

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

MODELO DE REGRESIÓN LINEAL MÚLTIPLE

Conjunto de datos: SWISS

Susana Chen, Carmen Martínez, Óscar Mesa y Eva de Vega

26 de marzo de 2023

1.Introducción

2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

1.Introducción

2. Estudio y evaluación del modelo completo

3.Selección del mejor modelo

4.Diagnóstico

5.Calcular el Error de test

6.Conclusión

1.Introducción

2. Estudio y evaluación del modelo completo
3. Selección del mejor modelo
4. Diagnóstico
5. Calcular el Error de test
6. Conclusión

1.INTRODUCCIÓN

1.Introducción

2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

1.Introducción

```
> head(conjunto_swiss)
```

	Fertility	Agriculture	Examination	Education	Catholic	Infant.Mortality
Courtellary	927.112	197	173	139	115	257
Delemont	13615.935	7390	983	1475	13901	3637
Franches-Mnt	7919.850	3399	428	428	7997	1730
Moutier	1990.560	847	278	162	783	471
Neuveville	1814.840	1027	401	354	122	486
Porrentruy	4906.928	2276	580	451	5840	1715

1.1. Descripción de las variables

- **Fertility:** Es la media estandarizada común de fertilidad.

```
> typeof(Fertility)
[1] "double"
```

- **Agriculture:** Número de hombres agricultores.

```
> typeof(Agriculture)
[1] "double"
```

- **Examination:** Número de reclutas que reciben la calificación más alta en el examen del ejército.

```
> typeof(Examination)
[1] "double"
```

- **Education:** Número de habitantes con estudios superiores.

```
> typeof(Education)
[1] "double"
```

1.Introducción

2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

- **Catholic:** Número de católicos.

```
> typeof(Catholic)
[1] "double"
```

- **Infant.Mortality:** Niños que viven menos de un año.

```
> typeof(Infant)
[1] "double"
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

2. ESTUDIO Y EVALUACIÓN DEL MODELO COMPLETO

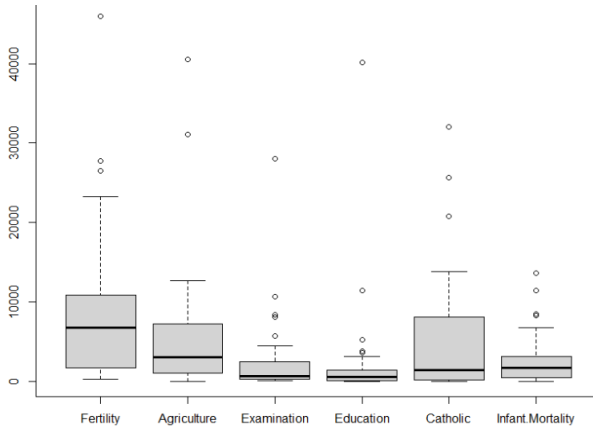
- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

2.Estudio y evaluación del modelo completo

```
> summary(conjunto_swiss)
```

Fertility	Agriculture	Examination	Education	Catholic	Infant.Mortality
Min. : 308.4	Min. : 86	Min. : 173	Min. : 21	Min. : 12	Min. : 61.0
1st Qu.: 1755.0	1st Qu.: 1117	1st Qu.: 330	1st Qu.: 175	1st Qu.: 201	1st Qu.: 478.5
Median : 6763.7	Median : 3095	Median : 708	Median : 586	Median : 1463	Median : 1715.0
Mean : 8386.6	Mean : 5508	Mean : 2379	Mean : 2008	Mean : 5130	Mean : 2493.6
3rd Qu.: 10927.5	3rd Qu.: 7287	3rd Qu.: 2516	3rd Qu.: 1484	3rd Qu.: 8168	3rd Qu.: 3152.0
Max. : 45908.6	Max. : 40489	Max. : 28012	Max. : 40126	Max. : 32055	Max. : 13628.0

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión



1.Introducción

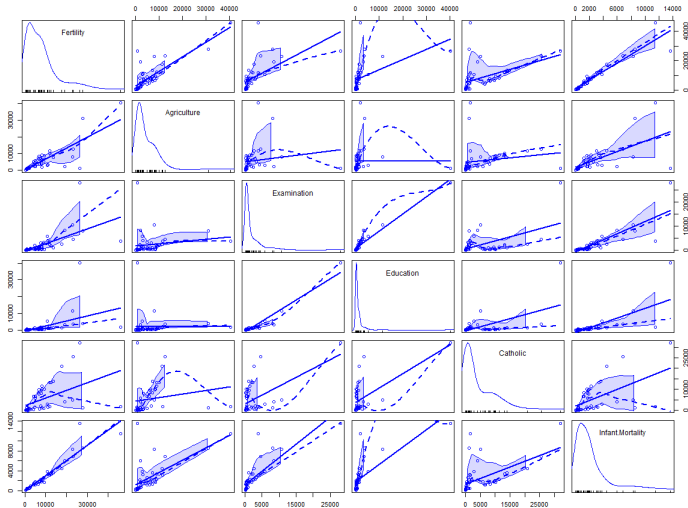
2. Estudio y evaluación del modelo completo

3.Selección del mejor modelo

4.Diagnóstico

5.Calcular el Error de test

6.Conclusión



- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3. SELECCIÓN DEL MEJOR MODELO

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3.Selección del mejor modelo

Definimos las variables explicativas y la variable respuesta:

