

1 - CORRELAÇÃO LINEAR SIMPLES – r_{xy}

Em pesquisas, freqüentemente, procura-se verificar se existe relação entre duas ou mais variáveis, isto é, saber se as alterações sofridas por uma das variáveis são acompanhadas por alterações nas outras. Por exemplo, peso vs. idade, consumo vs. renda, altura vs. peso, de um indivíduo.

O termo correlação significa relação em dois sentidos (co + relação), e é usado em estatística para designar a força que mantém unidos dois conjuntos de valores. A verificação da existência e do grau de relação entre as variáveis é o objeto de estudo da correlação.

Uma vez caracterizada esta relação, procura-se descrevê-la sob forma matemática, através de uma função. A estimação dos parâmetros dessa função matemática é o objeto da regressão.

Os pares de valores das duas variáveis poderão ser colocados num diagrama cartesiano chamado “diagrama de dispersão”. A vantagem de construir um diagrama de dispersão está em que, muitas vezes sua simples observação já nos dá uma idéia bastante boa de como as duas variáveis se relacionam.

Uma medida do grau e do sinal da correlação é dada pela covariância entre as duas variáveis aleatórias X e Y que é uma medida numérica de associação linear existente entre elas, e definida por:

$$\text{Cov}(X, Y) = \frac{1}{n} \left[\sum x.y - \frac{\sum x. \sum y}{n} \right].$$

É mais conveniente usar para medida de correlação, o coeficiente de correlação linear de Pearson, como estimador de ρ_{xy} , definido por:

$$r_{xy} = \frac{\text{Cov}(x, y)}{\sqrt{\sigma_x \sigma_y}} = \frac{S_{xy}}{\sqrt{S_{xx} \cdot S_{yy}}}$$

$$r_{xy} = \frac{\sum xy - \frac{\sum x \cdot \sum y}{n}}{\left[\left[\sum x^2 - \frac{(\sum x)^2}{n} \right] \cdot \left[\sum y^2 - \frac{(\sum y)^2}{n} \right] \right]^{\frac{1}{2}}} = \frac{S_{xy}}{(S_{xx} \cdot S_{yy})^{\frac{1}{2}}} = \sqrt{\frac{S_{xy} \cdot S_{xy}}{S_{xx} \cdot S_{yy}}} = \sqrt{\frac{b \cdot S_{xy}}{S_{yy}}}$$

onde: as somas de quadrados são:

$$S_{xy} = \sum x \cdot y - \frac{\sum x \cdot \sum y}{n} \quad S_{xx} = \sum x^2 - \frac{(\sum x)^2}{n} \quad S_{yy} = \sum y^2 - \frac{(\sum y)^2}{n}$$

n = número de pares das observações.

A partir de X e Y são determinadas todas as somas necessárias para este cálculo:

Y	X	X ²	Y ²	X . Y
:	:	:	:	:
:	:	;	:	:
Σy	Σx	Σx^2	Σy^2	$\Sigma(x.y)$

O coeficiente de correlação r_{xy} linear é um número puro que varia de -1 a $+1$ e sua interpretação dependerá do valor numérico e do sinal, como segue:

$r_{xy} = -1$	\Rightarrow	correlação perfeita negativa
$-1 < r_{xy} < 0$	\Rightarrow	correlação negativa
$r_{xy} = 0$	\Rightarrow	correlação nula
$0 < r_{xy} < 1$	\Rightarrow	correlação positiva
$r_{xy} = 1$	\Rightarrow	correlação perfeita positiva
$0,2 < r_{xy} < 0,4$	\Rightarrow	correlação fraca*
$0,4 < r_{xy} < 0,7$	\Rightarrow	correlação moderada*
$0,7 < r_{xy} < 0,9$	\Rightarrow	correlação forte*

*possui o mesmo significado para os casos negativos ou positivos.

Análise do Diagrama de Dispersão

O diagrama de dispersão mostrará que a correlação será tanto mais forte quanto mais próximo estiver o coeficiente de -1 ou $+1$, e será tanto mais fraca quanto mais próximo o coeficiente estiver de zero.

a) Correlação perfeita negativa ($r_{xy} = -1$): Quando os pontos estiverem perfeitamente alinhados, mas em sentido contrário, a correlação é denominada perfeita negativa.

b) Correlação negativa ($-1 < r_{xy} < 0$): A correlação é considerada negativa quando valores crescentes da variável X estiverem associados a valores decrescentes da variável Y, ou valores decrescentes de X associados a valores crescentes de Y.

