

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*

O Estoque de Moeda ($M1$) está relacionado com a variação dos preços. Verifique se existe correlação entre o IPC americano com a oferta monetária, considerando dados do período de 1960 a 2003.



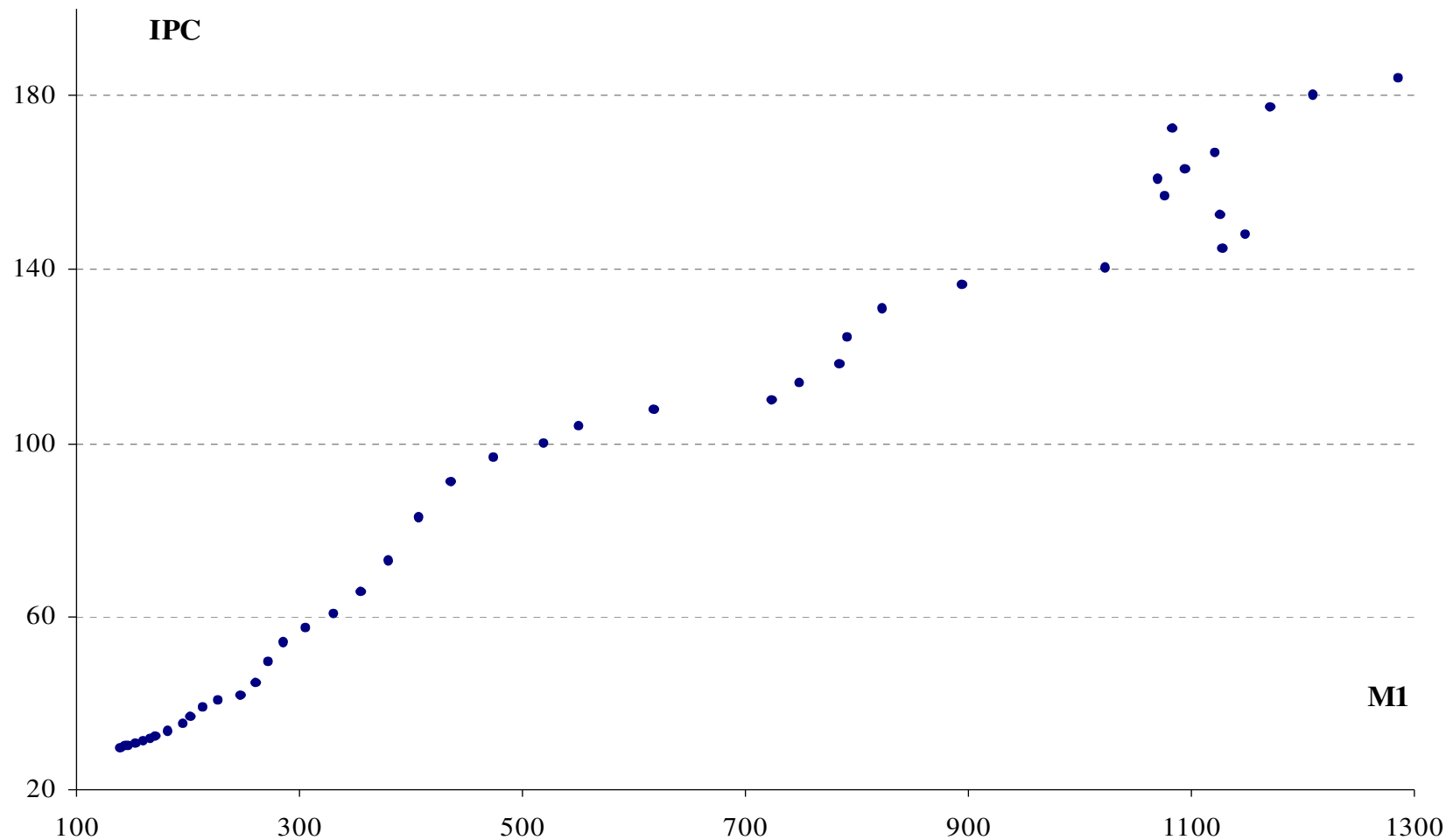
<i>Ano</i>	<i>$Y = M1$</i>	<i>$X = IPC$</i>
1960	140,7	29,6
1961	145,2	29,9
1962	147,8	30,2
1963	153,3	30,6
1964	160,3	31,5
1965	167,8	32,4
...
2000	1172,9	177,1
2002	1210,4	179,9
2003	1287,1	184,0



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 ideia 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 é mais conveniente utilizar o coeficiente de correlação linear de Pearson (momento produto).



*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 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

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

a

z

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

e

n

d

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

o

Tem - se :

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



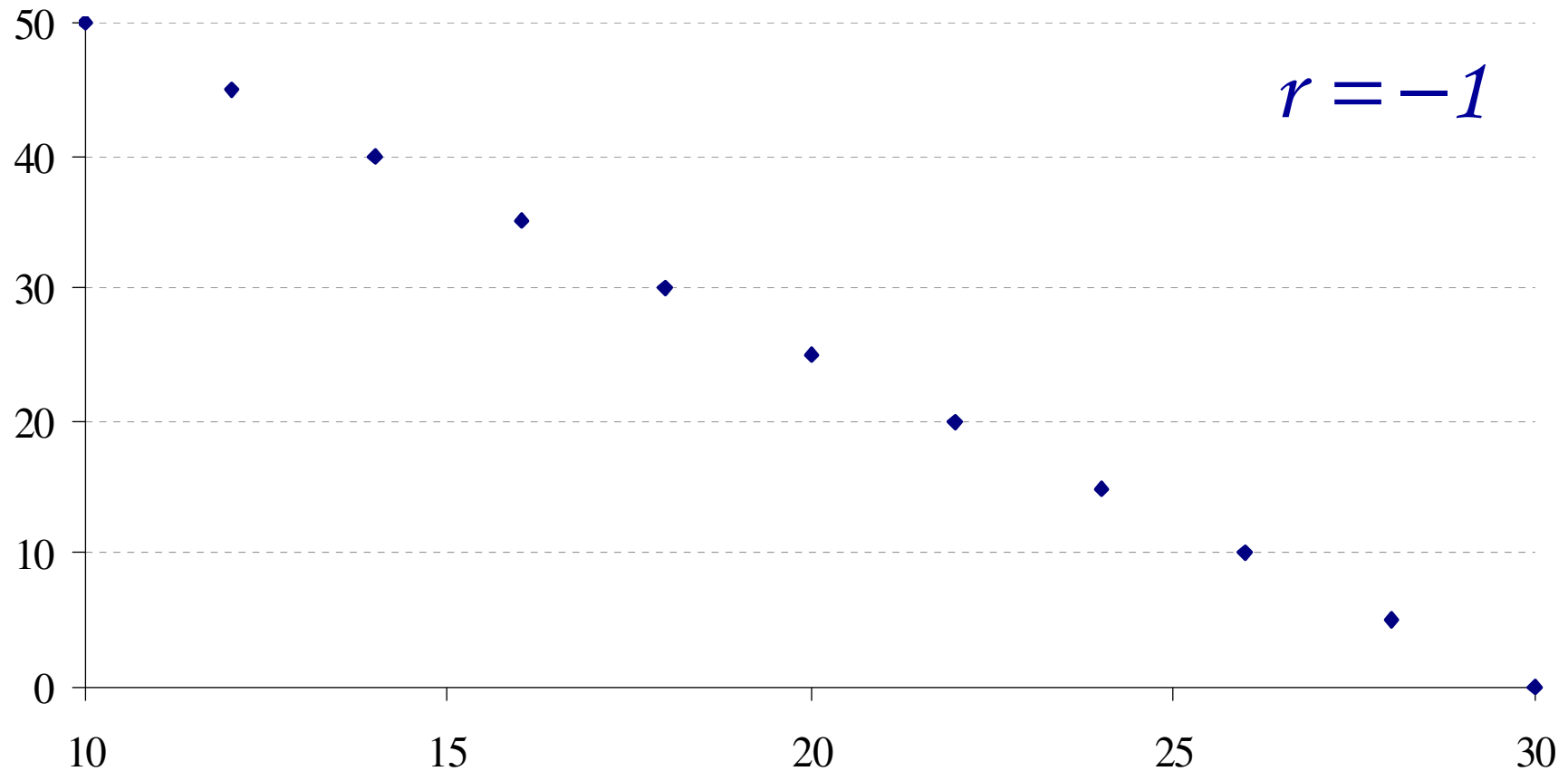
A vantagem do coeficiente de correlação (de Pearson) é ser adimensional e variar de -1 a $+1$, que o torna de fácil interpretação.



