

Correlação

Prof. Lorí Viali, Dr.

viali@mat.ufrgs.br

<http://www.mat.ufrgs.br/~viali/>

*É o grau de associação entre
duas ou mais variáveis. Pode ser:*

correlacional

ou

experimental.



Numa relação *experimental* os valores de uma das variáveis são controlados.

No relacionamento *correlacional*, por outro lado, não se tem nenhum controle sobre as variáveis sendo estudadas.



Indicadores de Associação



Um engenheiro químico está investigando o efeito da temperatura de operação do processo no rendimento do produto. O estudo resultou nos dados da tabela seguinte:



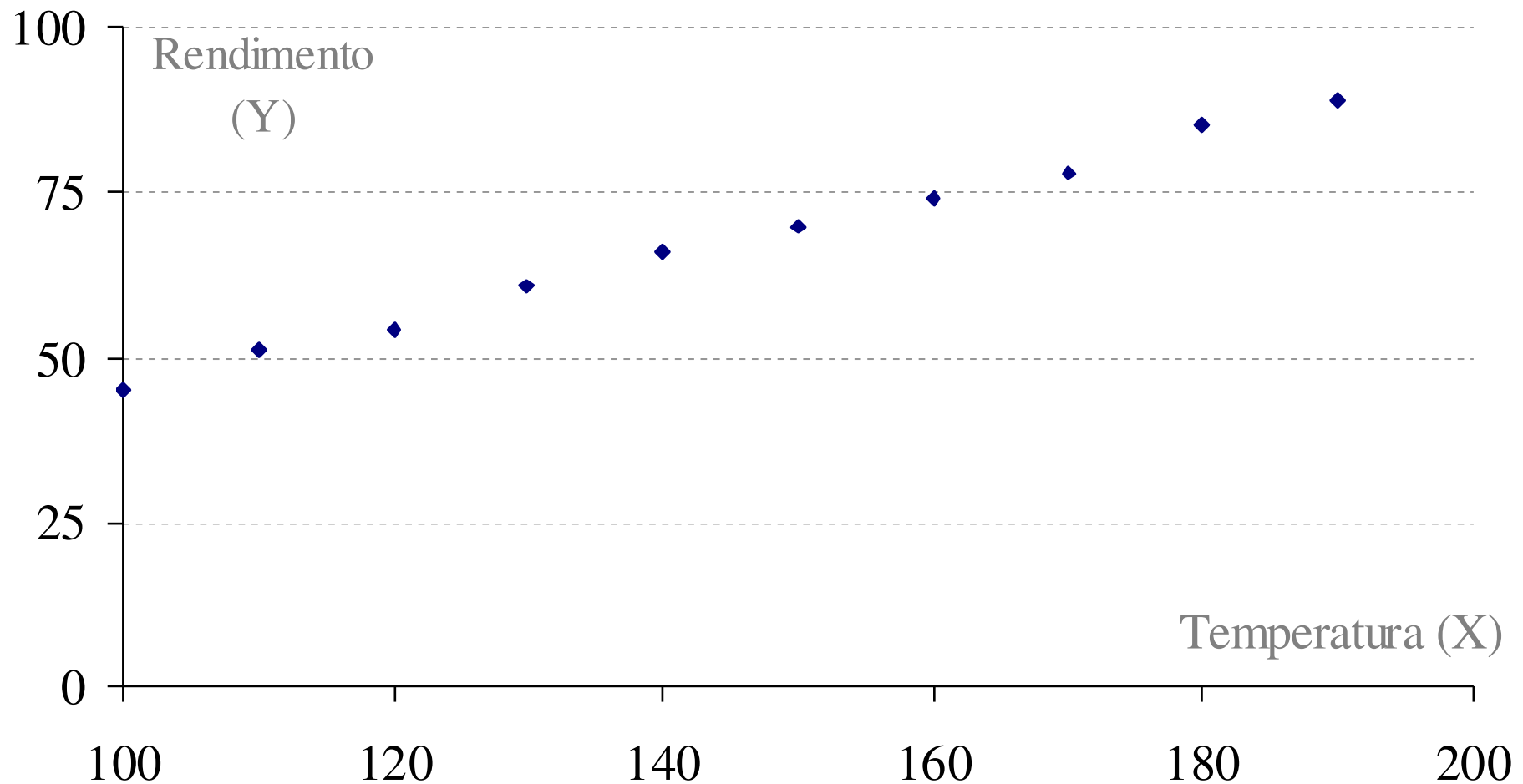
<i>Temperatura, C° (X)</i>	<i>Rendimento (Y)</i>
100	45
110	51
120	54
130	61
140	66
150	70
160	74
170	78
180	85
190	89



O primeiro passo para determinar se existe relacionamento entre as duas variáveis é obter o diagrama de dispersão (scatter diagram).



Diagrama de Dispersão



O diagrama de dispersão fornece uma idéia do tipo de relacionamento entre as duas variáveis. Neste caso, percebe-se que existe um relacionamento linear.



*Quando o relacionamento
entre duas variáveis
quantitativas for do tipo linear,
ele pode ser medido através do:*



Coeficiente de Correlação



Observado um relacionamento linear entre as duas variáveis é possível determinar a intensidade deste relacionamento. O coeficiente que mede este relacionamento é denominado de Coeficiente de Correlação (linear).



Quando se está trabalhando com amostras o coeficiente de correlação é indicado pela letra “r” e é uma estimativa do coeficiente de correlação populacional que é representado por “ ρ ” (rho).



*Determinação
do
Coeficiente
de
Correlação*



Para determinar o coeficiente de correlação (grau de relacionamento linear entre duas variáveis) vamos determinar inicialmente a variação conjunta entre elas, isto é, a covariância.



A covariância entre duas variáveis X e Y , é representada por “ $Cov(X; Y)$ ” e calculada por:

$$Cov(X, Y) = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$



Mas

$$\begin{aligned}\sum (x_i - \bar{x})(y_i - \bar{y}) &= \\ &= \sum [x_i y_i - \bar{x} y_i - x_i \bar{y} + \bar{x} \bar{y}] = \\ &= \sum x_i y_i - \sum \bar{x} y_i - \sum x_i \bar{y} + \sum \bar{x} \bar{y} = \\ &= \sum x_i y_i - \bar{x} \sum y_i - \bar{y} \sum x_i + \sum \bar{x} \bar{y} = \\ &= \sum x_i y_i - n \bar{x} \bar{y} - n \bar{x} \bar{y} + n \bar{x} \bar{y} = \\ &= \sum x_i y_i - n \bar{x} \bar{y}\end{aligned}$$



Então:

$$\begin{aligned} \text{Cov}(X, Y) &= \frac{\sum (x_i - \bar{X})(y_i - \bar{Y})}{n - 1} = \\ &= \frac{\sum x_i y_i - n\bar{X}\bar{Y}}{n - 1} \end{aligned}$$



A covariância poderia ser utilizada para medir o grau e o sinal do relacionamento entre as duas variáveis, mas ela é difícil de interpretar por variar de $-\infty$ a $+\infty$. Assim vamos utilizar o coeficiente de correlação linear de Pearson.



