



escola  
britânica de  
artes criativas  
& tecnologia

**Profissão: Cientista de Dados**

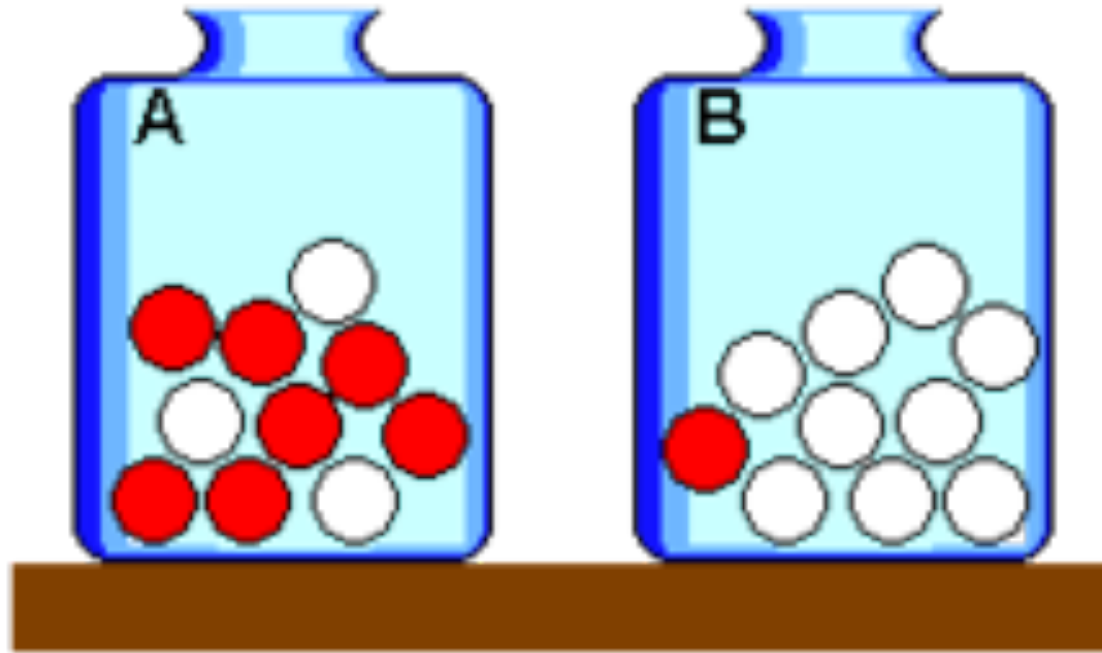
Probabilidade

Introdução:

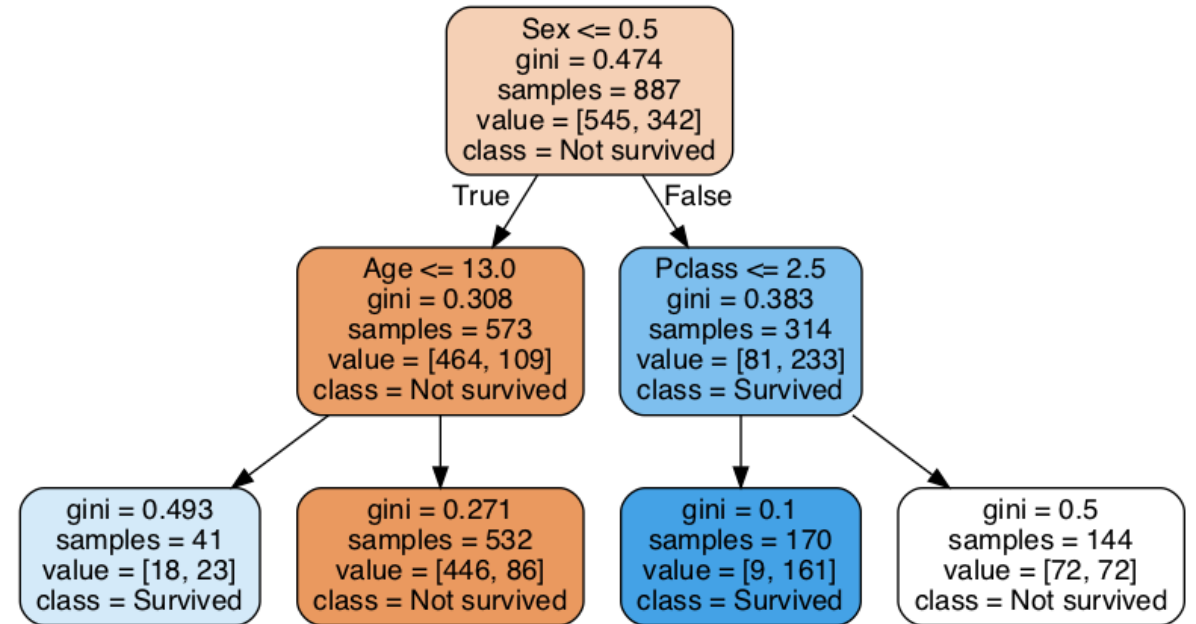
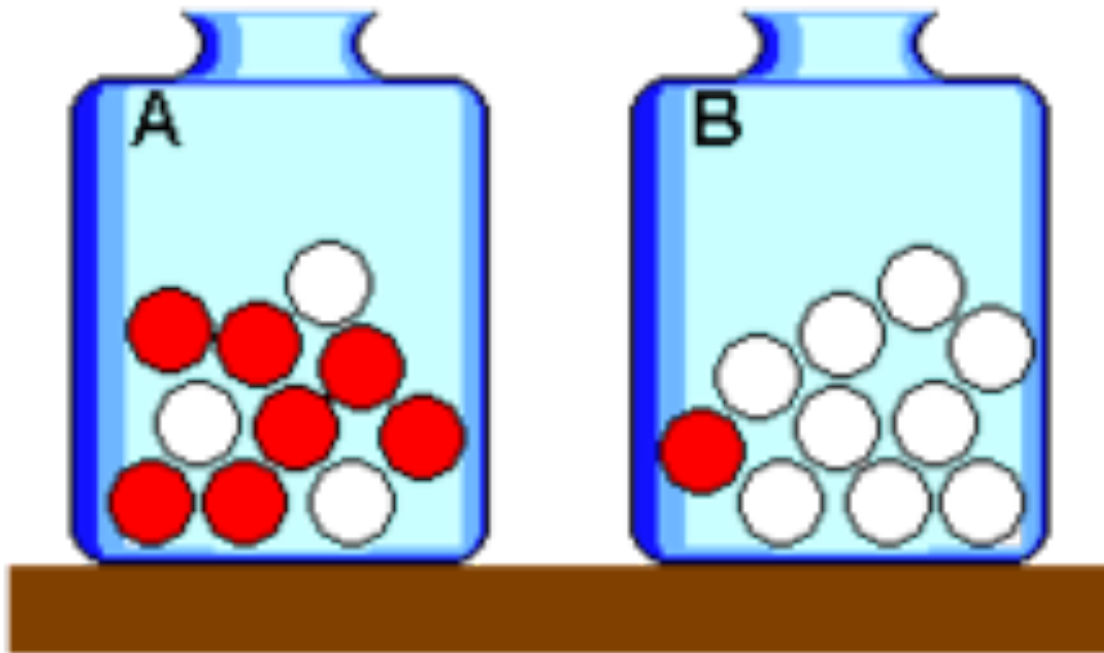
# Variáveis Aleatórias e Amostragem

# Probabilidade

---

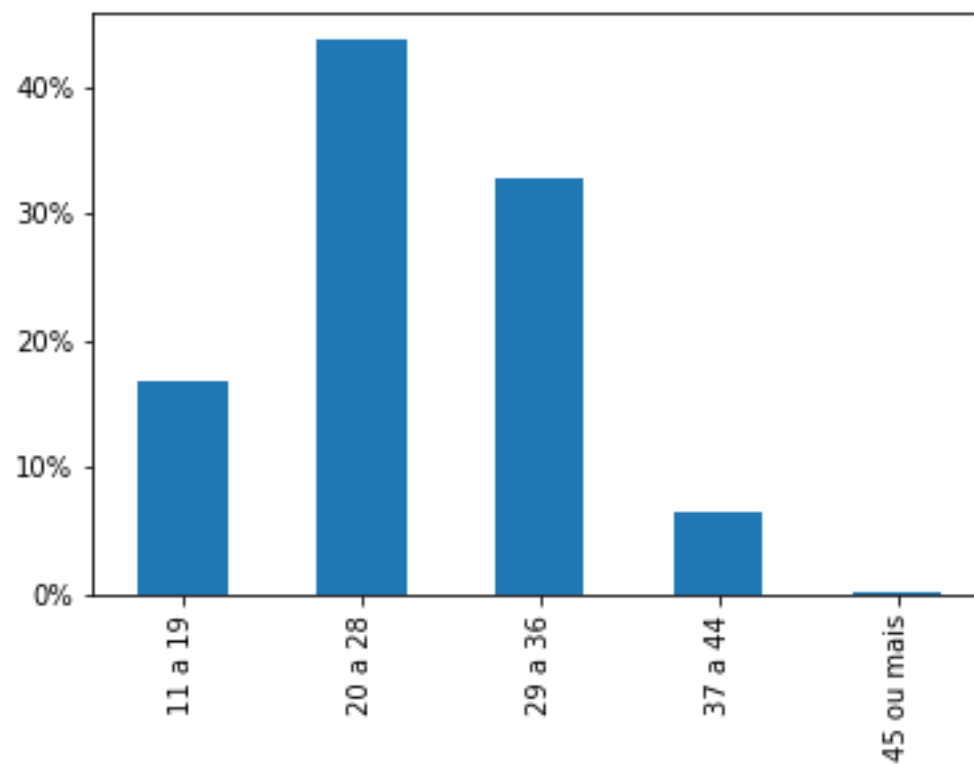
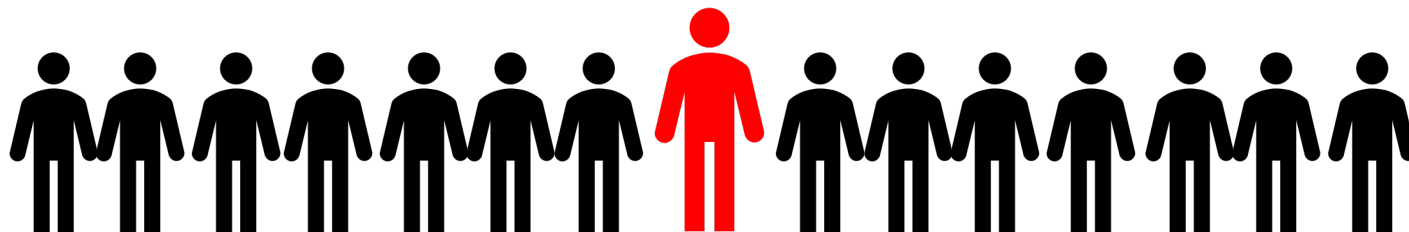


