

# 1 Regressão e Correlação

Nas unidades anteriores, descrevemos a distribuição de valores de uma única variável, com esse objetivo aprendemos a calcular medidas de tendência central e variabilidade. Porém, se considerarmos duas ou mais variáveis surge um novo problema: as relações que podem existir entre as variáveis estudadas.

Vamos verificar as relações entre as seguintes variáveis:

- Altura e peso - espera-se que quanto mais alto mais pesado é o indivíduo;
- Quantidade de memória RAM e tempo de processamento - espera-se que com mais memória RAM tenha-se um tempo menor de processamento;
- Temperatura e Umidade do ar - não se pode associar a temperatura a uma menor ou maior umidade do ar.

Para estudar a relação entre duas variáveis quantitativas na utilizamos a análise de regressão e correlação destas variáveis.

Correlação é um número entre -1 e 1 que mede o grau relacionamento entre duas variáveis quantitativas

Regressão é o estudo que busca ajustar uma equação a um conjunto de dados de forma que a relação entre duas variáveis quantitativas possa ser expressa matematicamente.

Definimos um conjunto de variáveis  $(x, y)$ , sendo  $x$  a variável independente e  $y$  a variável dependente. A primeira forma de verificar a relação de duas variáveis é traçar o gráfico de dispersão do dados.

O gráfico de dispersão contém uma variável independente representada no eixo horizontal e a variável dependente representada no eixo vertical.

O gráfico de dispersão da um idéia da existência de correlação, entretanto não apresenta qual a magnitude da correlação. Para determinar a magnitude da correlação utilizamos o coeficiente de correlação populacional ( $\rho$ ). Em geral trabalhamos com amostras, e para estimar o coeficiente de correlação populacional pode-se utilizar o coeficiente de correlação amostral.

$$r = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}}$$

sendo que:

- $r > 0$  - correlação positiva;
- $r < 0$  - correlação negativa;
- $r = 0$  - ausência de correlação.

O valor obtido para o coeficiente de correlação amostral tem como finalidade estimar o populacional, ou seja, verificar se na população existe uma associação entre as variáveis em estudo.

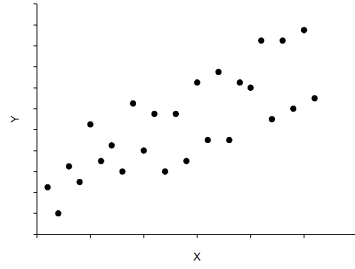


Figura 1: indícios de correlação positiva, aumentando x, y também aumenta

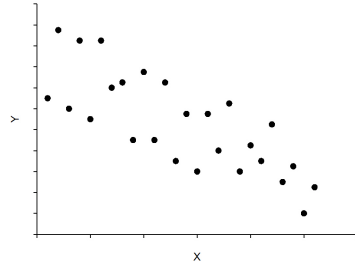


Figura 2: indícios de correlação negativa, aumentando x, y diminui

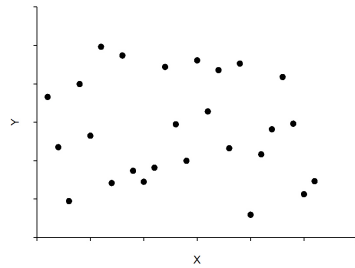


Figura 3: indícios de ausência correlação

Desta forma, deve ser realizado um teste de hipótese sobre o coeficiente populacional, com base no resultado obtido na amostra, que pode ser definido da seguinte maneira:

$$\begin{cases} H_0 : \rho = 0 \\ H_1 : \rho \neq 0 \end{cases}$$

Rejeita-se  $H_0$  se  $|t_c| > t_{\frac{\alpha}{2}}$ , em que

$$t_c = \frac{r - 0}{\sqrt{\frac{1-r^2}{n-2}}}$$

nesse caso  $v = n - 2$  graus de liberdade

Exemplo: Numa pesquisa feita com 7 famílias com renda bruta mensal entre 10 e 25 salários mínimos mediram-se:

- X: renda bruta mensal (em salários mínimos)
- Y: porcentagem da renda bruta anual gasta com assistência médica

x	10	12	14	16	18	20	22
y	11,8	10,2	12,1	13,2	15,1	15,4	15,6

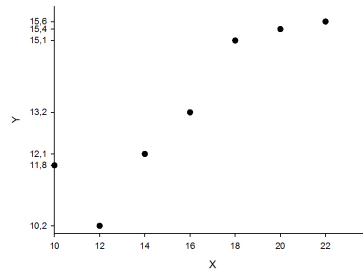


Figura 4: Gráfico de dispersão

$$\begin{aligned}\bar{x} &= \frac{\sum_{i=1}^n x_i}{n} = \frac{112}{7} = 16 \\ \bar{y} &= \frac{\sum_{i=1}^n y_i}{n} = \frac{93,4}{7} = 13,3 \\ r &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= \frac{49,6}{\sqrt{112 \times 26,25}} = 0,9148\end{aligned}$$

Verificou que o valor da correlação é  $r=0,9148$ . Vamos testar a hipótese se este valor é diferente de zero.

$$\begin{cases} H_0 : \rho = 0 \\ H_1 : \rho \neq 0 \end{cases}$$

Temos  $v = n - 2 = 7 - 2 = 5$  graus de liberdade

$$t_c = \frac{r - 0}{\sqrt{\frac{1-r^2}{n-2}}} = \frac{0,9148}{\sqrt{\frac{1-0,9148^2}{5}}} = 5,06$$

Tomando-se  $\alpha = 0,05$ , temos  $t_{0,025;5} = 2,571$ .

