(n_2 - 1)s_2^2}{\sigma_2^2}} \sim F_{n_1-1, n_2-1}$$

$$\frac{\cancel{(n_1 - 1)}s_1^2}{\cancel{(n_1 - 1)}\sigma_1^2} = \frac{s_1^2}{\sigma_1^2} \frac{\sigma_2^2}{s_2^2}$$

$$\frac{\cancel{(n_2 - 1)}s_2^2}{\cancel{(n_2 - 1)}\sigma_2^2}$$

$$\frac{s_1^2 \sigma_2^2}{s_2^2 \sigma_1^2} \sim F_{n_1-1, n_2-1}$$

# Distribuição amostral relacionada com $\hat{p}$

$$Y \sim \text{Binomial}(n, p) \quad Y = \sum_{i=1}^n X_i \quad X_i \sim \text{Bernoulli} \quad p = P(X_i = 1) \Leftrightarrow P(\text{sucesso})$$

$$\frac{Y}{n} = \hat{p} \quad \text{Proporção Amostral} \quad E(\hat{p}) = p \quad \text{Var}(\hat{p}) = \frac{pq}{n}$$

Qual a distribuição de  $\hat{p}$  ?

Se  $n$  for grande (ou seja, adotando-se o TLC):

$$\hat{p} \sim N\left(p, \frac{pq}{n}\right)$$

$$\frac{\hat{p} - p}{\sqrt{\frac{pq}{n}}} \sim N(0,1)$$

# Distribuição amostral relacionada com $\hat{p}_1$ e $\hat{p}_2$

$$Y_1 \sim \text{Binomial}(n_1, p_1)$$

$$Y_2 \sim \text{Binomial}(n_2, p_2)$$

$$\frac{Y_1}{n_1} = \hat{p}_1 \quad E(\hat{p}_1) = p_1 \quad \text{Var}(\hat{p}_1) = \frac{p_1 q_1}{n_1}$$

$$\frac{Y_2}{n_2} = \hat{p}_2 \quad E(\hat{p}_2) = p_2 \quad \text{Var}(\hat{p}_2) = \frac{p_2 q_2}{n_2}$$

Se  $n_1$  e  $n_2$  forem grandes (ou seja, adotando-se o TLC):

$$\hat{p}_1 \sim N\left(p_1, \frac{p_1 q_1}{n_1}\right) \quad \hat{p}_2 \sim N\left(p_2, \frac{p_2 q_2}{n_2}\right)$$

$$\frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} \sim N(0, 1)$$

# Distribuições amostrais (Resumo)

$$\text{para } \bar{X} \begin{cases} N(0,1) & \text{se } \sigma^2 \text{ é conhecida} \\ t_{n-1} & \text{se } \sigma^2 \text{ é desconhecida} \end{cases}$$

$$\text{para } s^2 \begin{cases} \chi_{n-1}^2 \end{cases}$$

$$\text{para } \bar{X}_1 - \bar{X}_2 \begin{cases} N(0,1) & \text{se } \sigma_1^2 \text{ e } \sigma_2^2 \text{ são conhecidas} \\ t_{n_1+n_2-2} & \text{se } \sigma_1^2 \text{ e } \sigma_2^2 \text{ são desconhecidas, mas } \sigma_1^2 = \sigma_2^2 \\ t_g & \text{se } \sigma_1^2 \text{ e } \sigma_2^2 \text{ são desconhecidas, mas } \sigma_1^2 \neq \sigma_2^2 \end{cases}$$

$$\text{para } \frac{s_1^2}{s_2^2} \begin{cases} F_{n_1-1, n_2-1} \end{cases}$$

$$\text{para } \hat{p} \begin{cases} N(0,1) \end{cases}$$

$$\text{para } \hat{p}_1 - \hat{p}_2 \begin{cases} N(0,1) \end{cases}$$