O coeficiente de correlação linear (de Pearson) é definido por:

$$r = \frac{\text{Cov}(X, Y)}{S_X S_Y}$$



Onde:

$$\text{Cov}(X, Y) = \frac{\sum X_i Y_i - n \bar{X} \bar{Y}}{n - 1}$$

$$S_X = \sqrt{\frac{\sum X_i^2 - n \bar{X}^2}{n - 1}}$$

$$S_Y = \sqrt{\frac{\sum Y_i^2 - n \bar{Y}^2}{n - 1}}$$



Esta expressão não é muito prática para calcular manualmente o coeficiente de correlação. Pode-se obter uma expressão mais conveniente para o cálculo manual e o cálculo de outras medidas necessárias mais tarde.



Tem-se:

$$\begin{aligned} r &= \frac{\text{Cov}(X, Y)}{S_X S_Y} = \\ &= \frac{\frac{\sum X_i Y_i - n \bar{X} \bar{Y}}{n-1}}{\sqrt{\frac{\sum X_i^2 - n \bar{X}^2}{n-1}} \sqrt{\frac{\sum Y_i^2 - n \bar{Y}^2}{n-1}}} = \\ &= \frac{\sum X_i Y_i - n \bar{X} \bar{Y}}{\sqrt{(\sum X_i^2 - n \bar{X}^2)(\sum Y_i^2 - n \bar{Y}^2)}} \end{aligned}$$



F
a

$$S_{XY} = \sum X_i Y_i - n \bar{X} \bar{Y}$$

z
e

$$S_{XX} = \sum X_i^2 - n \bar{X}^2$$

n
d
o

$$S_{YY} = \sum Y_i^2 - n \bar{Y}^2$$

Tem - se :

$$r = \frac{S_{XY}}{\sqrt{S_{XX} \cdot S_{YY}}}$$

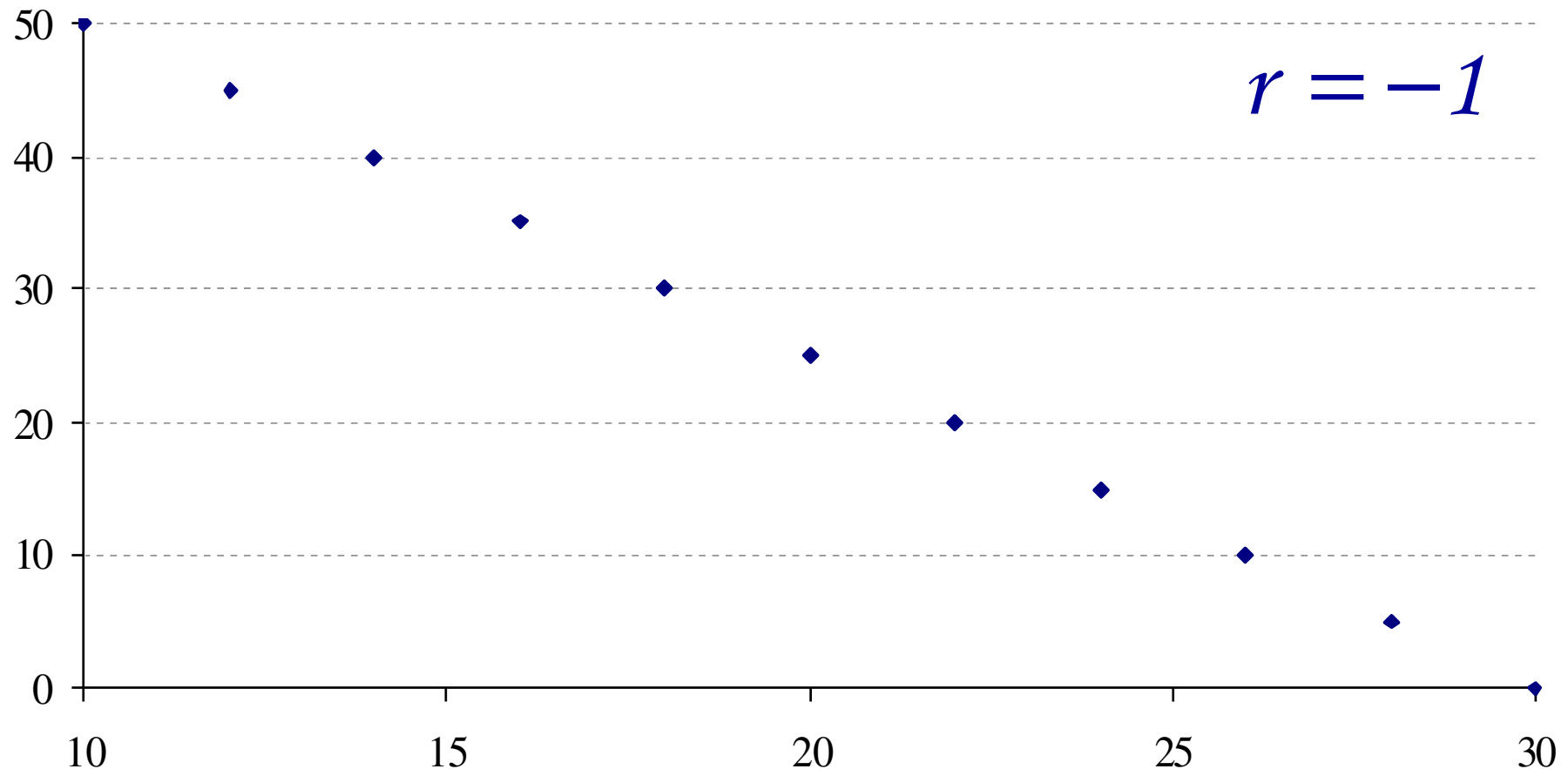
A vantagem do coeficiente de correlação (de Pearson) é ser adimensional e variar de -1 a $+1$, que o torna de fácil interpretação.



Assim se $r = -1$, temos um relacionamento linear negativo perfeito, isto é, os pontos estão todos alinhados e quando X aumenta Y decresce e vice-versa.



