

\*

Unusual )

(Multiple Regression)

(Observation

(Outlier)

(studentized residuals)

(Boxplot)

(leverage observation)

(Influence)

(Hat matrix)

.(dfbeta)

(Multicollinearity)

.(Ridge Regression)

.(SAS.9)

## **Types of Unusual Observations in Multiple Regression & Some Methods of it's Diagnostic with Application**

### **ABSTRACT**

The work in this paper is a diagnosis of three types of unusual observations in multiple regression, the outlier observation is diagnostic by using Boxplot & studentized

---

تاريخ التسلم : 2008/ 3/ 24 تاريخ القبول : 2008/ 7/ 6

\*

residuals, Diagnostic leverage points by using (Hat matrix) , & diagnostic of the influence observations by using (dfbeta).

In practice the work is comparing the effects of omitting outlier observations in the normal distribution of the residuals to the equation which is building to the Thalassaemia disease and Treating the multicollinearity by omitting some variables and using ridge regression, getting a good model agrees with the viewing of medicine by using (SAS.9) package.

-:

(1774) (Bernoulli)

.(Evans, 1999)

...

(Bacon,1995)

(Outlier)

(Leverage)

(Influence)

(Hoagolin & Welsch,1978)

.(Chatterjee,et al.,1986)

-:

-:

**-1**

(Unusual Observation)

( )

(Explanatory Variables)

Y

.(Beckman & Cook,1983)

Simple )

(Sinha,1997)

(Regression

X

(Scatter plot)

(Residuals)

Y

(Dependent Variable)

( )

( )

X-Y

(Boxplot)

(normal probability plot)

high-leverag )

(observation

Ordinary Least )

(OLS)(Square

(bias the result)

-:

...

-(Outlier) -a

(Residual)

.(Barnett & Lewis,1994)

.(studentized residuals)

-(Boxplot) 1-a

Box )

.(and Whisker Plot

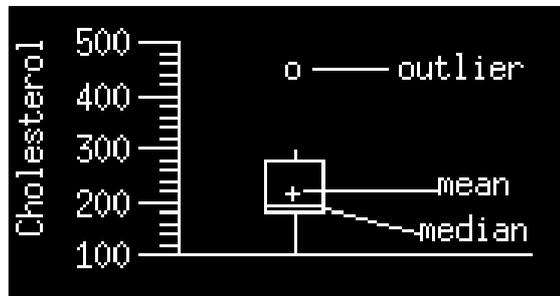
(Pattern)

(1)

(+)

.(o)

:(1)



(+)

.(skewed)

heavy-)

.( Michael,et al(1989)) (tailed

-:(studentized residuals) 2-a

\*(Standardized Residuals)

-(i)

$$e_i^s = \frac{e_i}{\sqrt{\hat{s}_{(i)}^2(1-h_i)}} \dots(1)$$

-:

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{v^2} \dots(2)$$

$S_e$  (\*)

1

$$v^2 = \sum_{i=1}^n (x_i - \bar{x})^2 \dots(3)$$

-:

-:  $e_i^s$

-(i)  $(y_i - \hat{y}_i)$  -:  $e_i$

-(i) -:  $\hat{s}_{(i)}$

-: n

-(i) -:  $x_i$

-:  $\bar{x}$

Hair et al.,(1998) ) ( $e_i^s > \pm 2$ )

.(Iglewicz, B. y Hoaglin, D.C. (1993);

-(Leverage observation ) -b  
 (Lever)  
 .(Leverage)

.(Hat Matrix)

-: Hat Matrix 1-b

-(Hat Matrix)

$$\text{Hat Matrix} = X(X'X)^{-1}X' \dots(4)$$

$$H = X(X'X)^{-1}X' \dots(5)$$

( h<sub>ii</sub> )

. i

n k ((2k+2)/n)

.(Iglewicz, B. y Hoaglin, D.C. (1993))

-:(Influence) -c

(Y)

( )

.(Hair,et al.,1998;Tsay et al.,2000;Hordo,2004)

.(DFBETA )

$$\frac{-(Dfbetas) \quad 1-c}{Dfbetas}$$

(Hair et al.,1998)

(i)

$$DFBETAS_{j,i} = \frac{\hat{B}_j - \hat{B}_{j(i)}}{s_{(i)}^2 \sqrt{(X'X)^{-1}_{jj}}} \dots (6) \quad -:$$

