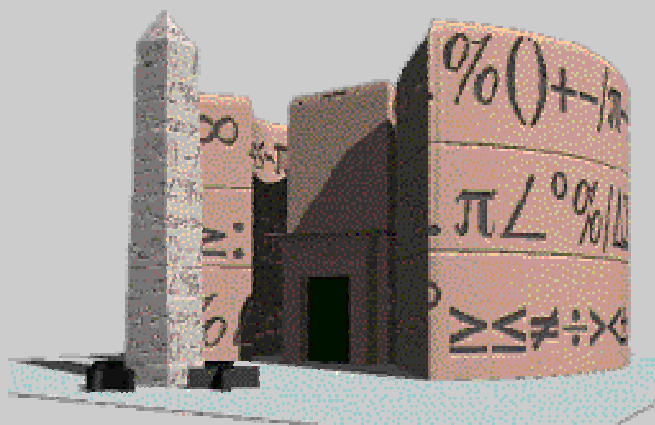


Material Didático

Série

Estatística Básica



Texto VIII

Correlação e Regressão

Prof. Lorí Viali, Dr.



SUMÁRIO

1. CORRELAÇÃO.....	3
1.1. INTRODUÇÃO.....	3
1.2. PADRÕES DE ASSOCIAÇÃO.....	4
1.3. INDICADORES DE ASSOCIAÇÃO.....	4
1.4. O COEFICIENTE DE CORRELAÇÃO.....	7
1.5. HIPÓTESES BÁSICAS.....	8
1.6. DEFINIÇÃO.....	8
1.7. PROPRIEDADES DE R.....	9
2. REGRESSÃO.....	10
2.1. ESTIMATIVA DOS PARÂMETROS DE REGRESSÃO.....	12
2.2. ESTIMATIVA DA VARIÂNCIA DO TERMO ERRO.....	14
2.3. DECOMPOSIÇÃO DA SOMA DOS QUADRADOS.....	17
2.3.1. <i>Decomposição dos desvios.....</i>	<i>17</i>
2.3.2. <i>Cálculo das variações.....</i>	<i>18</i>
2.4. COEFICIENTE DE DETERMINAÇÃO OU DE EXPLICAÇÃO.....	19
3. EXERCÍCIOS.....	20
4. RESPOSTAS.....	24
5. REFERÊNCIAS.....	28



CORRELAÇÃO E REGRESSÃO

1. CORRELAÇÃO

1.1. INTRODUÇÃO

Ao se estudar uma variável o interesse eram as medidas de tendência central, dispersão, assimetria, etc. Com duas ou mais variáveis além destas medidas individuais também é de interesse conhecer se elas têm algum relacionamento entre si, isto é, se valores altos (baixos) de uma das variáveis implicam em valores altos (ou baixos) da outra variável. Por exemplo, pode-se verificar se existe associação entre a taxa de desemprego e a taxa de criminalidade em uma grande cidade, entre verba investida em propaganda e retorno nas vendas, etc.

A associação entre duas variáveis poder ser de dois tipos: **correlacional** e **experimental**. Numa relação experimental os valores de uma das variáveis são controlados pela atribuição ao acaso do objeto sendo estudado e observando o que acontece com os valores da outra variável. Por exemplo, pode-se atribuir dosagens casuais de uma certa droga e observar a resposta do organismo; pode-se atribuir níveis de fertilizante ao acaso e observar as diferenças na produção de uma determinada cultura.

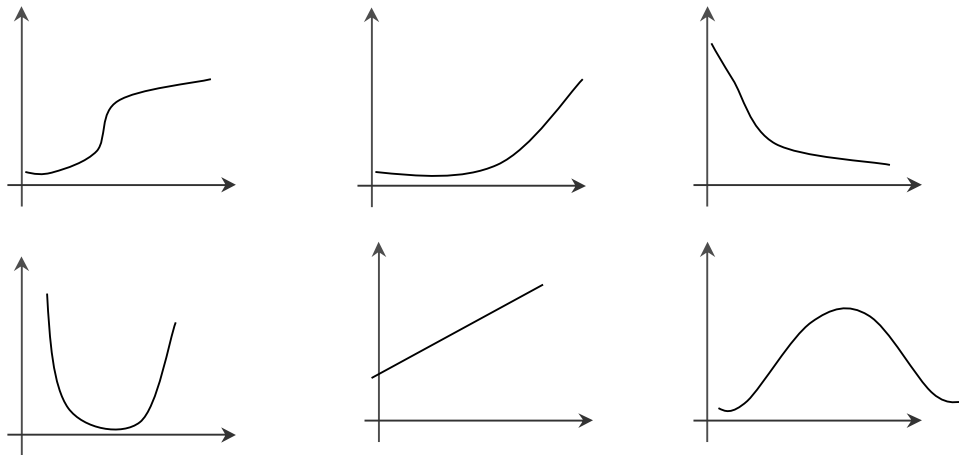
No relacionamento correlacional, por outro lado, não se tem nenhum controle sobre as variáveis sendo estudadas. Elas são observadas como ocorrem no ambiente natural, sem nenhuma interferência, isto é, as duas variáveis são aleatórias. Assim a diferença entre as duas situações é que na experimental nós atribuímos valores ao acaso de uma forma não tendenciosa e na outra a atribuição é feita pela natureza.

Freqüentemente é necessário estudar o relacionamento entre duas ou mais variáveis. Ao estudo do relacionamento entre duas ou mais variáveis denominamos de **correlação** e **regressão**. Se o estudo tratar apenas de duas variáveis tem-se a correlação e a regressão simples, se envolver mais do que duas variáveis, tem-se a correlação e a regressão múltiplas. A regressão e a correlação tratam apenas do relacionamento do tipo linear entre duas variáveis.

A análise de correlação fornece um número que resume o *grau de relacionamento linear* entre as duas variáveis. Já a análise de regressão fornece uma equação que descreve o comportamento de uma das variáveis em função do comportamento da outra variável.



Figura 1.1 - Vários tipos de relacionamento entre as variáveis X e Y



1.2. PADRÕES DE ASSOCIAÇÃO

Independente do tipo (correlacional ou experimental) a relação entre as variáveis pode ser resumida através de uma equação indicando o padrão de associação entre as duas variáveis. As relações mais comuns encontradas estão ilustradas na figura 1.1.

Quando não é possível perceber uma relação sistemática entre as variáveis é dito que as variáveis são **não correlacionadas**, são **independentes** ou ainda que são **ortogonais**.

1.3. INDICADORES DE ASSOCIAÇÃO

Suponha-se que queiramos determinar se duas variáveis aleatórias estão de alguma forma correlacionadas. Por exemplo, suponha-se que se queira determinar se o desempenho dos empregados no trabalho está de alguma forma associado ao escore obtido num teste vocacional.

Tabela de contingência 2x2. Uma vez que a correlação entre duas variáveis aleatórias reflete o quanto os altos escores de uma delas implicam em altos escores da outra e baixos escores de uma implicam em baixos escores da outra e vice-versa, no caso de uma relação negativa, pode-se começar a análise identificando, justamente quantos elementos de uma das variáveis são altos e quantos são baixos. Para determinar se um escore ou valor é alto ou baixo, pode-se convencionar que qualquer valor acima da mediana é alto e qualquer valor abaixo da mediana é baixo. Classificando desta forma pode-se ter então, para o exemplo, 4 possíveis resultados:

- ✓ Tanto o desempenho no trabalho quanto no teste estão acima da mediana (+ +)
- ✓ O desempenho no trabalho está acima mas o do teste está abaixo da mediana (+ -)



- ✓ Tanto o desempenho no trabalho quanto o do teste estão abaixo da mediana (– –)
- ✓ O desempenho no trabalho está abaixo da mediana, mas o teste não (– +)

Estas quatro possibilidades podem ser arranjadas em uma tabela de contingência 2x2, como a mostrada abaixo:

Tabela 1.1 – Desempenho no trabalho e no teste

Desempenho no trabalho	Escore no teste vocacional	
	Abaixo da mediana (–)	Acima da mediana (+)
Acima da mediana (+)	(–, +) 10 empregados	(+, +) 40 empregados
Abaixo da mediana (–)	(–, –) 40 empregados	(+, –) 10 empregados

Observe-se que se não existir relação entre as duas variáveis deve-se esperar número idêntico de empregados em cada uma das células da tabela, isto é, se a pessoa o escore da pessoa no teste vocacional está acima ou abaixo da mediana não tem nada a ver com o seu escore no desempenho no trabalho estar acima ou abaixo da mediana.

