

# Inferência Estatística

## Correlação e Regressão

- Coeficiente de Correlação
- **Regressão Linear Simples**

# Regressão Linear Simples

- Em muitos problemas há duas ou mais variáveis que são relacionadas, e pode ser importante modelar essa relação.
- Por exemplo, pode-se ter interesse em prever
  - as vendas futuras de um produto em função do seu preço,
  - a perda de peso de uma pessoa em decorrência do número de dias que se submete a uma determinada dieta,
  - a produção de uma determinada cultura em função da quantidade de nutriente aplicada no solo.

- Outro exemplo, as vendas de um produto podem estar relacionadas ao valor gasto em marketing com esse produto. Assim, é possível construir um modelo relacionando vendas a gastos com marketing, e então pode-se usar esse modelo para fins previsão de vendas.
- Em geral vamos supor que há uma **variável dependente** (ou variável de resposta) **Y** que depende de uma **variável preditora** (ou variável explicativa) **X**.
- A regressão linear simples estima uma equação matemática (ou modelo) que, dado o valor de X (variável preditora), prevê o valor de Y (variável dependente).
- É dito regressão **linear simples**, pois supõe-se tendência linear entre as variáveis e simples por ser uma única variável preditora.

➤ Modelo de regressão linear simples

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

↗ erro aleatório

- O coeficiente  $\beta_0$  é a *interseção* (valor de  $Y$  para  $X = 0$ ).
- O coeficiente  $\beta_1$  é a *inclinação* da reta, que pode ser positiva, negativa ou nula.
- Se há  $n$  pares de dados  $(y_1, x_1), \dots, (y_n, x_n)$  é possível estimar os parâmetros  $\beta_0$  e  $\beta_1$  usando o Método dos Mínimos Quadrados.
- Temos então  $b_0$  e  $b_1$ , estimativas amostrais de  $\beta_0$  e  $\beta_1$ .

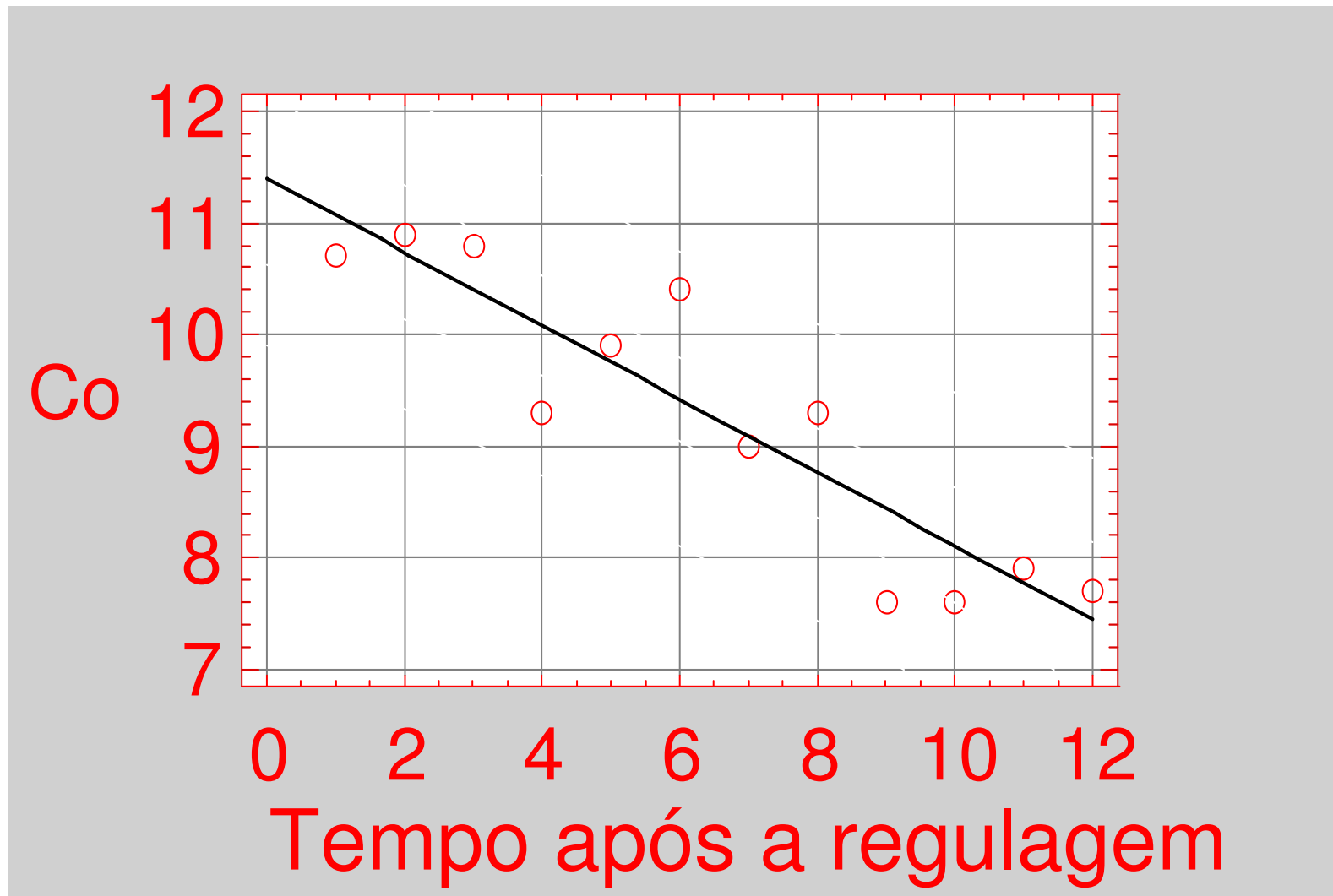
$$b_1 = \frac{\sum x_i y_i - (\sum x_i)(\sum y_i)/n}{\sum x_i^2 - (\sum x_i)^2/n}$$