-:

$$\dots (i) \quad -: \hat{B}_J$$

$$\dots (i) \quad -: \hat{B}_{j(i)}$$

$$\dots (i) \quad -: s_{(i)}^2$$

$$\dots (X'X)^{-1} \quad -: (X'X)^{-1}_{jj}$$

$$|Dfbetas_{J,I}| > \frac{2}{\sqrt{n}}$$

-: -2

(Y)

(2005 )

$$Y = a_0 + a_1X_1 + a_2X_2 + a_3X_3 + a_4X_4 + a_5X_5 + a_6X_6 + a_7X_7 + a_8X_8 + a_9X_9 + a_{10}X_{10} + e \dots (7)$$

$$\dots ( ) \quad -: X_1 \quad -:$$

$$\dots ( ) \quad -: X_2$$

$$\dots -: X_4 \quad \dots ( ) \quad -: X_3$$

$$\dots -: X_6 \quad \dots ( ) \quad -: X_5$$

$$\dots -: X_8 \quad \dots \quad -: X_7$$

$$\dots ( ) \quad -: X_{10} \quad \dots \quad -: X_9$$

...

(7)

(2005 )

(2005 )

-:

$$Y = -0.00137 + 0.789X_1 - 0.126X_4 + 0.0162X_7 \dots (8)$$

(X4)

(OLS)

(OLS)

(2)

(scatter plot matrix)

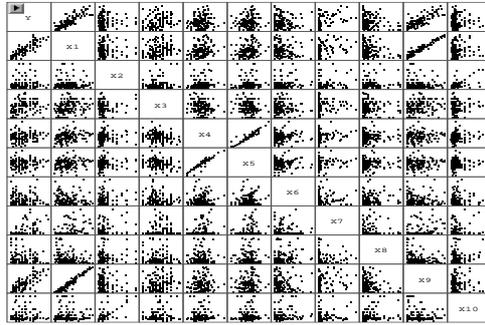
)

((a1-a150)

(1-150)

(1)

(2)



)

(

-:

:

:

(W)(Shapiro-Wilk)

(W)

(1)

(5%)

)

(3)

(Kernal density plot) (( )

(\*)**(1)**

Variable	DF	Parameter Estimate	Standard Error	t Value	Pr >  t	Variance
						Inflation
						Factor
Intercept	1	5.13662	6.0365	0.85	0.3963	0
X1	1	2.40727	0.21524	11.18	<.0001	41.35083
X2	1	-0.3269	0.24222	-1.35	0.1793	2.72755
X3	1	0.64232	0.59046	1.09	0.2786	1.23351
X4	1	-13.42824	2.56653	-5.23	<.0001	14.96864
X5	1	1.76712	0.70535	2.51	0.0134	11.70736
X6	1	0.16335	0.37435	0.44	0.6633	1.15793
X7	1	0.02645	0.01253	2.11	0.0366	1.33558
X8	1	-0.04441	0.06963	-0.64	0.5247	1.2398
X9	1	-1.28947	0.17118	-7.53	<.0001	37.85468
X10	1	-0.07048	0.2232	-0.32	0.7526	2.98425

Source	DF	Sum of Squares	Analysis of Variance	F Value	Pr > F
Model	10	184752	18475	71.97	<.0001
Error	139	35683	256.71456		
Corrected Total	149	220436			
Root MSE		16.02231	R-square	<b>0.8381</b>	
Dependent Mean		<b>63.72667</b>	Adj R-sq	<b>0.8265</b>	

