

# Introdução à Estatística usando o R com Aplicação em Análises Laboratoriais

Profa Carolina & Prof Gilberto

Instituto de Matemática e Estatística  
Universidade Federal da Bahia

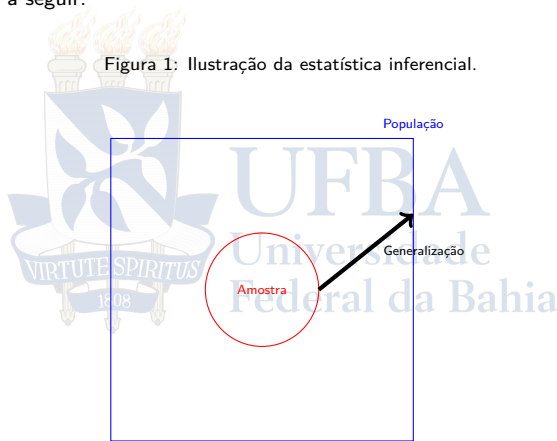
05 de outubro de 2019



# Probabilidade: Motivação

Com estatística descritiva podemos fazer afirmações válidas para amostra, mas queremos fazer afirmações válidas para toda a população. Com esse objetivo, vamos usar inferência estatística (ou estatística inferencial) para fazer generalizações da amostra para a população, conforme ilustrado na Figura 1. As técnicas de inferência estatística, usam probabilidade para fazer as generalizações como apresentado a seguir.

Figura 1: Ilustração da estatística inferencial.



# O que faremos nesse curso?

- **Estimação pontual:** Aproximar um parâmetro.

Exemplo: Estimar o teor alcóolico de uma bebida.

- **Intervalo de confiança:** Encontrar uma estimativa intervalar para um parâmetro.

Exemplo: Encontrar números  $a$  e  $b$  tal que o teor alcóolico verdadeiro está entre  $a$  e  $b$  com uma probabilidade estabelecida pelo pesquisador.

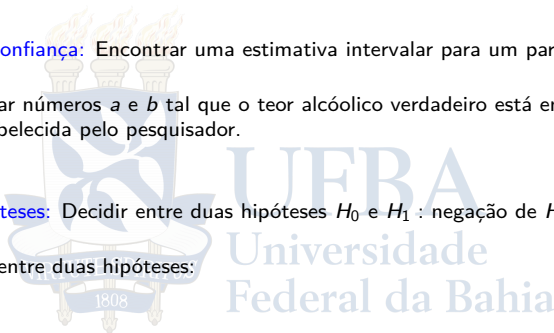
- **Teste de hipóteses:** Decidir entre duas hipóteses  $H_0$  e  $H_1$  : negação de  $H_0$ .

Exemplo: Decidir entre duas hipóteses:

$H_0$  : O teor alcóolico da bebida é 10%,

$H_1$  : O teor alcóolico da bebida não é 10%.

Em todos esses casos, precisamos usar probabilidade.



# Probabilidade

## Fenômeno Aleatório

Procedimento ou evento cujo resultado não é possível antecipar de forma determinística. Por exemplo:

- Teremos uma guerra total na Venezuela envolvendo o Brasil, Colômbia e Estados Unidos da América?
- Qual o resultado do lançamento de um dado “justo”?

## Notação e nomes

- **Espaço amostral:** O conjunto de todos os resultados de um fenômeno aleatório.

**Notação:**  $\Omega$

- **Evento:** Subconjunto de um espaço amostral.

**Notação:**  $A, B, C, \dots$

- **Ponto amostral:** Um resultado possível de um fenômeno aleatório.

**Notação:**  $\omega$ .

- **Probabilidade:** A plausibilidade de um ponto amostral  $\omega$  de  $A$  ser o resultado do fenômenos aleatório.

**Notação:**  $P(A)$ .

## Classificação de variáveis aleatórias

# Variável aleatória discreta

Seja  $X$  uma variável aleatória discreta. Então, podemos definir

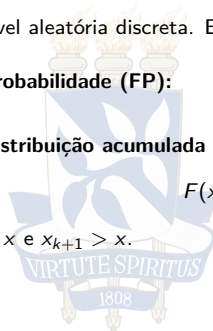
- **Função de probabilidade (FP):**

$$f(x) = P(X = x)$$

- **Função de distribuição acumulada (FDA):**

$$F(x) = f(x_1) + \cdots + f(x_k),$$

em que  $x_k \leq x$  e  $x_{k+1} > x$ .