# Probabilidade

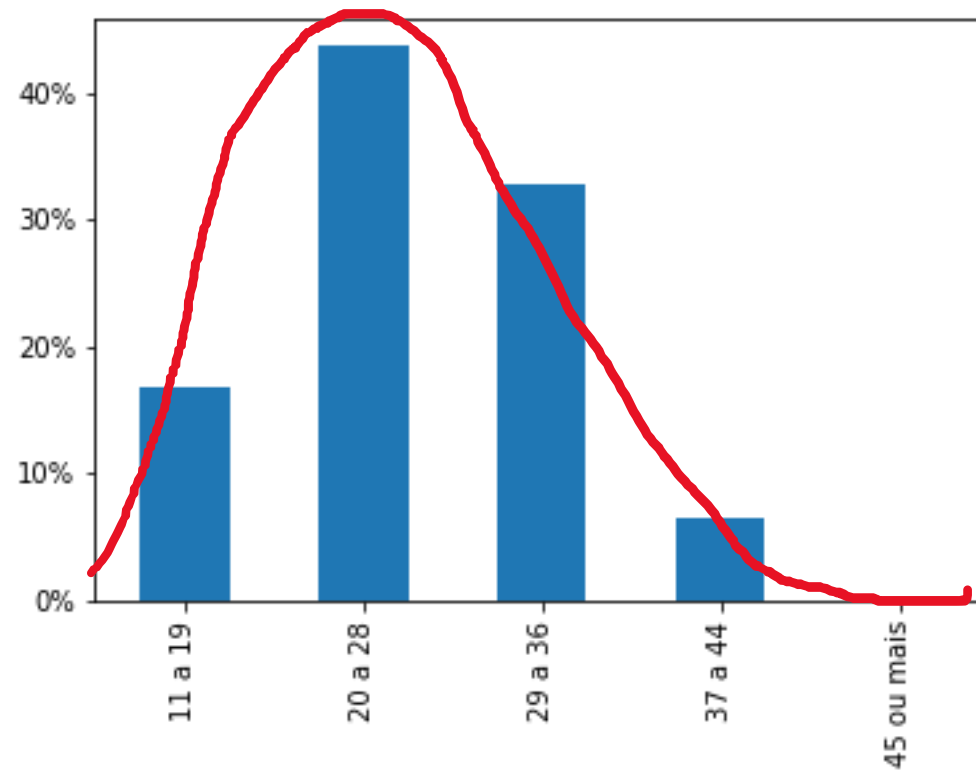
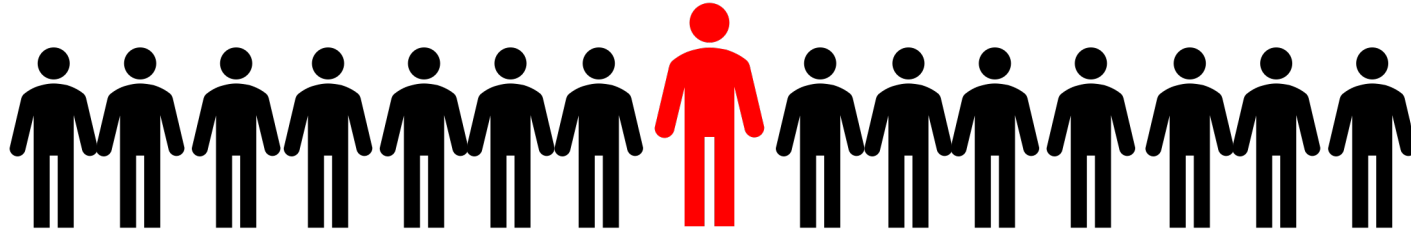


Amostragem tipo  
*cross section*

# Probabilidade



# Probabilidade



# Amostragem de Processo Estocástico





# Modelos discretos

# Modelos discretos - Bernoulli



$$X = \begin{cases} 1 & \text{se o resultado do ensaio é "sucesso"} \\ 0 & \text{se o resultado do ensaio é "fracasso"} \end{cases}$$

Definimos: "cara" = "sucesso", "coroa" = "fracasso"

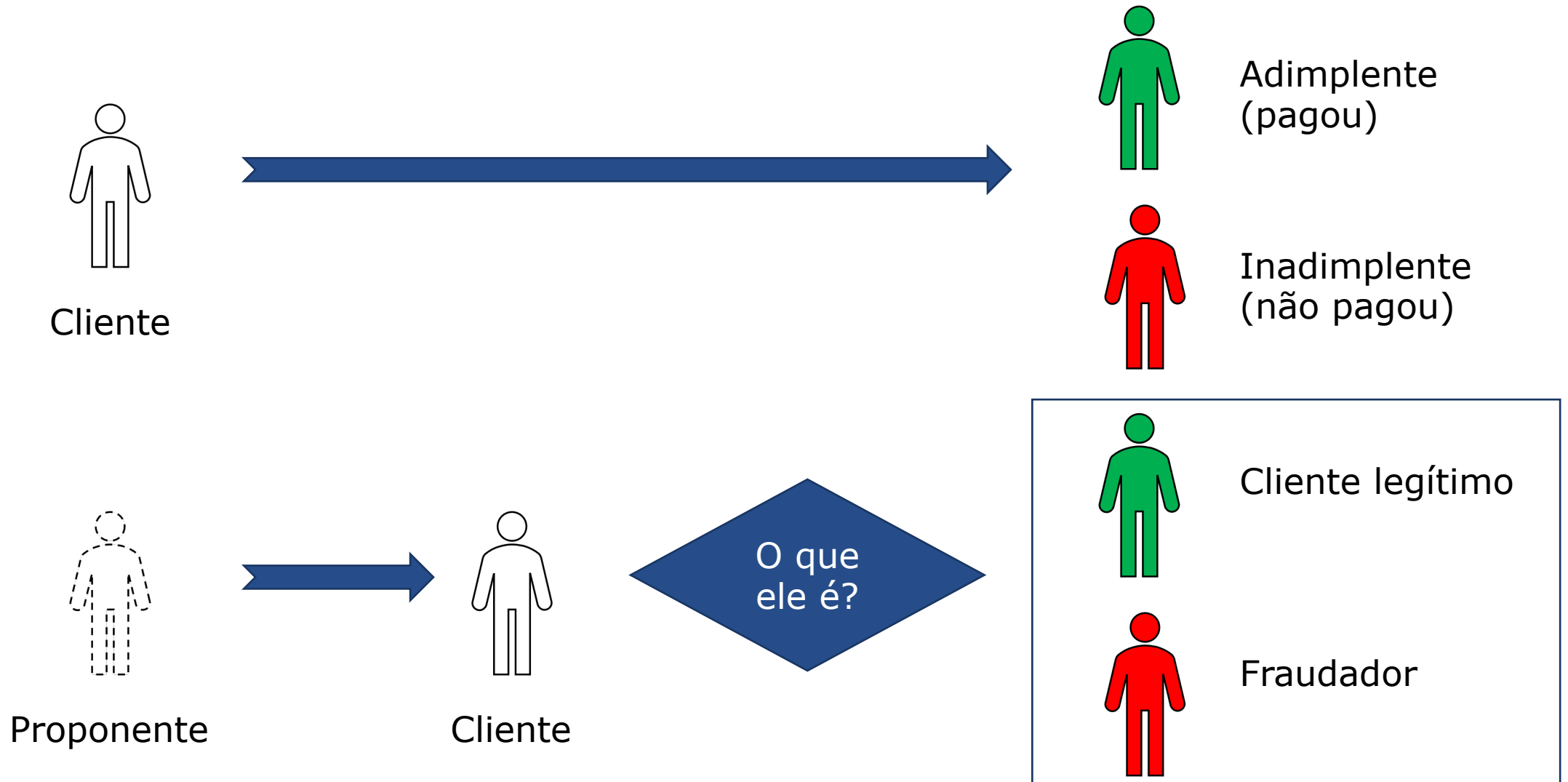
$$P(X = 1) = 0,5$$

$$P(X = 0) = 0,5$$

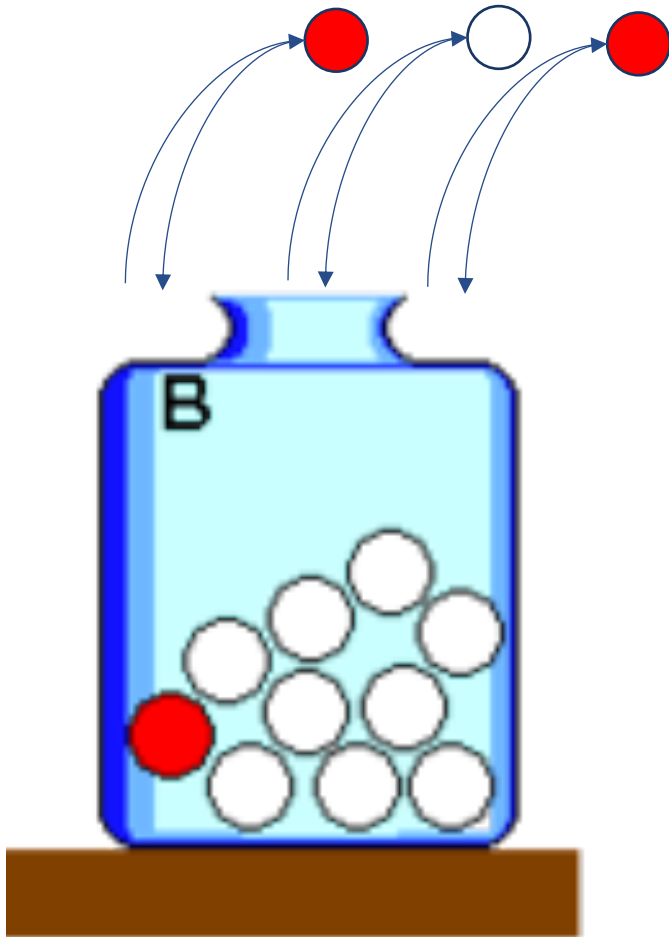
$$P(X = 1) = p$$

$$P(X = 0) = 1 - p$$

# Modelos discretos - Bernoulli

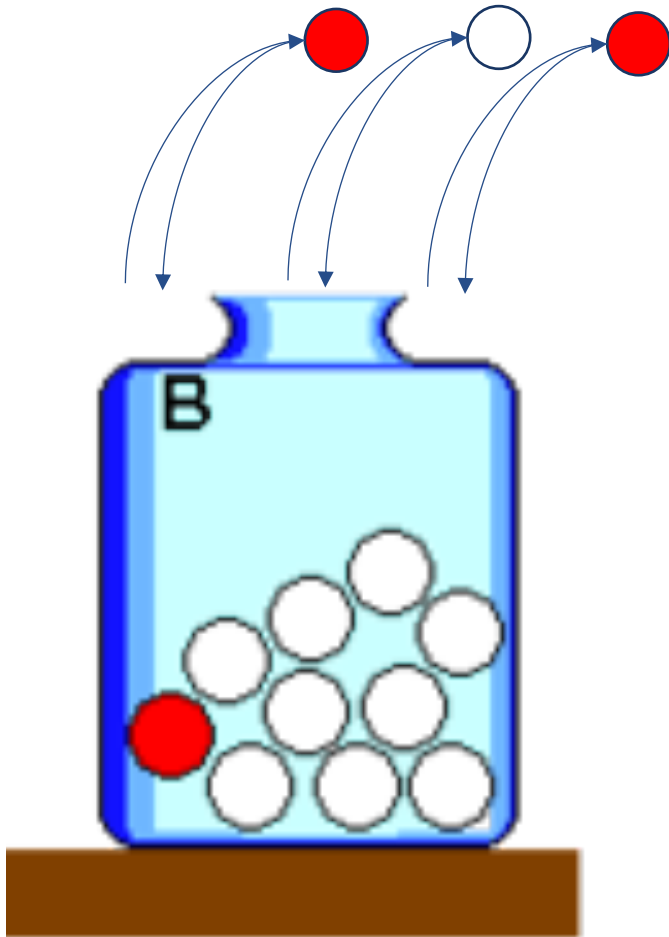


# Modelos discretos - Binomial



Se retirarmos  $N$  bolas, com reposição, qual a probabilidade de retirarmos  $x$  bolas brancas?