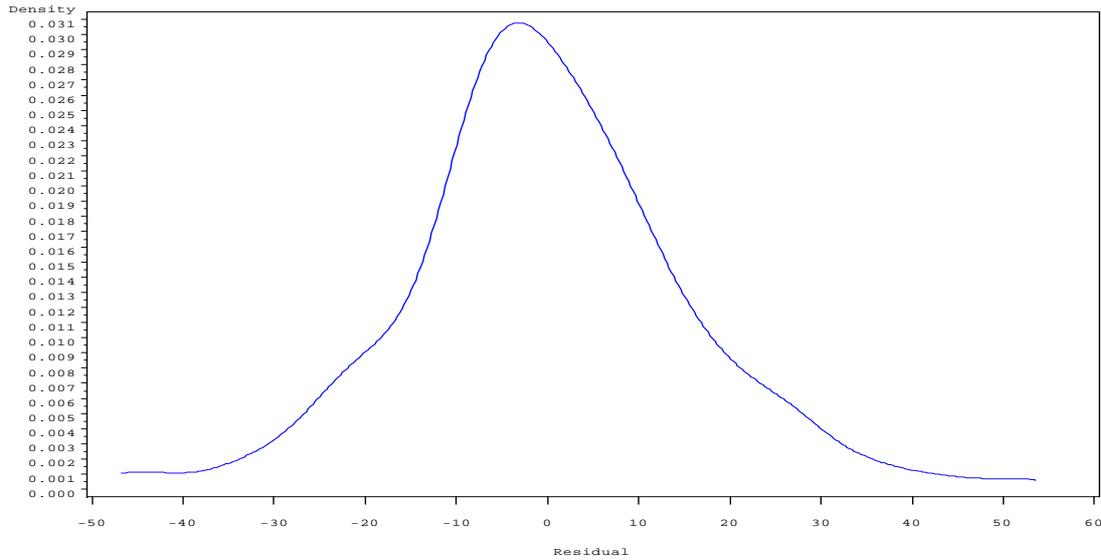
Tests for Normality			
Test		Statistic---	-----p Value-----
Shapiro-Wilk	W	0.981634	Pr < W 0.0426

(1)

(\*)

.(2005 )

(\*) (3)  
 ( )  
 (Kernal density plot)



-: 1-2

.(studentized residuals)

(4 ) (stem & leaf)

(± 2)

(a127) (a142) (a96))

(2)

((a3) (a99) (a75) (a104) (a109)

.(a96)

(o)



:(2)

Highest		Lowest	
Obs	Value	Obs	Value
40	1.90524	142	- 3.1265
75	2.30187	127	- 3.0108
104	2.42619	99	- 2.2285
109	2.9294	3	- 2.1901
96	3.53914	132	- 1.8775

2-2 تشخيص المشاهدات ذات القوة الرافعة (Leverage):-

$$(k) \quad (n) \quad ((2k+2)/n)$$

$$(0.1467)$$

$$(0.1467) \quad (3)$$

$$(0.1467)$$

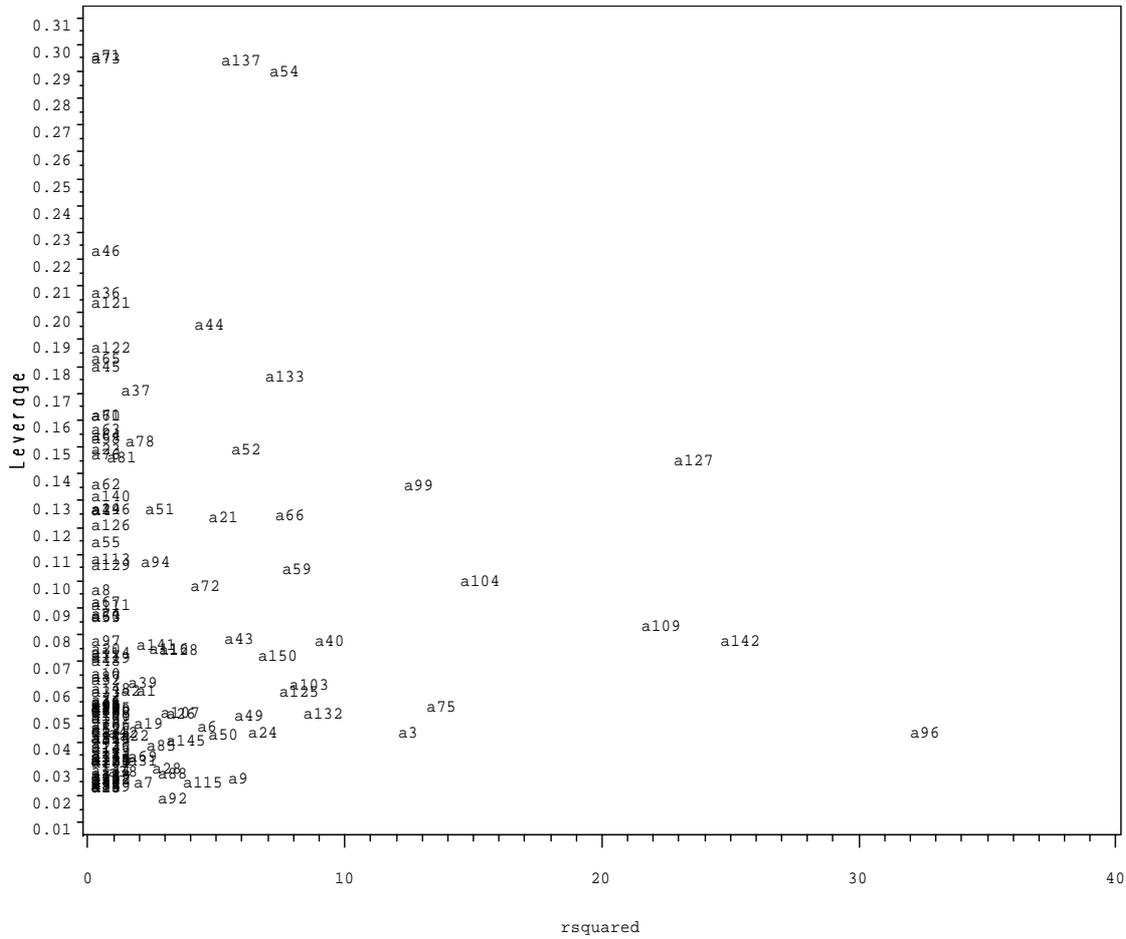
(5)