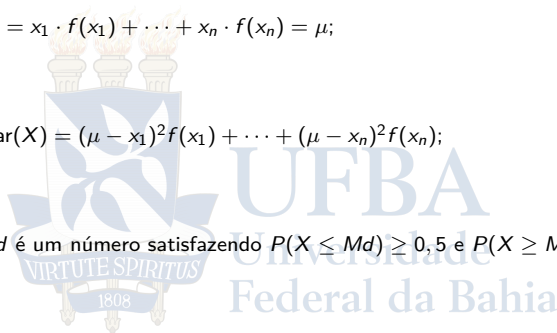


UFBA  
Universidade  
Federal da Bahia

# Medidas de resumo para variável aleatória discreta

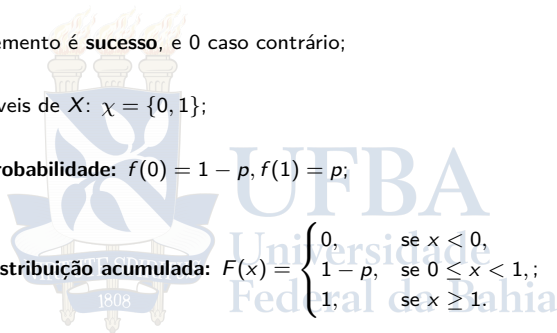
Seja  $X$  uma variável aleatória discreta com suporte  $\chi = \{x_1, \dots, x_n\}$  e função de probabilidade  $f(x)$ . Então

- **Média:**  $E(X) = x_1 \cdot f(x_1) + \dots + x_n \cdot f(x_n) = \mu$ ;
- **Variância:**  $\text{Var}(X) = (\mu - x_1)^2 f(x_1) + \dots + (\mu - x_n)^2 f(x_n)$ ;
- **Mediana:**  $Md$  é um número satisfazendo  $P(X \leq Md) \geq 0,5$  e  $P(X \geq Md) \geq 0,5$ ;
- **Desvio Padrão:**  $\text{DP}(X) = \sqrt{\text{Var}(X)}$ .



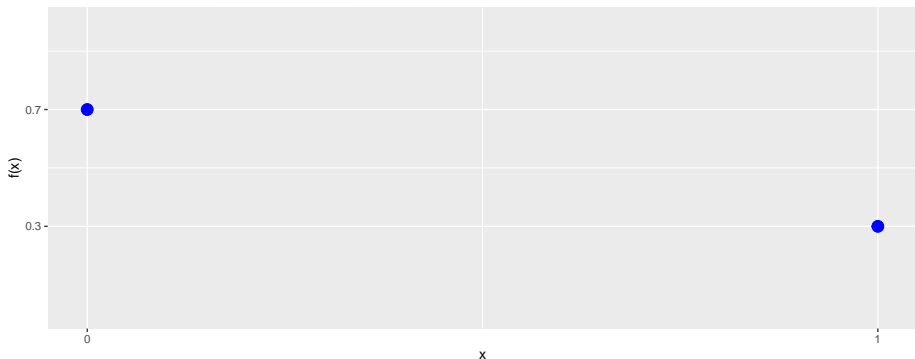
# Distribuição Bernoulli

- Cada elemento da população pode ser **sucesso** ou **fracasso**;
- $P(\text{sucesso}) = p$  e  $P(\text{fracasso}) = 1 - p$ ;
- $X : 1$  se o elemento é **sucesso**, e  $0$  caso contrário;
- Valores possíveis de  $X$ :  $\chi = \{0, 1\}$ ;
- **Função de probabilidade:**  $f(0) = 1 - p, f(1) = p$ ;
- **Função de distribuição acumulada:**  $F(x) = \begin{cases} 0, & \text{se } x < 0, \\ 1 - p, & \text{se } 0 \leq x < 1, \\ 1, & \text{se } x \geq 1. \end{cases}$
- $E(X) = n \cdot p$ ;
- $\text{Var}(X) = n \cdot p \cdot (1 - p)$ .



