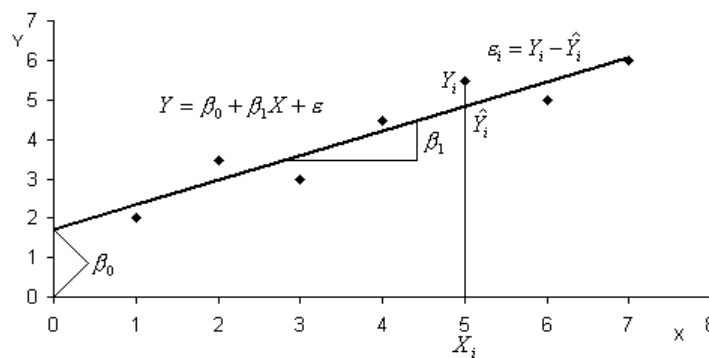


REGRESSÃO E CORRELAÇÃO

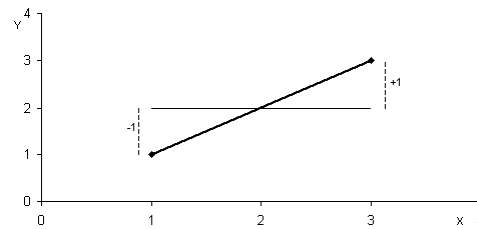


AJUSTE DE UMA RETA





MINIMIZAÇÃO DOS DESVIOS

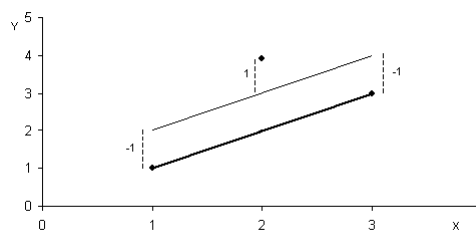


$$\sum (y_i - \hat{y}_i)$$

3



MINIMIZAÇÃO DOS DESVIOS ABSOLUTOS



$$\sum |y_i - \hat{y}_i|$$

4



EXEMPLO 1

- Considere o seguinte conjunto de pontos

X	Y
1	1
2	1
3	2
4	2
5	4

5



RETAS DE AJUSTE

R1 $Y = -0.1 + 0.7X$

R2 $Y = 0.5 + 0.5X$

R3 $Y = -0.7 + 0.9X$

6

RETAS



R1	R2	R3
0.6	1	0.2
1.3	1.5	1.1
2	2	2
2.7	2.5	2.9
3.4	3	3.8

7

DESVIOS



Desv1	Desv2	Desv3
0.4	0	0.8
-0.3	-0.5	-0.1
0	0	0
-0.7	-0.5	-0.9
0.6	1	0.2
0	0	0

8

DESVIOS ABSOLUTOS



$ \text{Desv1} $	$ \text{Desv2} $	$ \text{Desv3} $
0.4	0	0.8
0.3	0.5	0.1
0	0	0
0.7	0.5	0.9
0.6	1	0.2
2	2	2

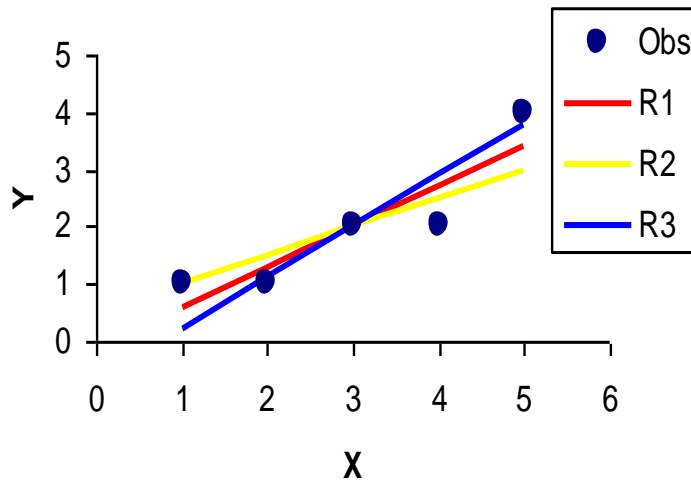
9

QUADRADO DOS DESVIOS



$(\text{Desv1})^2$	$(\text{Desv2})^2$	$(\text{Desv3})^2$
0.16	0	0.64
0.09	0.25	0.01
0	0	0
0.49	0.25	0.81
0.36	1	0.04
1.10	1.50	1.50

10



11



EXEMPLO 2



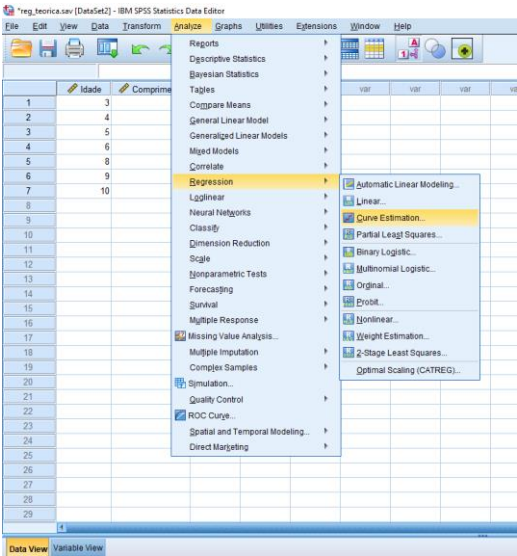
Comprimento alar
(cm) em função da
idade (dias) para
andorinhas