(a96,a142,a127,a109) (5)  
(a71,a73,a137,a54)

(3)

Leverage	name
0.30768	a71
0.30598	a137
0.30174	a54
0.28871	a73
0.21695	a46
0.20131	a36
0.19795	a121
0.18931	a44
0.18093	a122
0.17669	a65
0.17384	a45
0.17043	a133
0.16533	a37
0.15583	a70
0.15563	a61
0.1504	a63
0.14813	a64

(5)



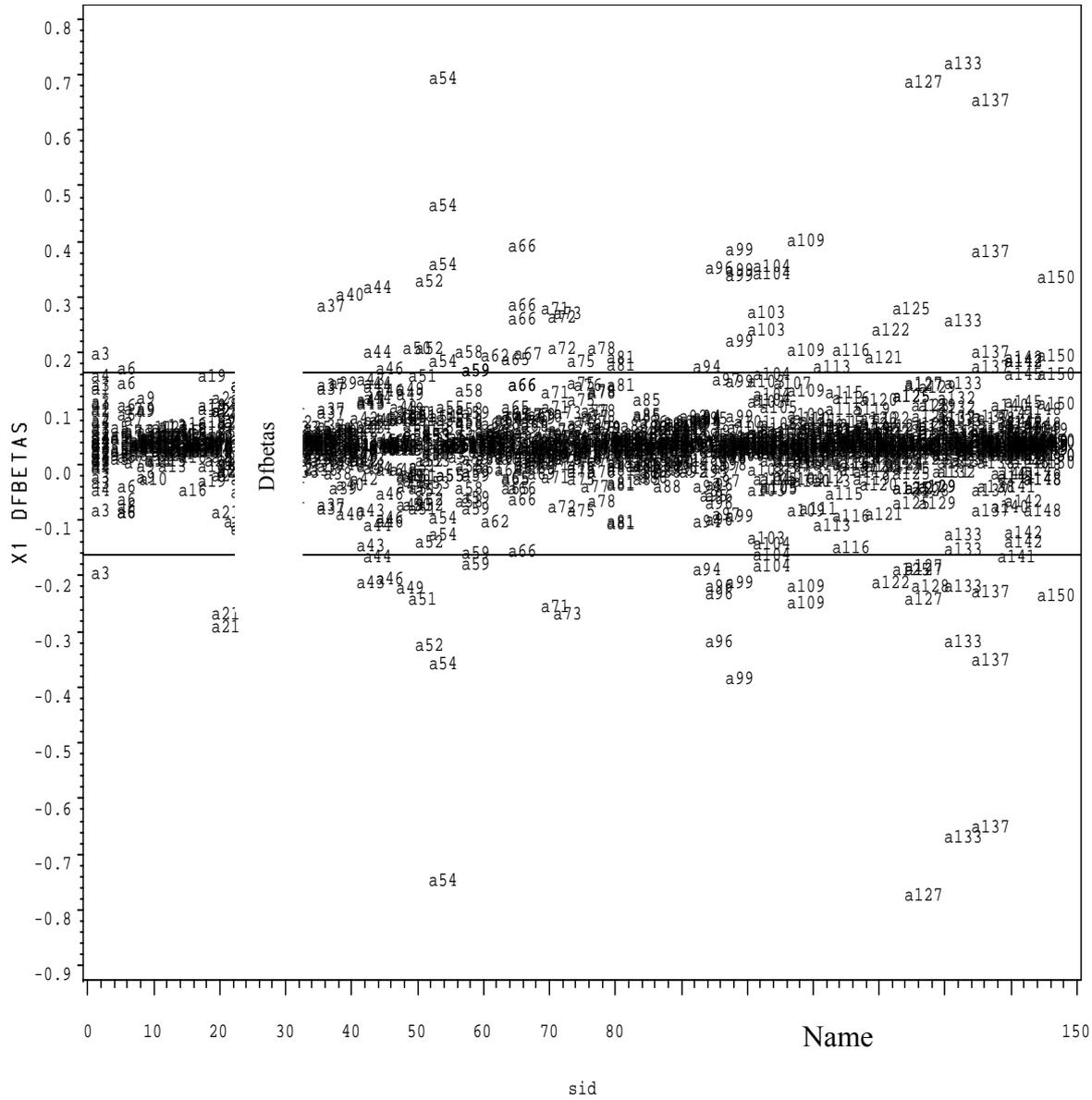
...