# Distribuição Bernoulli – função de probabilidade

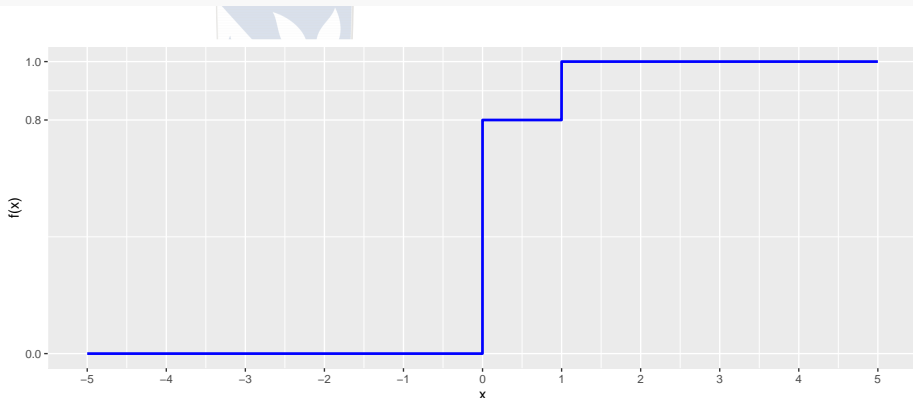
```
# gráfico da função de probabilidade
p <- 0.3 # probabilidade de sucesso
x <- c(0,1)
y <- dbinom(x, 1, p)
tibble(x = x, `f(x)`=y) %>%
  ggplot(aes(x, `f(x)`)) +
  geom_point(colour = 'blue', size=4) +
  scale_x_continuous(breaks = c(0,1)) +
  scale_y_continuous(breaks = c(1-p,p), limits = c(0,1)) +
  labs(y = 'f(x)')
```





# Distribuição Bernouli – função de distribuição acumulada

```
# Função de distribuição acumulada
p <- 0.2
x <- seq(from = -5, to = 5, by = 0.001)
y <- pbinom(x, 1, p)
# gráfico -- FDA
tibble(x=x, `f(x)`=y) %>% ggplot() +
  geom_line(aes(x, `f(x)`), color = 'blue', size = 1) +
  scale_x_continuous(breaks = seq(from = -5, to = 5, by = 1)) +
  scale_y_continuous(breaks = c(0, 1-p, 1))
```



# Distribuição Bernoulli

```
# Variável Bernoulli:  $X \sim \text{Bernoulli}(p)$   
p <- 0.3  
# Função densidade  
(dbinom(c(0,1),1,p))
```

```
## [1] 0.7 0.3
```



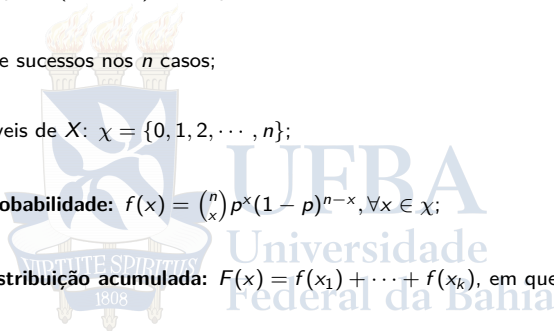
```
# simular valores da variável Bernoulli  
n <- 1000 # tamanho da amostra  
amostra <- rbinom(n, 1, p)  
tibble(x = amostra) %>%  
  summarise(media = mean(x), mediana = median(x),  
             dp = sd(x), cv = sd(x) * 100 / mean(x),  
             q1 = quantile(x, probs = 0.25),  
             q3 = quantile(x, probs = 0.75))
```

```
## # A tibble: 1 x 6  
##   media mediana    dp    cv    q1    q3  
##   <dbl>   <dbl> <dbl> <dbl> <dbl> <dbl>  
## 1 0.291     0 0.454 156.     0     1
```