Dias	Comp.
3	1,4
4	1,5
5	2,1
6	2,4
8	3,1
9	3,2
10	3,3

DPS

12

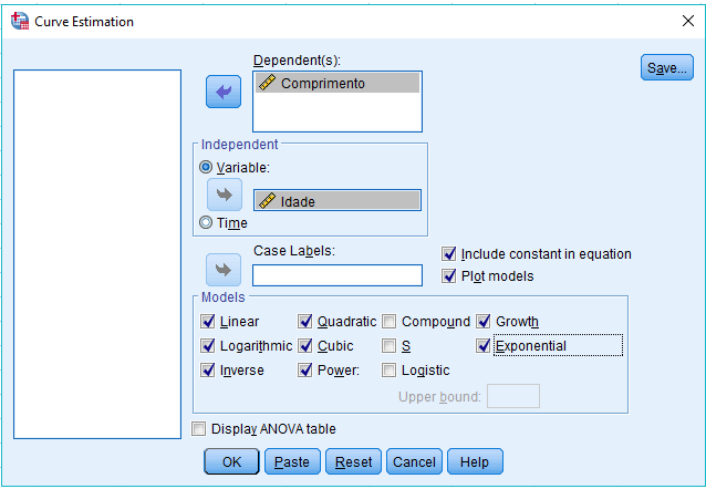
EXEMPLO 2 (SPSS)



DPS

13

EXEMPLO 2



DPS

14



Models (curve estimation algorithms)

Previous Next

CURVEFIT allows the user to specify a model with or without a constant term designated by β_0 . If this constant term is excluded, simply set it zero or one depending upon whether it appears in an additive or multiplicative manner in the models listed below.

- (1) Linear $E(Y_i) = \beta_0 + \beta_1 t$
- (2) Logarithmic $E(Y_i) = \beta_0 + \beta_1 \ln(t)$
- (3) Inverse $E(Y_i) = \beta_0 + \beta_1 / t$
- (4) Quadratic $E(Y_i) = \beta_0 + \beta_1 t + \beta_2 t^2$
- (5) Cubic $E(Y_i) = \beta_0 + \beta_1 t + \beta_2 t^2 + \beta_3 t^3$
- (6) Compound $E(Y_i) = \beta_0 \beta_1^t$
- (7) Power $E(Y_i) = \beta_0 t^{\beta_1}$
- (8) S $E(Y_i) = \exp(\beta_0 + \beta_1 / t)$
- (9) Growth $E(Y_i) = \exp(\beta_0 + \beta_1 t)$
- (10) Exponential $E(Y_i) = \beta_0 e^{\beta_1 t}$
- (11) Logistic $E(Y_i) = \left(\frac{1}{u} + \beta_0 \beta_1^t \right)^{-1}$

DPS

15

EXEMPLO 2



*Output1 [Document1] - IBM SPSS Statistics Viewer

File Edit View Data Transform Insert Format Analyze Graphs Utilities Extensions Window Help

Output Log Curve Fit

Curve Fit

[DataSet2] C:\Users\ACB\OneDrive\Aulas2017_18\reg_teorica.sav

Model Description

Model Name	MOD_1
Dependent Variable	1
Equation	1
	2
	3
	4
	5
	6
	7
	8
Independent Variable	Idade
Constant	Included
Variable Whose Values Label Observations in Plots	Unspecified
Tolerance for Entering Terms in Equations	0,0001

a. The model requires all non-missing values to be positive.

DPS

16



Output

Log

Curve Fit

Title

Notes

Active Dataset

Model Description

Case Processing

Variable Processi

Model Summary a

Curvelfit for Compr

Variable Processing Summary

	Variables	
	Dependent Comprimento	Independent Idade
Number of Positive Values	7	7
Number of Zeros	0	0
Number of Negative Values	0	0
Number of Missing Values	User-Missing	0
	System-Missing	0

Model Summary and Parameter Estimates

Dependent Variable: Comprimento

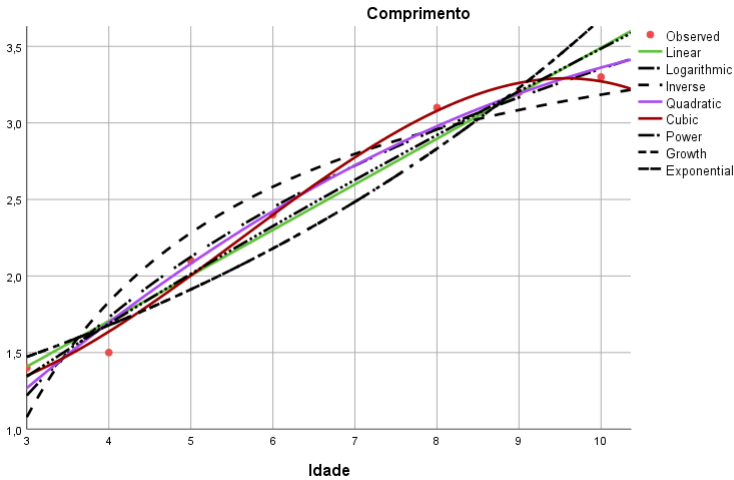
Equation	R Square	F	Model Summary			Parameter Estimates			
			df1	df2	Sig.	Constant	b1	b2	b3
Linear	0,964	132,174	1	5	0,000	0,515	0,298		
Logarithmic	0,971	165,753	1	5	0,000	-0,727	1,772		
Inverse	0,915	53,833	1	5	0,001	4,087	-9,026		
Quadratic	0,980	99,685	2	4	0,000	-0,274	0,579	-0,021	
Cubic	0,991	106,896	3	3	0,002	1,471	-0,387	0,141	-0,008
Power	0,968	149,638	1	5	0,000	0,563	0,792		
Growth	0,931	67,190	1	5	0,000	-0,006	0,131		
Exponential	0,931	67,190	1	5	0,000	0,994	0,131		