**3-2-تشخيص المشاهدات المؤثرة (influence):-**

	Dfbeta	(Dfbeta)
(4)	.(	
	(a127,a54,a137,a133)	(6)
	(-0.808,-0.7813,0.7212,0.701)	(Dfbeta)
		(a127)
		X <sub>10</sub>
	(127)	
t	X <sub>10</sub>	
	.(0.47)	(-0.32)

(6)

(DFbeta)







4-2 اثر حذف المشاهدات الشاردة فى معاملات الانحدار فى التوزيع الطبيعى:-

(5)

((a109) (a127) (a142) (a96))

(a96)

(W)

(f)

(MSE)

(VIF)

 $X_{10}$ 

5-2 الكشف عن مشكلة تعدد العلاقات الخطية ومعالجتها:-

 $(X_1, X_4, X_5, X_9)$ 

(6)

(eigen value)

(30)

(condition index)

variance )

(0.00624)

(VIF)

(0.5)

(proportion

( (7) )

(0.76%)







$$C \quad ) (C)^* \quad (7)$$

(0.005)      0.1

(VIF)      C

(RMSE) (Root Mean Square Error)

$$(X_1) \quad (C=0.02) \quad (X_4, X_5, X_9)$$

(C=0.08)

(8)

$$(5\%) \quad (10) \quad (VIF) \quad C=0.08$$

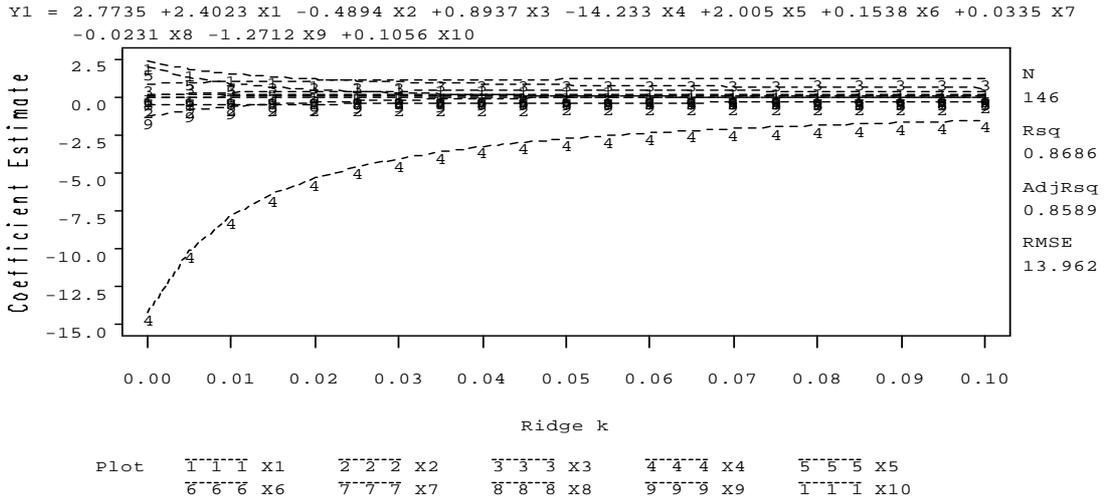
(X<sub>3</sub>, X<sub>7</sub>, X<sub>9</sub>, X<sub>10</sub>)

-:

$$y=7.97258+1.25253X_3+0.05396X_7+0.07576X_9+0.42331X_{10}$$

(7)

### the Ridge Regression



$(0 < C < 1)$

(\*)

$(X'X)$

(OLS) (Ordinary least square)

(C=0.08)

(9 8)

( )

(RMSE)

(8 )

(9 )

...