Como  $|t_c| > t_{\frac{\alpha}{2}}$ , rejeita-se  $H_0$  ao nível de 5% de significância. Logo a correlação é diferente de zero e é igual a 0,9148.

Pelo diagrama de dispersão e pelo coeficiente de correlação, verificamos que existe uma relação linear entre as variáveis  $X$  e  $Y$ , podemos determinar uma função que exprima esse relacionamento. A função que expressa a relação linear entre  $X$  e  $Y$  é dada por

$$y = a + bx + \epsilon$$

Tabela 1: Tabela auxiliar para o calculo da correlação

Observação	x	y	$(x - \bar{x})$	$(y - \bar{y})$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$	$(y - \bar{y})^2$
1	10	11,8	-6	-1,5	9	36	2,25
2	12	10,2	-4	-3,1	12,4	16	9,61
3	14	12,1	-2	-1,2	2,4	4	1,44
4	16	13,2	0	-0,1	0	0	0,01
5	18	15,1	2	1,8	3,6	4	3,24
6	20	15,4	4	2,1	8,4	16	4,41
7	22	15,6	6	2,3	13,8	36	5,29
Total	112	93,4			49,6	112	26,25

em que:

- $a$  é coeficiente linear, interpretado como o valor da variável de dependente quando a variável independente é igual a 0;
- $b$  é coeficiente de regressão, interpretado como acréscimo na variável dependente para a variação de uma unidade na variável.
- $\epsilon$  são os erros aleatórios de uma população normal, com média 0 e variância constante.

Os estimadores para os coeficientes são:

$$a = \bar{y} - b\bar{x} \quad b = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$$

A análise de variância é uma técnica utilizada para se testar o ajuste da equação como um todo, ou seja, um teste para verificar se a equação de regressão obtida é significativa ou não.

Tabela 2: Análise de Variância para Regressão Linear Simples

Fontes de Variação	GL	Soma de Quadrados (SQ)	Quadrado Médio (QM)	Fc
Regressão	1	SQRegressão	QMRegressão	QMRegressão/QMErro
Erro	n-2	SQErro	QMErro	
Total	n-1	SQTotal		

$$SQTotal = \sum_i (y_i - \bar{y})^2$$

$$SQRegressão = \frac{\left( \sum_i (x_i - \bar{x})(y_i - \bar{y}) \right)^2}{\sum_i (x_i - \bar{x})^2}$$

$$SQErro = SQTotal - SQRegressão$$

$$QMRegressão = \frac{SQRegressão}{1}$$

$$QMErro = \frac{SQErro}{n - 2}$$

$$\begin{cases} H_0 : a = 0 \text{ ou } b = 0 \\ H_1 : a \neq 0 \text{ e } b \neq 0 \end{cases}$$

O teste de hipótese para avaliar se o modelo de regressão é significativo é feito da seguinte forma:

- Estabelecer o nível de significância  $\alpha$ ;
- Obter o valor tabelado  $F_\alpha$ ;
- Rejeita-se a hipótese  $H_0$ , se  $F_c > F_\alpha$ .

O coeficiente de determinação  $r^2$ , é definido por:

$$r^2 = \frac{\text{SQRegressão}}{\text{SQTotal}} \quad 0 < r^2 < 1$$

ele representa a porcentagem da variação total que é explicada pela equação de regressão, quanto maior o seu valor melhor.

Após ter verificado o ajuste da equação de regressão pode-se utiliza-la para fazer previsões.

Exemplo: Utilizando o exemplo da renda bruta mensal (em salários mínimos) e a porcentagem da renda bruta anual gasta com assistência médica.

Vamos ajustar o modelo

$$y = a + bx$$

Utilizando os calculo da tabela 1

$$\begin{aligned} b &= \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} \\ &= \frac{49,6}{112} = 0,44 \\ a &= \bar{y} - b\bar{x} \\ &= 6,26 \end{aligned}$$

Assim a equação de regressão é igual a

$$y = 6,26 + 0,44x$$

Vamos verificar se a regressão é significativa

$$\begin{aligned}
 \text{SQTotal} &= \sum_i (y_i - \bar{y})^2 = 26,25 \\
 \text{SQRegressão} &= \frac{\left( \sum_i (x_i - \bar{x})(y_i - \bar{y}) \right)^2}{\sum_i (x_i - \bar{x})^2} \\
 &= \frac{(49,6)^2}{112} = 21,97 \\
 \text{SQErro} &= \text{SQTotal} - \text{SQRegressão} \\
 &= 26,25 - 21,97 = 4,28
 \end{aligned}$$

Tabela 3: Análise de Variância para Regressão Linear Simples

Fontes de Variação	GL	Soma de Quadrados (SQ)	Quadrado Médio (QM)	Fc	$F_\alpha$
Regressão	1	21,97	21,97	25,55	6,60
Erro	5	4,28	0,86		
Total	6	26,25			

Como o  $F_c > F_\alpha$ , rejeita-se  $H_0$ , logo o modelo de regressão linear é significativo.

Obtendo o  $r^2$

$$r^2 = \frac{\text{SQRegressão}}{\text{SQTotal}} = \frac{21,97}{26,25} = 0,8370 = 83,70\%$$

Assim verifica-se que é a renda bruta explica 83,70% da variação do gasto com assistência médica.