O que pode ser visto na tabela acima é que parece existir uma forte correlação entre as duas variáveis, pois ao invés de igual número em cada célula o que se tem é um número grande de ambas as variáveis acima da mediana e um número grande de escores de ambas as variáveis abaixo da mediana. Das 50 pessoas com escore acima da mediana no teste, 40 deles (80%) apresentaram escore acima da mediana no desempenho do trabalho. Da mesma forma dos 50 que tiveram classificações abaixo da mediana, 40 deles apresentaram escore abaixo da mediana no desempenho do trabalho. Se não houvesse correlação seria de se esperar que dos 50 que tiveram escores acima da mediana no teste 25 tivessem escores acima da mediana no desempenho do trabalho e 25 abaixo.

A tabela 1.2 mostra outras possíveis saídas para este tipo de esquema de classificação cruzada. Novamente 100 elementos são classificados em 4 células de acordo com o critério anterior. A parte (a) da tabela mostra uma associação positiva, a parte (b) uma negativa e a parte (c) que não deve existir associação entre duas variáveis X e Y.



Tabela 1.2 - Indicativos da presença de associação entre duas variáveis X e Y.

Valor de X	(a) Relação positiva		(b) Relação negativa			(c) Sem relação		
	Valor de Y		Valor de X	Valor de Y		Valor de Y		
	Abaixo da mediana	Acima da mediana		Abaixo da mediana	Acima da mediana	Valor de X	Abaixo da mediana	Acima da mediana
Acima da mediana	15	35	Acima da mediana	35	15	Acima da mediana	25	25
Abaixo da mediana	35	15	Abaixo da mediana	15	35	Abaixo da mediana	25	25

Diagramas de dispersão. As tabelas de contingência 2x2 fornecem somente a indicação grosseira da relação entre duas variáveis, a não ser o fato de que os valores estão situados acima e abaixo da mediana, qualquer outra informação é desperdiçada. Vamos considerar um exemplo, envolvendo duas variáveis contínuas.

Um comerciante de temperos está curioso sobre a grande variação nas vendas de loja para loja e acha que as vendas estão associadas com o espaço nas prateleiras dedicados a sua linha de produto em cada ponto de venda. Dez lojas foram selecionadas ao acaso através do país e as duas seguintes variáveis foram mensuradas: (1) total de espaço de frente (comprimento x altura em cm^2) dedicados a sua linha de produtos e (2) total das vendas dos produtos, em reais, no último mês. Os dados são apresentados na tabela 1.3.

Tabela 1.3 – Vendas x espaço dedicado aos produtos (em cm^2).

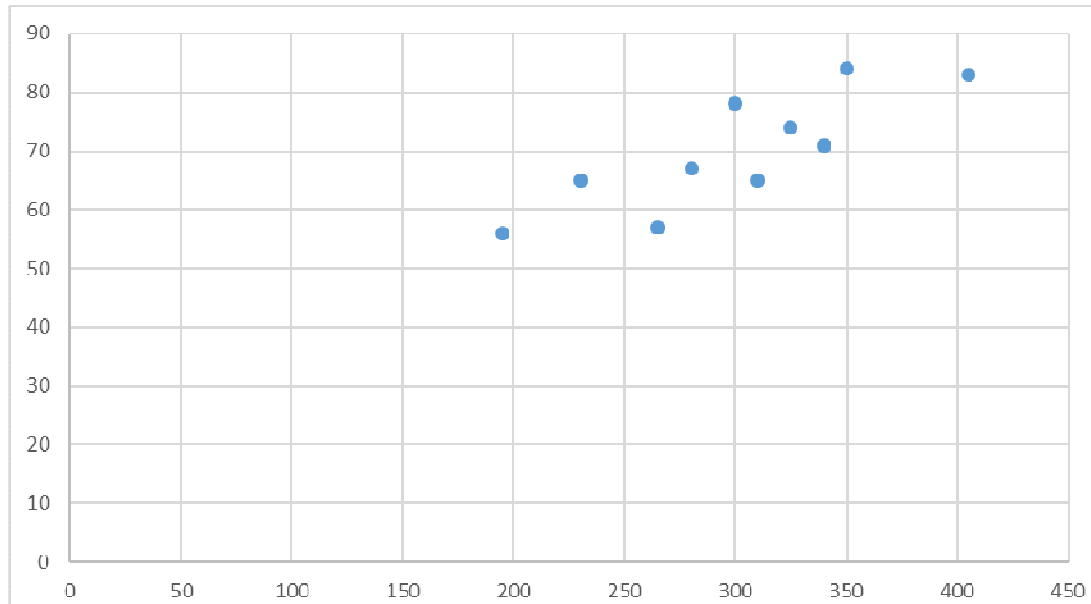
Local	Espaço	Vendas
1	340	71
2	230	65
3	405	83
4	325	74
5	280	67
6	195	56
7	265	57
8	300	78
9	350	84
10	310	65

Pela observação da tabela não é fácil perceber o tipo de relacionamento que possa existir entre as duas variáveis. Para ter uma idéia melhor, as variáveis são colocadas no que é denominado de



diagrama de dispersão. Uma das variáveis (X) é representada no eixo horizontal e a outra variável (Y) no eixo vertical, conforme figura 1.2.

Figura 1.2 – Diagrama de dispersão das variáveis apresentadas na tabela 1.3



Uma olhada rápida no diagrama de dispersão mostra a existência de um relacionamento entre as variáveis, com altos valores de uma das variáveis associados a altos valores da outra variável. Se não houvesse relacionamento entre elas, os pontos estariam distribuídos ao acaso no gráfico sem mostrarem alguma tendência.

1.4. O COEFICIENTE DE CORRELAÇÃO

Apesar do diagrama de dispersão nos fornecer uma idéia do tipo e extensão do relacionamento entre duas variáveis X e Y, seria altamente desejável ter um número que medisse esta relação. Esta medida existe e é denominada de **coeficiente de correlação**. Quando se está trabalhando com amostras o coeficiente de correlação é indicado pela letra **r** que é, por sua vez, uma estimativa do coeficiente de correlação populacional: **ρ** (rho).

O coeficiente de correlação pode variar de $-1,00$ a $+1,00$, com um coeficiente de $+1$, indicando uma correlação **linear** positiva perfeita. Neste caso, as duas variáveis serão exatamente iguais em termos de escores padronizados z, isto é, um elemento apresentando um escore padronizado de 1,5 em uma das variáveis vai apresentar o mesmo escore padronizado na outra variável. Um coeficiente de correlação de -1 , indica correlação linear perfeita negativa, com os escores padronizados exatamente iguais em valores absolutos, diferindo apenas no sinal.



Uma correlação de +1 ou -1 é raramente observado. O mais comum é que o coeficiente fique situado no intervalo entre estes dois valores. Um coeficiente de correlação “0”, significa que não existe um relacionamento **linear** entre as duas variáveis.