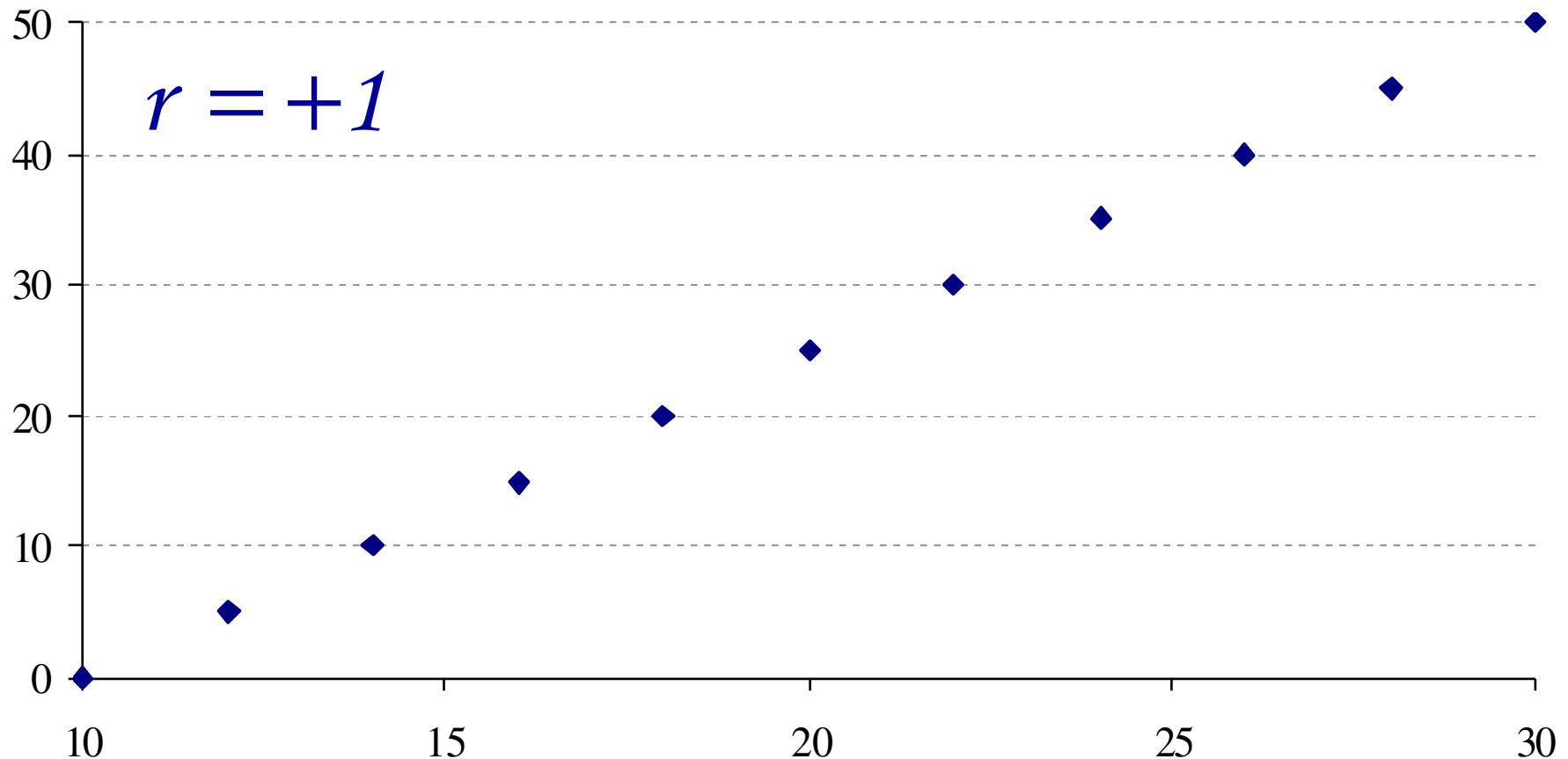
Correlação perfeita e negativa



*Se $r = +1$, temos uma
relacionamento linear positivo
perfeito, isto é, os pontos estão todos
alinhados e quando X aumenta Y
também aumenta.*