# Modelos discretos - Binomial



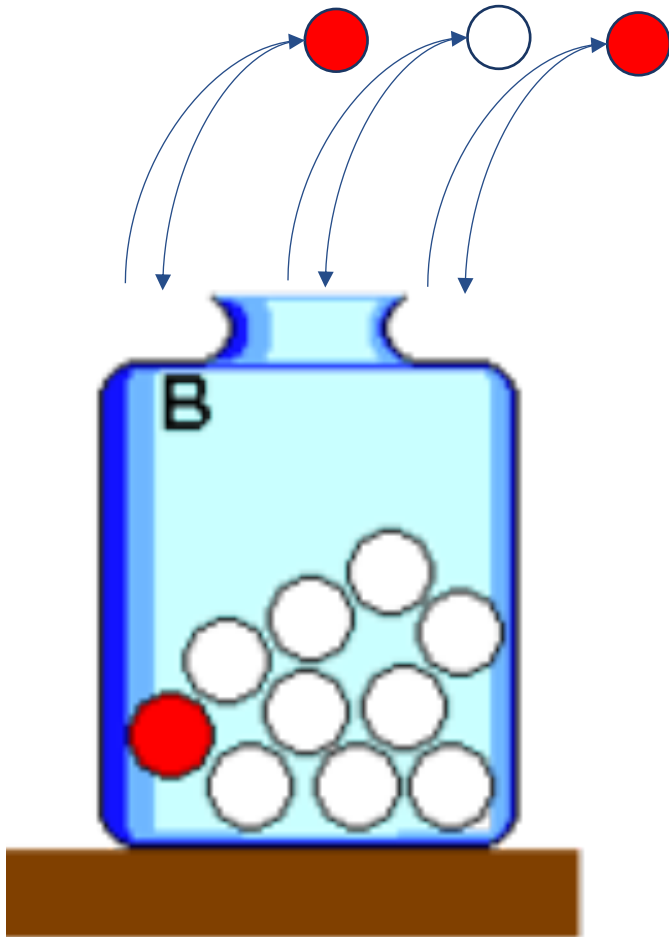
Qual a probabilidade de retirarmos 1 bola branca?



$$p * (1 - p) * (1 - p)$$

$$= p(1 - p)^2$$

# Modelos discretos - Binomial



Qual a probabilidade de retirarmos 1 bola branca?

OBS: Considerando a ordem, podemos fazer isso de 3 formas diferentes:

Forma 1:     ●     ○     ●

Forma 2:     ○     ●     ●

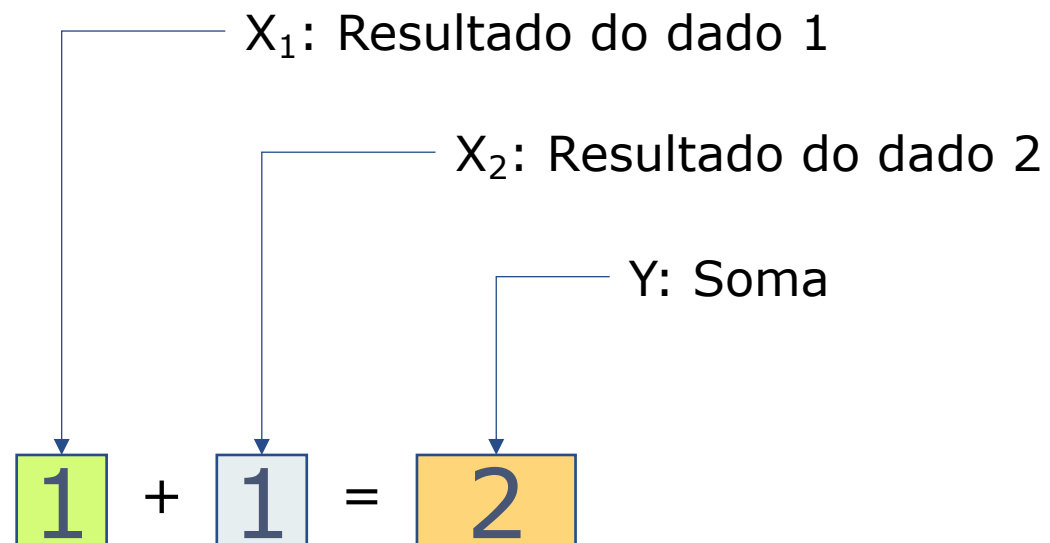
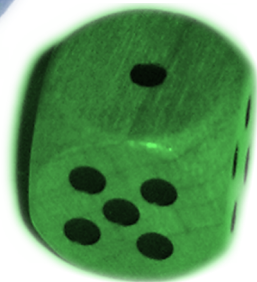
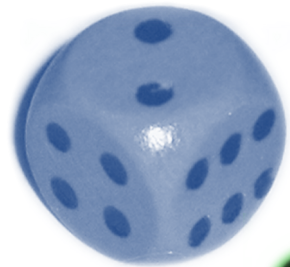
Forma 3:     ●     ●     ○

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

# Distribuição condicional e independência



# Conjunta

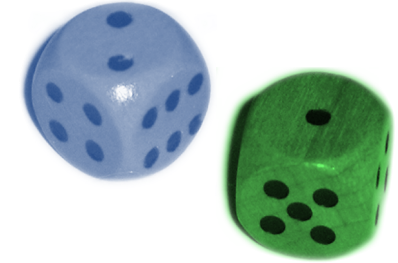


# Conjunta

1	+	1	=	2
1	+	2	=	3
1	+	3	=	4
1	+	4	=	5
1	+	5	=	6
1	+	6	=	7

2	+	1	=	3
2	+	2	=	4
2	+	3	=	5
2	+	4	=	6
2	+	5	=	7
2	+	6	=	8

3	+	1	=	4
3	+	2	=	5
3	+	3	=	6
3	+	4	=	7
3	+	5	=	8
3	+	6	=	9



	dado 1
	dado 2
	Soma

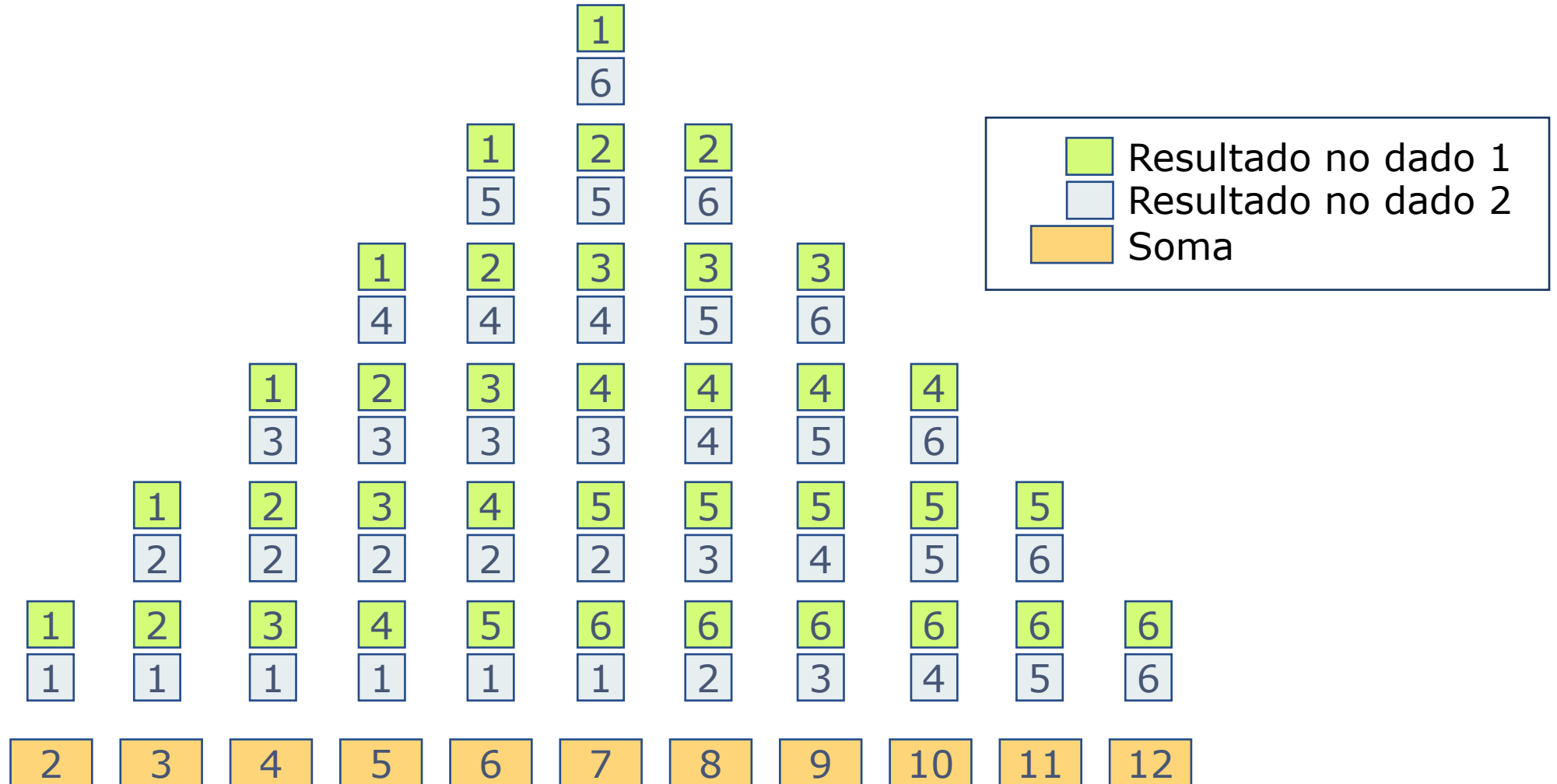
4	+	1	=	5
4	+	2	=	6
4	+	3	=	7
4	+	4	=	8
4	+	5	=	9
4	+	6	=	10

5	+	1	=	6
5	+	2	=	7
5	+	3	=	8
5	+	4	=	9
5	+	5	=	10
5	+	6	=	11

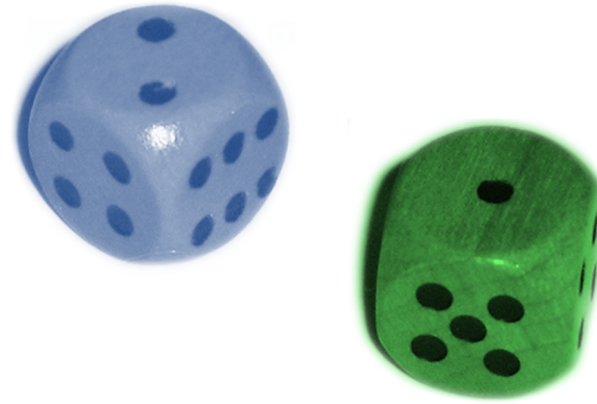
6	+	1	=	7
6	+	2	=	8
6	+	3	=	9
6	+	4	=	10
6	+	5	=	11
6	+	6	=	12

$$\begin{aligned} &P(X_1 = i, X_2 = j) \\ &= P(X_1 = i) \cdot P(X_2 = j) \\ &= \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{36} \end{aligned}$$