1.5. HIPÓTESES BÁSICAS

A suposição básica sobre o coeficiente de correlação é que o relacionamento entre as duas variáveis seja linear. Isto é, o coeficiente de correlação é adequado para avaliar somente o relacionamento linear. As duas variáveis podem estar perfeitamente relacionadas, mas se não for de forma linear o valor do coeficiente pode ser zero ou próximo de zero.

Uma segunda hipótese é que as variáveis envolvidas sejam aleatórias e que sejam medidas no mínimo em escala de intervalo. Ele não se aplica a variáveis em escala nominal ou ordinal ou quando uma das variáveis é manipulada experimentalmente, pois neste caso, a escolha dos valores experimentais vai influenciar o valor de r obtido.

Uma terceira hipótese é que as duas variáveis tenham uma distribuição conjunta normal bivariada. Isto é equivalente a dizer que para cada x dado a variável y é normalmente distribuída.

Suponha-se que existam apenas duas variáveis X e Y. Uma amostra da variável “X”, assumindo os valores particulares X_1, X_2, \dots, X_n e uma amostra da variável “Y” assumindo os valores particulares Y_1, Y_2, \dots, Y_n são obtidas e suponha-se ainda que o objetivo é saber se existe algum tipo de relacionamento linear entre estas duas variáveis. Isto poderá ser medido pelo **coeficiente de correlação** que fornece o grau de relacionamento linear entre duas variáveis.

1.6. DEFINIÇÃO

Na população o coeficiente de correlação é representado por ρ e na amostra por r . Assim dadas duas amostras, uma da variável X e outra da variável Y, o coeficiente de correlação amostral poderá ser calculado através da seguinte expressão:

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}} = \frac{n \sum x_i \cdot y_i - (\sum x_i)(\sum y_i)}{\sqrt{[n \sum x_i^2 - (\sum x_i)^2][n \sum y_i^2 - (\sum y_i)^2]}}$$

Uma população que tenha duas variáveis não correlacionadas linearmente pode produzir uma amostra com coeficiente de correlação diferente de zero. Para testar se a amostra foi ou não retirada de uma população de coeficiente de correlação não nulo entre duas variáveis, precisamos saber qual é a distribuição amostral da estatística r .



1.7. PROPRIEDADES DE R

As propriedades mais importantes do coeficiente de correlação são:

O intervalo de variação vai de -1 a +1.

O coeficiente de correlação é uma medida adimensional, isto é, ele é independente das unidades de medida das variáveis X e Y.

Quanto mais próximo de +1 for “r”, maior o grau de relacionamento linear positivo entre X e Y, ou seja, se X varia em uma direção Y variará na mesma direção.

Quanto mais próximo de -1 for “r”, maior o grau de relacionamento linear negativo entre X e Y, isto é, se X varia em um sentido Y variará no sentido inverso.

Quanto mais próximo de zero estiver “r” menor será o relacionamento linear entre X e Y. Um valor igual a zero, indicará ausência **apenas** de relacionamento linear. Isto não quer dizer que não existam outros tipos de relacionamento entre X e Y diferentes do relacionamento linear.



2. REGRESSÃO

Uma vez constatado que existe correlação linear entre duas variáveis, pode-se tentar prever o comportamento de uma delas em função da variação da outra.

Para tanto será suposto que existem apenas duas variáveis. A variável **X** (denominada variável controlada, explicativa ou independente) com valores observados X_1, X_2, \dots, X_n e a variável **Y** (denominada variável dependente ou explicada) com valores Y_1, Y_2, \dots, Y_n . Os valores de **Y** são aleatórios, pois eles dependem não apenas de **X**, mas também de outras variáveis que não estão sendo representadas no modelo. Estas variáveis são consideradas no modelo através de um termo aleatório denominado “erro”. A variável **X** pode ser **aleatória** ou então **controlada**.

Desta forma pode-se considerar que o modelo para o relacionamento linear entre as variáveis **X** e **Y** seja representado por uma equação do tipo:

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

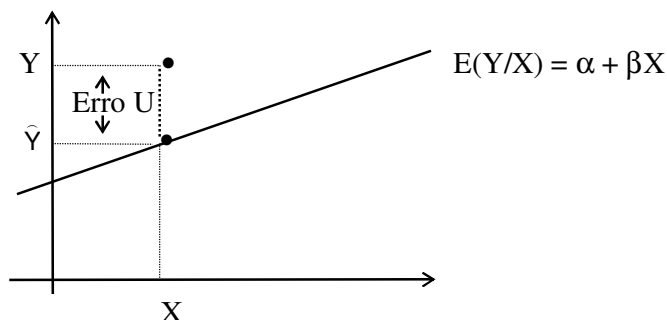
onde “**U**” é o termo erro, isto é, “**U**” representa as outras influências na variável **Y** além da exercida pela variável “**X**”.

Esta equação permite que **Y** seja maior ou menor do que $\alpha + \beta X$, dependendo de “**U**” ser positivo ou negativo. De forma ideal o termo “**U**” deve ser pequeno e independente de **X**, de modo que se possa modificar **X**, sem modificar “**U**”, e determinar o que ocorrerá, em média, a **Y**, isto é:

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

Os dados $\{(X_i, Y_i), i = 1, 2, \dots, n\}$ podem ser representados graficamente marcando-se cada par.

Figura 2.1 – O modelo de regressão linear





(X_i, Y_i) como um ponto de um plano. Os termos U_i são iguais a distância vertical entre os pontos observados (X_i, Y_i) , e os pontos calculados $(X_i, \alpha + \beta X_i)$. Isto está ilustrado na figura 2.1.

Um modelo de regressão consiste em um conjunto de hipóteses sobre a distribuição dos termos “erro” e as relações entre as variáveis X e Y .

Algumas destas hipóteses são:

(i) $E(U_i) = 0$;

(ii) $\text{Var}(U_i) = \sigma^2$

Na hipótese (i) o que se está supondo é que os U_i são variáveis aleatórias independentes com valor esperado igual a zero e na (ii) que a variância de cada U_i é a mesma e igual a σ^2 , para todos os valores de X .

Supõem-se ainda que a variável independente X , *permaneça fixa*, em observações sucessivas e que a variável dependente Y seja função linear de X . Os valores de Y devem ser independentes um do outro. Isto ocorre em geral, mas em alguns casos, como, por exemplo, observações diferentes são feitas no mesmo indivíduo em diferentes pontos no tempo esta suposição poderá não ocorrer.

Como o valor esperado de U_i é zero, o valor esperado da variável dependente Y , para um determinado valor de X , é dado pela função de regressão $\alpha + \beta X$ ou seja:

$$E(Y/X) = E(\alpha + \beta X + U) = \alpha + \beta X + E(U) = \alpha + \beta X \quad [1]$$

já que $\alpha + \beta X$ é constante para cada valor de X dado.

O símbolo $E(Y/X)$ é lido **valor esperado de Y , dado X** . A variância de Y , para determinado valor de X , é igual a:

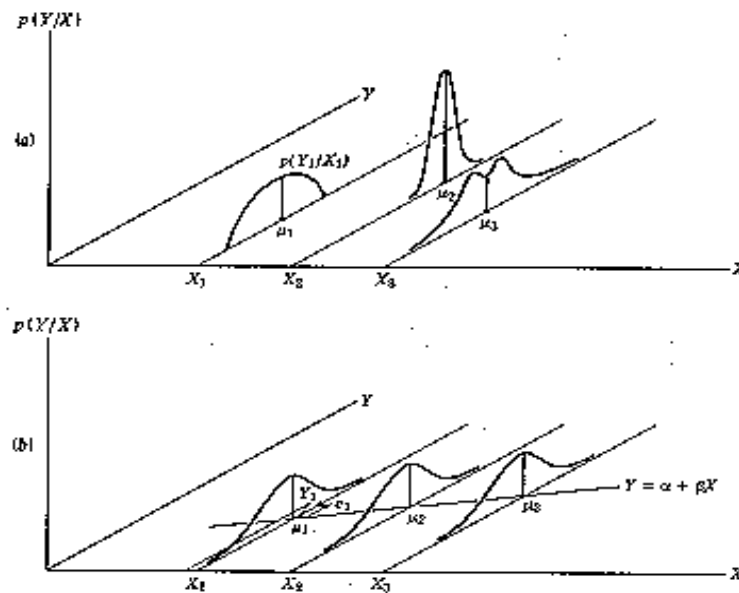
$$V(Y/X) = V(\alpha + \beta X + U) = V(U) = \sigma^2 \quad [2]$$

A hipótese de que $V(Y/X)$ é a mesma para todos os valores de X , denominada de homocedasticidade, é útil pois permite que se utilize cada uma das observações sobre X e Y para estimar σ^2 . O termo “homo” significa “o mesmo” e “cedasticidade” significa “disperso”.