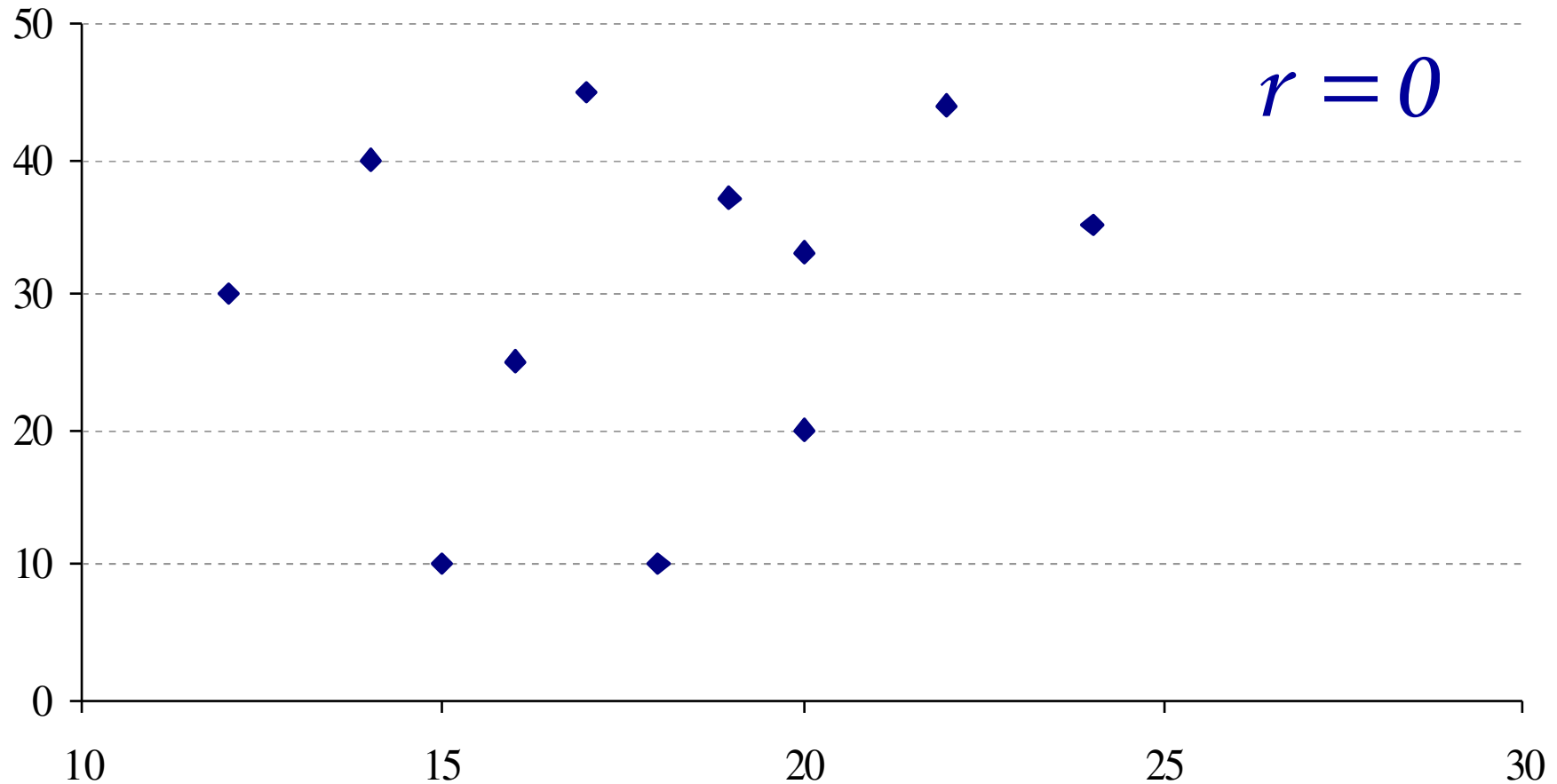
Correlação perfeita e positiva



Assim se $r = 0$, temos uma ausência de relacionamento linear, isto é, os pontos não mostram “alinhamento”.



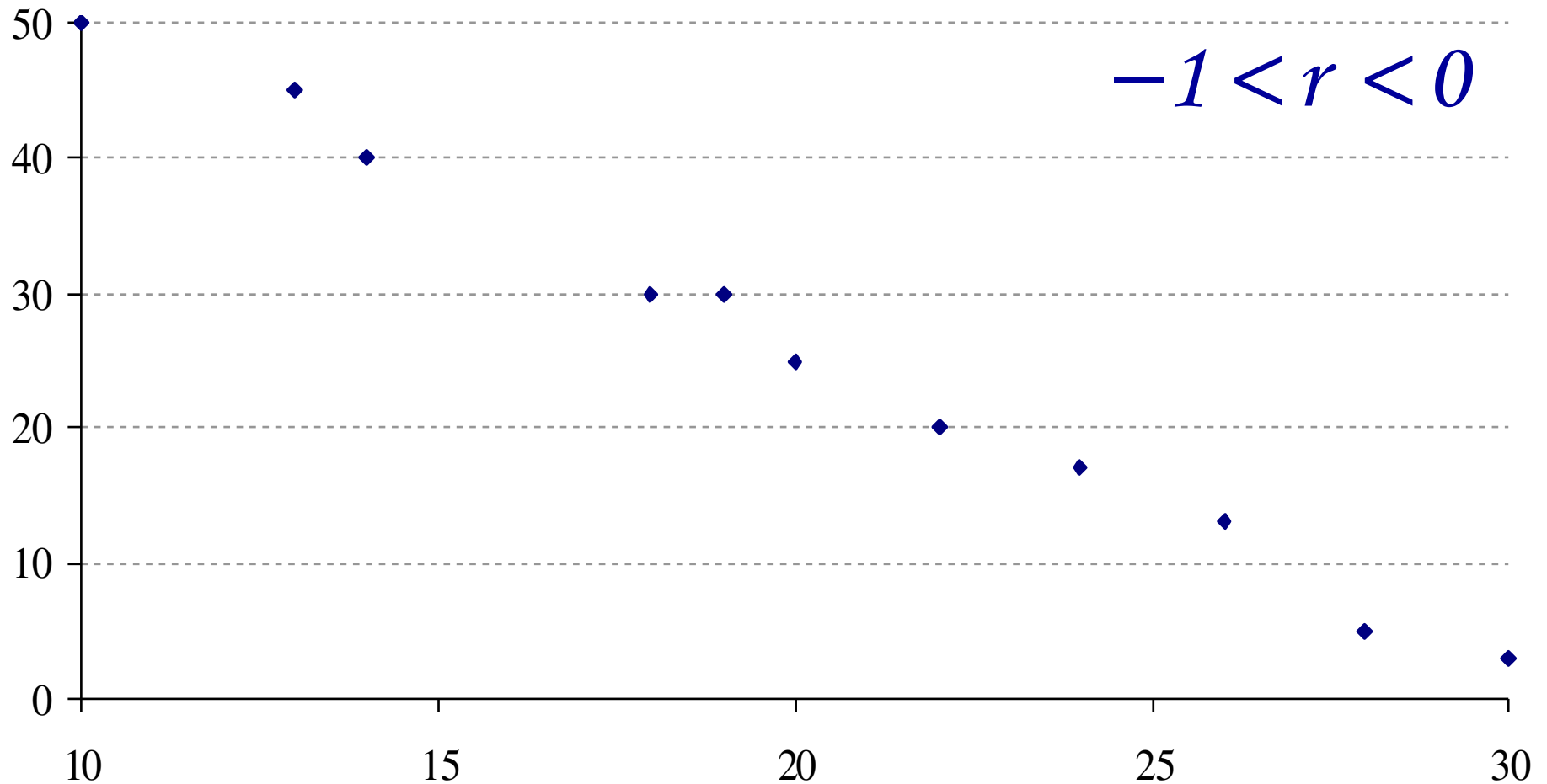
Correlação nula



Assim se $-1 < r < 0$, temos uma relacionamento linear negativo, isto é, os pontos estão mais ou menos alinhados e quando X aumenta Y decresce e vice-versa.