The independent variable is Idade.

17



EXEMPLO 2

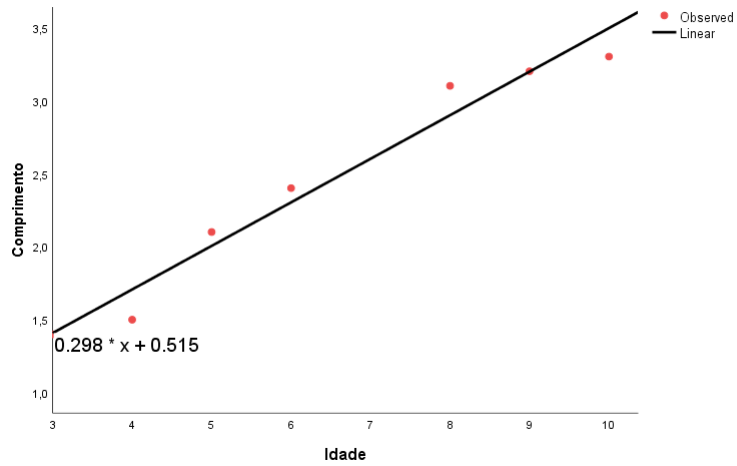


DPS

18



RECTA DE MÍNIMOS CUADRADOS

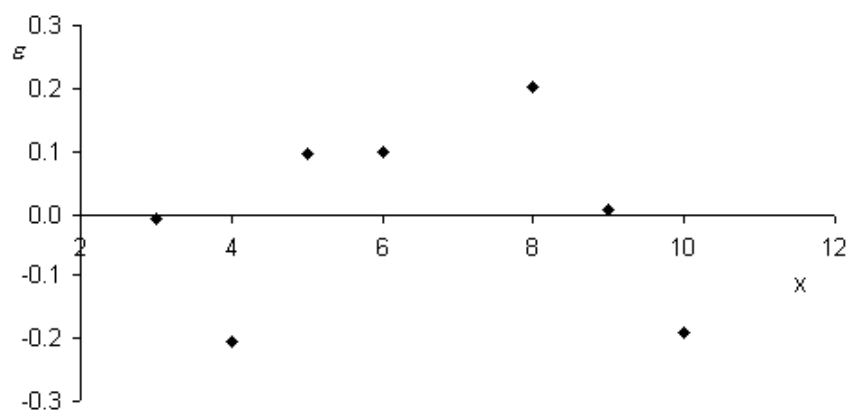


DPS

19



RESÍDUOS



20



Estimadores

$$Y_i = \beta_0 + \beta_1 \cdot (X_i - \bar{X}) + \varepsilon_i \quad i = 1, \dots, n$$

β_0

$$\hat{\beta}_0 = \frac{1}{n} \sum_i Y_i = \bar{Y}$$

β_1

$$\hat{\beta}_1 = \frac{\sum_i (X_i - \bar{X}) \cdot (Y_i - \bar{Y})}{\sum_i (X_i - \bar{X})^2} = \frac{s_{XY}}{s_{XX}}$$

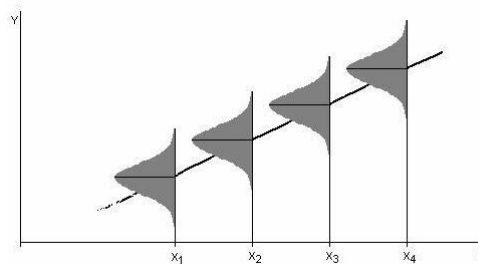
σ^2

$$s^2 = \frac{1}{n-2} \sum_i \hat{\varepsilon}_i^2 = \frac{1}{n-2} \sum_i \left\{ Y_i - \left[\hat{\beta}_0 + \hat{\beta}_1 \cdot (X_i - \bar{X}) \right] \right\}^2$$