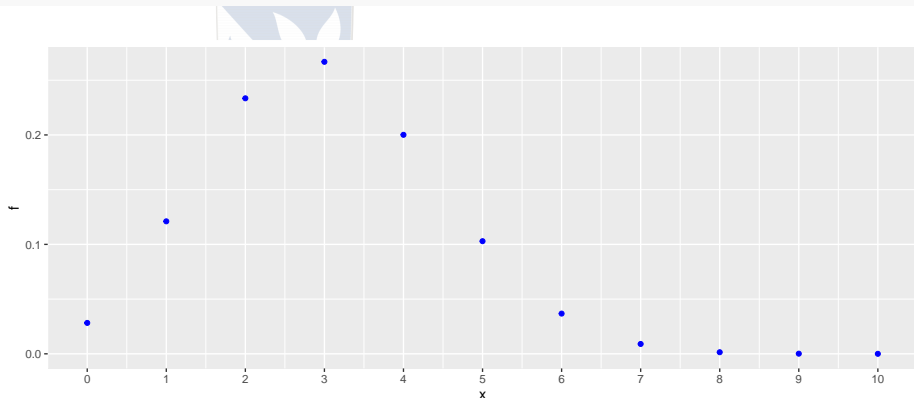
# Distribuição Binomial

- Temos  $n$  casos em que cada caso pode ser **sucesso** ou **fracasso**;
- $P(\text{sucesso}) = p$  e  $P(\text{fracasso}) = 1 - p$ ;
- $X$  : número de sucessos nos  $n$  casos;
- Valores possíveis de  $X$ :  $\chi = \{0, 1, 2, \dots, n\}$ ;
- **Função de probabilidade:**  $f(x) = \binom{n}{x} p^x (1 - p)^{n-x}, \forall x \in \chi$ ;
- **Função de distribuição acumulada:**  $F(x) = f(x_1) + \dots + f(x_k)$ , em que  $x \leq x_k$  e  $x_{k+1} > x$ ;
- $E(X) = n \cdot p$ ;
- $\text{Var}(X) = n \cdot p \cdot (1 - p)$ .



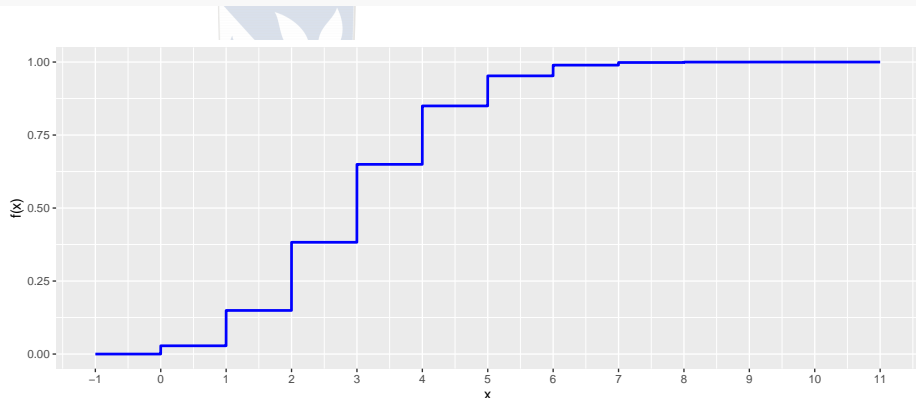
# Distribuição binomial – função de probabilidade

```
# gráfico da função de probabilidade  
n <- 10  
x <- 0:n  
p <- 0.3  
f <- dbinom(x, n, p)  
tibble(x=x, f = f) %>%  
  ggplot() +  
  geom_point(aes(x=x, y=f), color = 'blue') +  
  scale_x_continuous(breaks = 0:10)
```



# Distribuição binomial – função de distribuição acumulada

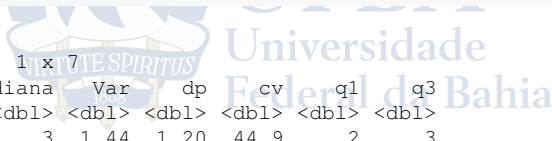
```
# gráfico da função de distribuição acumulada
n <- 10
p <- 0.3
x <- seq(from = -1, to = 11, by = 0.001)
y <- pbinom(x, n, p)
tibble(x = x, `f(x)` = y) %>%
  ggplot() +
  geom_line(aes(x, `f(x)`), stat = 'identity', size = 1, color = 'blue') +
  scale_x_continuous(breaks = seq(from = -1, to = 11, by = 1))
```



# Distribuição binomial – amostra aleatória

```
# Amostra da distribuição binomial
m <- 100 # tamanho da amostra
n <- 10 # número de casos
p <- 0.3
amostra <- rbinom(m,n,p)
tibble(x = amostra) %>%
  summarise(media = mean(x), mediana = median(x), Var = var(x),
            dp = sd(x), cv = sd(x) * 100 / mean(x),
            q1 = quantile(x, probs = 0.25),
            q3 = quantile(x, probs = 0.75))
```