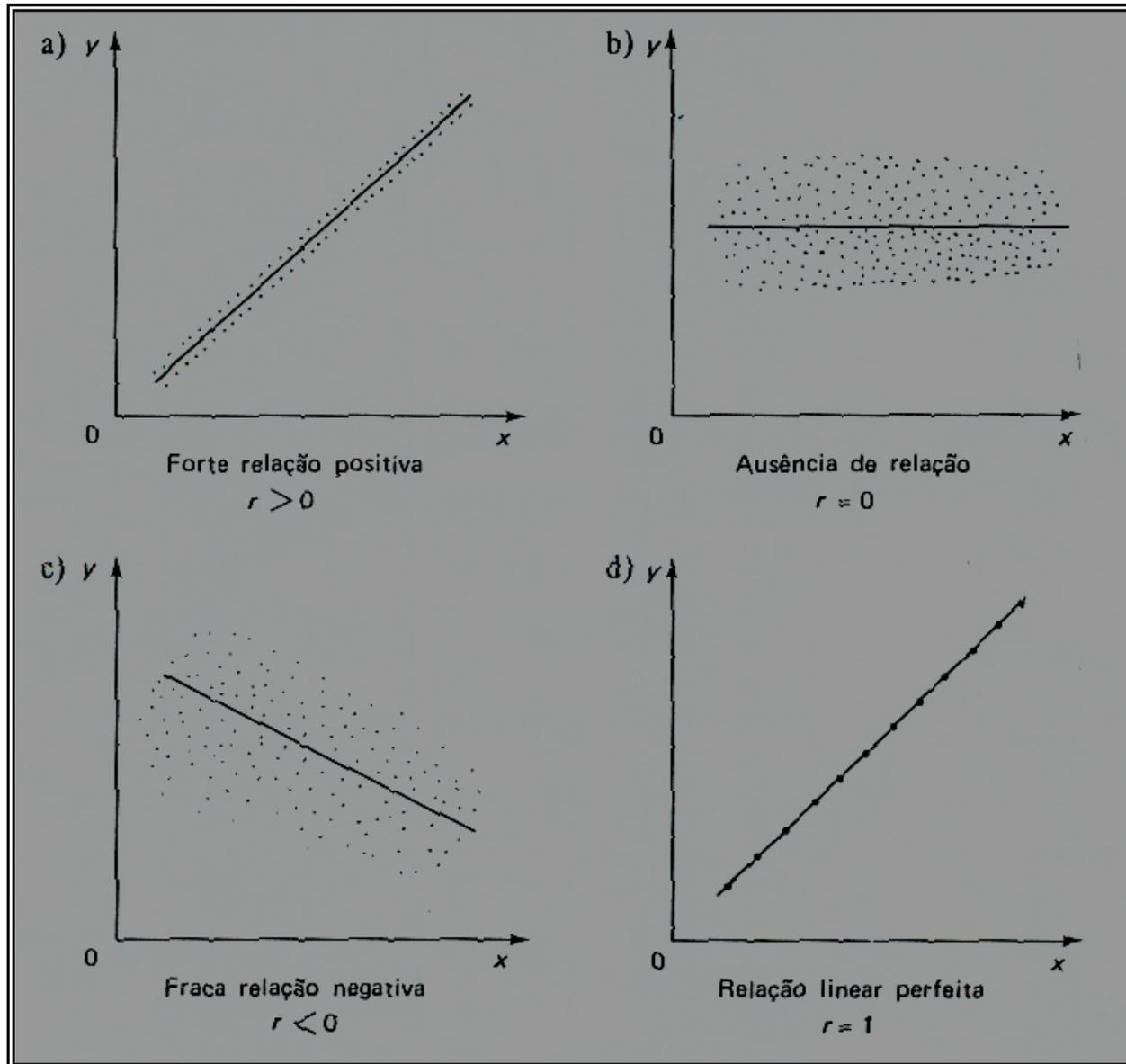
c) **Correlação nula ($r_{xy} = 0$):** Quando não houver relação entre as variáveis X e Y , ou seja, quando os valores de X e Y ocorrerem independentemente, não existe correlação entre elas.

d) **Correlação positiva ($0 < r_{xy} < 1$):** Será considerada positiva se os valores crescentes de X estiverem associados a valores crescentes de Y .

e) **Correlação perfeita positiva ($r_{xy} = 1$):** A correlação linear perfeita positiva corresponde ao caso anterior, só que os pontos (X, Y) estão perfeitamente alinhados.

f) **Correlação espúria:** Quando duas variáveis X e Y forem independentes, o coeficiente de correlação será nulo. Entretanto, algumas vezes, isto não ocorre, podendo, assim mesmo, o coeficiente apresentar um valor próximo de -1 ou $+1$. Neste caso a correlação é espúria.