21

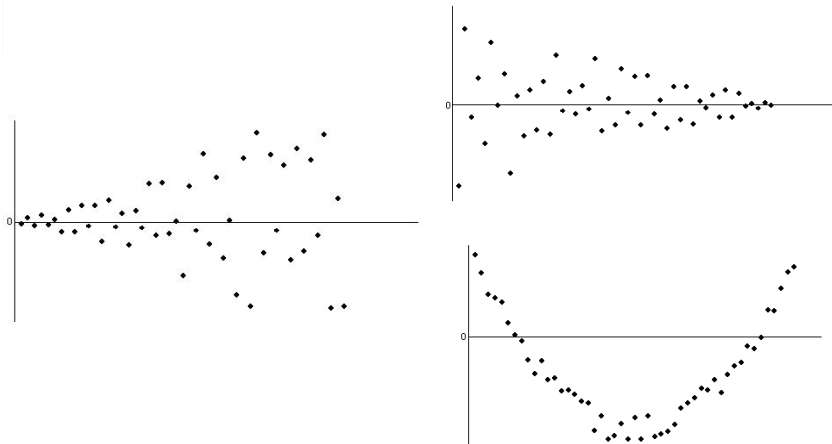


DISTRIBUIÇÃO DOS ERROS



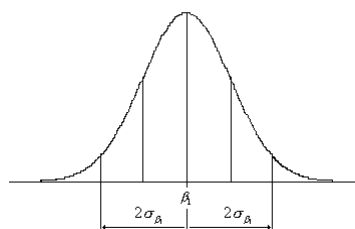
22

RESÍDUOS



23

DISTRIBUIÇÃO DO DECLIVE



24



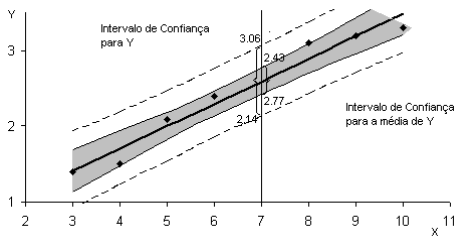
IC e Testes de hipóteses

	IC	TH
β_0	$\hat{\beta}_0 \pm t_{n-2,(\alpha/2)} \cdot \frac{s}{\sqrt{n}}$	$H_0 : \beta_0 = b_0$ $H_1 : \beta_0 \neq b_0, \beta_0 > b_0 \text{ ou } \beta_0 < b_0$ $ET = \frac{\hat{\beta}_0 - b_0}{s/\sqrt{n}}$ $H_0 \text{ verdadeira} \Rightarrow ET \sim t_{n-2}$
β_0'	$(\hat{\beta}_0 - \bar{X} \cdot \hat{\beta}_1) \pm t_{n-2,(\alpha/2)} \cdot s \cdot \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{s_{XX}}}$	$H_0 : \beta_0' = b_0'$ $H_1 : \beta_0' \neq b_0', \beta_0' > b_0' \text{ ou } \beta_0' < b_0'$ $ET = \frac{(\hat{\beta}_0 - \bar{X} \cdot \hat{\beta}_1) - b_0'}{s \cdot \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{s_{XX}}}}$ $H_0 \text{ verdadeira} \Rightarrow ET \sim t_{n-2}$
β_1	$\hat{\beta}_1 \pm t_{n-2,(\alpha/2)} \cdot \frac{s}{\sqrt{s_{XX}}}$	$H_0 : \beta_1 = b_{10}$ $H_1 : \beta_1 \neq b_{10}, \beta_1 > b_{10} \text{ ou } \beta_1 < b_{10}$ $ET = \frac{\hat{\beta}_1 - b_{10}}{\frac{s}{\sqrt{\sum_i (x_i - \bar{x})^2}}}$ $H_0 \text{ verdadeira} \Rightarrow ET \sim t_{n-2}$

25



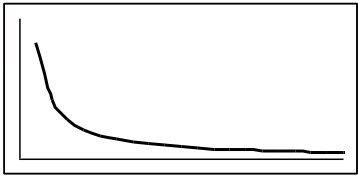
INTERVALO DE CONFIANÇA



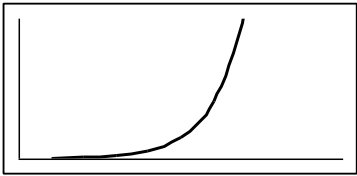
26



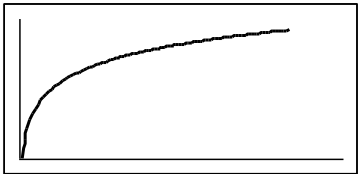
REGRESSÃO NÃO LINEAR



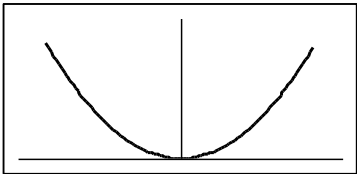
$$\hat{Y} = \beta_0 + \beta_1 \frac{1}{X}$$



$$\hat{Y} = \beta_0 + \beta_1 e^X$$



$$\hat{Y} = \beta_0 + \beta_1 \ln X$$



$$\hat{Y} = \beta_0 + \beta_1 X^2$$

27



REGRESSÃO NÃO LINEAR

Modelo	Transformação
<ul style="list-style-type: none">$Y_i = \alpha' + \frac{\beta}{X_i} + e_i$	$U_i = \frac{1}{X_i}$ $Y_i = \alpha' + \beta.U_i + e_i$
<ul style="list-style-type: none">$Y_i = e^{\alpha' + \beta.X_i + e_i}$	$Z_i = \ln Y_i$ $Z_i = \alpha' + \beta.X_i + e_i$
<ul style="list-style-type: none">$Y_i = e^{\alpha' + \frac{\beta}{X_i} + e_i}$ com $\alpha' > 0, \beta < 0$	$U_i = \frac{1}{X_i}$ $Z_i = \ln Y_i$ $Z_i = \alpha' + \beta.U_i + e_i$