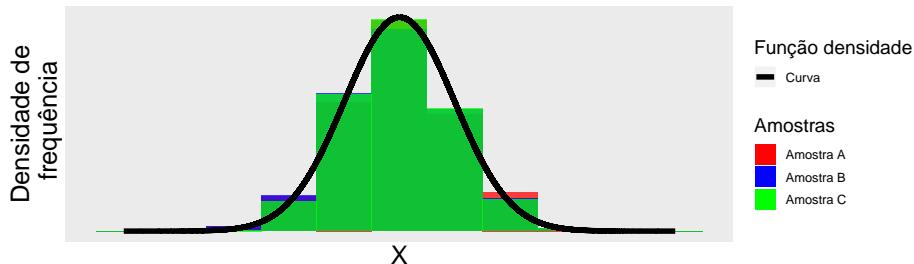
```
## # A tibble: 1 x 7
##   media mediana  Var    dp    cv    q1    q3
##   <dbl>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1  2.67     3  1.44  1.20  44.9     2     3
```



# Variável aleatória contínua

## Motivação

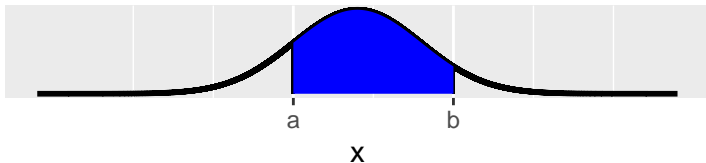
- Para cada amostra, temos um histograma;
- Queremos encontrar uma curva que aproxime bem todos os histogramas possíveis;



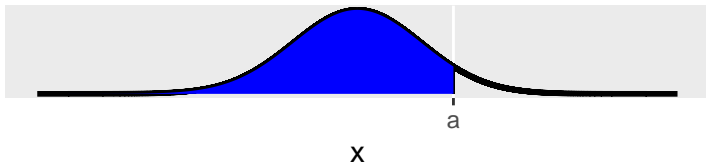
Chamamos a curva preta de **função densidade**.

# Propriedades de variável aleatória contínua

- $P(X = a) = 0$ ;
- Usamos a notação  $P(a < X < b)$ ;
- $P(a \leq X < b) = P(a < X \leq b) = P(a \leq X \leq b) = P(a < X < b)$ ;



- $F(a) = P(X \leq a)$





# Distribuição normal

- Valores da variável aleatória concentrados em torno da média populacional  $\mu$ ;
- Valores da variável aleatória afastados da média populacional  $\mu$  são pouco prováveis;
- Valores possíveis da variável: todos os números reais  $\chi = \mathbb{R}$ ;
- Função densidade (fd): curva em formato de sino;
- $\mu$  é a média da população e  $\sigma^2$  é a variância da população;
- Função de distribuição acumulada (fda):  $F(x) = P(X \leq x)$ ;
- Usamos a notação:  $X \sim N(\mu, \sigma^2)$ ;
- Seja  $\Phi(x)$  a fda de uma variável  $Z \sim N(0, 1)$ , então se  $X \sim N(\mu, \sigma^2)$  temos que

$$F(x) = \Phi\left(\frac{x - \mu}{\sigma}\right)$$

## Distribuição normal – continuação

- Se  $X \sim N(\mu, \sigma^2)$ , então

$$\begin{aligned}P(a < X < b) &= F(b) - F(a) \\&= \Phi\left(\frac{b - \mu}{\sigma}\right) - \Phi\left(\frac{a - \mu}{\sigma}\right)\end{aligned}$$

- Seja  $\Phi(x)$  a fda de uma variável  $Z \sim N(0, 1)$ , então

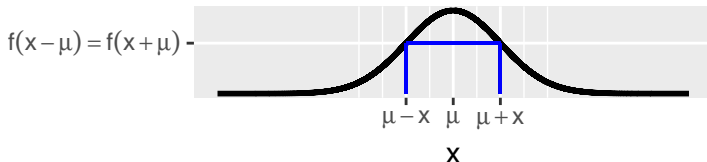
$$\Phi(a) = 1 - \Phi(-a), \text{ se } a < 0$$