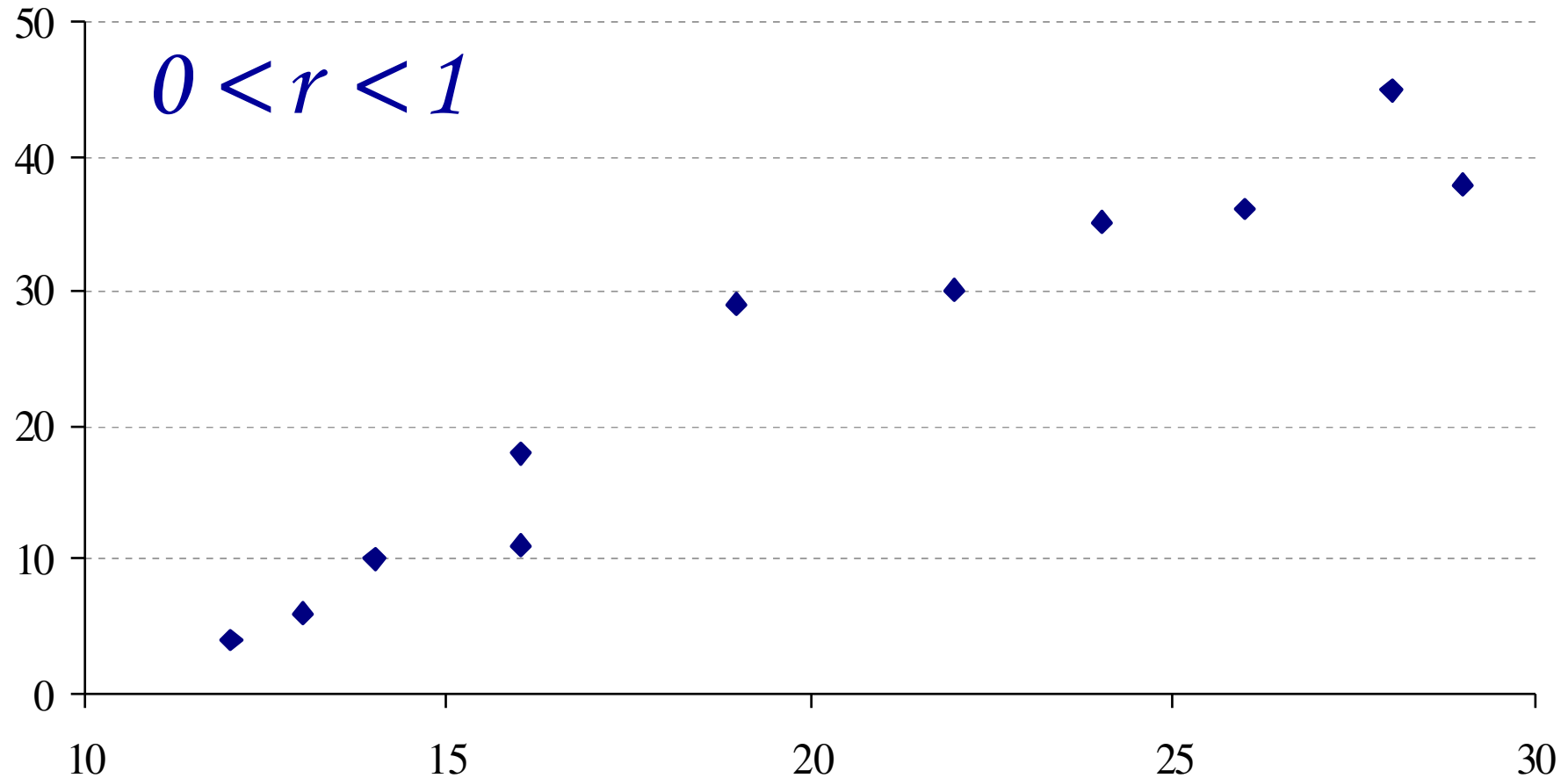
Correlação negativa



Assim se $0 < r < 1$, temos uma relacionamento linear positivo, isto é, os pontos estão mais ou menos alinhados e quando X aumenta Y também aumenta.



Correlação positiva



Observação:

Uma correlação amostral não significa necessariamente uma correlação populacional e vice-versa. É necessário testar o coeficiente de correlação para verificar se a correlação amostral é também populacional.