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



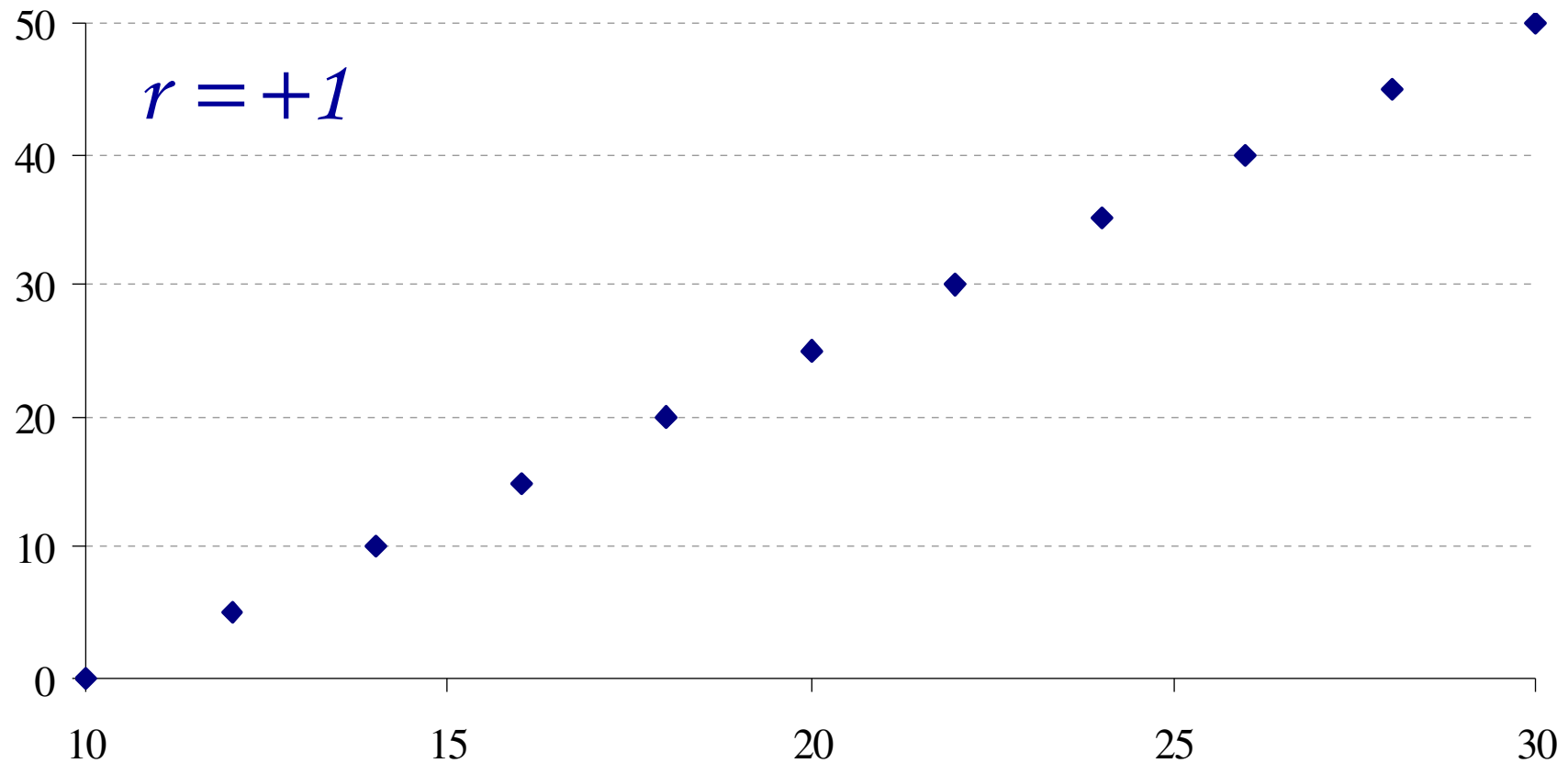
Correlação perfeita e negativa



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



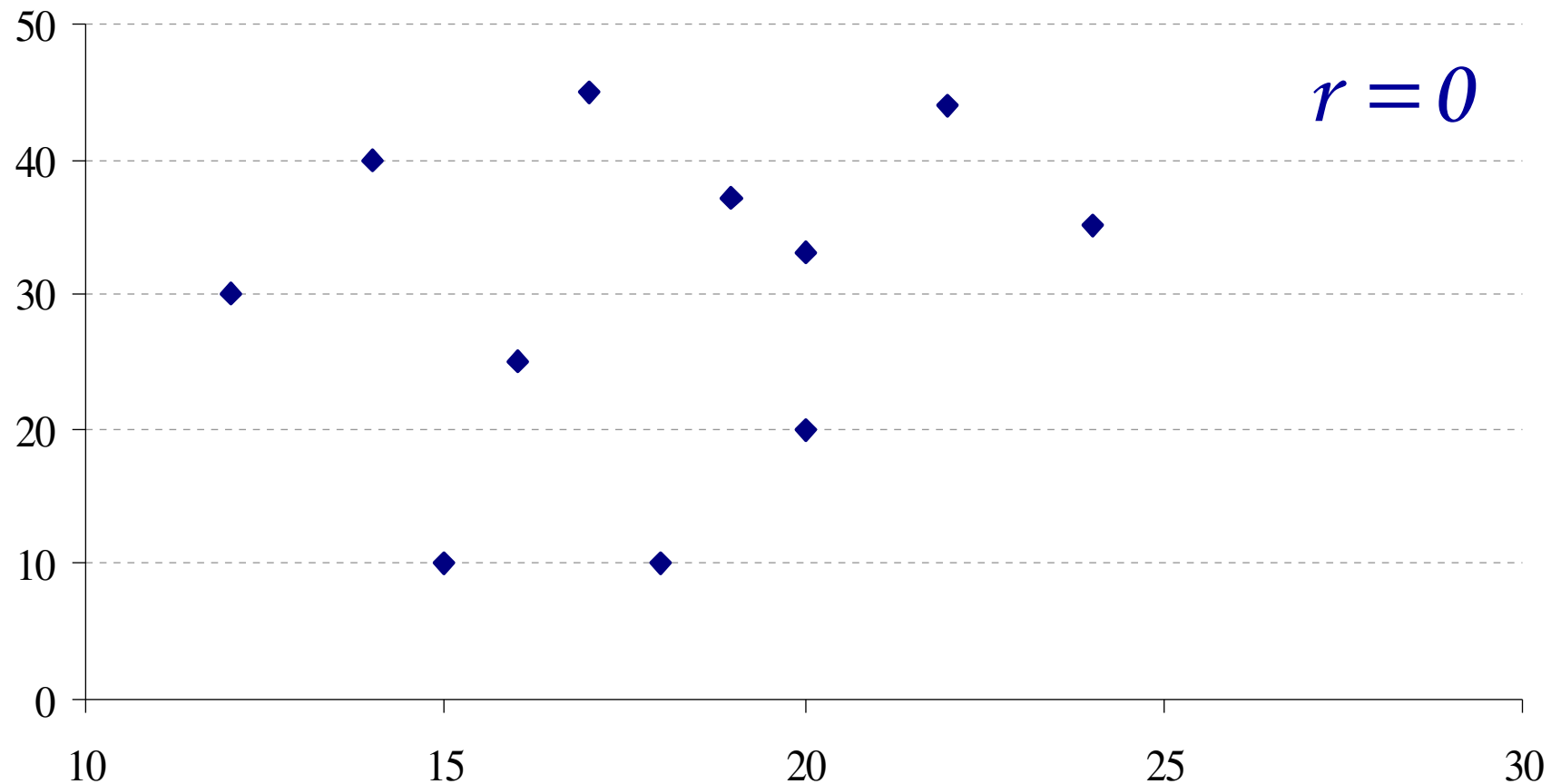