De [1] e [2] decorre que, para um dado valor de X , a variável dependente Y tem função densidade de probabilidade (condicional) com média $\alpha + \beta X$ e variância σ^2 . A figura 2.2, ilustra a função densidade. Na parte superior da figura é ilustrado o caso heterocedástico e na parte inferior o caso homocedástico.



Figura 2.2 – Função densidade de Y dado X



A posição da função densidade $f(Y/X)$ varia em função da variação do valor de X . Note-se que a média da função densidade se desloca ao longo da função de regressão $\alpha + \beta X$.

Em resumo, o modelo de regressão proposto consiste nas seguintes hipóteses:

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

$$E(Y/X) = \alpha + \beta X;$$

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

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

A variável X permanece fixa em observações sucessivas;

Os erros U são normalmente distribuídos.

2.1. ESTIMATIVA DOS PARÂMETROS DE REGRESSÃO

Se fosse conhecido toda a população de valores (X_i, Y_i) então seria possível determinar os valores exatos dos parâmetros α , β e σ^2 . Como, em geral, se trabalha com amostras se faz necessário, então, estimar estes parâmetros com base nos valores da amostra.

Existem alguns métodos para ajustar uma linha entre as variáveis X e Y o mais utilizado é o denominado **método dos mínimos quadrados (MMQ)**. A reta obtida através deste método, não é



necessariamente, o “melhor” ajustamento possível, mas possui muitas propriedades estatísticas que são desejáveis.

Sejam \mathbf{a} e \mathbf{b} estimadores de $\boldsymbol{\alpha}$ e $\boldsymbol{\beta}$ e $\mathbf{E}_i = Y_i - a - bX_i$ o desvio observado em relação a reta ajustada, isto é, \mathbf{E}_i é um estimador do termo \mathbf{U}_i . O método dos mínimos quadrados exige que os estimadores \mathbf{a} e \mathbf{b} sejam escolhidos de tal forma que a soma dos quadrados dos desvios dos mesmos em relação à reta de regressão ajustada seja mínima, isto é:

$$\phi = \sum_{i=1}^n \mathbf{E}_i^2 = \sum_{i=1}^n (Y_i - a - bX_i)^2 = \text{mínimo.}$$

Para tornar mínima esta soma em relação a \mathbf{a} e \mathbf{b} , é necessário diferenciar a expressão parcialmente em relação aos valores \mathbf{a} e \mathbf{b} . Após algumas simplificações vai-se obter:

$$\sum Y_i = na + b\sum X_i \quad (\text{i})$$

$$\sum X_i Y_i = a\sum X_i + b\sum (X_i)^2 \quad (\text{ii})$$

que são denominadas de equações normais da regressão, onde “n” é o número de pares de observações.

Obs.: Para simplificar a notação foram desconsiderados os índices nos somatórios.

Dividindo-se a equação (i) por “n” e isolando o valor de \mathbf{a} vem:

$$a = \frac{\sum y_i}{n} - b\left(\frac{\sum X_i}{n}\right) = \bar{Y} - b\bar{X}$$

levando-se este resultado na equação (ii) tem-se:

$$b = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2} = \frac{\sum X_i Y_i - \frac{\sum X_i \sum Y_i}{n}}{\sum X_i^2 - \frac{(\sum X_i)^2}{n}} = \frac{n\sum X_i Y_i - \sum X_i \sum Y_i}{n\sum X_i^2 - (\sum X_i)^2}$$

A reta estimada de regressão será então:

$$\hat{Y} = \mathbf{a} + \mathbf{bX}$$

com os valores de “a” e “b” obtidos através das seguintes expressões:

$$b = \frac{n\sum X_i Y_i - \sum X_i \sum Y_i}{n\sum X_i^2 - (\sum X_i)^2} \quad \text{e} \quad a = \bar{Y} - b\bar{X}$$

Utiliza-se o valor \hat{Y} , porque o valor de Y, obtido a partir da reta estimada de regressão, para um dado valor de X, é uma estimativa do valor $\mathbf{E}(Y/X)$, isto é, do valor esperado de Y dado X.

Exemplo:

São fornecidos 5 pares de valores, na tabela abaixo, correspondentes as variáveis **X** e **Y**. A estimativa da reta de regressão entre **X** e **Y**, é obtida utilizando as expressões de **a** e **b** acima e usando os resultados obtidos na tabela 2.1.

Tabela 2.1 - Valores para estimar a linha de regressão

X	Y	X ²	XY
1	3	1	3
2	3	4	6
4	7	16	28
5	6	25	30
8	12	64	96
20	31	110	163

$$\bar{X} = 20 / 5 = 4;$$

$$\bar{Y} = 31/5 = 6,2$$

$$b = (5.163 - 20.31) / (5.110 - 400) = 1,30$$

$$a = \bar{Y} - b\bar{X} = 6,20 - 1,30.4 = 1$$

Então a linha estimada será: $\hat{Y} = 1.3X + 1$

Esta reta é o “melhor” ajustamento para estes dados e seria diferente para cada amostra das variáveis X e Y, retiradas desta mesma população. Esta reta pode ser considerada uma estimativa da verdadeira linha de regressão onde 1,3 seria uma estimativa do valor β (parâmetro angular) e 1 uma estimativa do valor α (parâmetro linear), que são os verdadeiros coeficientes de regressão.

2.2. ESTIMATIVA DA VARIÂNCIA DO TERMO ERRO

O termo erro, U, é uma variável aleatória, supostamente com média zero e variância constante. Então, intuitivamente parece plausível usar os resíduos da reta de regressão pelos método dos mínimos quadrados para se estimar a variância σ^2 dos termos “erro”. A variância amostral desses resíduos é igual a:

$$\hat{\sigma}^2 = \frac{\sum (E - \bar{E})^2}{n}, \text{ onde } \bar{E} = \sum E / n. \text{ Observe-se, entretanto, que:}$$

$$\sum E = \sum (Y - a - bX) = \sum Y - na - b\sum X = 0, \text{ pela primeira equação normal (i).}$$

Portanto, $\hat{\sigma}^2$ pode ser escrito como: $\hat{\sigma}^2 = \sum E^2 / n$.



Mas $\hat{\sigma}^2$, neste caso, é um estimador tendencioso. Pode-se obter um estimador não tendencioso, multiplicando $\hat{\sigma}^2$ por $n / (n - 2)$. O novo estimador, não tendencioso, será representado S^2 e sua raiz quadrada:

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

é denominada de “erro-padrão da estimativa” ou “erro-padrão amostral da regressão”.

Obs.: A utilização de “ $n - 2$ ” é consequência do fato de que se deve estimar dois parâmetros, α e β , antes de obter os resíduos E . Como resultado, há somente “ $n - 2$ ” graus de liberdade associados à quantidade $\sum E^2$.