$$b_0 = \bar{Y} - b_1 \bar{X}$$

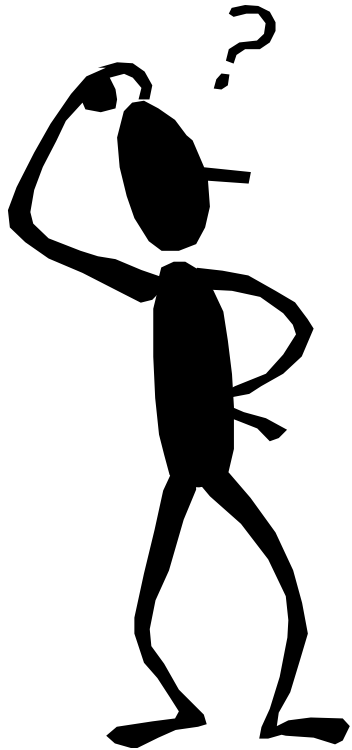
**Exemplo:** Após uma regulagem eletrônica um veículo apresenta um rendimento ideal no que tange a consumo de combustível. Contudo, com o passar do tempo esse rendimento vai se degradando. Os dados a seguir representam o rendimento medido mês a mês após a regulagem. Ajuste um modelo linear a esses dados.

X:meses após a regulagem	1	2	3	4	5	6
Y : rendimento	10,7	10,9	10,8	9,3	9,5	10,4
X:meses após a regulagem	7	8	9	10	11	12
Y : rendimento	9,0	9,3	7,6	7,6	7,9	7,7

# Rendimento de combustível



# Cálculos iniciais



Meses(X)	Rendimento(Y)	X^2	Y^2	X*Y
1	10,7	1	114,49	10,7
2	10,9	4	118,81	21,8
3	10,8	9	116,64	32,4
4	9,3	16	86,49	37,2
5	9,5	25	90,25	47,5
6	10,4	36	108,16	62,4
7	9	49	81	63
8	9,3	64	86,49	74,4
9	7,6	81	57,76	68,4
10	7,6	100	57,76	76
11	7,9	121	62,41	86,9
12	7,7	144	59,29	92,4
<b>78</b>	<b>110,7</b>	<b>650</b>	<b>1039,55</b>	<b>673,1</b>
<b>6,5</b>	<b>9,225</b>			