Correlação perfeita e positiva



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



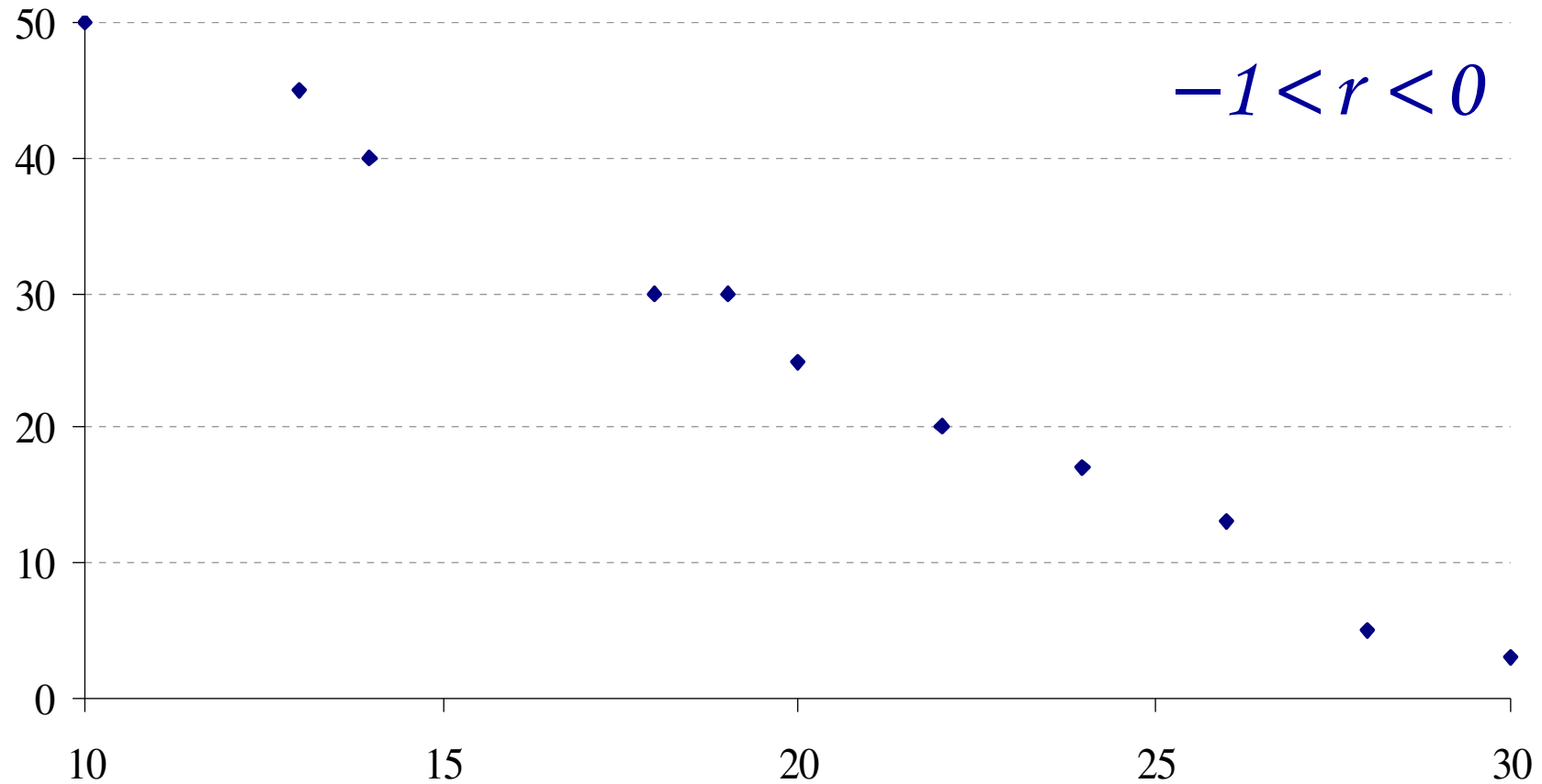
Correlação nula



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.