A expressão acima, para o cálculo do erro amostral da regressão, apresenta o inconveniente de exigir o cálculo de cada valor previsto de Y , através da linha de regressão, tornando sua obtenção muito trabalhosa. Existe, entretanto, uma alternativa para se obter este valor (erro padrão da estimativa) sem a necessidade de calcular todos os valores previstos.

Observe-se que:

$$\sum E^2 = \sum (Y - \hat{Y})^2 = \sum (Y - a - bX)^2 = \sum [Y - \bar{Y} + b(\bar{X} - bX)]^2 = \sum (Y - \bar{Y})^2 - 2b \sum (X - \bar{X})(Y - \bar{Y}) + \sum b^2(\bar{X} - X)^2.$$

Fazendo:

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

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

$$\sum (X - \bar{X})(Y - \bar{Y}) = \sum XY - \frac{\sum X \sum Y}{n} = S_{XY}$$

Lembrando que:

$$b = \frac{n \sum X_i Y_i - \sum X_i \sum Y_i}{n \sum X_i^2 - (\sum X_i)^2} = \frac{\sum X_i Y_i - \frac{\sum X_i \sum Y_i}{n}}{\sum X_i^2 - \frac{(\sum X_i)^2}{n}}, \text{ segue que } b = S_{XY}/S_{XX} \text{ e que } S_{XY} = bS_{XX}$$

Então vem:

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

Assim:



$$S^2 = \frac{\sum E^2}{n-2} = \frac{\sum (Y-a-bX)^2}{n-2} = \frac{S_{YY} - b^2 S_{XX}}{n-2} = \frac{S_{YY} - b S_{XY}}{n-2}$$

Pode-se verificar que S^2 definido desta maneira é um estimador não-tendencioso de σ^2 , isto é, $E(S^2) = \sigma^2$.

O erro padrão da regressão será dado, então, por:

$$s = \sqrt{\frac{S_{YY} - b^2 S_{XX}}{n-2}} = \sqrt{\frac{S_{YY} - b S_{XY}}{n-2}}$$

Exemplo:

Considerando as variáveis X e Y acima e a linha de regressão anterior determinar uma estimativa do erro padrão da regressão.

Os cálculos necessários estão na tabela 2.2.

Tabela 2.2 – Determinação do erro padrão da regressão

X	Y	Y _c	E = Y - Y _c	E ²
1	3	2,3	0,7	0,49
2	3	3,6	-0,6	0,36
4	7	6,2	0,8	0,64
5	6	7,5	-1,5	2,25
8	12	11,40	0,6	0,36
20	31	31	0	4,10

O erro padrão da regressão será então:

$$s = \sqrt{\frac{\sum E^2}{n-2}} = \sqrt{\frac{\sum (Y-a-bX)^2}{n-2}} = \sqrt{\frac{4,10}{5-3}} = \sqrt{1,3667} = 1,17$$

Este mesmo cálculo poderá ser efetuado pela expressão definida acima, sem a necessidade de se obter os valores estimados.

Tabela 2.3 – Determinação do erro padrão da regressão

X	Y	X ²	Y ²	XY
1	3	1	9	3
2	3	4	9	6
4	7	16	49	28
5	6	25	36	30
8	12	64	144	96
20	31	110	247	163



Neste caso, tem-se:

$$S_{XX} = \sum X^2 - \frac{(\sum X)^2}{n} = 110 - 20^2/5 = 30$$

$$S_{YY} = \sum Y^2 - \frac{(\sum Y)^2}{n} = 247 - 31^2/5 = 54,80$$

$$S_{XY} = \sum XY - \frac{\sum X \sum Y}{n} = 163 - (20 \cdot 31)/5 = 39$$

O valor de “b” será:

$$b = S_{XY}/S_{XX} = 39/30 = 1,30$$

Portanto o erro padrão da regressão será:

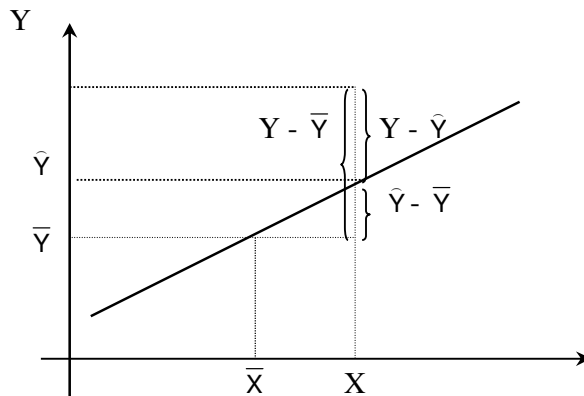
$$s = \sqrt{\frac{S_{YY} - b^2 S_{XX}}{n-2}} = \sqrt{\frac{S_{YY} - b S_{XY}}{n-2}} = \sqrt{\frac{54,80 - 13,39}{5-2}} = \sqrt{\frac{4,10}{3}} = \sqrt{1,3667} = 1,1690 = 1,17$$

2.3. DECOMPOSIÇÃO DA SOMA DOS QUADRADOS

2.3.1. DECOMPOSIÇÃO DOS DESVIOS

Pelo figura 2.3, pode-se perceber que o desvio em relação a Y (desvio total), isto é, $Y - \bar{Y}$ pode ser decomposto em dois outros desvios:

Figura 2.3 – Desvios na regressão



- * O desvio explicado pela linha de regressão, isto é, $\hat{Y} - \bar{Y}$ e
- * O desvio não-explicado (resíduos) pela linha de regressão, isto é, $Y - \hat{Y}$.



É fácil perceber que a variação total, $\sum(Y - \bar{Y})$, é a soma da variação explicada, $\sum(\hat{Y} - \bar{Y})$, e a não-explicada, $\sum(Y - \hat{Y})$, pois:

$$Y - \bar{Y} = Y - \hat{Y} + \hat{Y} - \bar{Y}, \text{ então:}$$

Aplicando somatório a ambos os membros vem:

$$\sum(Y - \bar{Y}) = \sum(Y - \hat{Y}) + \sum(\hat{Y} - \bar{Y})$$

Pode-se verificar também que a propriedade aditiva dos desvios é extensiva à soma dos quadrados desses desvios, ou seja:

$$\sum(Y - \bar{Y})^2 = \sum(Y - \hat{Y})^2 + \sum(\hat{Y} - \bar{Y})^2$$

De fato:

$$\sum(Y - \bar{Y})^2 = \sum(Y - \hat{Y} + \hat{Y} - \bar{Y})^2 = \sum[(Y - \hat{Y}) + (\hat{Y} - \bar{Y})]^2 = \sum(Y - \hat{Y})^2 + \sum(\hat{Y} - \bar{Y})^2 - 2\sum(Y - \hat{Y})(\hat{Y} - \bar{Y})$$

Mas

$$\sum(Y - \hat{Y})(\hat{Y} - \bar{Y}) = \sum(Y - \hat{Y})(a + bX - a - b\bar{X}) = b\sum X(Y - \hat{Y}) - b\bar{X}\sum X(Y - \hat{Y})$$

Pelas condições do método dos mínimos quadrados, tem-se:

$$\sum(\hat{Y} - \bar{Y}) = 0 \text{ e } \sum X(Y - \hat{Y}) = 0, \text{ em consequência}$$

$$\sum(Y - \hat{Y})(\hat{Y} - \bar{Y}) = 0, \text{ logo, segue que:}$$

$$\sum(Y - \bar{Y})^2 = \sum(Y - \hat{Y})^2 + \sum(\hat{Y} - \bar{Y})^2,$$

isto é, que a soma dos quadrados dos desvios calculados em torno da média de Y (variação total = VT) é igual à soma dos quadrados dos desvios em torno da linha de regressão (variação residual = VR) mais a soma dos quadrados dos desvios da linha de regressão em torno da média (variação explicada = VE).