# Conjunta

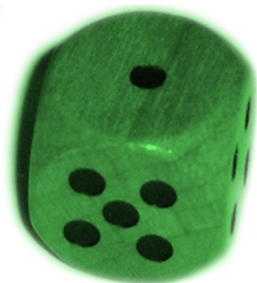


# Independência



$$P(X_1 = i, |X_2 = j) = P(X_1 = i)$$

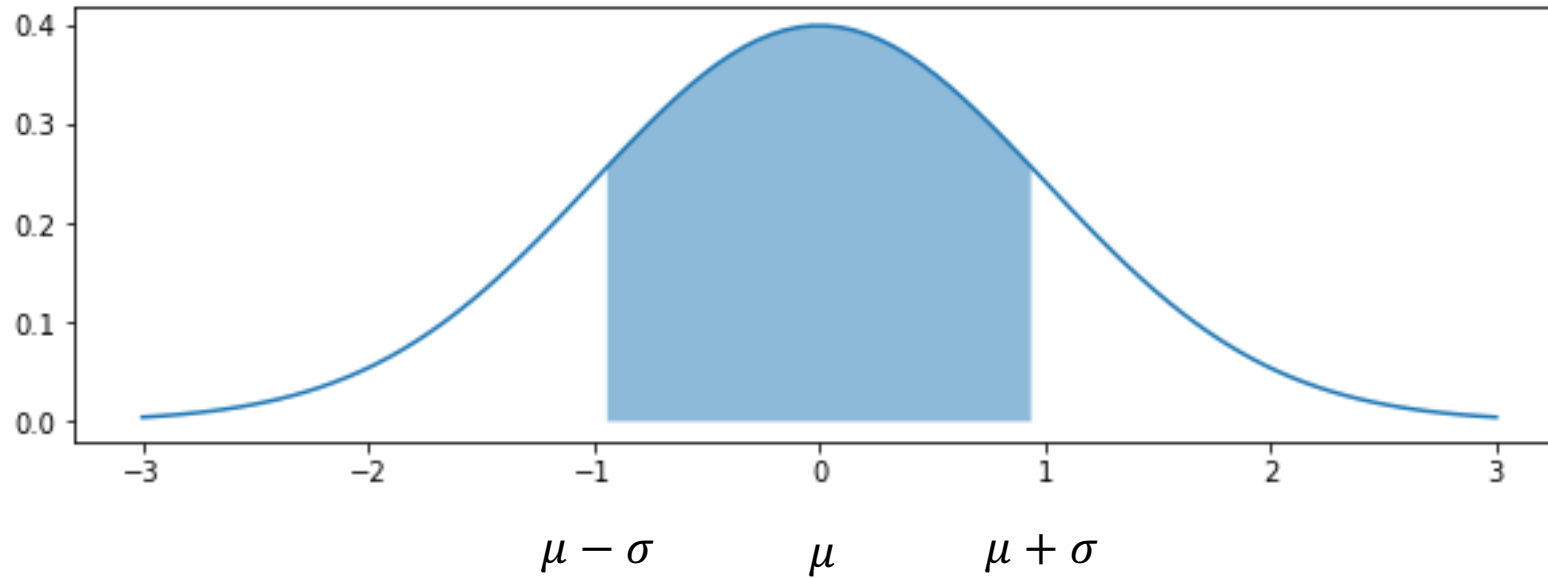
# Condicional



$$P(Y = i | X_1 = 6) = \begin{cases} \frac{1}{6} & \text{se } i = 7 \\ \frac{1}{6} & \text{se } i = 8 \\ \frac{1}{6} & \text{se } i = 9 \\ \frac{1}{6} & \text{se } i = 10 \\ \frac{1}{6} & \text{se } i = 11 \\ \frac{1}{6} & \text{se } i = 12 \\ 0 & \text{caso contrário} \end{cases}$$

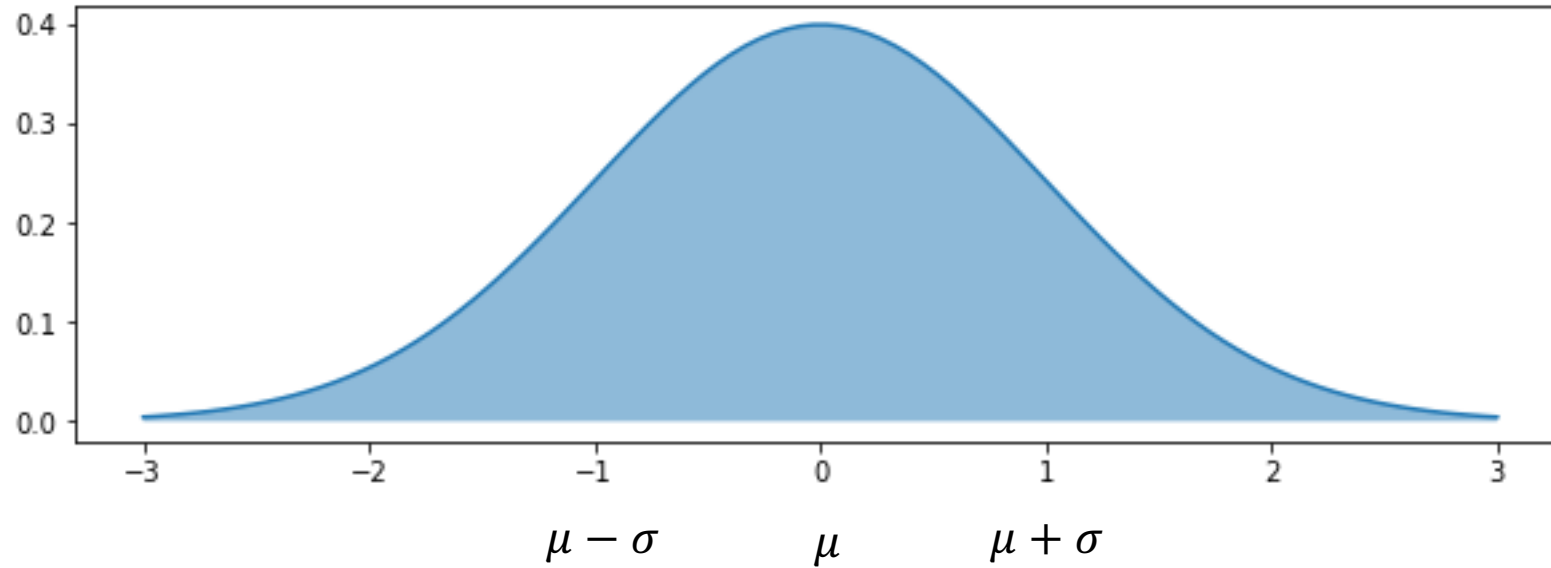
# Função Densidade de Probabilidade

# Função Densidade de Probabilidade



$$P(a \leq X \leq b) = \int_a^b f(x)dx$$

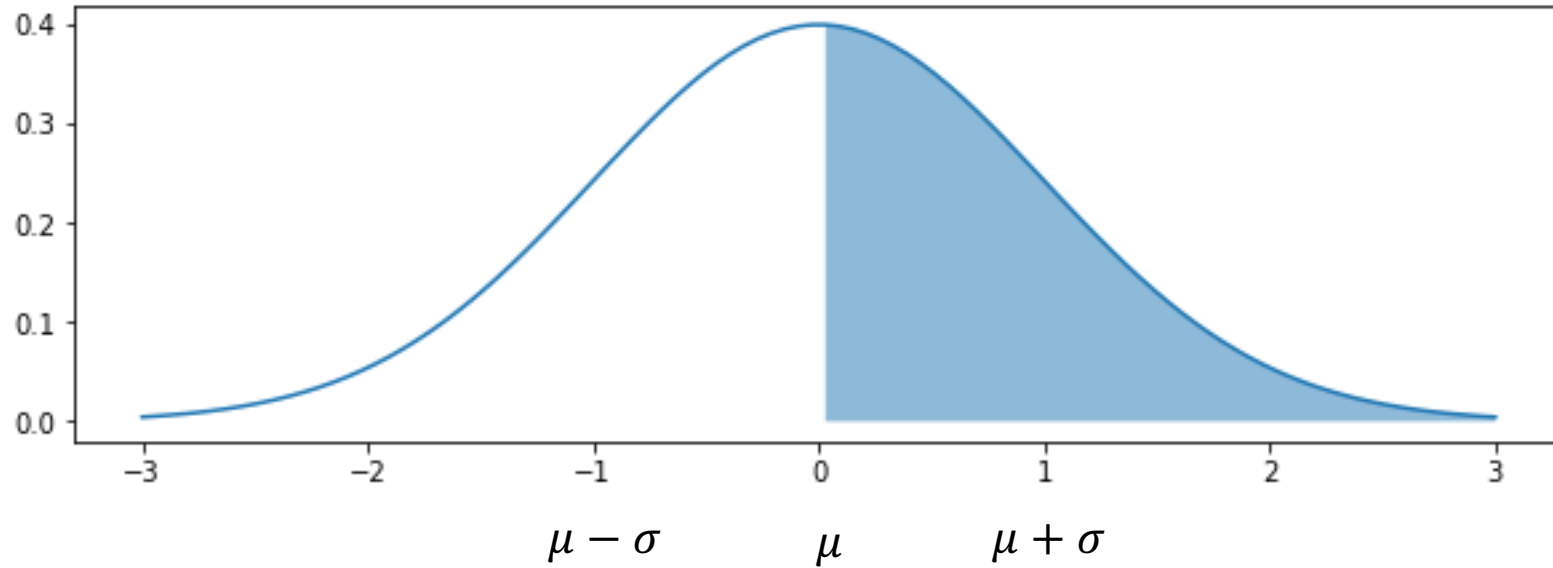
# Função Densidade de Probabilidade



$$P(-\infty \leq X \leq +\infty) = \int_{-\infty}^{+\infty} f(x)dx = 1$$



# Função Densidade de Probabilidade

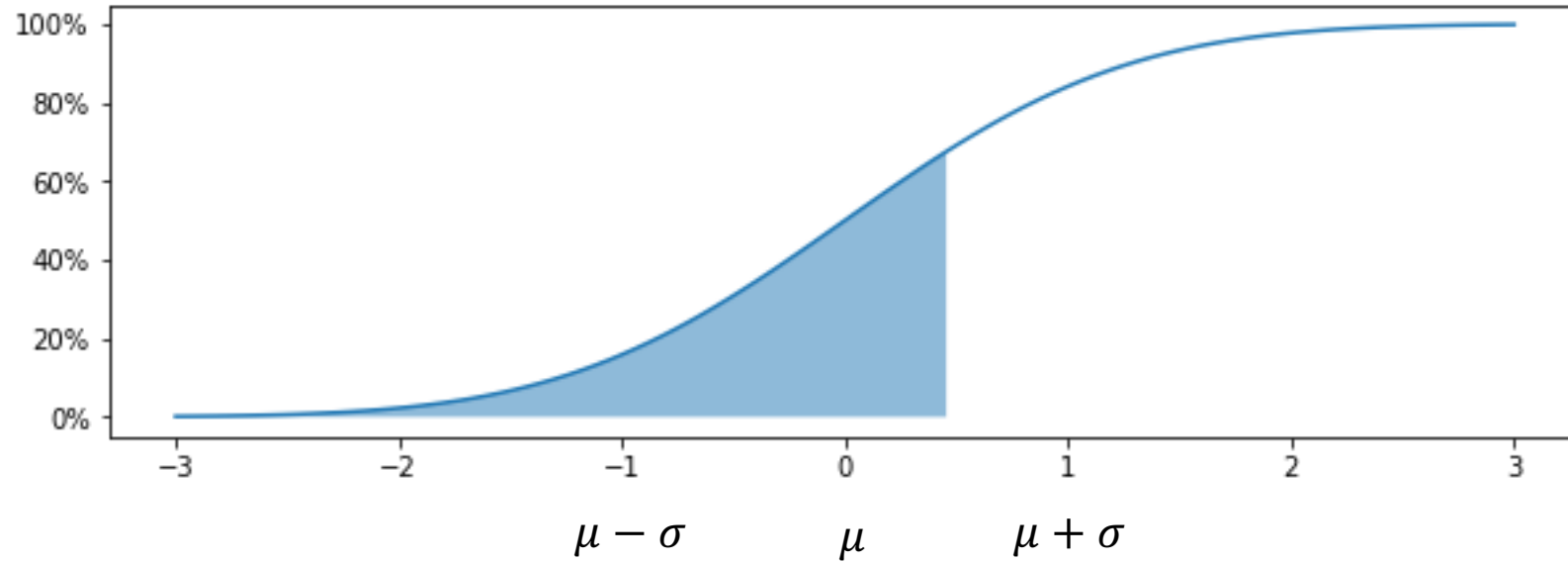


$$P(0 \leq X \leq +\infty) = \int_0^{+\infty} f(x)dx = 50\%$$

Quando a função é simétrica, a probabilidade para valores maiores que média é de 50%