$$\Sigma x = 78$$

$$\Sigma y = 110,7$$

$$\Sigma x y = 673,1$$

$$\bar{X} = 6,50$$

$$\bar{Y} = 9,225$$

$$\Sigma x^2 = 650$$

$$\Sigma y^2 = 1039,55$$

# Cálculos

$$\begin{array}{lll} \Sigma x = 78 & \bar{X} = 6,50 & \Sigma x^2 = 650 \\ \Sigma y = 110,7 & \bar{Y} = 9,225 & \Sigma xy = 637,1 \end{array}$$

$$S_{xy} = \sum x_i y_i - (\sum x_i)(\sum y_i)/n = 673,1 - (78 \times 110,70)/12 = -46,45$$

$$S_{xx} = \sum x_i^2 - (\sum x_i)^2 / n = 650 - (78)^2 / 12 = 143$$

Estimativa dos parâmetros:

$$b_1 = -46,45 / 143,00 = -0,325$$

$$b_0 = 9,225 - (-0,325) 6,50 = 11,34$$

$$b_1 = \frac{\sum x_i y_i - (\sum x_i)(\sum y_i)/n}{\sum x_i^2 - (\sum x_i)^2 / n}$$

$$b_0 = \bar{Y} - b_1 \bar{X}$$

Equação de regressão

$$Y = 11,34 - 0,325 X$$



Interpretar e utilizar!!!



# Coeficiente de Determinação

$r^2$  é conhecido como Coeficiente de Determinação

$r^2$  = quadrado do coeficiente de correlação  $r$

*$r^2$  equivale a proporção da variância dos valores de  $Y$  que pode ser atribuída à regressão com a variável  $X$ .*

➤ Para o exemplo, resultou  $r^2 = (-0,907)^2 = 0,82$ , ou seja, 82% da variabilidade nos resultados de rendimento de combustível pode ser devida ao tempo decorrido após a regulagem.

18% da variabilidade total é devido a outros fatores que não foram investigados.

# Intervalos de Confiança e Testes de Hipótese

Como os resíduos de  $Y$  supostamente seguem a distribuição normal, e como os valores de  $b_0$  e  $b_1$  são funções lineares de  $Y$ :

$$b_0 \rightarrow N(\beta_0, \sigma_{b_0}^2) \quad b_1 \rightarrow N(\beta_1, \sigma_{b_1}^2)$$

Esses resultados podem ser usados em testes de hipótese. Por exemplo, se a hipótese é:

$$H_0 : \beta_1 = 0$$

$$H_A : \beta_1 \neq 0$$

então calcula-se:

$$t = b_1 / S_{b_1}$$

$\Rightarrow H_0$  é rejeitada se  $|t| > t_{\alpha/2, n-2}$  .

O intervalo de confiança para  $\beta_1$  virá dado por

$$b_1 - t_{\alpha/2} S_{b_1} < \beta_1 < b_1 + t_{\alpha/2} S_{b_1}$$

➤ Usando os dados do **problema do consumo de combustível**, vamos construir um intervalo de confiança para a inclinação  $b_1$  e verificar a hipótese.

$$S_{xx} = 143$$

$$S_{yy} = 18,34$$

$$S_{xy} = -46,45$$

$$SQR = S_{yy} - b_1 S_{xy} = 3,24$$

$$S^2 = \frac{SQR}{n-2} = 0,324$$

$$S_{b_1}^2 = \frac{S^2}{S_{xx}} = 0,00227 \Rightarrow S_{b_1} = 0,0476$$

Intervalo de confiança para  $\beta_1$

$$b_1 - t_{\alpha/2} S_{b1} < \beta_1 < b_1 + t_{\alpha/2} S_{b1}$$

$$t_{0,025;10} = 2,228$$

$$\begin{aligned} -0,325 - 2,228 (0,0476) < \beta_1 < -0,325 + 2,228 (0,0476) \\ -0,431 < \beta_1 < -0,219 \end{aligned}$$

Como esse intervalo não inclui o zero, a hipótese  $\beta_1 = 0$  é rejeitada, ou seja, existe uma relação entre o consumo de combustível e o tempo decorrido após a regulagem.