2.3.2. CÁLCULO DAS VARIAÇÕES

(a) Variação Total: VT ou s_Y^2

$$VT = \sum(Y - \bar{Y})^2 = S_{YY}, \text{ onde } S_{YY} = \sum Y^2 - (\sum Y)^2 / n$$

(b) Variação Explicada: VE ou $s_{\hat{Y}}^2$

$$VE = \sum(\hat{Y} - \bar{Y})^2 = \sum(a + bX - \bar{Y})^2 = \sum(\bar{Y} - b\bar{X} + bX - \bar{Y})^2 = \sum[(b(X - \bar{X}))]^2 = b^2\sum(X - \bar{X})^2 = b^2S_{XX}$$



Logo:

$$VE = b^2 S_{XX} \text{ ou } VE = \left(\frac{S_{XY}}{S_{XX}} \right)^2 S_{XX} = b S_{XY}$$

(c) **Varição Residual: VR ou $s_{\hat{Y}/X}^2$**

De acordo com a propriedade aditiva das variações, pode-se calcular VR por diferença.

Assim:

$$VR = \sum(Y - \hat{Y})^2 = VT - VE \text{ ou } VR = S_{YY} - b S_{XY}$$

2.4. COEFICIENTE DE DETERMINAÇÃO OU DE EXPLICAÇÃO

Além dos testes de hipóteses e dos intervalos de confiança, outro indicador que fornece elementos para a análise do modelo adotado é o coeficiente de determinação ou de explicação, definido por:

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

O coeficiente de determinação indica quantos por cento a variação explicada pela regressão representa sobre a variação total. Deve-se ter:

$$0 \leq R^2 \leq 1$$

Se R^2 for igual a 1, isto significa que todos os pontos observados se situam “exatamente” sobre a reta de regressão. Tendo-se, neste caso, um ajuste perfeito. As variações da variável Y são 100% explicadas pelas variações da variável X, não ocorrendo desvios em torno da função estimada.

Por outro lado, se $R^2 = 0$, isto quer dizer que as variações de Y são exclusivamente aleatórias e explicadas pelas variações de outros fatores que não X.



3. EXERCÍCIOS

(01) Para cada uma das situações abaixo, diga o que é mais adequado: a análise de regressão ou a análise de correlação. Por quê?

- (01.1) Uma equipe de pesquisadores deseja determinar se o rendimento na Universidade sugere êxito na profissão escolhida.
- (01.2) Deseja-se estimar o número de quilômetros que um pneu radial pode rodar antes de ser substituído.
- (01.3) Deseja-se prever quanto tempo será necessário para executar uma determinada tarefa por uma pessoa, com base no tempo de treinamento.
- (01.4) Deseja-se verificar se o tempo de treinamento é importante para avaliar o desempenho na execução de uma dada tarefa.
- (01.5) Um gerente deseja estimar as vendas semanais com base nas vendas das segundas e terças-feiras.

(02) Suponha que uma cadeia de supermercados tenha financiado um estudos dos gastos com mercadorias para famílias de 4 pessoas. O estudo se limitou a famílias com renda líquida entre 8 e 20 salários mínimos. Obteve-se a seguinte equação:

$\hat{Y} = -1,20 + 0,40X$, onde \hat{Y} = despesa mensal estimada com mercadorias e X = renda líquida mensal.

- (02.1) Estimar a despesa de uma família com renda mensal líquida de 15 s.m.
- (02.2) Um dois diretores da empresa ficou intrigado com o fato de que a equação sugerir que uma família com renda de 3 s.m. líquidos mensais não gaste nada em mercadorias. Qual a explicação?
- (02.3) Explique por que a equação acima não poderia ser utilizada para estimar
- (a) As despesas com mercadorias de famílias de 5 pessoas.
- (b) As despesas com mercadorias de famílias com renda de 20 a 40 s.m. líquidos mensais.

(03) Utilize os valores abaixo para estimar as equações de regressão:

(03.1) $\sum X = 200$, $\sum Y = 300$, $\sum XY = 6200$, $\sum X^2 = 3600$ e $n = 20$

(03.2) $\sum X = 7,2$, $\sum Y = 37$, $\sum XY = 3100$, $\sum X^2 = 620$ e $n = 36$

(04) Para cada uma das situações abaixo, grafe os valores em um diagrama e se uma equação linear parecer apropriada para explicar os dados, determine os seus parâmetros.

(04.1)

Tamanho do pedido(X)	25	20	40	45	22	63	70	60	55	50	30
Custo Total (Y)	2000	3500	1000	800	3000	1300	1500	1100	950	900	1600

(04.2)

Vendas em mil (X)	201	225	305	380	560	600	685	735	510	725	450	370	150
Lucro em mil (Y)	17	20	21	23	25	24	27	27	22	30	21	19	15



(05) Suponha que uma população se constitua dos seis pontos seguintes:

(1, 2), (4, 6), (2, 4), (2, 3), (3, 5) e (5, 10)

(05.1) Grafite os pontos em um diagrama de dispersão.

(05.2) Determine a equação de regressão: $Y = \alpha + \beta X + u$.

(05.3) Os termos-erro verificam a condição $E(u) = 0$?

(05.4) Selecione uma amostra de tamanho $n = 4$, da população acima e estime a equação de regressão determinada no item 5.2. Grafite o resultado no mesmo diagrama construído em 5.1.

(06) Verifique que a reta de regressão $\hat{Y} = a + bX$, sempre passa pelo ponto (\bar{X}, \bar{Y}) .

(07) Os dados abaixo foram colhidos de cinco fábricas diferentes de uma determinada indústria:

Custo total (Y)	80	44	51	70	61
Produção (X)	12	4	6	11	8

(07.1) Estime uma função linear da forma $\hat{Y} = a + bX$ para o custo total dessa indústria.

(07.2) Qual o significado econômico das estimativas “a” e “b”?

(07.3) Determine o erro padrão da regressão.

(08) Em uma amostra aleatória de 1990, 50 homens americanos entre 35 e 54 anos de idade acusaram a seguinte relação entre renda anual Y (em dólares) e a escolaridade X (em anos). $\hat{Y} = 1200 + 800X$. A renda média foi de 10000 dólares e a escolaridade média foi de 11,0 anos. Sabendo, ainda, que $\sum X^2 = 9000$ e que o desvio padrão residual em relação à reta ajustada foi de 7300 dólares, determine:

(08.1) A renda de uma pessoa que tenha completado 2 anos de educação secundária ($x = 10$ anos).

(08.2) Se é válida a afirmação que cada ano de escolaridade custa 800 dólares?

(09) Uma pesquisa foi realizada com o objetivo de determinar os efeitos da falta de sono sobre a capacidade de as pessoas resolverem problemas simples. Foram testadas 10 pessoas, mantendo-se cada grupo de 2 pessoas sem dormir por um determinado número de horas. Após cada um destes períodos, cada pessoa teve de resolver um teste com adições simples, anotando-se então os erros cometidos. Os dados resultantes estão na tabela abaixo:

Número de erros (Y)	6, 8	6, 10	8, 14	12, 14	12, 16
Número de horas sem dormir (X)	8	12	16	20	24

(9.1) Determine a estimativa da linha de regressão do número de erros em função do número de horas sem dormir.

(9.2) Determine a dispersão dos termos erro em torno da linha de regressão.

(10) Realizou-se uma pesquisa de mercado com o objetivo de estudar a relação entre o tempo necessário para um consumidor tomar uma decisão (sobre o que comprar) e o número de embalagens alternativas do mesmo produto apresentadas a esse consumidor. Eliminaram-se as marcas das embalagens, a fim de reduzir o efeito da preferência por uma ou outra marca. Os consumidores fizeram suas escolhas somente com base na descrição do produto, anotada nas embalagens pelos fabricantes. O tempo necessário, Y, para que cada um tomasse sua decisão foi anotado para 15 participantes, resultando nos seguintes dados:



Tempo para decisão, Y (em segundos)	5, 7, 8, 8, 9	7, 8, 9, 9, 10	9, 10, 10, 11, 12
Número de alternativas (X)	2	3	4

(10.1) Determine a reta dos mínimos quadrados de Y em função de X.