28



REGRESSÃO LINEAR E MÚLTIPLA

Um modelo de regressão linear múltipla descreve uma relação entre várias variáveis quantitativas **independentes**, X_1, X_2, \dots, X_J , e uma variável quantitativa **dependente**, Y , nos termos seguintes:

$$Y_i = \beta_0 + \beta_1 \cdot (X_{1i} - \bar{X}_1) + \beta_2 \cdot (X_{2i} - \bar{X}_2) + \dots + \beta_J \cdot (X_{ji} - \bar{X}_j) + \varepsilon_i \quad \begin{matrix} i = 1, \dots, n \\ j = 1, \dots, J \end{matrix}$$

onde:

- $(X_{1i}, X_{2i}, \dots, X_{ji}, Y_i)$ i-ésima observação das variáveis $X_{1i}, X_{2i}, \dots, X_{ji}$ e Y .
- \bar{X}_j média aritmética das observações X_{ji}
- $\beta_0, \beta_1, \beta_2, \dots, \beta_J$ parâmetros fixos da relação linear entre $X_{1i}, X_{2i}, \dots, X_{ji}$ e Y
- ε_i erro aleatório associado ao valor observado Y_i



RESÍDUOS

Pressupostos para ε_i

- Têm valor esperado nulo e variância constante, σ^2 ;
 - São mutuamente independentes;
 - São normalmente distribuídos.
- $$\left. \begin{matrix} \bullet \text{ Têm valor esperado nulo e variância constante, } \sigma^2; \\ \bullet \text{ São mutuamente independentes;} \\ \bullet \text{ São normalmente distribuídos.} \end{matrix} \right\} \varepsilon_i \sim IN(0, \sigma^2)$$

Se estas hipóteses se verificarem então: $Y_i \sim IN(\mu_{Y_i}, \sigma^2)$

ESTIMADORES MÍNIMOS QUADRADOS



$$\hat{\beta}_0 = \frac{1}{n} \sum_i Y_i = \bar{Y}$$

$$\begin{cases} \hat{\beta}_1.S_{X_1X_1} + \hat{\beta}_2.S_{X_1X_2} + \dots + \hat{\beta}_J.S_{X_1X_J} = S_{X_1Y} \\ \hat{\beta}_1.S_{X_2X_1} + \hat{\beta}_2.S_{X_2X_2} + \dots + \hat{\beta}_J.S_{X_2X_J} = S_{X_2Y} \\ (...) \\ \hat{\beta}_1.S_{X_JX_1} + \hat{\beta}_2.S_{X_JX_2} + \dots + \hat{\beta}_J.S_{X_JX_J} = S_{X_JY} \end{cases}$$

$$S_{X_{j1}X_{j2}} = \sum_n (X_{ji} - \bar{X}_{j1})(X_{ji} - \bar{X}_{j2})$$

$$S_{X_jY} = \sum_n (X_{ji} - \bar{X}_j)(Y_i - \bar{Y})$$

$$\begin{aligned} s^2 &= \frac{1}{n-J-1} \sum_i \hat{\epsilon}_i^2 = \\ &= \frac{1}{n-J-1} \sum_i \left\{ Y_i - \left[\hat{\beta}_0 + \hat{\beta}_1.(X_{1i} - \bar{X}_1) + \hat{\beta}_2.(X_{2i} - \bar{X}_2) + \dots + \hat{\beta}_J.(X_{Ji} - \bar{X}_J) \right] \right\}^2 \end{aligned}$$

EXEMPLO 3



Determine a relação existente entre o calor envolvido no endurecimento, representado pela variável Y e os pesos de duas substâncias X₁ e X₂, tendo em consideração os seguintes valores obtidos numa experiência:

Y	78.5	74.3	104.3	87.6	95.6	109.2	102.7	72.5	93.1	115.9
X1	7	1	11	11	7	11	3	1	2	21
X2	26	29	59	31	52	55	71	31	54	47

The screenshot shows the IBM SPSS Statistics Data Editor window with the file 'reg_multipla.sav'. The 'Analyze' menu is open, and the path 'Analyze > Regression > Linear...' is highlighted. The data table has 16 rows and 2 columns: 'Y' and an unlabeled column. The values for 'Y' are: 78,50, 74,30, 104,30, 87,60, 95,60, 109,20, 102,70, 72,50, 93,10, 115,90, and then 11 empty rows.

	Y
1	78,50
2	74,30
3	104,30
4	87,60
5	95,60
6	109,20
7	102,70
8	72,50
9	93,10
10	115,90
11	
12	
13	
14	
15	
16	

33

The 'Linear Regression' dialog box is shown. The 'Dependent' variable is 'Y'. The 'Independent(s)' variables are 'X1' and 'X2'. The 'Method' is set to 'Enter'. The 'Statistics...', 'Plots...', 'Save...', 'Options...', 'Style...', and 'Bootstrap...' buttons are visible on the right. The 'OK', 'Paste', 'Reset', 'Cancel', and 'Help' buttons are at the bottom.