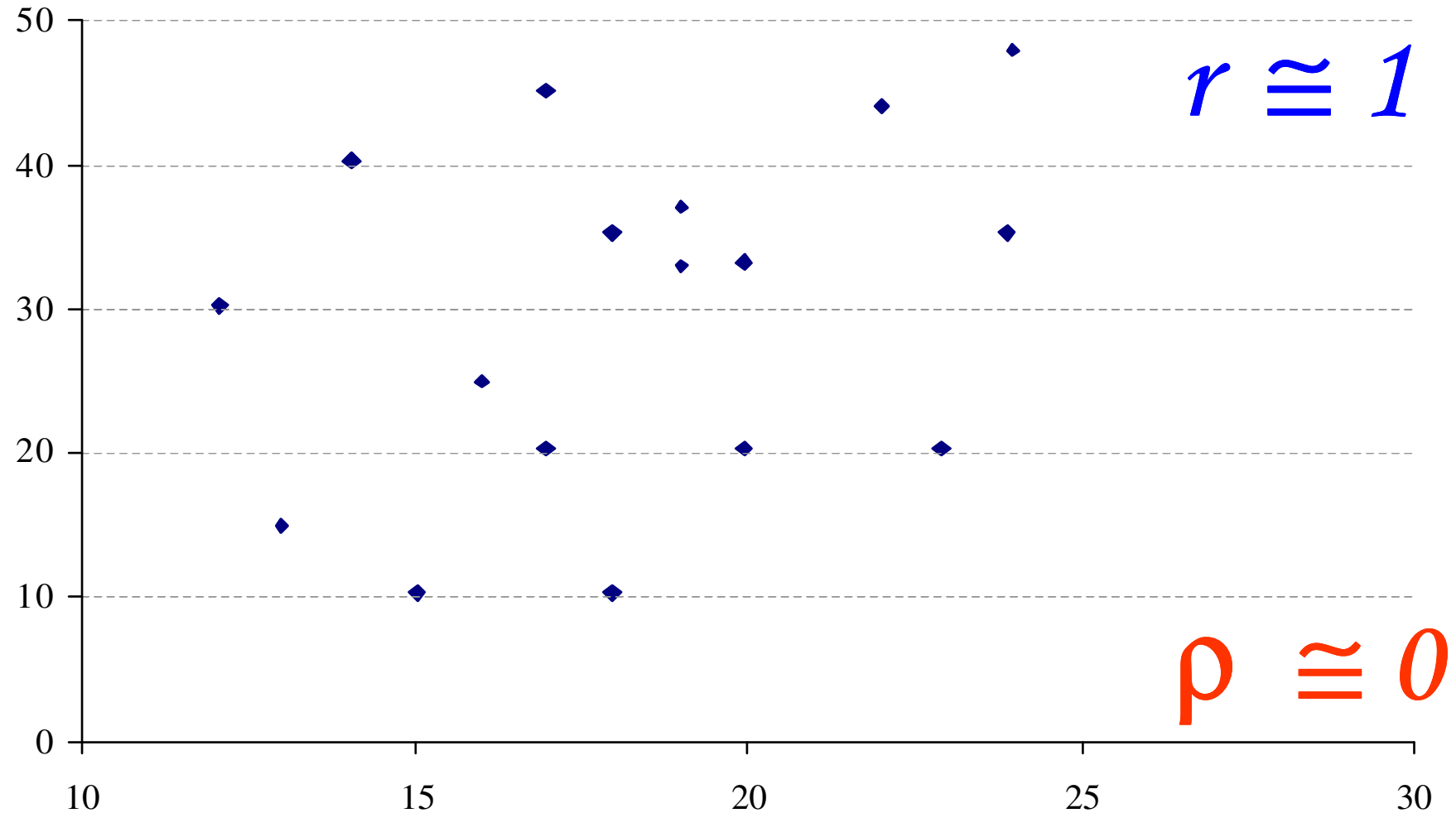


Ilustração

Observada uma amostra de seis pares, pode-se perceber que a correlação é quase um, isto é, $r \cong 1$. No entanto, observe o que ocorre quando mais pontos são acrescentados, isto é, quando se observa a população!



Correlação amostral χ populacional



Exemplo



Determinar o “grau de relacionamento linear” entre as variáveis $X =$ temperatura de operação do processo versus $Y =$ rendimento do produto, conforme tabela.



<i>X</i>	<i>Y</i>	<i>XY</i>	<i>X</i>	<i>Y</i>
100	45	4500	10000	2025
110	51	5610	12100	2601
120	54	6480	14400	2916
130	61	7930	16900	3721
140	66	9240	19600	4356
150	70	10500	22500	4900
160	74	11840	25600	5476
170	78	13260	28900	6084
180	85	15300	32400	7225
190	89	16910	36100	7921
1450	673	101570	218500	47225

Vamos calcular “ r ”
utilizando a expressão em
destaque vista anteriormente,
isto é, através das quantidades,
 S_{xy} , S_{xx} e S_{yy} .



Tem-se: $n = 10 \quad \sum X = 1450 \quad \sum Y = 673$

$$\bar{X} = 145 \quad \bar{Y} = 67,3 \quad \sum XY = 101570$$

$$\sum X^2 = 218500 \quad \sum Y^2 = 47225$$

Então: $S_{XY} = \sum X_i Y_i - n \bar{X} \bar{Y} =$
 $= 101570 - 10 \cdot 145 \cdot 67,3 =$
 $= 3985$



$$\begin{aligned} S_{XX} &= \sum X_i^2 - n \bar{X}^2 = \\ &= 218500 - 10.145^2 = \\ &= 8250 \end{aligned}$$

$$\begin{aligned} S_{YY} &= \sum Y_i^2 - n \bar{Y}^2 = \\ &= 47225 - 10.67,3^2 = \\ &= 1932,10 \end{aligned}$$



$$\begin{aligned} r &= \frac{S_{xy}}{\sqrt{S_{xx} \cdot S_{yy}}} = \\ &= \frac{3985}{\sqrt{8250 \cdot 1932,10}} = \\ &= 0,9981 \end{aligned}$$



Apesar de “ r ” ser um valor adimensional, ele não é uma taxa. Assim o resultado não deve ser expresso em percentagem.



Teste para o Coeficiente de Correlação



O valor de “r” é obtido com base em uma amostra. Ele é portanto, uma estimativa do verdadeiro valor da correlação populacional (ρ).



A teoria dos testes de hipóteses pode ser utilizada para verificar se com base na estimativa “r” é possível concluir se existe ou não correlação populacional, isto é, desejamos testar:



$$\mathcal{H}_0: \rho = 0$$

$$\mathcal{H}_1: \rho > 0$$

(teste unilateral/unicaudal à direita)

$$\rho < 0$$

(teste unilateral/unicaudal à esquerda)

$$\rho \neq 0$$

(teste bilateral/bicaudal).



O teste para a existência de correlação linear entre duas variáveis é realizado por:

$$t_{n-2} = \frac{r - \mu_r}{\hat{\sigma}_r} = \frac{r - 0}{\sqrt{\frac{1 - r^2}{n - 2}}} =$$

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



Rejeita-se a Hipótese nula se:

$$t_{n-2} > t_c$$

(teste unilaterial/unicaudal à direita)

$$t_{n-2} < t_c$$

(teste unilaterial/unicaudal à esquerda)

$$|t_{n-2}| > t_c$$

(teste bilateral/bicaudal).



Onde t_c é tal que:

$$\mathcal{P}(t < t_c) = 1 - \alpha$$

(teste unilateral/unicaudal à direita)

$$\mathcal{P}(t < t_c) = \alpha$$

(teste unilateral/unicaudal à esquerda)

$$\mathcal{P}(t < t_c) = \alpha/2 \text{ ou } \mathcal{P}(t > t_c) = \alpha/2$$

(teste bilateral/bicaudal).



Exemplo



Suponha que uma amostra de $n = 12$, alunos forneceu um coeficiente de correlação amostral de $r = 0,66$, entre $X =$ “nota em cálculo” e $Y =$ “nota em Probabilidade e Estatística”. Verifique se é possível afirmar que uma nota boa em Cálculo está relacionada com uma nota boa em Probabilidade e Estatística a 1% de significância.



Solução:

Hipóteses:

$$\mathcal{H}_0: \rho = 0$$

$$\mathcal{H}_1: \rho > 0$$

Dados:

$$n = 12$$

$$r = 0,66$$

$$\alpha = 1\%$$

Trata-se de um teste unilateral à direita para o coeficiente de correlação.



A variável teste é:

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

Então:

$$t_{10} = r \sqrt{\frac{n-2}{1-r^2}} = 0,66 \sqrt{\frac{12-2}{1-0,66^2}} = 2,778$$