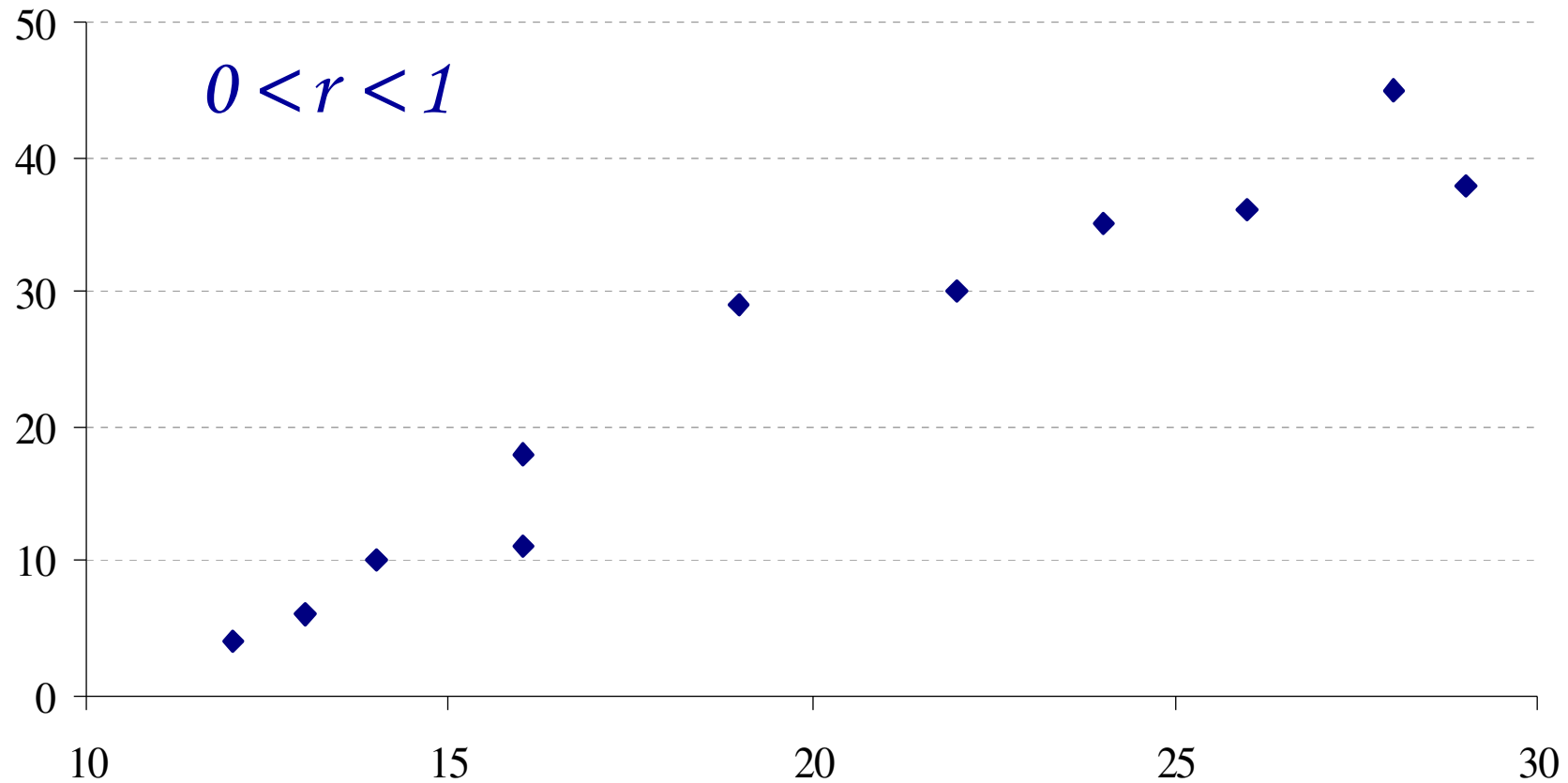
Correlação negativa



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.