34



Regression>Linear>Statistics

Linear Regression: Statistics

Regression Coefficients

☒ Estimates
☒ Confidence intervals
Level(%): 95
☐ Covariance matrix

☒ Model fit
☐ R squared change
☐ Descriptives
☐ Part and partial correlations
☐ Collinearity diagnostics

Residuals

☐ Durbin-Watson
☐ Casewise diagnostics
Outliers outside: 3 standard deviations
☒ All cases

Continue Cancel Help

Regression>Linear>Save

Linear Regression: Save

Predicted Values

☒ Unstandardized
☐ Standardized
☐ Adjusted
☐ S.E. of mean predictions

Residuals

☐ Unstandardized
☒ Standardized
☐ Studentized
☐ Deleted
☐ Studentized deleted

Distances

☐ Mahalanobis
☐ Cook's
☐ Leverage values

Influence Statistics

☐ DfBeta(s)
☐ Standardized DfBeta(s)
☐ DfFit
☐ Standardized DfFit
☐ Covariance ratio

Prediction Intervals

☐ Mean ☐ Individual
Confidence interval: 95 %

Coefficient statistics

☐ Create coefficient statistics
☒ Create a new dataset
Dataset name:
☒ Write a new data file
File:

Export model information to XML file

☒ Include the covariance matrix

Continue Cancel Help

35



Variables Entered/Removed^a

Model	Variables Entered	Variables Removed	Method
1	X2, X1 ^b	.	Enter

a. Dependent Variable: Y
b. All requested variables entered.

Model Summary^b

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	,988 ^a	,977	,970	2,57617

a. Predictors: (Constant), X2, X1
b. Dependent Variable: Y

36



ANOVA (Modelo)

H_0 : O modelo de regressão considerado não serve

Decisão: Como valor $p < 0,05$, rejeita-se a H_0 , pelo que o modelo de regressão considerado é estatisticamente significativo

ANOVA ^a						
Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	1976,924	2	988,462	148,940	,000 ^b
	Residual	46,457	7	6,637		
	Total	2023,381	9			

a. Dependent Variable: Y

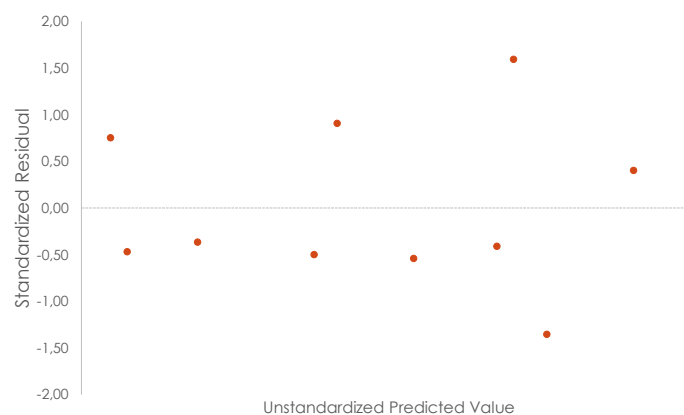
b. Predictors: (Constant), X2, X1

Complementos de Estatística, Prof^a
Ana Cristina Braga

37



RESÍDUOS (homoscedasticidade)



Complementos de Estatística, Prof^a
Ana Cristina Braga

38

Resíduos (Normalidade)



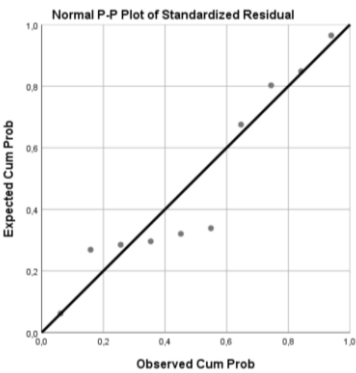
- Teste analítico (KS com correção de Lilliefors)
- Método gráfico (P-P ou Q-Q plot)

NPar Tests

One-Sample Kolmogorov-Smirnov Test

		Standardized Residual
N		10
Normal Parameters ^{a,b}	Mean	,0000000
	Std. Deviation	,88191710
Most Extreme Differences	Absolute	,261
	Positive	,261
	Negative	-,169
Test Statistic		,261
Asymp. Sig. (2-tailed)		,051 ^c

- a. Test distribution is Normal.
- b. Calculated from data.
- c. Lilliefors Significance Correction.



Complementos de Estatística, Prof^a
Ana Cristina Braga

39

Resíduos (média zero)



T-Test

One-Sample Statistics

	N	Mean	Std. Deviation	Std. Error Mean
Standardized Residual	10	,0000000	,88191710	,27888668

One-Sample Test

Test Value = 0

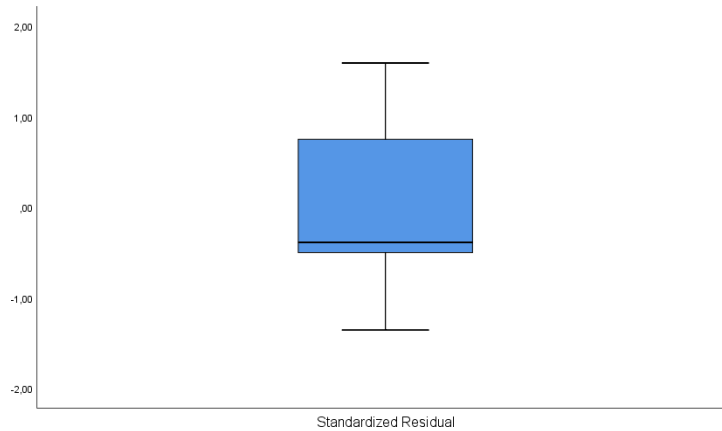
	t	df	Sig. (2-tailed)	Mean Difference	95% Confidence Interval of the Difference	
					Lower	Upper
Standardized Residual	1,19E-014	9	1,000	3,32E-015	-,631	,631