(10.2) Determine o erro padrão da estimativa, ou seja, o desvio padrão amostral da regressão.

(11) Mediu-se a altura de uma amostra de 5 meninos (em polegadas) na idade de 4 anos e novamente na idade de 18 anos. Os resultados obtidos estão abaixo:

Na idade de 4 anos	40	43	40	40	42
Na idade de 18 anos	68	74	70	68	70

(11.1) Determine o coeficiente de correlação entre as duas categorias de alturas.

(11.2) Qual é o percentual das variações da variável “Altura aos 18 anos” que não é explicada pela variação na altura aos quatro anos.

(12) A equação de regressão estimada abaixo resume um estudo da relação entre o uso do fumo e a incidência de câncer pulmonar, relacionando o número X de anos que uma pessoa fumou com a percentagem Y de incidência de câncer pulmonar em cada grupo.

$$\hat{Y} = -2 + 1,70.X \quad e \quad r = 0,60.$$

(12.1) Explique o significado das estimativas “-2” e “1,70” na equação de regressão.

(12.2) Qual a taxa de incidência de câncer pulmonar para as pessoas que fumam há 20 anos?

(12.3) Se “r” fosse igual a “um” seria possível concluir que o fumo é a única causa de câncer pulmonar?

(13) Explique se concorda ou não com as seguintes afirmativas:

(13.1) Um coeficiente de correlação de +1,0 entre duas variáveis X e Y indica que X causa Y, mas um coeficiente de correlação de -1,0 significa que X não causa Y.

(13.2) Se o coeficiente de regressão é zero, o coeficiente de correlação é também zero.

(13.3) Se o coeficiente angular é 1 (um), isto significa que existe perfeita correlação entre X e Y.

(13.4) É possível que o coeficiente de correlação amostral seja positivo, quando não existe, de fato, nenhuma correlação entre as variáveis X e Y.

(13.5) Não se pode utilizar a técnica da regressão pelo método dos mínimos quadrados quando a relação básica entre X e Y não for linear.

(14) Um estudo de duas safras forneceu as seguintes informações:

Safra A: $\hat{Y} = 200 + 0,8X$, $r = 0,70$ e $S = 30$ Safra B: $\hat{Y} = 50 + 1,20X$, $r = 0,9$ e $S = 20$, onde Y é a produção por alqueire e X é a quantidade de chuva (em polegadas) no período da safra.

(14.1) Se não houvesse chuva, estas duas equações poderiam ser usadas para predizer a quantidade produzida nas duas safras? Por quê?

(14.2) Qual das duas safras tira mais proveito do aumento das chuvas? Por quê?

(14.3) Para qual das duas safras é possível predizer a produção com melhor aproximação? Por quê?

(15) Examine os cinco pares de pontos dados na tabela



X	-2	-1	0	1	2
Y	4	1	0	1	4

(15.1) Qual é a relação matemática entre X e Y?

(15.2) Determine o valor de r.

(15.3) Mostre que calculando-se a linha de regressão de Y em relação a X tem-se $b = 0$.

(15.4) Por que, aparentemente, não existe relação entre X e Y como estão indicando b e r?



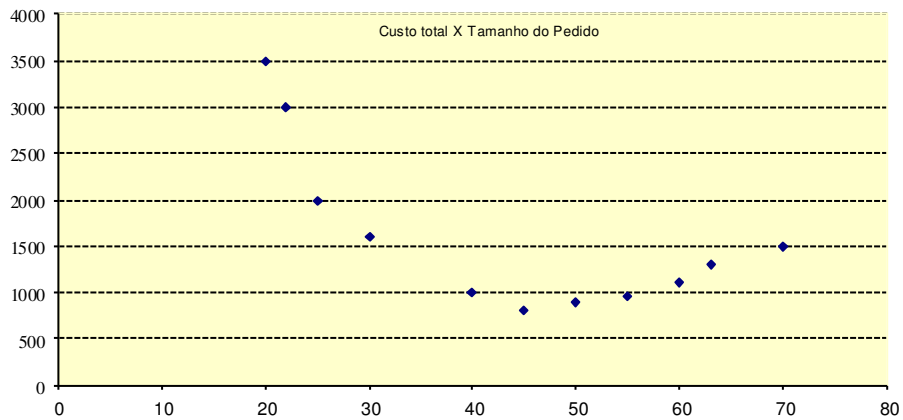
4. RESPOSTAS

(01) (01.1) Correlação (01.2) Regressão (01.3) Regressão
(01.4) Correlação (01.5) Regressão

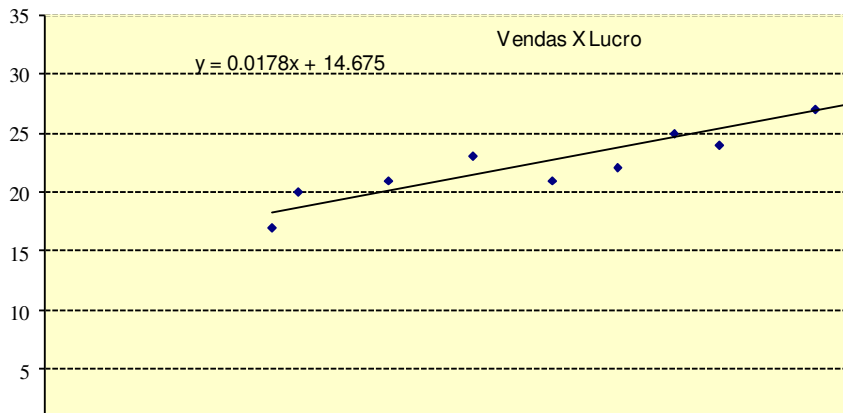
(02) (02.1) 4,80 s.m.
(02.2) A equação é para rendas entre 8 e 20 sm.
(02.3) (a) A equação foi determinada para famílias de 4 pessoas.
(b) Os dados utilizados são para famílias entre 8 e 20 sm.

(03) (03.1) $\hat{Y} = -5 + 2.X$ (03.2) $\hat{Y} = -35 + 5.X$

(04) (04.1) Neste caso, com base no diagrama, uma linha reta não é adequada.

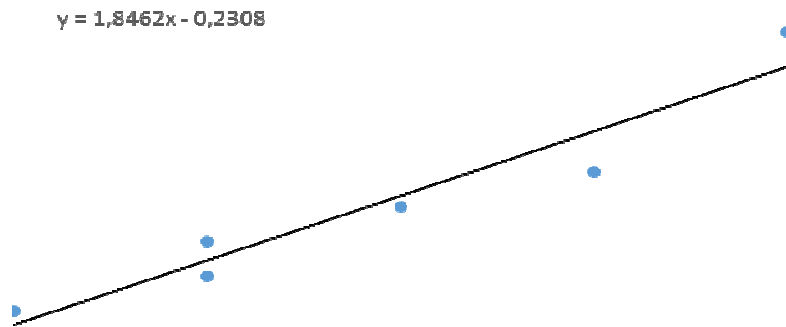


(04.2) Neste caso, uma linha é adequada e sua equação está sobre o gráfico abaixo.