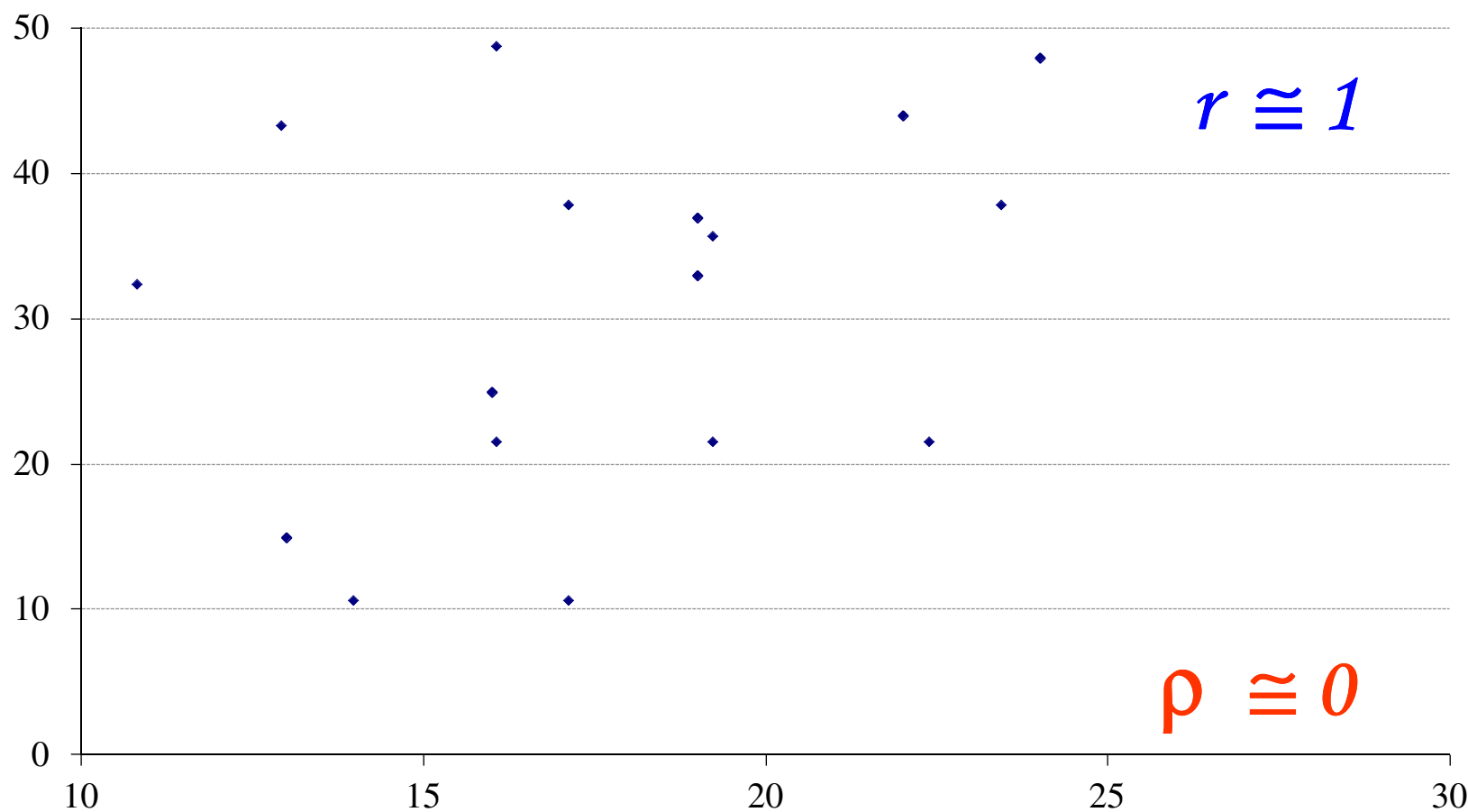


Ilustração

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 X populacional



Exemplo



Determinar o “grau de relacionamento linear” entre as variáveis $X = \text{Índice de Preços ao Consumidor}$ versus $Y = \text{Estoque de Moeda}$, para os valores da Economia Americana de 1960 a 2003.



<i>Ano</i>	<i>X</i>	<i>Y</i>	<i>XY</i>	<i>X²</i>	<i>Y²</i>
1960	140,7	29,6			
1961	145,2	29,9			
1962	147,8	30,2			
1963	153,3	30,6			
1964	160,3	31,5			
1965	167,8	32,4			
...			
2000	1172,9	177,1			
2002	1210,4	179,9			
2003	1287,1	184,0			
<i>Total</i>	<i>25894,5</i>	<i>4102,9</i>	<i>3295760,69</i>	<i>21856837,21</i>	<i>503187,97</i>

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



Tem-se:

$$n = 44 \quad \sum X = 25894,50 \quad \sum Y = 4102,90$$

$$\bar{X} = 588,5114 \quad \bar{Y} = 93,2477 \quad \sum XY = 13295760,69$$

$$\sum X^2 = 21856837,21 \quad \sum Y^2 = 503187,97$$

Então:

$$\begin{aligned} S_{XY} &= \sum X_i Y_i - n \bar{X} \bar{Y} = \\ &= 881157,4161 \end{aligned}$$



$$\begin{aligned} S_{XX} &= \sum X_i^2 - n \bar{X}^2 = \\ &= 6617629,7043 \end{aligned}$$

$$\begin{aligned} S_{YY} &= \sum Y_i^2 - n \bar{Y}^2 = \\ &= 120601,8698 \end{aligned}$$



$$\begin{aligned} r &= \frac{S_{XY}}{\sqrt{S_{XX} \cdot S_{YY}}} = \\ &= \frac{881157,4161}{\sqrt{6617629,7043 \cdot 120601,8698}} = \\ &= 0,9863 \end{aligned}$$



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



*Regressão
Linear Simples*

Em muitas situações duas ou mais variáveis estão relacionadas e surge então a necessidade de determinar a natureza deste relacionamento.



A análise de regressão é uma técnica estatística para modelar e investigar o relacionamento entre duas ou mais variáveis.



De fato a regressão pode ser dividida em dois problemas:

(i) o da especificação e

(ii) o da determinação.



A especificação

O problema da especificação é descobrir dentre os possíveis modelos (linear, quadrático, exponencial, etc.) qual o mais adequado.



A determinação

O problema da determinação é uma vez definido o modelo (linear, quadrático, exponencial, etc.) estimar os parâmetros da equação.



O modelo

Normalmente é suposto que exista uma variável Y (dependente ou resposta), que está relacionada a “ k ” variáveis (independentes ou regressoras) X_i ($i = 1, 2, \dots, k$).



A variável resposta Y é aleatória, enquanto que as variáveis regressoras X_i são normalmente controladas. O relacionamento entre elas é caracterizado por uma equação denominada de “equação de regressão”.



Quando existir apenas uma variável regressora (X) tem-se a regressão simples, se Y depender de duas ou mais variáveis regressoras, então tem-se a “regressão múltipla”.