# Função Distribuição Acumulada de Probabilidade

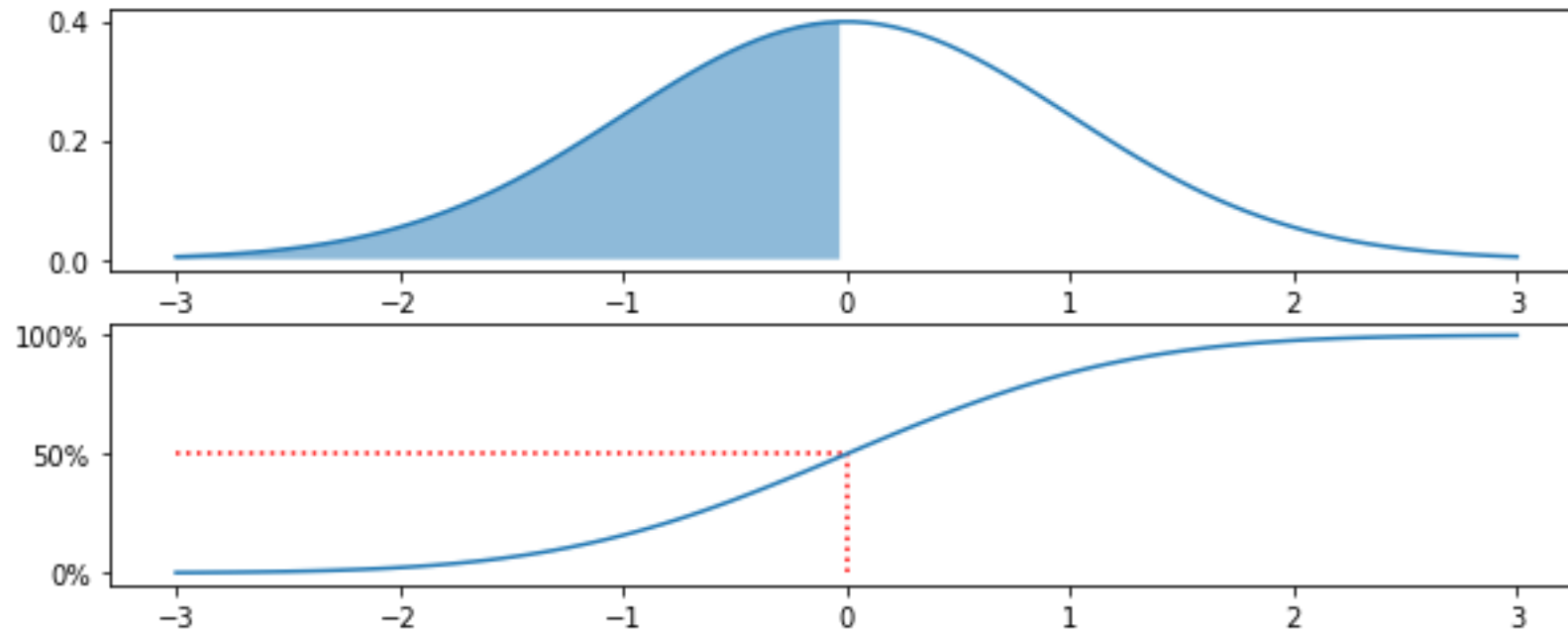
# Distribuição acumulada



$$F(x) = P(X < x)$$

$$F(+\infty) = 1$$

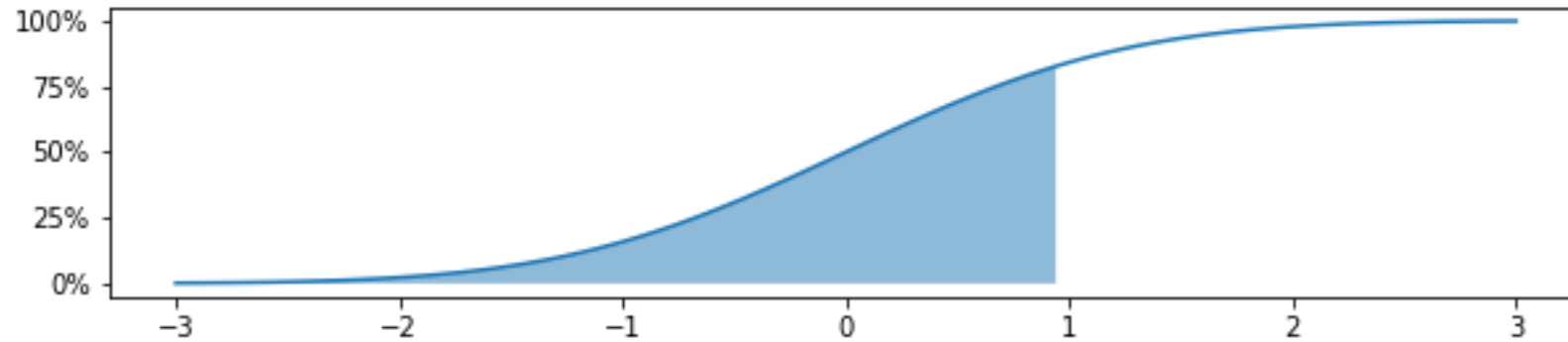
# Distribuição acumulada



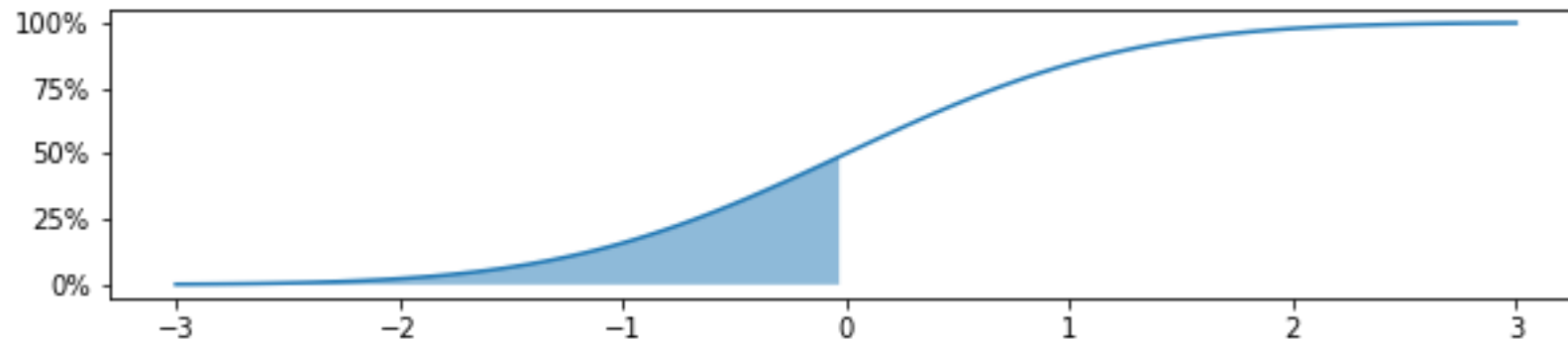
$$F(x) = \int_{-\infty}^x f(x)dx = P(X \leq x)$$