Complementos de Estatística, Prof^a
Ana Cristina Braga

40



Verificação de outliers



Complementos de Estatística, Profª
Ana Cristina Braga

41



COEFICIENTE DE CORRELAÇÃO

Coeficiente de correlação de Pearson

$$R = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}} = \frac{s_{XY}}{\sqrt{s_{xx}} \cdot \sqrt{s_{YY}}}$$

42



TESTES DE ASSOCIAÇÃO

Unilateral à direita

$$H_0 : \rho = 0$$

$$H_1 : \rho > 0$$

Unilateral à esquerda

$$H_0 : \rho = 0$$

$$H_1 : \rho < 0$$

Bilateral

$$H_0 : \rho = 0$$

$$H_1 : \rho \neq 0$$

Estatística de teste

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

Região de Rejeição:

$$t > t_{n-2,(\alpha)}$$

$$t < -t_{n-2,(\alpha)}$$

$$|t| > t_{n-2,(\alpha/2)}$$

43

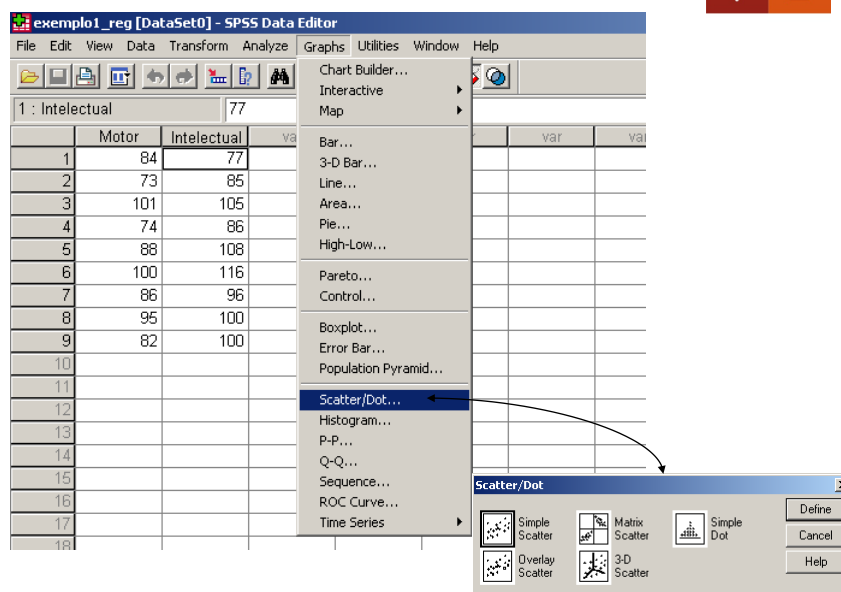


EXEMPLO

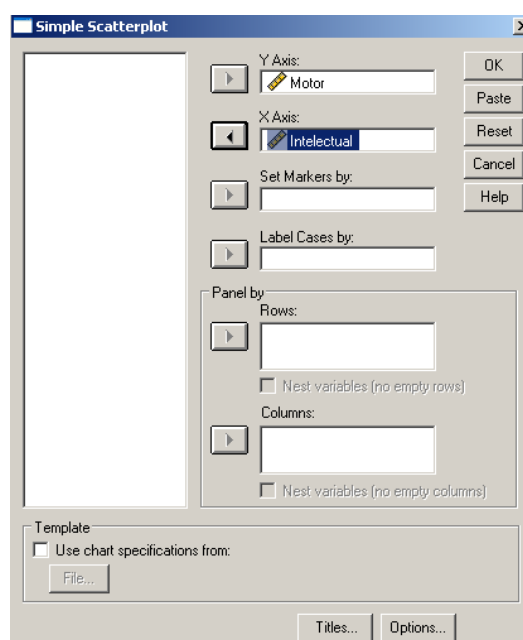
- Índice de Desenvolvimento de Griffiths
- avaliações motora e intelectual para 9 crianças com a idade de 4 anos