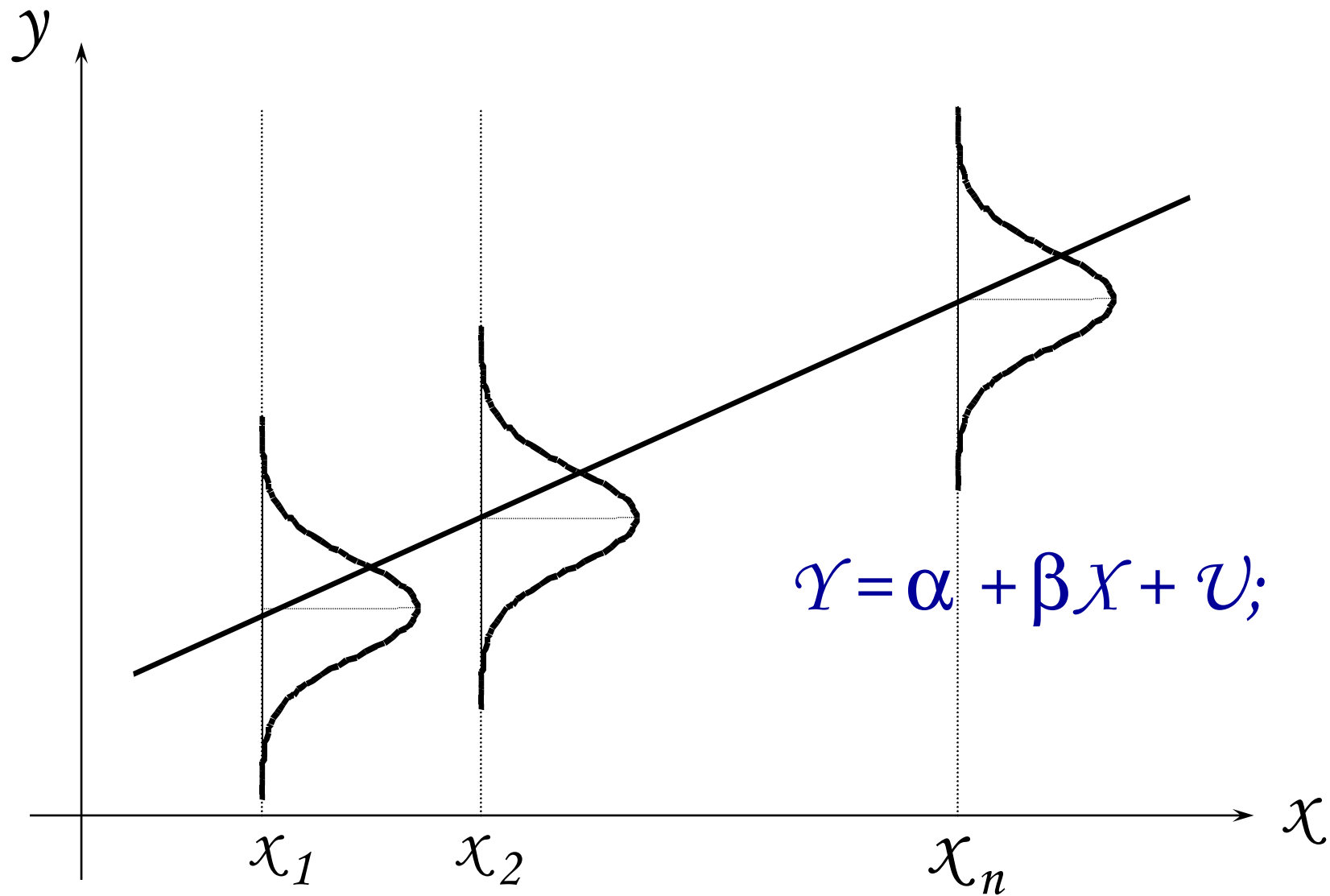


O modelo considerado

Vamos supor que a regressão é do tipo simples e que o modelo seja linear, isto é, vamos supor que a equação de regressão seja do tipo: $Y = \alpha + \beta X + \mathcal{U}$.



O modelo linear simples



O termo “ \mathcal{U} ” é o termo erro, isto é, “ \mathcal{U} ” representa outras influências sobre a variável Y , além da exercida pela variável “ X ”. A variação residual (termo \mathcal{U}) é suposto de média zero e desvio constante e igual a σ .



Ou ainda pode-se admitir que o modelo fornece o valor médio de Y , para um dado “ x ”, isto é:

$$E(Y/x) = \alpha + \beta x$$



Em resumo, as hipóteses são:

$$Y = \alpha + \beta X + U;$$

$$E(Y/x) = \alpha + \beta X, \text{ isto é, } E(U) = 0$$

$$V(Y/x) = \sigma^2;$$

$$\text{Cov}(U_i, U_j) = 0, \text{ para } i \neq j;$$

A variável X permanece fixa em observações sucessivas e os erros U são normalmente distribuídos.



A equação de regressão

O modelo suposto $E(Y/x) = \alpha + \beta X$ é populacional.

Vamos supor que se tenha n pares de observações, digamos: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ e que através deles queremos estimar o modelo acima.



A reta estimada será representada por:

$$\hat{Y} = a + bX \quad \text{ou} \quad Y = a + bX + E$$

Onde “a” é um estimador de α e “b” é um estimador de β , sendo \hat{Y} um estimador de $E(Y/x)$.



O método utilizado

Existem diversos métodos para a determinação da reta desejada. Um deles, denominado de MMQ (Métodos dos Mínimos Quadrados), consiste em minimizar a “soma dos quadrados das distâncias da reta aos pontos”.



Tem-se:

$$Y_i = a + bx_i + E_i$$

Então:

$$E_i = Y_i - (a + bx_i)$$

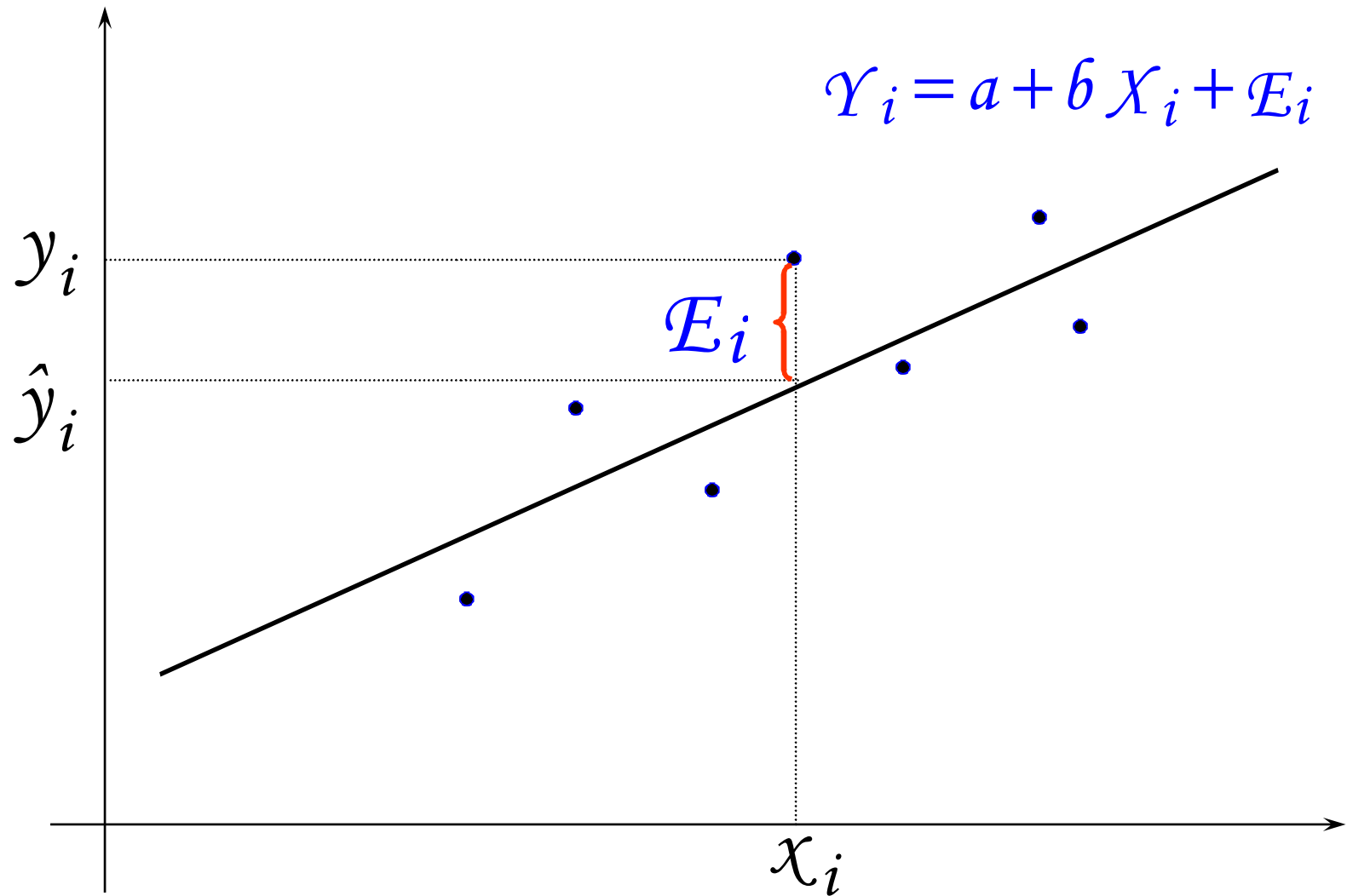


Deve-se minimizar:

$$\begin{aligned}\Phi &= \sum_{i=1}^n \mathcal{E}_i^2 = \sum_{i=1}^n (\mathcal{Y}_i - \hat{\mathcal{Y}}_i)^2 = \\ &= \sum_{i=1}^n (\mathcal{Y}_i - a - bX_i)^2\end{aligned}$$