```
> Fertility<-conjunto_swiss$Fertility  
> Agriculture<-conjunto_swiss$Agriculture  
> Examination<-conjunto_swiss$Examination  
> Education<-conjunto_swiss$Education  
> Catholic<-conjunto_swiss$Catholic  
> Infant<-conjunto_swiss$Infant.Mortality
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3.Selección del mejor modelo

En forma matricial nos queda la siguiente ecuación de regresión:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} \text{ con } i \in \{1, \dots, 47\}$$
$$Y_{47 \times 1} = X_{47 \times 6} \beta_{6 \times 1} + \varepsilon_{47 \times 1}$$

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3.Selección del mejor modelo

Nuestro modelo y sus coeficientes son:

```
> model<-lm(Fertility~Agriculture+Examination+Education+Catholic+Infant, data=conjunto_swiss)
> model$coefficients
```

(Intercept)	Agriculture	Examination	Education	Catholic	Infant
117.18755122	0.06906232	-0.44990220	-0.35570581	0.01914249	3.83994829

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3.Selección del mejor modelo

Realizamos el siguiente contraste de hipótesis:

$$H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4 = \beta_5 = 0$$

$$H_1 : \beta_i \neq 0$$

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3.Selección del mejor modelo

```
> summary(model)
```

Call:

```
lm(formula = Fertility ~ Agriculture + Examination + Education +  
    Catholic + Infant, data = conjunto_swiss)
```

Residuals:

Min	1Q	Median	3Q	Max
-2269.2	-476.5	-67.7	406.1	2346.0

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	117.18755	195.25119	0.600	0.551684
Agriculture	0.06906	0.04651	1.485	0.145189
Examination	-0.44990	0.17573	-2.560	0.014244 *
Education	-0.35571	0.09509	-3.741	0.000562 ***
Catholic	0.01914	0.02647	0.723	0.473683
Infant	3.83995	0.21926	17.513	< 2e-16 ***

signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 902 on 41 degrees of freedom

Multiple R-squared: 0.9911, Adjusted R-squared: 0.99

F-statistic: 915.4 on 5 and 41 DF, p-value: < 2.2e-16

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3.1.Métodos basados en pruebas

3.1.1) Método Forward

En cada iteración, vamos añadiendo la **variable explicativa** que cumpla:

$$p\text{-valor} < \alpha_{critic} = 0,05$$

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

Primera Iteración

Infant.Mortality: $Pr(> t) < 2 \cdot 10^{-16} < \alpha_{critic} = 0,05$

	Df	Sum of Sq	RSS	AIC	F value	Pr>F
<none>			3756931492	857,25		
Agriculture	1	2490455546	1266475947	808,14	88,490	$3,406 \cdot 10^{-12}$
Examination	1	1397162087	2359769405	837,39	26,643	$5,359 \cdot 10^{-6}$
Education	1	785611536	2971319956	848,22	11,898	0,001231
Catholic	1	820867437	2936064056	847,66	12,581	0,000924
Infant.Mortality	1	3341249575	415681917	755,78	361,710	$< 2,2 \cdot 10^{-16}$

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

Segunda Iteración

Examination: $Pr(> t) < 2,2 \cdot 10^{-16} < \alpha_{critic} = 0,05$

	Df	Sum of Sq	RSS	AIC	F value	Pr>F
<none>			41568197	755,78		
Agriculture	1	298539564	117142353	698,25	112,1349	$1,109 \cdot 10^{-13}$
Examination	1	367054893	48627024	643,93	332,1284	$< 2,2 \cdot 10^{-16}$
Education	1	365564388	50117529	658,35	320,9423	$< 2,2 \cdot 10^{-16}$
Catholic	1	17595127	398086790	755,75	1,9448	0,1702

Observación

Aunque había otra variable explicativa con, aparentemente el mismo *p-valor* (en este caso era *Education*) nos hemos decantado por *Examination* ya que su *valor F* era más grande y sabemos que cuanto más grande es, más pequeño es el *p-valor*.