Algumas situações que podem se apresentar os diagramas de dispersão



OBSERVAÇÕES:

- ⇒ Correlação não é o mesmo que causa e efeito. Duas variáveis podem estar altamente correlacionadas e, no entanto, não haver relação de causa e efeito entre elas.
- ⇒ Se duas variáveis estiverem amarradas por uma relação de causa e efeito elas estarão, obrigatoriamente, correlacionadas.
- ⇒ O estudo de correlação pressupõe que as variáveis X e Y tenham uma distribuição normal.
- ⇒ A palavra simples que compõe o nome correlação linear simples, indica que estão envolvidas no cálculo somente duas variáveis.
- ⇒ O coeficiente de correlação linear de Pearson mede a correlação em estatística paramétrica.
- ⇒ Coeficiente de correlação de Spearman (correlação por postos) é o correspondente à área não paramétrica.

$$\Rightarrow \text{Var}(x) = \frac{S_{xx}}{n} \quad \text{Cov}(x,y) = \frac{S_{xy}}{n} \quad \text{Var}(y) = \frac{S_{yy}}{n}$$

- ⇒ Testar $\rho = 0$ é equivalente a testar $\beta = 0$ na equação de regressão, pois

$$\hat{\rho}^2 = \hat{\beta}^2 \cdot \frac{S_{xx}}{S_{yy}} = \frac{SQ_{\text{Reg}}}{SQ_{\text{Total}}}$$

TESTES DO COEFICIENTE DE CORRELAÇÃO – SIGNIFICÂNCIA DE r_{xy}

O coeficiente de correlação r_{xy} é apenas uma estimativa do coeficiente de correlação populacional ρ_{xy} e não devemos esquecer que o valor de r_{xy} é calculado com base em de “n” pares de dados constituindo amostras aleatórias.

Muitas vezes os pontos da amostra podem apresentar uma correlação e, no entanto a população não, neste caso, estamos diante de um problema de inferência, pois $r_{xy} \neq 0$ não é garantia de que $\rho_{xy} \neq 0$.

Podemos resolver o problema aplicando um teste de hipóteses para verificarmos se o valor de r_{xy} é coerente com o tamanho da amostra n , a um nível de significância α , que realmente existe correlação linear entre as variáveis.

H0: $\rho = 0$ (não existe correlação entre X e Y)

H1: $\rho \neq 0$ (existe correlação entre X e Y).

$$t_c = \frac{r_{xy} \cdot \sqrt{n-2}}{\sqrt{1-r_{xy}^2}} = \frac{r_{xy}}{S_r} \approx \text{distribuição "t" de Student com } n-2 \text{ graus de liberdade.}$$

Onde, $S_r = \sqrt{\frac{1-r^2}{n-2}}$, é o erro padrão do coeficiente de correlação.

COEFICIENTE DE DETERMINAÇÃO r_{xy}^2

Indica a proporção de variação da variável independente que é explicada pela variável dependente, ou seja, é uma ferramenta que avalia a qualidade do ajuste.

$$R^2 = r_{xy}^2, 0 \leq R^2 \leq 1$$

Quanto mais próximo da unidade o R^2 estiver, melhor a qualidade do ajuste. O seu valor fornece a proporção da variável Y explicada pela variável X através da função ajustada.

Exemplo: $R^2 = r_{xy}^2 = (0,9929)^2 = 0,9858 = 98,50 \%$.

É a proporção que Y é explicada por X ; ou seja; 98,50% da variação do número de livros é explicado pelo tempo que freqüentou a escola.

CORRELAÇÃO LINEAR POR POSTOS OU SPEARMAN - r_s

De todas as estatísticas baseadas em postos, o coeficiente de correlação por postos de Spearman, foi a que surgiu primeiro, e é talvez a mais conhecida hoje. É uma medida de associação que exige que ambas as variáveis se apresentem em escala de mensuração pelo menos ordinal, de modo que os elementos em estudo possam dispor-se por postos em duas séries ordenadas.

Este teste não-paramétrico destina-se a determinar o grau de associação entre duas variáveis X e Y, dispostas em pontos ordenados, o objetivo é estudar a correlação entre duas classificações.

Resumo do Procedimento

1º) Dispor em postos as duas variáveis X e Y de 1 a n (n=número de pares de dados);

2º) Relacionar os n elementos, dar o posto de cada elemento;

3º) Determinar $d_i = (\text{posto } x - \text{posto } y)$, d_i^2 e $\sum d_i^2$;

4º) Se a proporção de empates de ambas as variáveis X ou Y é grande então calcula-se r_s pela fórmula:

$$r_s = \frac{\sum x^2 + \sum y^2 - \sum d_i^2}{2 \cdot \sqrt{\sum x^2 \cdot \sum y^2}}$$

Onde: $\sum x^2 = \frac{n^3 - n}{12} - \sum Tx$ $\sum y^2 = \frac{n^3 - n}{12} - \sum Ty$

$T = \frac{t^3 - t}{12}$, onde t, corresponde ao número de empates, usado para corrigir a soma de quadrados.

Caso contrário se aplica a fórmula: $r_s = 1 - \frac{6 \sum di^2}{n^3 - n}$

5º) A significância de r_s é testada com $t_c = \frac{r_s \sqrt{n-2}}{\sqrt{1-r_s^2}}$; com n-2 graus de liberdade, que é o mesmo teste anterior (Pearson).