-----  
-:-----

-:

.1

)

.2

(X<sub>7</sub>)

( X<sub>3</sub>)( )

(

)( )

(X<sub>9</sub>)

.(X<sub>10</sub>)

(1.25243)

(0.07576)

(0.42331)

.(0.05396)

.3

(17.8257) (RMSE)

(19.5796)

.4

-:\_\_\_\_\_

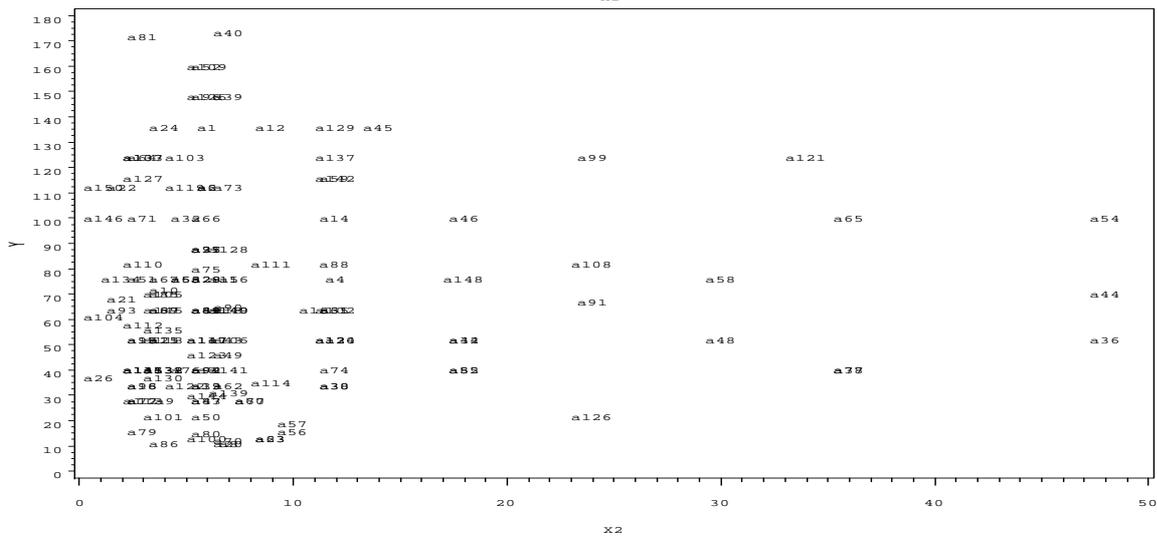
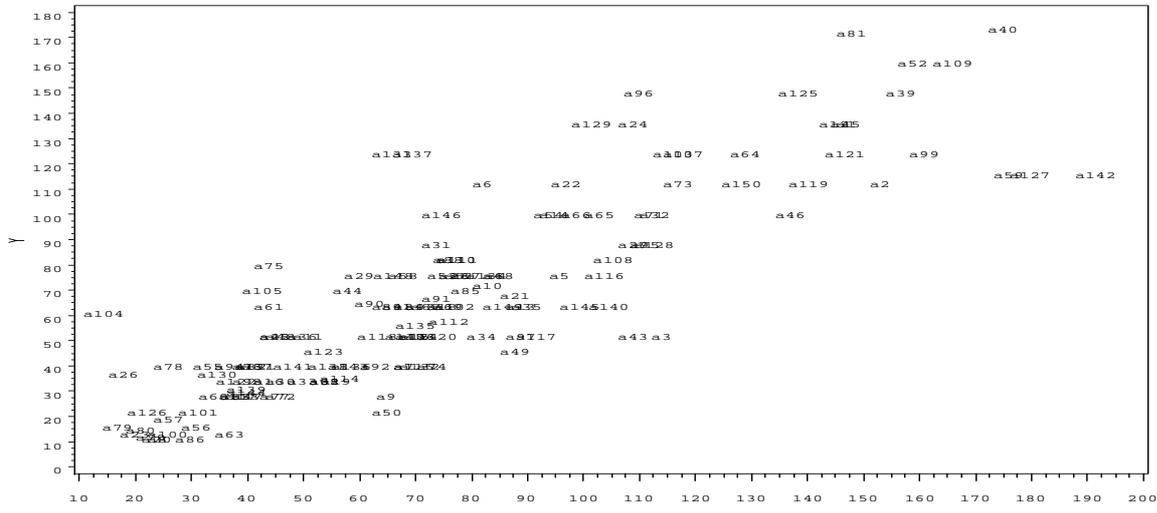
":(2005)