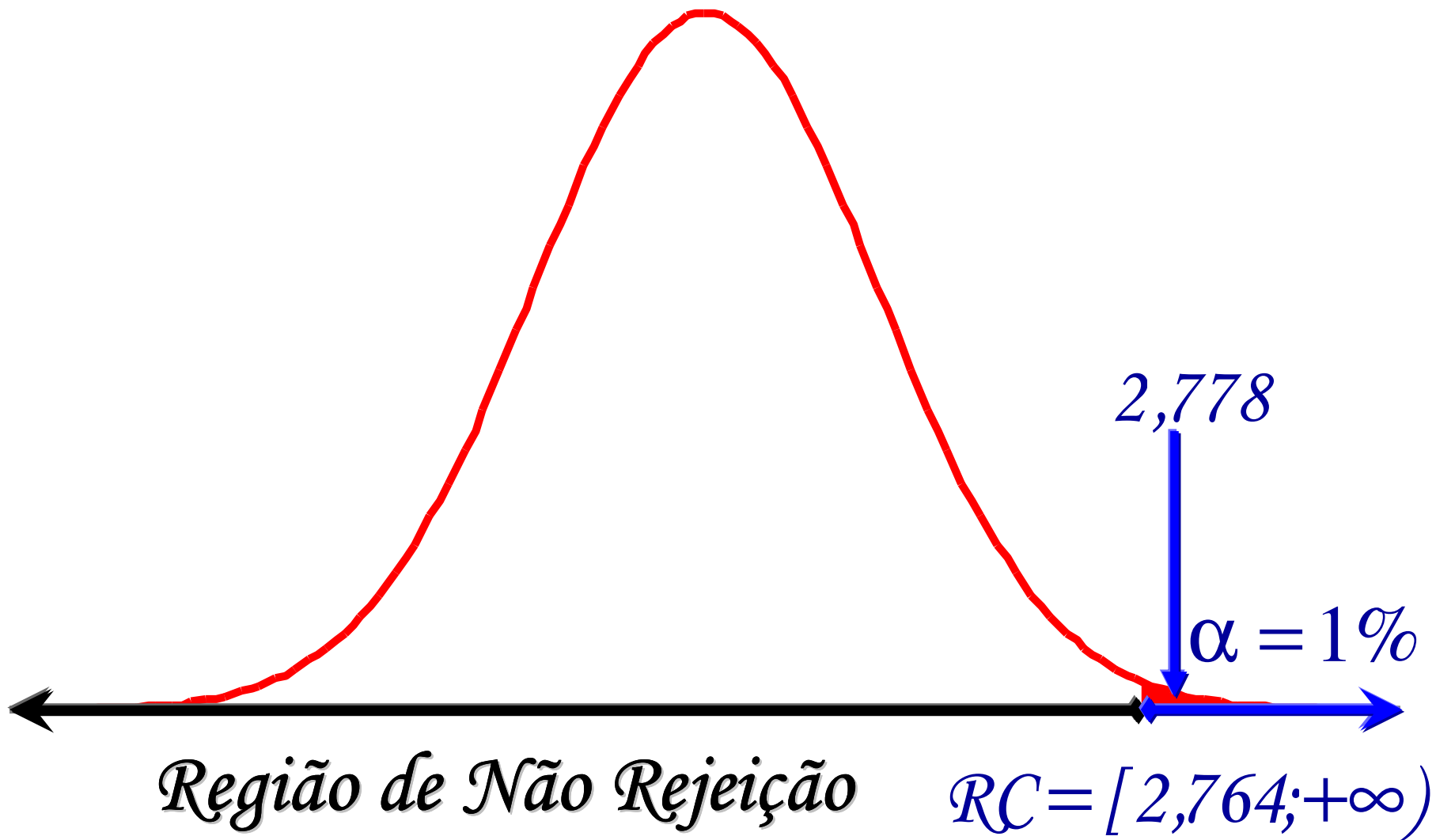
O valor crítico t_c é tal que: $P(T > t_c) = 1 - \alpha$

Então $t_c = 2,764$. Assim $\mathcal{R}C = [2,764; \infty)$

DECISÃO e CONCLUSÃO:

Como $t_{10} = 2,778 \in \mathcal{R}C$ ou $2,778 > 2,764$, Rejeito H_0 , isto é, a 1% de significância, pode-se afirmar que a nota de Cálculo está relacionada com a de Probabilidade e Estatística.





OPÇÃO:

Trabalhar com a significância do resultado obtido (2,778), isto é, o valor-p. Para isto, deve-se calcular $P(T_{10} > 2,778)$. Utilizando o Excel, tem-se:



DISTT

X	2,778	=	2,778
Graus_liberdade	10	=	10
Caudas	1	=	1

= 0,009758761

Retorna a distribuição t de Student.

Caudas especifica o número de caudas da distribuição a ser retornado:
distribuição uni-caudal = 1; distribuição bi-caudal = 2.

Resultado da fórmula = 0,009758761

OK Cancelar

Como a significância do resultado (0,98%) é menor que a significância do teste (1%) é possível rejeitar a hipótese nula.



A transformada de Fisher



O procedimento realizado para testar o coeficiente de correlação só é válido para testar a hipótese nula de que não existe correlação, isto é, $\rho = 0$. Outros tipos de testes só podem ser realizados através da transformada “zeta” de Fisher.



A transformada “ ζ ” é dada por:

$$\zeta = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right)$$

*O que equivale a considerar “ r ”
como a tangente hiperbólica de “ ζ ”*



A vantagem desta transformação é que os valores de “ ζ ” estão distribuídos aproximadamente de acordo com uma normal de média:

$$\mu_{\zeta} = \frac{1}{2} \ln \left(\frac{1 + \rho}{1 - \rho} \right)$$

E desvio:

$$\sigma_{\zeta} = \sqrt{\frac{1}{n - 3}}$$



Esta transformação permite, realizar, testes de hipóteses e construir intervalos de confiança para o coeficiente de correlação, através de ζ e da distribuição normal.



$$\mathcal{H}_0: \rho = \rho_0$$

$$\mathcal{H}_1: \rho > \rho_0$$

(teste unilateral/unicaudal à direita)

$$\rho < \rho_0$$

(teste unilateral/unicaudal à esquerda)

$$\rho \neq \rho_0$$

(teste bilateral/bicaudal).



O teste para a existência de correlação linear populacional entre duas variáveis X e Y é realizado por:

$$z = \frac{\zeta - \mu_{\zeta}}{\sigma_{\zeta}} = \frac{\zeta - \frac{1}{2} \ln \left(\frac{1 + \rho}{1 - \rho} \right)}{\sqrt{\frac{1}{n - 3}}}$$



Rejeita-se a Hipótese nula se:

$$z > z_c$$

(teste unilateral/unicaudal à direita)

$$z < z_c$$

(teste unilateral/unicaudal à esquerda)

$$|z| > z_c$$

(teste bilateral/bicaudal).



Onde z_c é tal que:

$$\Phi(z_c) = 1 - \alpha$$

(teste unilateral/unicaudal à direita)

$$\Phi(z_c) = \alpha$$

(teste unilateral/unicaudal à esquerda)

$$\Phi(z_c) = \alpha/2 \text{ ou } \Phi(z_c) = 1 - \alpha/2$$

(teste bilateral/bicaudal).