$$F(+\infty) = 1$$

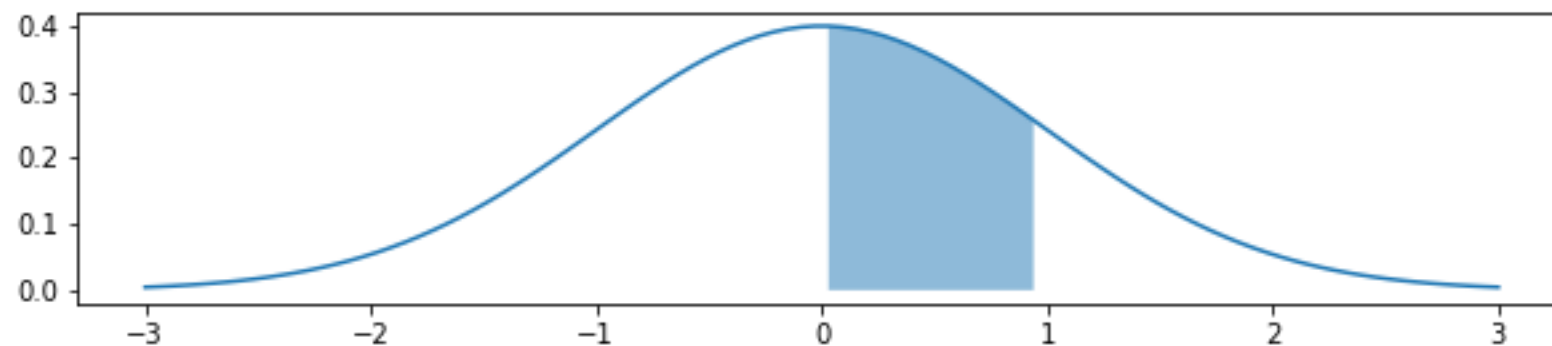
# Função Densidade de Probabilidade



$F(1)$



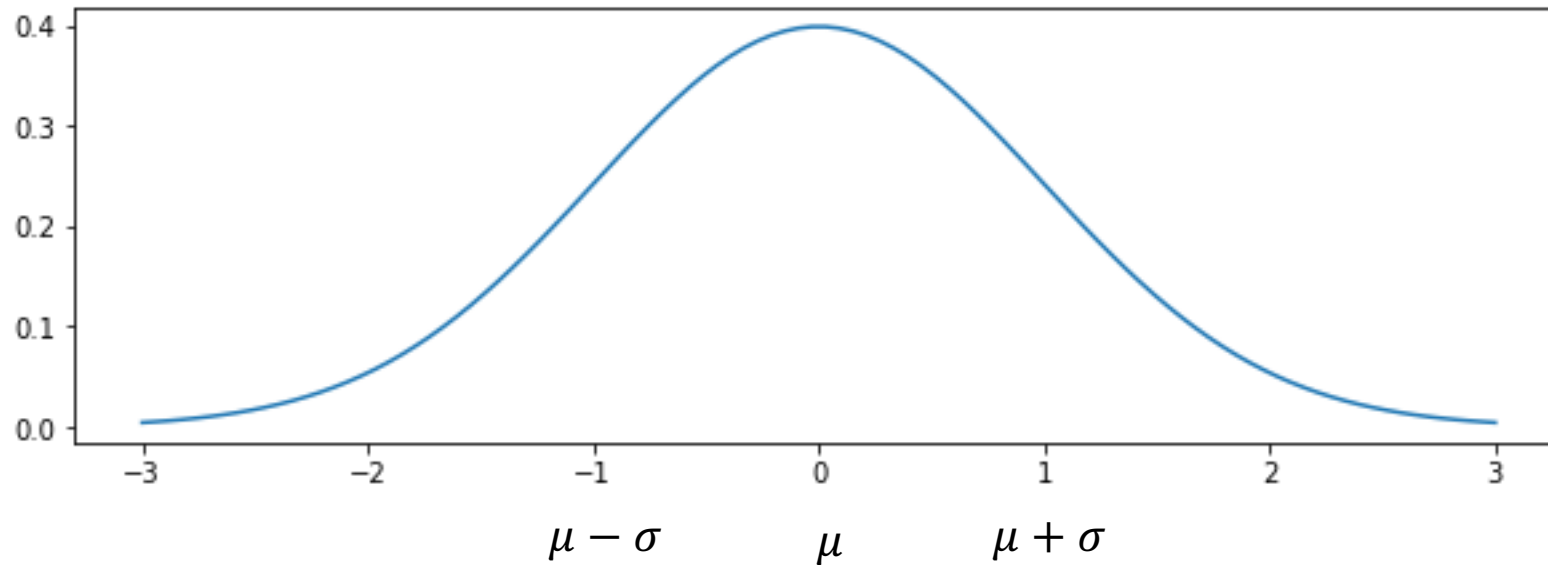
$F(0)$



$$P(0 < X < 1) = F(1) - F(0)$$

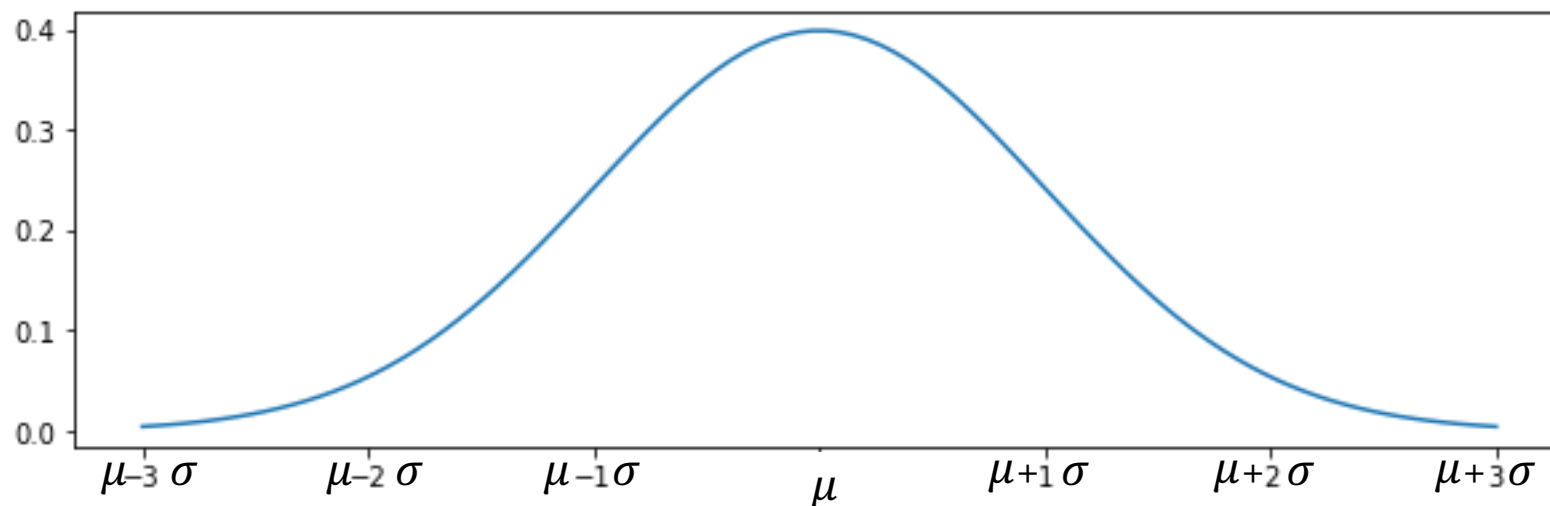
# Distribuição Normal

# Distribuição normal

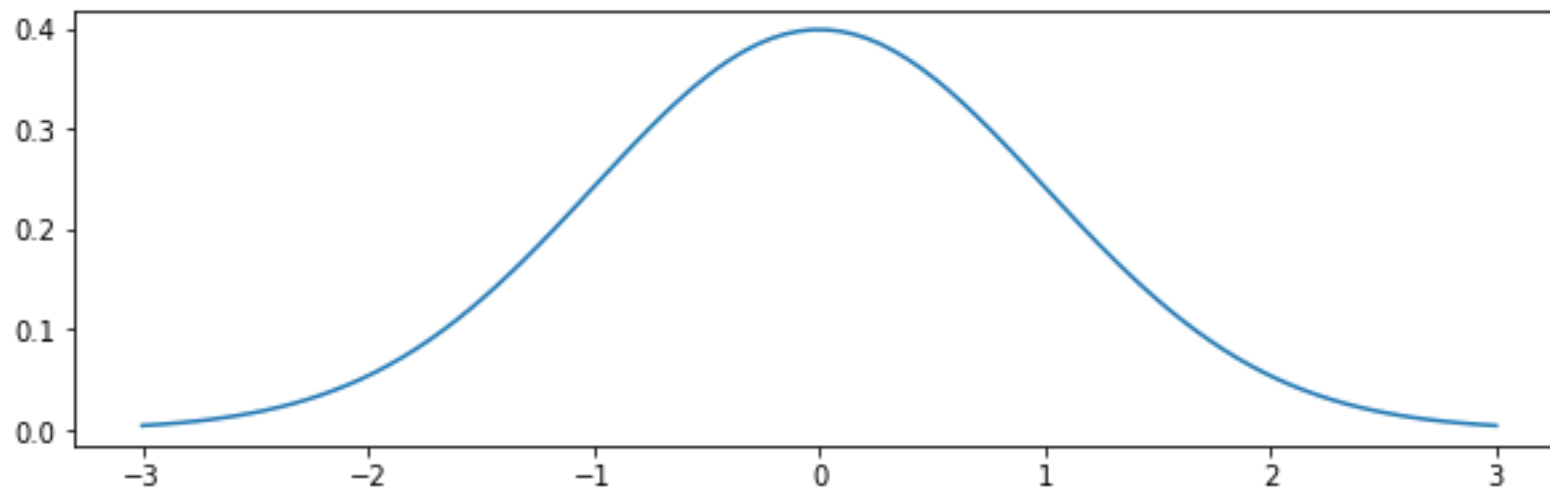


$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

# Distribuição normal



$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

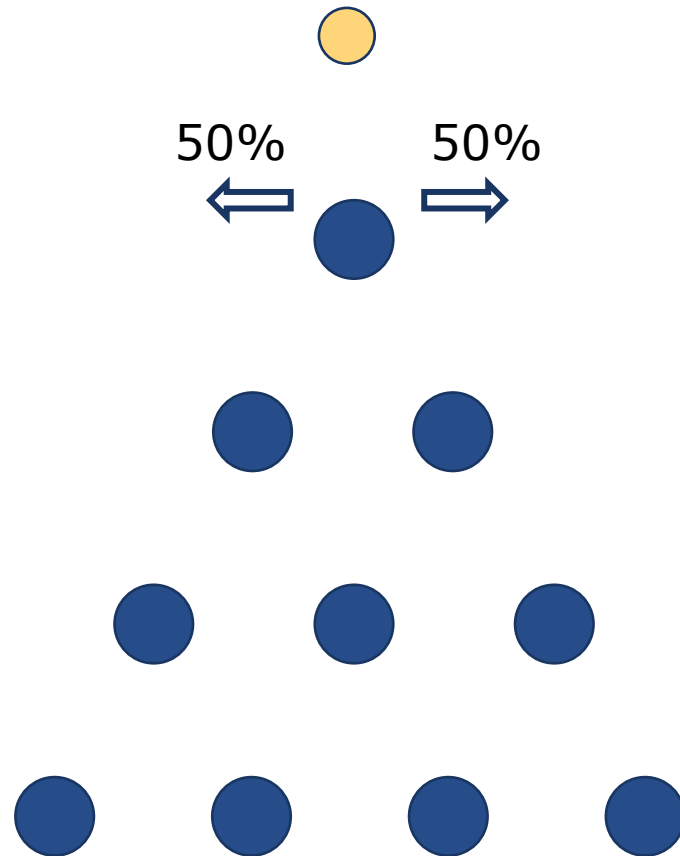


$$x^* = \frac{x - \mu}{\sigma}$$

$$f(x^*) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x^*)^2}$$

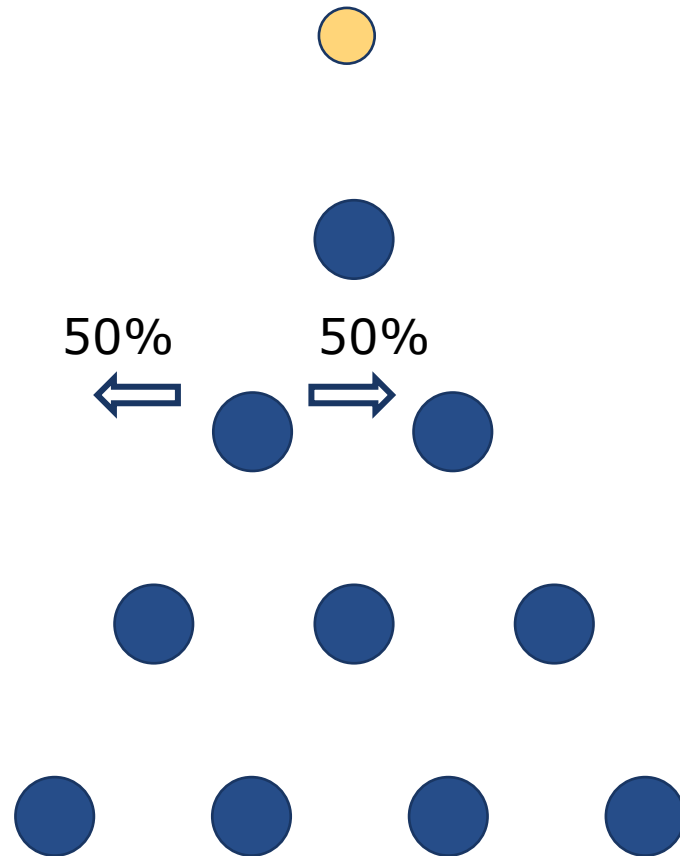


# Por que a Normal é tão importante?



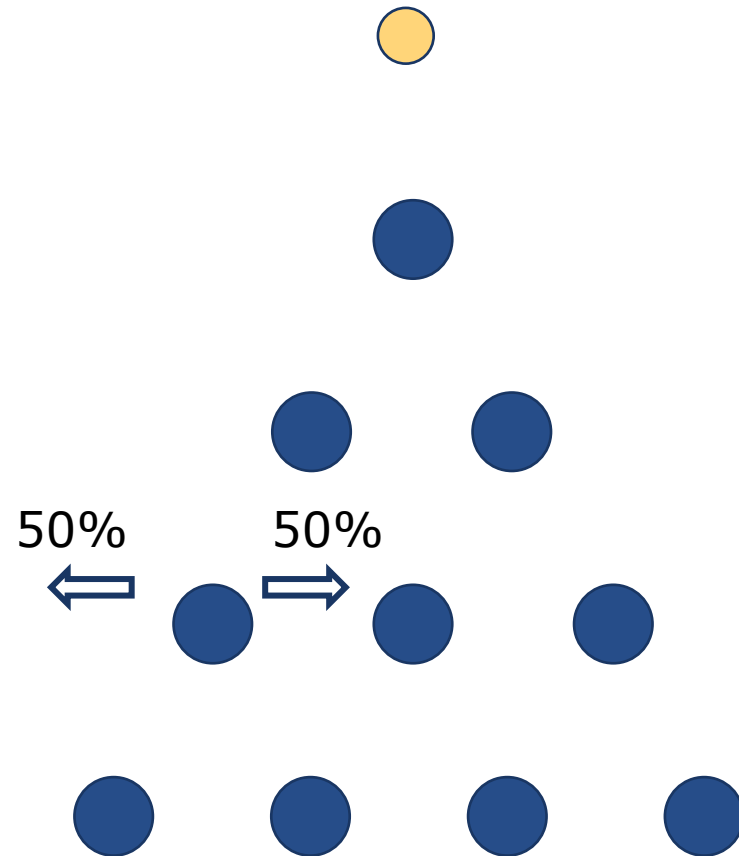