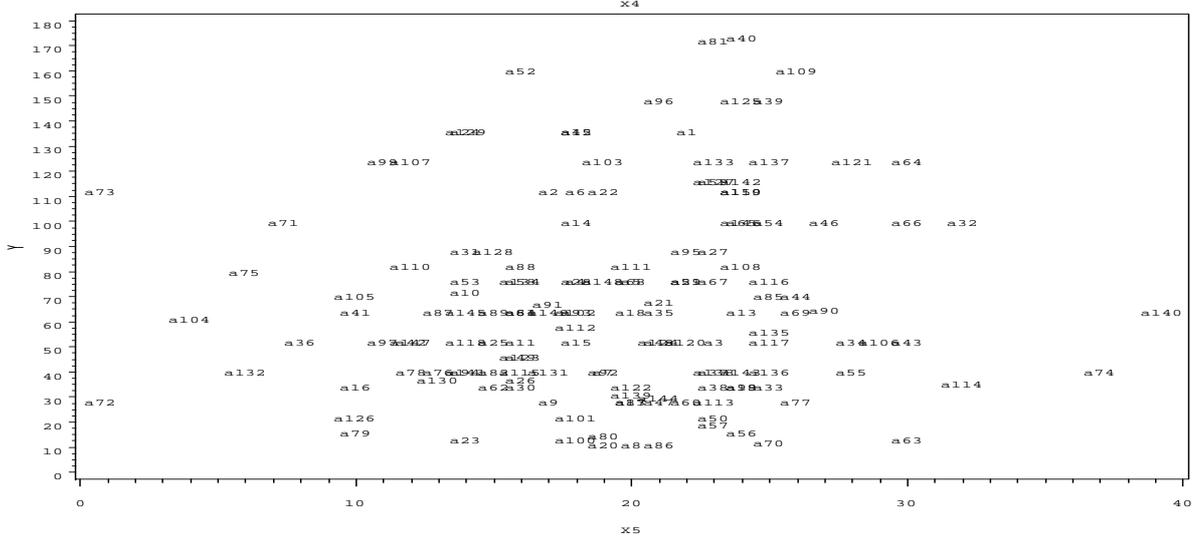
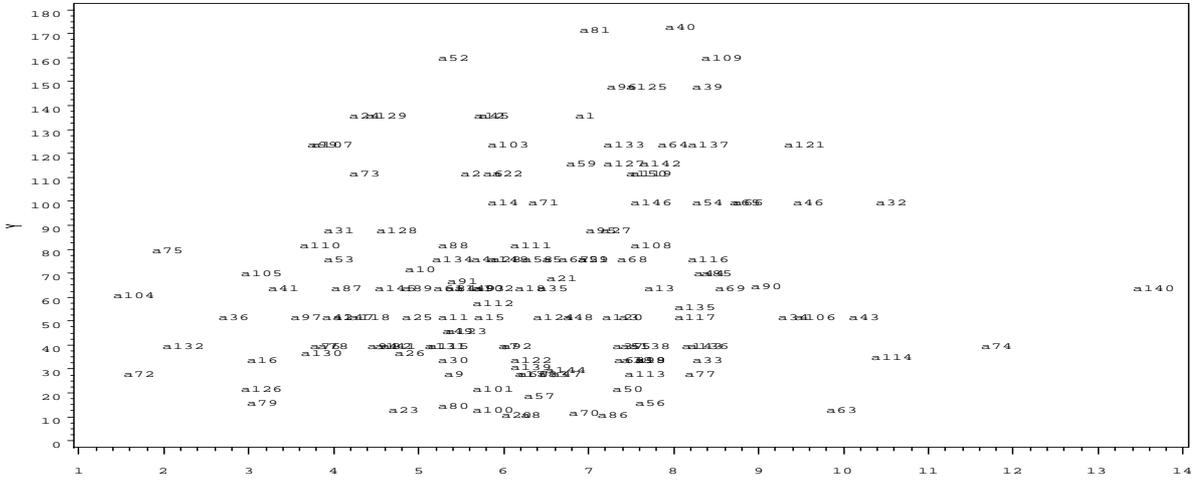
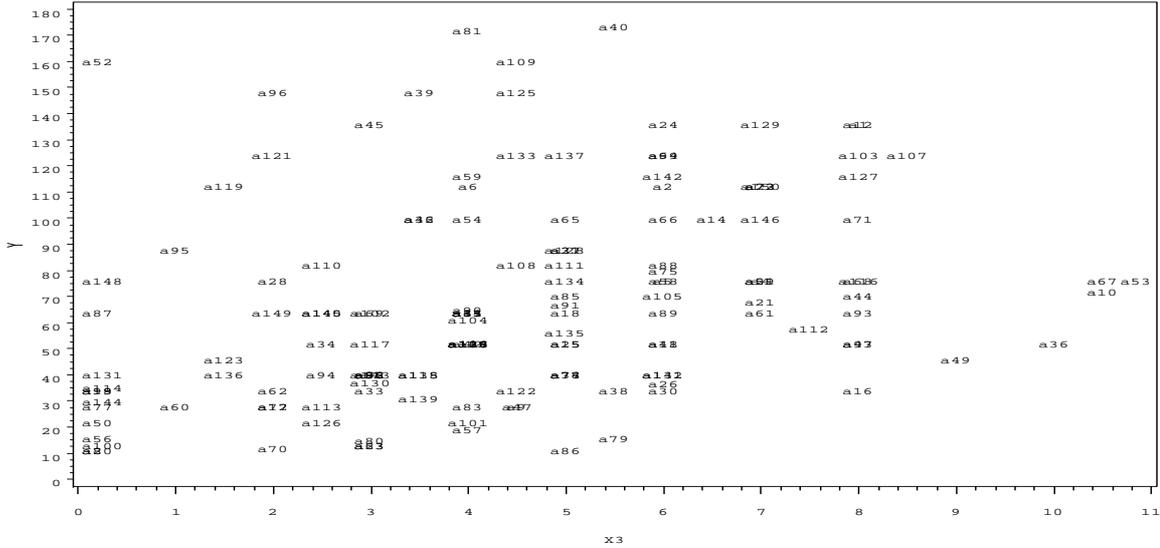
.1