O método dos mínimos quadrados



Derivando parcialmente tem-se:

$$\frac{\partial \phi}{\partial a} = -2 \sum_{i=1}^n (\gamma_i - a - b x_i)$$

$$\frac{\partial \phi}{\partial b} = -2 \sum_{i=1}^n x_i (\gamma_i - a - b x_i)$$



Igualando as derivadas parciais a

zero vem:

$$\sum_{i=1}^n (\mathcal{Y}_i - a - b X_i) = 0$$

$$\sum_{i=1}^n x_i (\mathcal{Y}_i - a - b X_i) = 0$$



Isolando as incógnitas, tem-se:

$$\sum \mathcal{Y}_i = na + b \sum \mathcal{X}_i$$

$$\sum \mathcal{X}_i \mathcal{Y}_i = n \sum \mathcal{X}_i + b \sum \mathcal{X}_i^2$$



Resolvendo para “a” e “b”, segue:

$$b = \frac{\sum X_i y_i - n \bar{X} \bar{Y}}{\sum X_i^2 - n \bar{X}^2} = \frac{S_{XY}}{S_{XX}}$$
$$a = \bar{Y} - b \bar{X}$$



Lembrando que:

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

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

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



Exemplo



Considerando os valores das variáveis “Oferta Monetária” e “Índice de Preços ao Consumidor”, consideradas anteriormente, determinar uma equação de regressão linear para prever o IPC dado um determinado nível de Oferta Monetária.



<i>Ano</i>	<i>$Y = IPC$</i>	<i>$X = M1$</i>
1960	29,6	140,7
1961	29,9	145,2
1962	30,2	147,8
1963	30,6	153,3
1964	31,5	160,3
1965	32,4	167,8
...
2000	177,1	1172,9
2002	179,9	1210,4
2003	184,0	1287,1



Da mesma forma que para calcular o coeficiente de correlação é necessário a construção de três novas colunas. Uma para X^2 , uma para Y^2 e outra para XY .



<i>Ano</i>	<i>X</i>	<i>Y</i>	<i>XY</i>	<i>X²</i>	<i>Y²</i>
1960	140,7	29,6			
1961	145,2	29,9			
1962	147,8	30,2			
1963	153,3	30,6			
1964	160,3	31,5			
1965	167,8	32,4			
...			
2000	1172,9	177,1			
2002	1210,4	179,9			
2003	1287,1	184,0			
<i>Total</i>	<i>25894,5</i>	<i>4102,9</i>	<i>3295760,69</i>	<i>21856837,21</i>	<i>503187,97</i>



Tem-se:

$$n = 44 \quad \sum X = 25894,50 \quad \sum Y = 4102,90$$

$$\bar{X} = 588,5114 \quad \bar{Y} = 93,2477 \quad \sum XY = 13295760,69$$

$$\sum X^2 = 21856837,21 \quad \sum Y^2 = 503187,97$$

Então:

$$\begin{aligned} S_{XY} &= \sum X_i Y_i - n \bar{X} \bar{Y} = \\ &= 881157,4161 \end{aligned}$$



$$\begin{aligned} S_{XX} &= \sum X_i^2 - n \bar{X}^2 = \\ &= 6617629,7043 \end{aligned}$$

$$\begin{aligned} S_{YY} &= \sum Y_i^2 - n \bar{Y}^2 = \\ &= 120601,8698 \end{aligned}$$



A equação de regressão, será, então:

$$b = \frac{S_{XY}}{S_{XX}} = \frac{881157,4161}{6617629,7043} = 0,1332 \cong 0,13$$

$$a = \bar{Y} - b\bar{X} = 93,2477 - 0,1332 \cdot 588,5114 = \\ = 14,8857 \cong 14,89$$

$$\hat{Y} = 14,89 + 0,13x$$



A pergunta que cabe agora é: este modelo representa bem os pontos dados? A resposta é dada através do erro padrão da regressão.



Variância Residual

e

Erro Padrão da Regressão



O objetivo do MMQ é minimizar a variação residual em torno da reta de regressão. Uma avaliação desta variação é dada por:

$$S^2 = \frac{\sum E^2}{n - 2} = \frac{\sum (\gamma - a - bX)^2}{n - 2}$$



O cálculo da variância residual, por esta expressão, é muito trabalhoso, pois é necessário primeiro determinar os valores previstos.

Entretanto é possível obter uma expressão que não requeira o cálculo dos valores previstos, isto

é, de $\hat{Y} = a + bX$.



Desenvolvendo o numerador da expressão,

vem:

$$\begin{aligned}\sum(\mathcal{Y}-a-bX)^2 &= \sum[\mathcal{Y}-(\bar{\mathcal{Y}}-b\bar{X})-bX]^2 = \\ &= \sum[\mathcal{Y}-\bar{\mathcal{Y}}+b\bar{X}-bX]^2 = \sum[\mathcal{Y}-\bar{\mathcal{Y}}-b(X-\bar{X})]^2 = \\ &= \sum(\mathcal{Y}-\bar{\mathcal{Y}})^2 - 2b\sum(X-\bar{X})(\mathcal{Y}-\bar{\mathcal{Y}}) + b^2\sum(X-\bar{X})^2 = \\ &= S_{\mathcal{Y}\mathcal{Y}} - 2bS_{X\mathcal{Y}} + b^2S_{XX}\end{aligned}$$



Uma vez que:

$$\begin{aligned}\sum (X - \bar{X})(Y - \bar{Y}) &= \\ &= \sum X_i Y_i - n\bar{X}\bar{Y} = S_{XY}\end{aligned}$$

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

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



Deste modo, tem-se:

$$\sum(Y - a - bX)^2 = S_{YY} - 2b S_{XY} + b^2 S_{XX}$$

Mas:
$$b = \frac{S_{XY}}{S_{XX}} \Rightarrow S_{XY} = b S_{XX}$$

Então:

$$\begin{aligned} \sum(Y - a - bX)^2 &= S_{YY} - 2b S_{XY} + b^2 S_{XX} = \\ &= S_{YY} - 2b^2 S_{XX} + b^2 S_{XX} = S_{YY} - b^2 S_{XX} \end{aligned}$$



Assim:

$$s = \sqrt{\frac{\sum E^2}{n-2}} = \sqrt{\frac{\sum (\gamma - a - bX)^2}{n-2}}$$

Será, finalmente:

$$s = \sqrt{\frac{S_{\gamma\gamma} - b^2 S_{XX}}{n-2}} = \sqrt{\frac{S_{\gamma\gamma} - b S_{X\gamma}}{n-2}}$$



Exemplo



Considerando os valores do exemplo anterior, determinar o erro padrão da regressão.

$$\text{Tem-se: } S_{XY} = 881157,4161$$

$$S_{XX} = 6617629,7043$$

$$b = \frac{S_{XY}}{S_{XX}} = \frac{881157,4161}{6617629,7043} = 0,1332$$



Então:

$$\begin{aligned} s &= \sqrt{\frac{S_{yy} - b S_{xy}}{n - 2}} = \\ &= \sqrt{\frac{120601,8698 - 0,1332 \cdot 881157,4161}{44 - 2}} = \\ &= 8,8278 \cong 8,83 \end{aligned}$$



A pergunta, agora, é: este erro é razoável?, quer dizer, ele não é muito grande?