Galton Board ou Tabuleiro de Galton:

# Por que a Normal é tão importante?



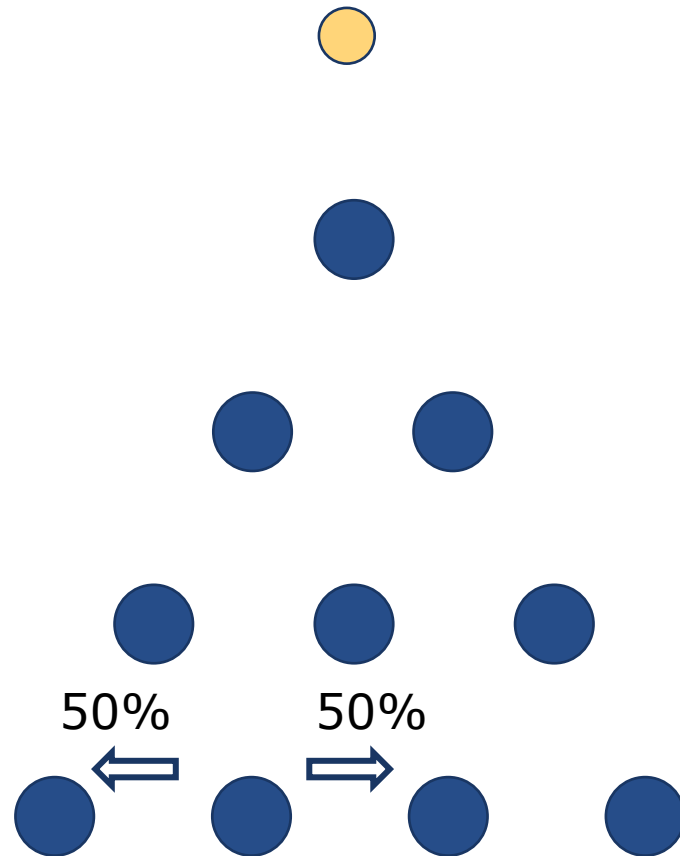
Galton Board ou Tabuleiro de Galton:

# Por que a Normal é tão importante?



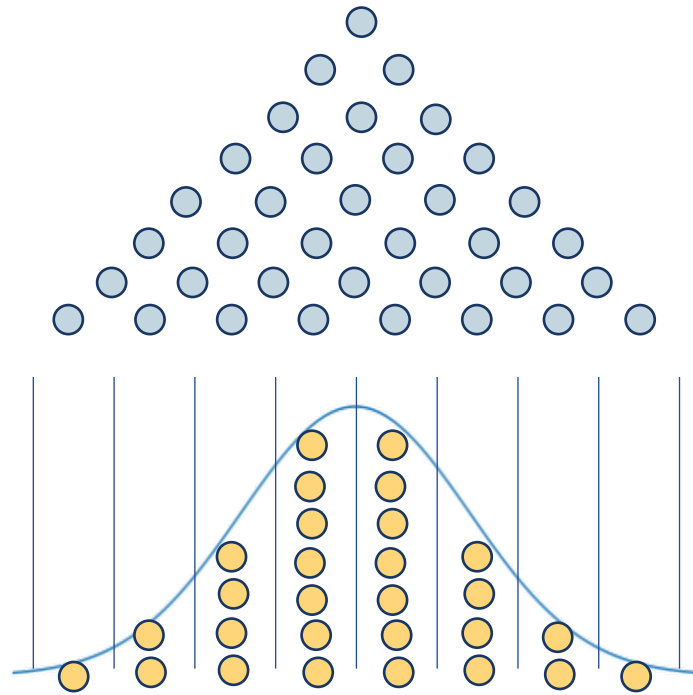
Galton Board ou Tabuleiro de Galton:

# Por que a Normal é tão importante?



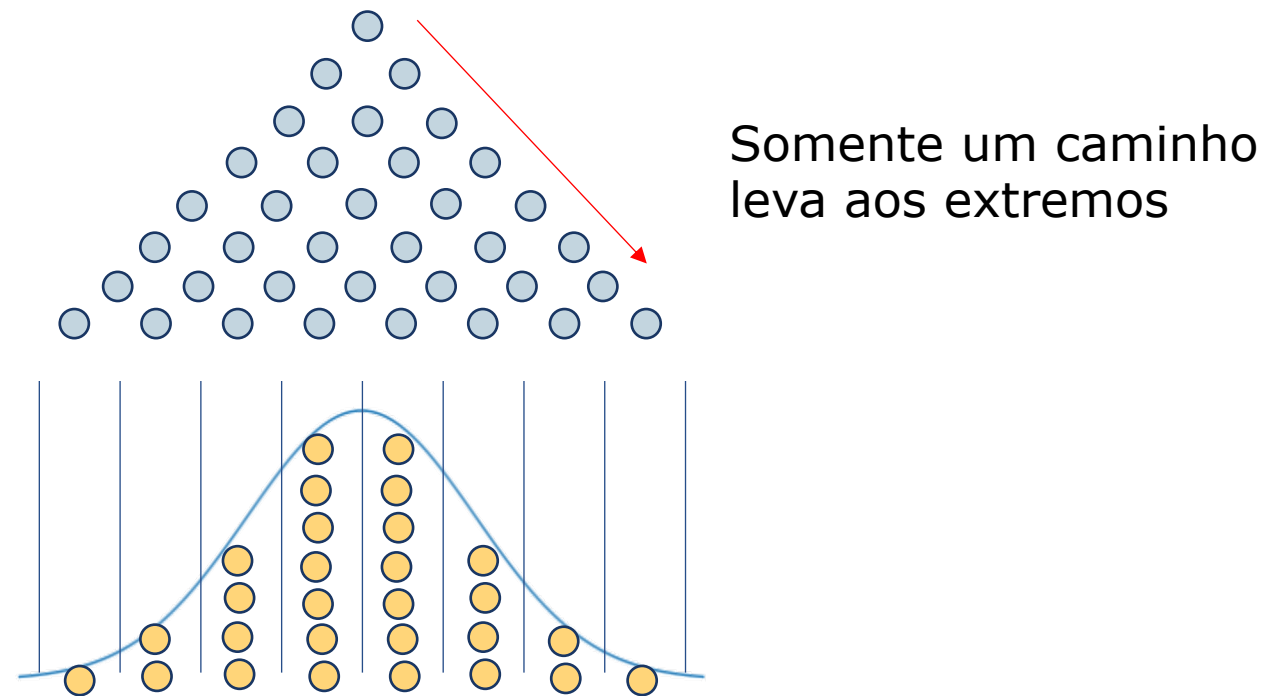
Galton Board ou Tabuleiro de Galton:

# Por que a Normal é tão importante?



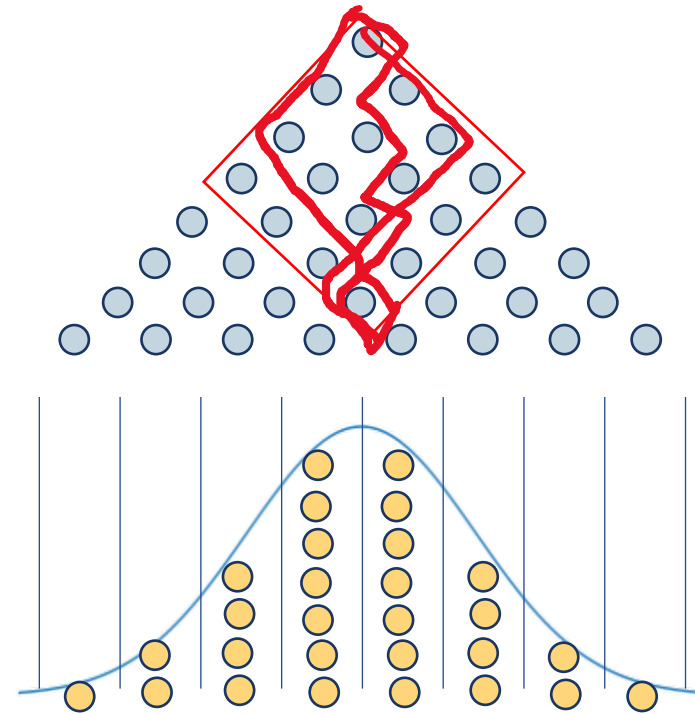
Galton Board ou Tabuleiro de Galton:

# Por que a Normal é tão importante?



Galton Board ou Tabuleiro de Galton:

# Por que a Normal é tão importante?



Vários caminhos levam ao meio

Galton Board ou Tabuleiro de Galton:

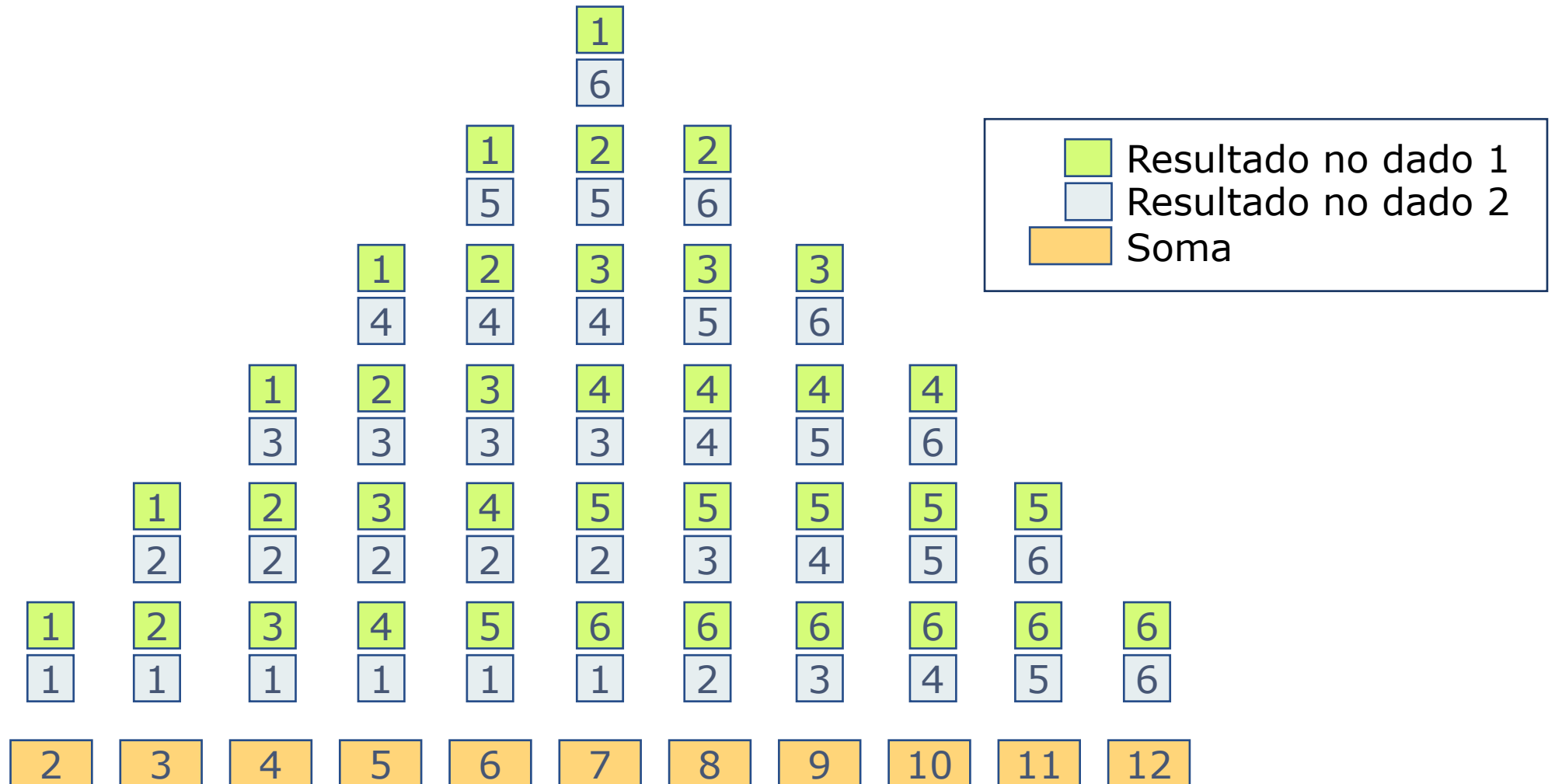
# Por que a Normal é tão importante?



Galton Board ou Tabuleiro de Galton:



# Conjunta, condicional e independência



# Combinação linear

- Se  $X_1 \sim N(\mu, \sigma^2)$
- $X_1^* = a + bX_1 \sim N(a + \mu, b^2\sigma^2)$  com a e b constantes.
- Ou seja:
  - A média de  $X_1^*$  é  $a + \mu$
  - A variância de  $X_1^*$  é  $b^2 \cdot \sigma^2$
  - O que significa que o desvio padrão de  $X_1^*$  é  $b \cdot \sigma$  (que é  $\sqrt{b^2\sigma^2}$ )
- POR ISSO QUE SE  $X_1 \sim N(\mu, \sigma^2)$  :
  - $X_1^* = \frac{X_1 - \mu}{\sigma} \Rightarrow X_1^* \sim N(0,1)$

# Combinação linear

- Se  $X_1 \sim N(\mu, \sigma^2)$  e  $X_2 \sim N(\mu, \sigma^2)$
- E se  $X_1$  e  $X_2$  são independentes,

$$X_1 + X_2 \sim N(2\mu, 2\sigma^2)$$

# Distribuição da média

- Dos dois resultados anteriores, temos a distribuição da média:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

Ou seja:

$$\text{Desvio Padrão}(\bar{X}) = \frac{\sigma}{\sqrt{n}}$$

Ou ainda:

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

# Distribuição da média

# Teorema Central do Limite

Se temos uma amostra  $X_1, X_2, \dots, X_n$  iid (independentes, identicamente distribuídas), então:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \approx N\left(\mu, \frac{\sigma^2}{n}\right)$$

Ou ainda:

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

# Teorema do Limite Central

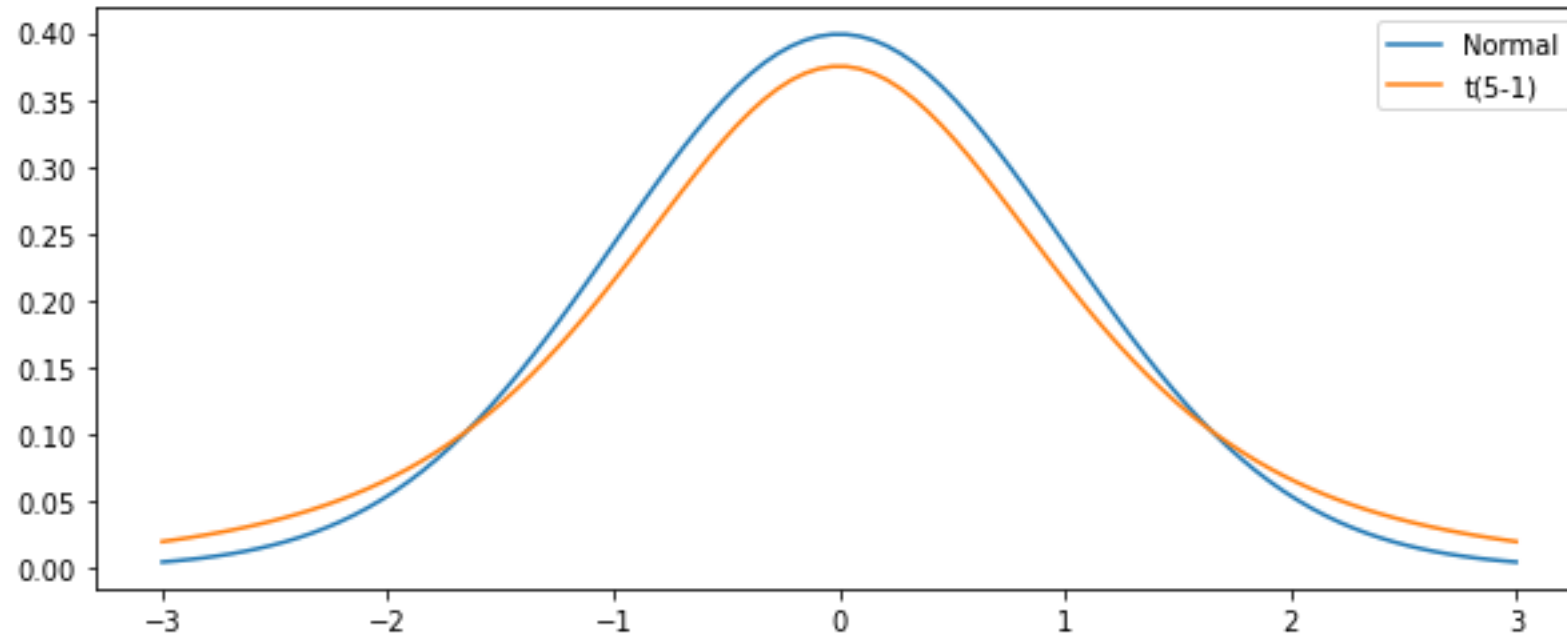
Mas se o parâmetro de variância  $\sigma^2$  não é conhecido, nos resta substituí-lo pela sua estimativa  $S^2$ .

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \approx N \left( \mu, \frac{S^2}{n} \right)$$

Ou ainda:

$$\frac{\bar{X} - \mu}{S / \sqrt{n-1}} \approx t(n-1)$$

# Média com Variância não conhecida

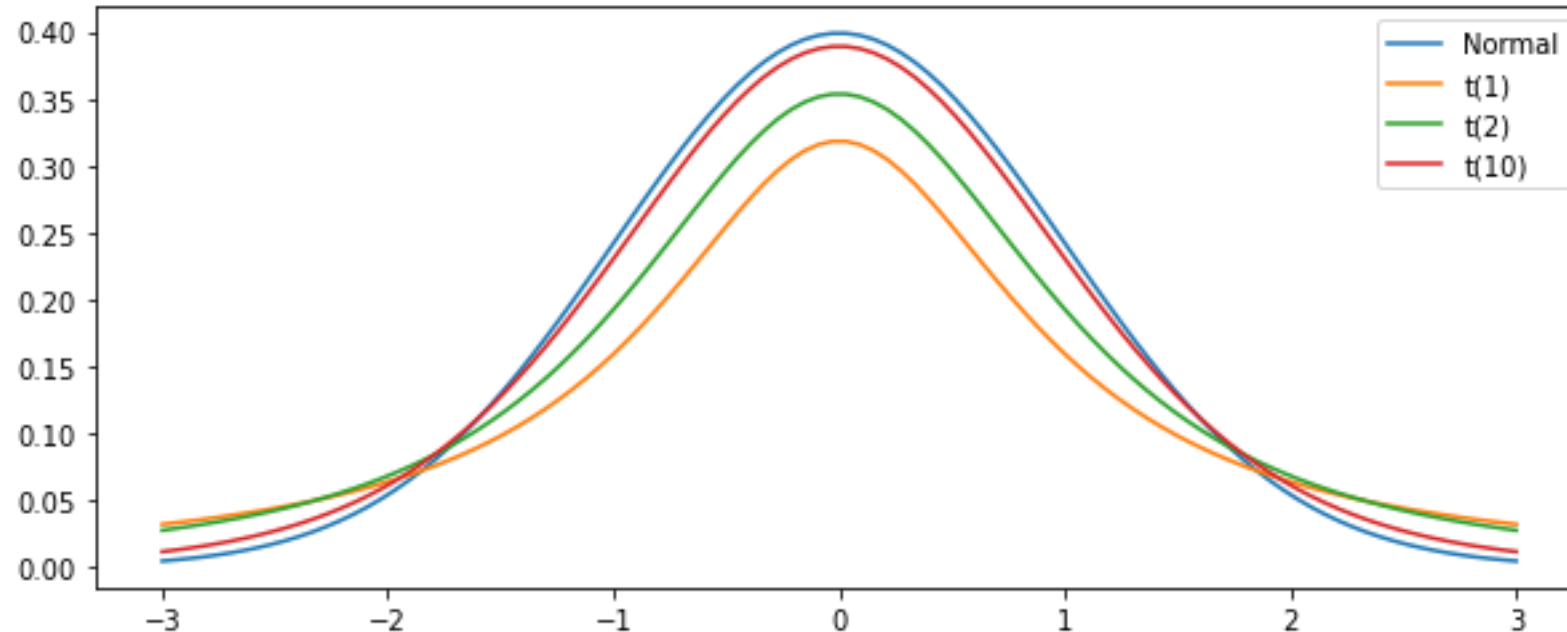


A distribuição t-student possui 'caudas mais pesadas'.

Isso significa que, com mesma média e desvio padrão, valores mais distantes da média são mais frequentes (têm maior probabilidade)



# Média com Variância não conhecida



Um fato interessante da distribuição t-student é que ela se aproxima da distribuição normal conforme o número de graus de liberdade aumenta. Na prática, a partir de 20 a aproximação pela normal já é bem razoável.