- Média, moda, mediana, variância para  $X \sim N(\mu, \sigma^2)$ :

$$E(X) = \mu, \quad Mo(X) = \mu, \quad Md(X) = \mu, \quad Var(X) = \sigma^2.$$

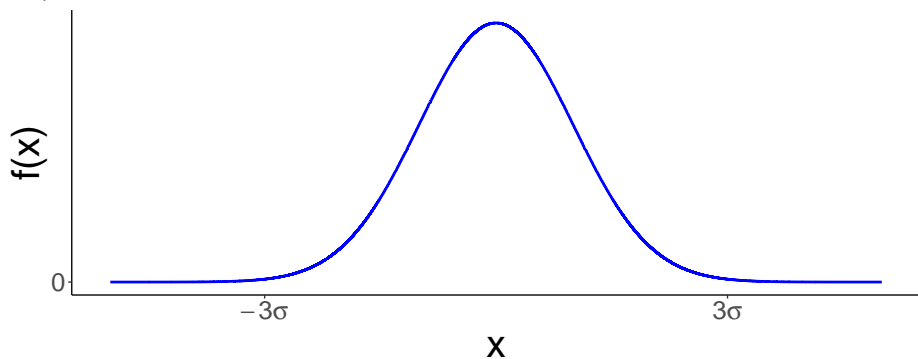
- A função densidade é simétrica em torno de  $\mu$ .

- Se  $X \sim N(\mu, \sigma^2)$ , então  $f(\mu - x) = f(\mu + x)$ .

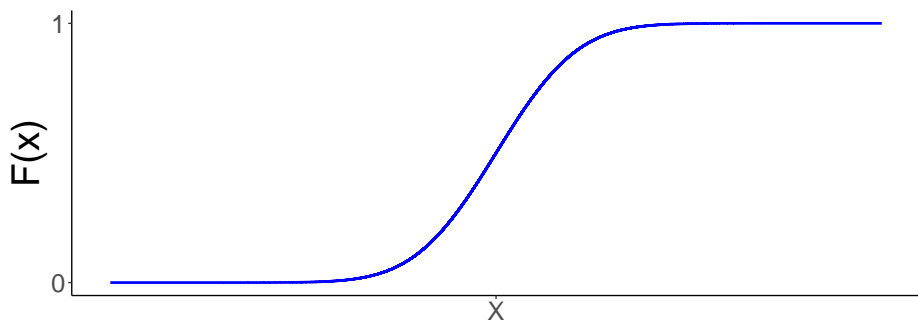


# Gráfico da função densidade

Função densidade tem formato de sino.



## Gráfico da função de distribuição acumulada



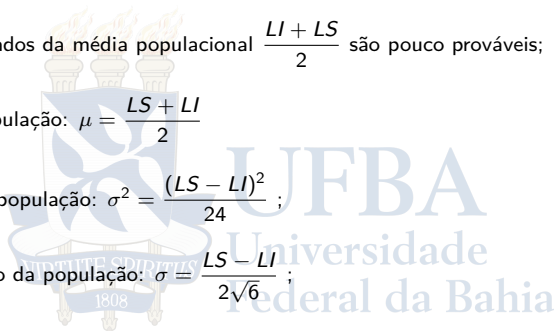
# Distribuição normal – amostra aleatória

```
media <- 2
s2 <- 1
n <- 1e+3 #tamanho da amostra
amostra <- rnorm(n, mean = media, sd = sqrt(s2))
tibble(x = amostra) %>%
  summarise(media = mean(x), mediana = median(x),
            Var = var(x), dp = sd(x), cv = sd(x) * 100 / mean(x),
            q1 = quantile(x, probs = 0.25),
            q3 = quantile(x, probs = 0.75))
```

```
## # A tibble: 1 x 7
##   media mediana   Var    dp    cv    q1    q3
##   <dbl>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1  1.99    2.01  1.03  1.02  50.9  1.28  2.70
```