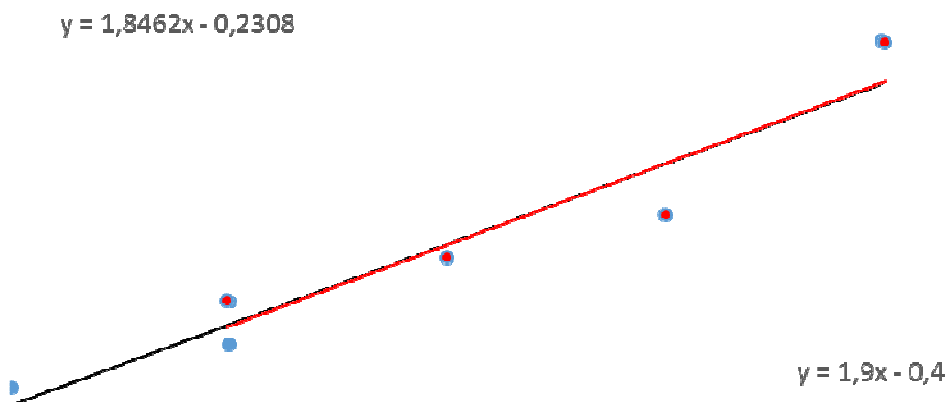
05. (05.1) e (05.2)



(05.3) e (05.4)

População				Amostra	
X	Y	Yc	Erro	X	Y
1	2	1.62	0.38	4	6
4	6	7.15	-1.15	2	4
2	4	3.46	0.54	3	5
2	3	3.46	-0.46	5	10
3	5	5.31	-0.31		
5	10	9.00	1.00		
17	30	30.00	0.00		

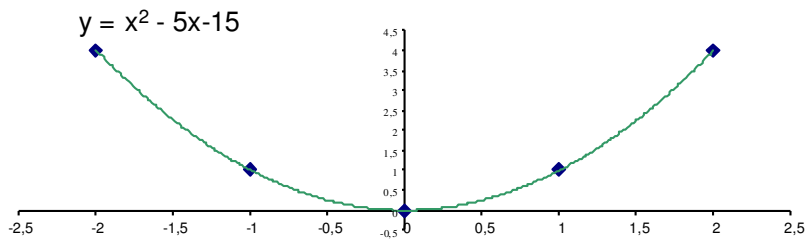
(05) (05.4)



Obs. Em virtude ter sido sorteado 4 pontos apenas, ou seja, foram eliminados dois que são simétricos em torno da linha tem-se esta situação particular onde a linha amostral coincidiu com a linha populacional, pois os pontos que ficaram fora do sorteio são simétricos em relação a linha.

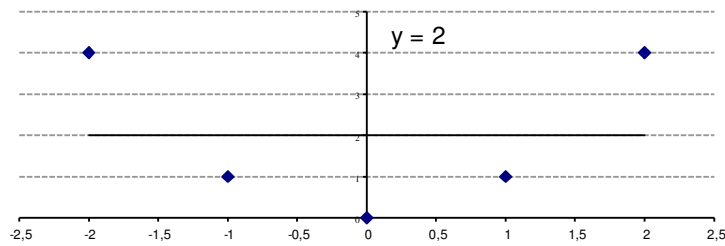


- (06) Basta mostrar que o ponto (\bar{X}, \bar{Y}) satisfaz a equação de regressão $\hat{Y} = a + bX$. Se substituirmos X por \bar{X} na equação o resultado deverá ser \bar{Y} . Mas $a + bX = a + b\bar{X} = \bar{Y} - b\bar{X} + b\bar{X} = \bar{Y}$. Uma vez que $a = \bar{Y} - b\bar{X}$.
- (07) (07.1) $\hat{Y} = 4,2589 + 26,2770.X$
(07.2) $a =$ Custo fixo $b =$ Custo marginal.
(07.3) $s = 0,37$.
- (08) (08.1) $\hat{Y} = 9200$ (08.2) Não, de fato cada ano aumenta a renda nesse valor.
- (09) (09.1) $\hat{Y} = 3 + 0,48X$ (09.2) 2,24
- (10) (11.1) $\hat{Y} = 4,30 + 1,50X$ ($r = 0,73$) (11.2) $s = 1,24$
- (11) (11.1) $r = 0,87$ (11.2) $1 - R^2 = 24,31\%$
- (12) (12.1) “-2” seria a taxa de incidência de câncer pulmonar que não está relacionada ao hábito de fumar, ou de quem nunca fumou. “1,70” é a variação na taxa de câncer pulmonar para cada ano que a pessoa fumou.
(12.2) $\hat{Y} = -2 + 1,70.20 = 32$.
(12.3) Não, pois “r” indica associação na amostra e pode não ser o mesmo na população. Além disso um coeficiente de correlação alto não implica necessariamente em relação de causa e efeito.
- (13) (13.1) Tanto um coeficiente de “+1” quanto um de “-1” indicam correlação perfeita entre as variáveis, mas não que exista necessariamente relação de causa e efeito.
(13.2) Coeficiente de regressão igual a zero implica em correlação também zero e vice-versa
(13.3) Não necessariamente, pois neste caso “1” é o valor de inclinação da linha e não grau de associação linear entre as duas variáveis.
(13.4) Sim é possível.
(13.5) A técnica dos mínimos quadrados pode ser utilizada para ajustar vários tipos de equação e não apenas uma linha reta.
- (14) (14.1) Neste caso, a interpretação deve ser mais cuidadosa, pois tanto o excesso de chuvas quanto a falta vão distorcer os dados e estas equações podem não ser mais válidas.
(14.2) A safra B tira mais proveito, provavelmente por ser uma cultura que precisa de mais chuvas.
(14.3) Para a safra B, pois existe uma melhor aderência dos dados a equação.
- (15) (15.1)



(15.2) $r = 0$

(15.3)



(15.4) Porque a correlação mostra apenas o relacionamento linear e, neste caso, o relacionamento é do tipo parábola (equação do segundo grau).



5. REFERÊNCIAS

- BUSSAB, Wilton O, MORETTIN, Pedro A. *Estatística Básica*. 3ª ed. São Paulo, Atual, 1986.
- DOWNING, Douglas, CLARK, Jeff. *Statistics the Easy Way*. Barron's Educational Series, Inc. New York, 1989.
- FONSECA, Jairo Simon da, MARTINS, Gilberto de Andrade, TOLEDO, Geraldo Luciano. *Estatística Aplicada*. São Paulo: Editora Atlas, 1976.
- FONSECA, Jairo Simon da, MARTINS, Gilberto de Andrade. *Curso de Estatística*. São Paulo: Editora Atlas S. A., 1980.
- HOFFMAN, Rodolfo. *Estatística para Economistas*. São Paulo. Livraria Pioneira Editora, 1980.
- KLEIBAUM, David G., KUPPER, Lawrence L. *Applied Regression Analysis and Other Multivariable Methods*. North Scituate, Massachusetts: Duxbury Press, 1978.
- MARKLAND, Robert E., SWEIGART, James R. *Quantitative Methods: Applications to Managerial Decision Making*. New York: John Wiley & Sons, 1987. 827p.
- MASON, Robert D., DOUGLAS, Lind A. *Statistical Techniques in Business And Economics*. IRWIN, Boston, 1990.
- MEYER, Paul L. *Probabilidade: aplicações à Estatística*. Tradução do Prof. Ruy C. B. Lourenço Filho. Rio de Janeiro, Livros Técnicos e Científicos Editora S.A., 1978.
- [MILLER, Charles D., HEEREN, Vern E., HORNSBY Jr., E. John. *Mathematical Ideas*. USA: Harper Collins Publishers, 1990.
- ROTHENBERG, Ronald I. *Probability and Statistics*. Hartcourt Brace Jovanovich, Publishers, Orlando, Florida, 1991.
- SALVATORE, Dominick. *Estatística e Econometria*. Tradução Newton Boer, revisão técnica Marco Antônio S. de Vasconcelos. São Paulo: McGraw-Hill do Brasil, 1982.