Exemplo



Suponha que uma amostra de $n = 35$, alunos forneceu um coeficiente de correlação amostral de $r = 0,75$, entre $X =$ “número de horas de estudo” e $Y =$ “nota em Probabilidade e Estatística”. Verifique se é possível afirmar que o “o número de horas de estudo” apresenta uma correlação de pelo menos $0,5$ na população com a “nota em Probabilidade e Estatística”, a 1% de significância.



Solução:

Hipóteses:

$$H_0: \rho = 0,5$$

$$H_1: \rho > 0,5$$

Dados:

$$n = 35$$

$$r = 0,75$$

$$\alpha = 1\%$$

Trata-se de um teste unilateral à direita para o coeficiente de correlação.



A variável teste é:

$$z = \frac{\zeta - \mu_\zeta}{\sigma_\zeta} = \frac{\zeta - \frac{1}{2} \ln \left(\frac{1 + \rho}{1 - \rho} \right)}{\sqrt{\frac{1}{n - 3}}}$$

Então:

$$\zeta = \frac{1}{2} \ln \left(\frac{1 + 0,75}{1 - 0,75} \right) = 0,9730$$



A média vale:

$$\mu_{\zeta} = \frac{1}{2} \ln \left(\frac{1+\rho}{1-\rho} \right) = \frac{1}{2} \ln \left(\frac{1+0,5}{1-0,5} \right) = 0,5493$$

E o desvio padrão vale:

$$\sigma_{\zeta} = \sqrt{\frac{1}{n-3}} = \sqrt{\frac{1}{35-3}} = \sqrt{\frac{1}{32}} = 0,1768$$



Padronizando, tem-se:

$$z = \frac{\zeta - \mu_\zeta}{\sigma_\zeta} = \frac{\zeta - \frac{1}{2} \ln \left(\frac{1 + \rho}{1 - \rho} \right)}{\sqrt{\frac{1}{n - 3}}} = \frac{0,9730 - 0,5493}{0,1768} = 2,40$$



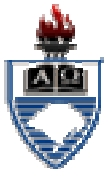
O valor crítico z_c é tal que:

$$P(Z > z_c) = \alpha = 1\%.$$

Ou $\Phi(z_c) = 99\%$.

Então $z_c = 2,33$.

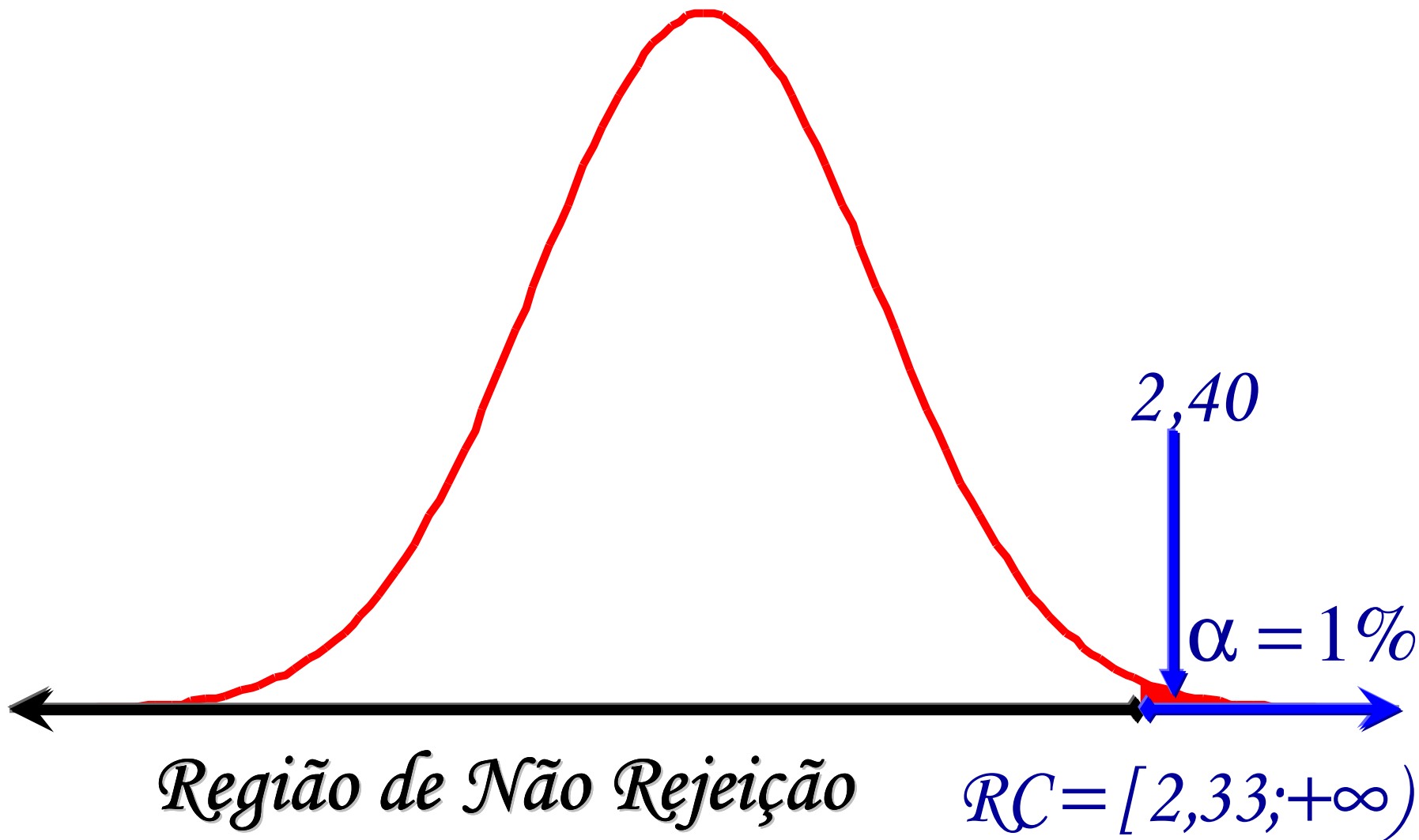
Assim $RC = [2,33; \infty)$



DECISÃO e CONCLUSÃO:

Como $z = 2,40 \in \mathcal{R}_C$ ou $2,40 > 2,33$, Rejeito H_0 , isto é, a 1% de significância, pode-se afirmar que “o número de horas de estudo” apresenta pelo menos 0,50 de correlação com a “nota em Probabilidade e Estatística”.





OPÇÃO:

Trabalhar com a significância do resultado obtido (2,40), isto é, o valor- p . Para isto, deve-se calcular $P(Z > 2,40)$, isto é, $\Phi(-2,40) = 0,82\%$. Como $p = 0,82\% < \alpha = 1\%$. Rejeito H_0 .