Motor	Intelectual
84	77
73	85
101	105
74	86
88	108
100	116
86	96
95	100
82	100

44



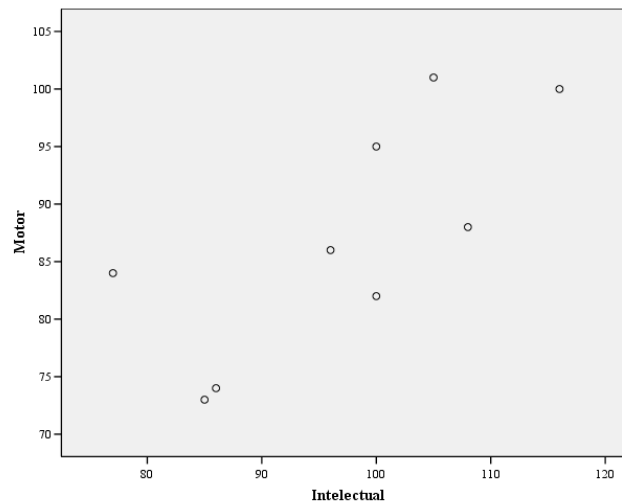
45



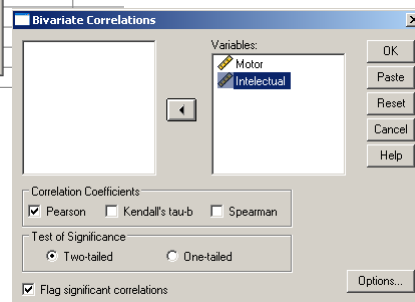
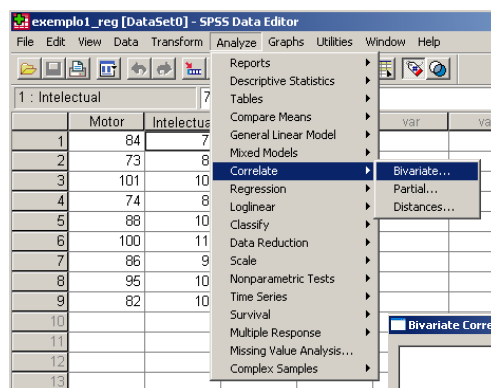
46



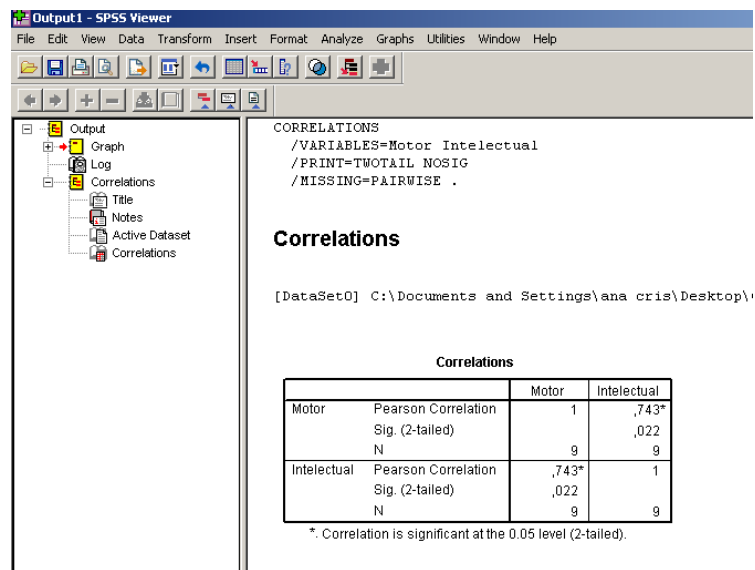
DIAGRAMA DE DISPERSÃO



47

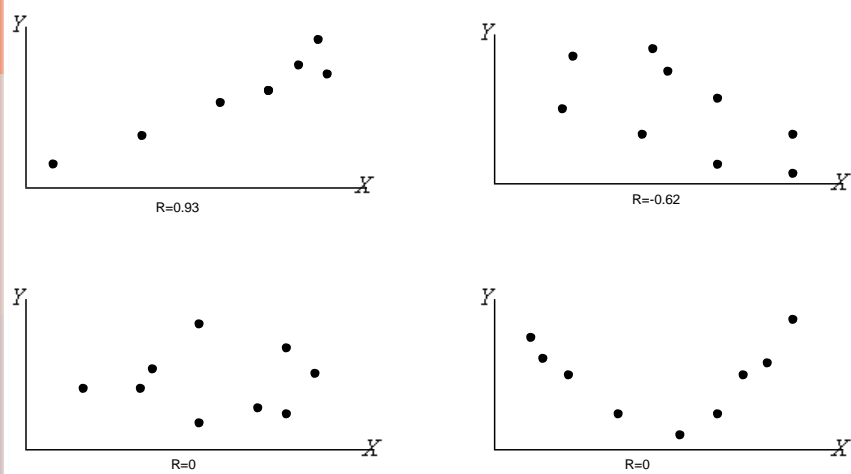


48



49

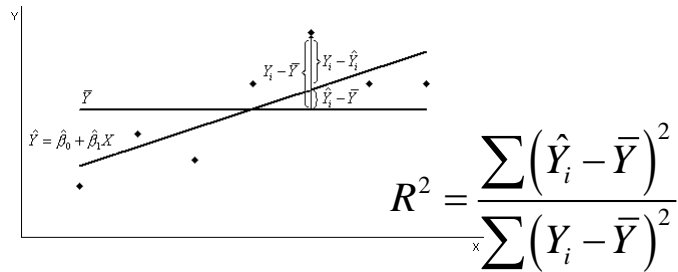
CORRELAÇÃO



50



COEFICIENTE DE DETERMINAÇÃO



51



Coeficiente de determinação (r^2), representa a proporção da variação de Y que é explicada pela regressão

$$r^2 = \frac{\hat{\beta}_1^2 \cdot s_{XX}}{s_{YY}} = \frac{\hat{\beta}_1^2 \cdot \sum_i (X_i - \bar{X})^2}{\sum_i (Y_i - \bar{Y})^2} = \frac{\text{variação de } Y \text{ explicada pela regressão}}{\text{variação total de } Y}$$

52