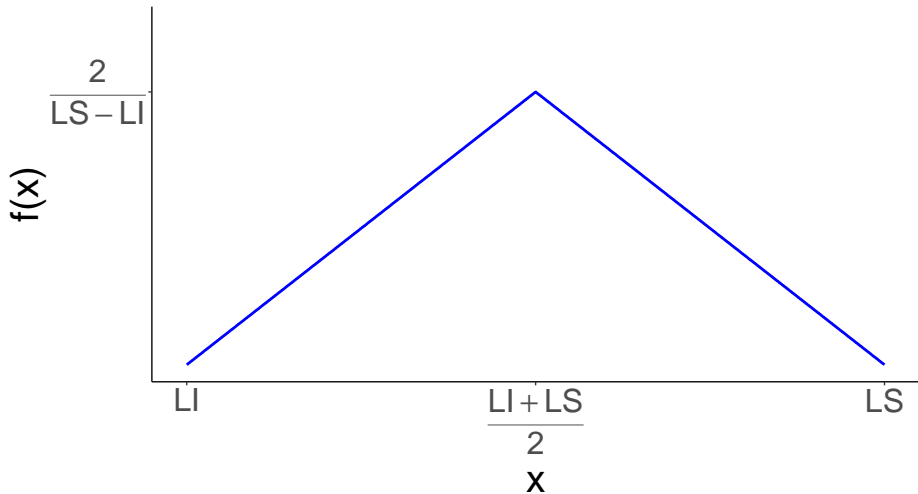
# Distribuição triangular

- Valores da variável aleatória sempre estão entre  $LI$  e  $LS$ ;
- Valores da variável aleatória estão concentrados em torno do valor  $\frac{LI + LS}{2}$ ;
- Valores afastados da média populacional  $\frac{LI + LS}{2}$  são pouco prováveis;
- Média da população:  $\mu = \frac{LS + LI}{2}$
- Variância da população:  $\sigma^2 = \frac{(LS - LI)^2}{24}$  ;
- Desvio padrão da população:  $\sigma = \frac{LS - LI}{2\sqrt{6}}$  ;
- Função densidade (fd): curva em formato de um triângulo;
- Função de distribuição acumulada (fda):  $F(x) = P(X \leq x)$ ;
- Usamos a notação:  $X \sim \text{triangular}(LI, LS)$ ;

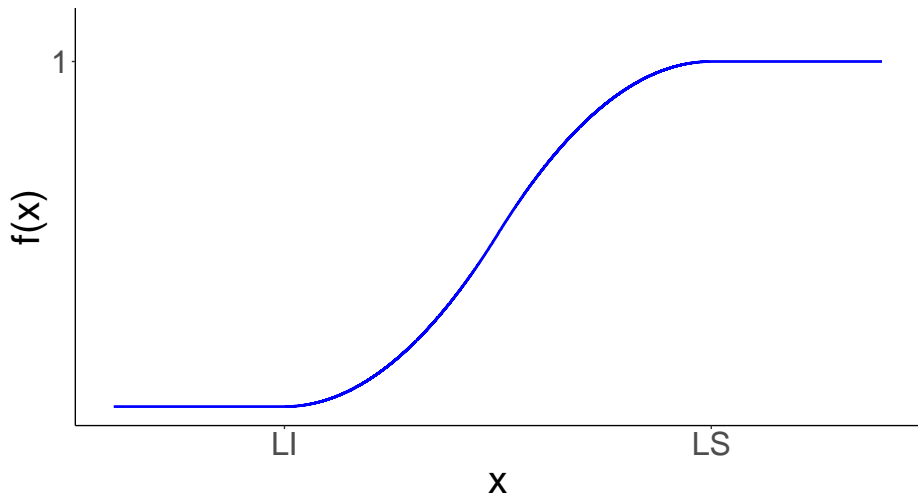


# Gráfico da função densidade

Função densidade em formato triangular.



## Gráfico da função de distribuição acumulada





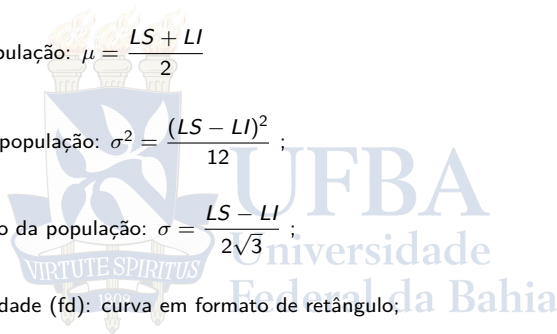
# Distribuição triangular: amostra aleatória

```
n <- 1000
min <- 0
max <- 10
amostra <- rtri(n, min = min, max = max, mode = (min + max) / 2)
tibble(x = amostra) %>%
  summarise(media = mean(x), mediana = median(x), Var = var(x),
            dp = sd(x), cv = sd(x) * 100 / mean(x),
            q1 = quantile(x, probs = 0.25),
            q3 = quantile(x, probs = 0.75))
```

```
## # A tibble: 1 x 7
##   media mediana Var dp cv q1 q3
##   <dbl>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1  5.00     4.97  4.21  2.05  41.1  3.49  6.58
```