A resposta envolve o cálculo do erro relativo, isto é, devemos comparar este resultado com a variável de interesse.



A variável envolvida aqui é a Υ , isto é, a base monetária, então, o erro relativo, será:

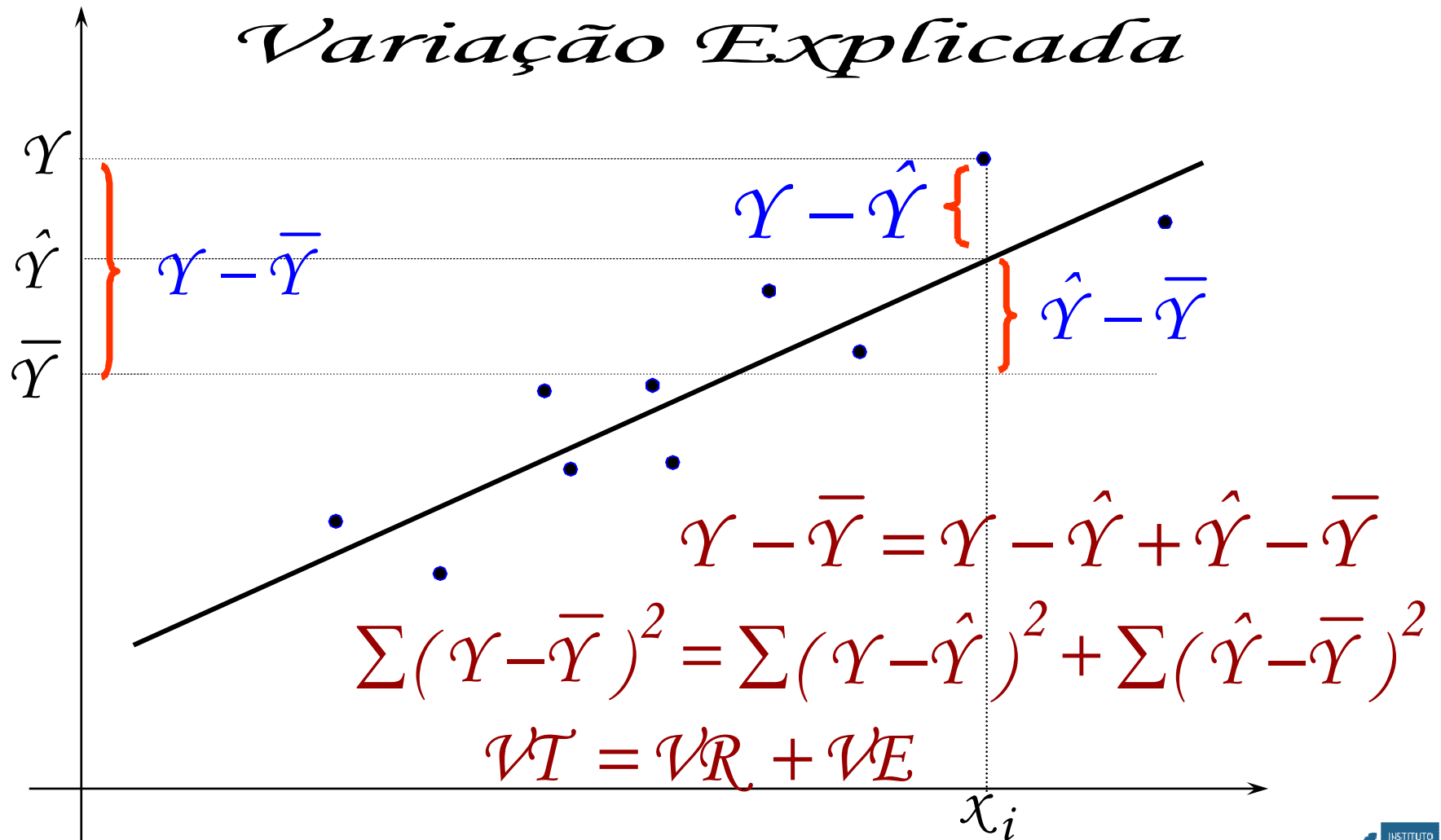
$$g_s = \frac{s}{\Upsilon} = \frac{8,8278}{93,2477} = 9,47\%$$



Decomposição da Variação



*Varição Total =
Varição Não-Explicada
+
Varição Explicada*



(a) Variação Total: $\mathcal{V}T$

$$\mathcal{V}T = \sum (\gamma - \bar{\gamma})^2 = S_{\gamma\gamma}$$

(b) Variação Residual: $\mathcal{V}R$

$$\mathcal{V}R = \sum (\gamma - \hat{\gamma})^2 = S_{\gamma\gamma} - b^2 S_{XX} = \mathcal{V}T - \mathcal{V}E$$

(c) Variação Explicada: $\mathcal{V}E$

$$\mathcal{V}E = \sum (\hat{\gamma} - \bar{\gamma})^2 = b^2 S_{XX}$$



Uma maneira de medir o grau de aderência (adequação) de um modelo é verificar o quanto da variação total de Y é explicada pela reta de regressão.



Para isto, toma-se o quociente entre a variação explicada, VE e a variação total, VT :

$$R^2 = VE / VT$$

Este resultado é denominado de “Coeficiente de Determinação”.



$$R^2 = \frac{VE}{VT} = \frac{b^2 S_{XX}}{S_{YY}} = \frac{b S_{XY}}{S_{YY}} = \frac{S_{XY}^2}{S_{XX} S_{YY}}$$

Este resultado mede o quanto as variações de uma das variáveis são explicadas pelas variações da outra variável.



Ou ainda, ele mede a parcela da variação total que é explicada pela reta de regressão, isto é:

$$VE = b^2 S_{XX} = R^2 S_{YY}$$

A variação residual corresponde a:

$$VR = (1 - R^2) S_{YY}$$

Assim $1 - R^2$ é o Coeficiente de

Indeterminação.