"

2. Bacon,D. (1995).”A maximum likelihood approach to correlational outlier identification”.Multivariate Behavioral Research,30,125-148.
3. Beckman,R.J. & Cook, R. D. (1983),”Outlier ...s”,Technometrics,25,119-163.
4. Chatterjee, S. & H. (1986).”Influential observations,High leverage points,and outlier in linear regression”,Statistical Science,1(3),379-393.
5. Evans,Victoria(1999).”Strategies for Detecting Outlier in Regression Analysis: An Introductory Primer”Paper presented at the annual meeting of the Southwest Educational Research Association ,San Antonio.
6. Hair, J. F. Jr., Anderson, R. E., Tatham, R. L. and Black, W. C. (1998) “Multivariate Data Analysis”, 5th edn. Prentice-Hall, New Jersey.
7. Hoaglin,D.C. & Welsch, R.E. (1978).”The Hat matrix in regression and ANOVA”.The American Statistician,32,17-22.
8. Hordo, M. (2004) ,”Outlier diagnostics on permanent sample plot network in Estonia”. – Research for Rural Development 2004. International Scientific Conference Proceedings. Latvia University of Agriculture, Jelgava, pp. 181-186.
9. Iglewicz, B. y Hoaglin, D.C. (1993). “How to Detect and Handle Outliers”. Quality Press, American Society for Quality Control, Milwaukee, Wisconsin.
10. Michael Frigge and David C. Hoaglin and Boris Iglewicz.( 1989) "Some Implementations of the Boxplot". The American Statistician. Vol. 43 (1),pp 50–54.
11. Sinha,Sanjoy Kumar,(1997)”Sequential Application of Multivariate Outlier Test:A Robust Approach”,Master Thesis,Dalhousie University,Halifax,Nova Scotia,pp.1-8 .

(1)





...