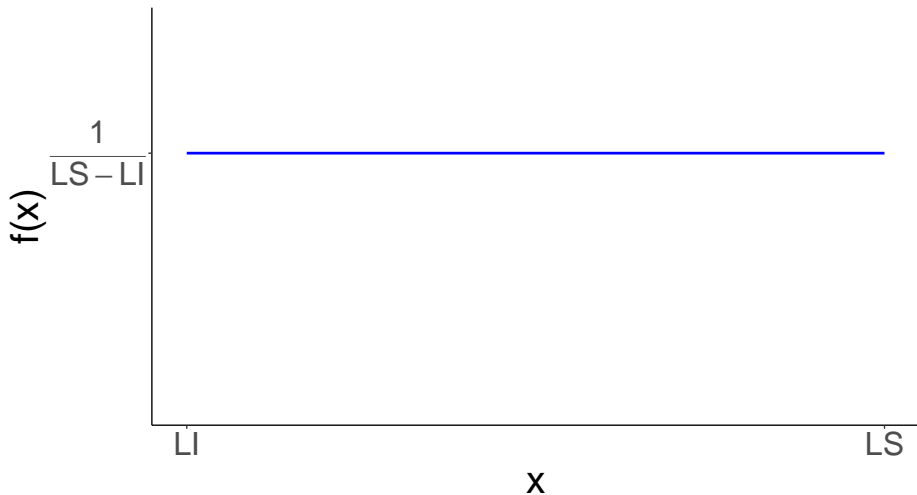
# Distribuição retangular ou uniforme

- Valores da variável aleatória sempre estão entre  $LI$  e  $LS$ ;
- Todos os valores entre  $LI$  e  $LS$  são igualmente prováveis;
- Média da população:  $\mu = \frac{LS + LI}{2}$
- Variância da população:  $\sigma^2 = \frac{(LS - LI)^2}{12}$  ;
- Desvio padrão da população:  $\sigma = \frac{LS - LI}{2\sqrt{3}}$  ;
- Função densidade (fd): curva em formato de retângulo;
- Função de distribuição acumulada (fda):  $F(x) = P(X \leq x)$ ;
- Usamos a notação:  $X \sim \text{retangular}(LI, LS)$ ;

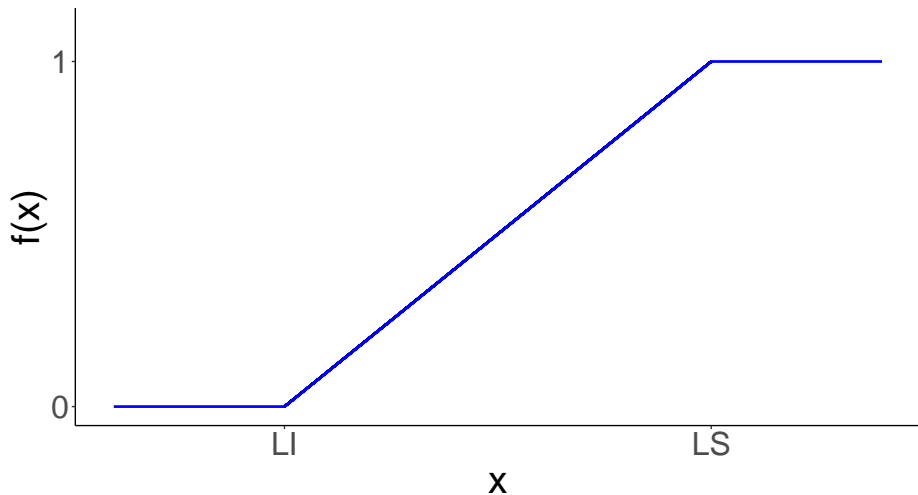


# Grafico da função densidade

Função densidade em formato retangular.



## Gráfico da função de distribuição acumulada



# Distribuição retangular: amostra aleatória

```
n <- 1000
min <- 0
max <- 10
amostra <- runif(n, min = min, max = max)
tibble(x = amostra) %>%
  summarise(media = mean(x), mediana = median(x), Var = var(x),
            dp = sd(x), cv = sd(x) * 100 / mean(x),
            q1 = quantile(x, probs = 0.25),
            q3 = quantile(x, probs = 0.75))
```

```
## # A tibble: 1 x 7
##   media mediana  Var    dp    cv    q1    q3
##   <dbl>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1  5.03     5.14  8.42  2.90  57.7  2.48  7.49
```