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

Tercera Iteración

Education: $Pr(> t) < 0,0002196 < \alpha_{critic} = 0,05$

	Df	Sum of Sq	RSS	AIC	F value	Pr>F
<none>			48627024	656,93		
Agriculture	1	2961590	45665434	655,98	2,7887	0,1021947
Education	1	13358152	35268872	643,83	16,2863	0,0002196
Catholic	1	1941217	46685807	657,01	1,7880	0,1882059

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3.2.Métodos basados en criterios

Los posibles modelos son:

```
> summary(model)
```

```
call:
lm(formula = Fertility ~ Agriculture + Examination + Education +
  Catholic + Infant, data = conjunto_swiss)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-2269.2  -476.5   -67.7    406.1   2346.0
```

```
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	117.18755	195.25119	0.600	0.551684
Agriculture	0.06906	0.04651	1.485	0.145189
Examination	-0.44990	0.17573	-2.560	0.014244 *
Education	-0.35571	0.09509	-3.741	0.000562 ***
Catholic	0.01914	0.02647	0.723	0.473683
Infant	3.83995	0.21926	17.513	< 2e-16 ***

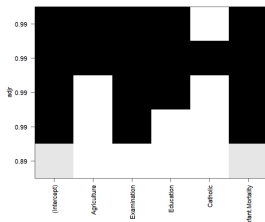
```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 902 on 41 degrees of freedom
Multiple R-squared:  0.9911,    Adjusted R-squared:  0.99
F-statistic: 915.4 on 5 and 41 DF,  p-value: < 2.2e-16
```

- 1 $Fertility \sim Infant$
- 2 $Fertility \sim Examination + Infant$
- 3 $Fertility \sim Examination + Education + Infant$
- 4 $Fertility \sim Agriculture + Education + Examination + Infant$
- 5 $Fertility \sim Agriculture + Examination + Education + Catholic + Infant$

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

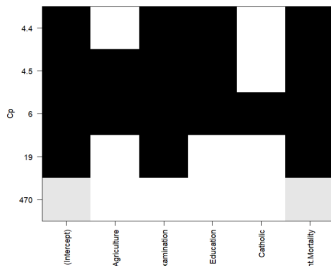
3.2.1.R cuadrado ajustado



```
> MR2adj<-summary(models)$adjr2
> MR2adj
[1] 0.8868972 0.9864684 0.9899572 0.9901522 0.990039
> which.max(MR2adj)
[1] 4
> plot(models,scale="adjr")
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

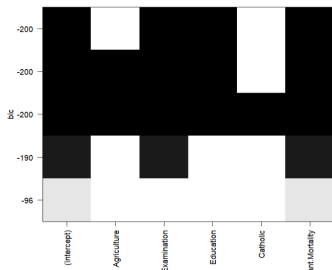
3.2.2.Cp de Mallows



```
> MCp<-summary(models)$cp
> MCp
[1] 467.953565 18.772028 4.353017 4.522777 6.00000
> which.min(MCp)
[1] 3
> plot(models,scale="Cp")
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3.2.3.Criterio de Información de Bayes (BIC)



```
> MBIC<-summary(models)$bic  
> MBIC  
[1] -95.76727 -192.76694 -204.01138 -202.18854 -198.93389  
> which.min(MBIC)  
[1] 3  
> plot(models,scale="bic")  
> install.packages("MASS")
```


- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3.2.4.Criterio de Información de Akaike (AIC)

```
> #AIC
> SCOPE<-(~.)
> stepAIC(modelo_completo, scope=SCOPE,k=2)
Start: AIC=645.21
Fertility ~ Agriculture + Examination + Education + Catholic +
  Infant
```

	Df	Sum of Sq	RSS	AIC
- Catholic	1	425301	33780501	643.81
<none>			33355200	645.21
- Agriculture	1	1793764	35148964	645.67
- Examination	1	5332289	38687490	650.18
- Education	1	11383630	44738830	657.01
- Infant	1	249523972	282879173	743.69

```
Step: AIC=643.81
Fertility ~ Agriculture + Examination + Education + Infant
```

	Df	Sum of Sq	RSS	AIC
<none>			33779724	643.81
- Agriculture	1	1489148	35268872	643.83
+ Catholic	1	425447	33354278	645.21
- Examination	1	8432237	42211961	652.28
- Education	1	11885710	45665434	655.98
- Infant	1	310151705	343931429	750.87

```
Call:
lm(formula = Fertility ~ Agriculture + Examination + Education +
  Infant, data = conjunto_swiss)
```

```
Coefficients:
(Intercept)  Agriculture  Examination  Education  Infant
 169.64188      0.06116     -0.50652     -0.32488     3.90494
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

3. Selección del mejor modelo

Como hemos obtenido dos modelos candidatos a mejor modelo, realizamos un contraste de hipótesis:

```
> modelo_mejor<-lm(Fertility~Examination+Education+Infant)
> modelo_mejor2<-lm(Fertility~Agriculture+Examination+Education+Infant)
> anova(modelo_mejor, modelo_mejor2)
Analysis of Variance Table
```

```
Model 1: Fertility ~ Examination + Education + Infant
Model 2: Fertility ~ Agriculture + Examination + Education + Infant
   Res.Df    RSS Df Sum of Sq    F Pr(>F)
1      43 35268872
2      42 33779724  1   1489148  1.8515 0.1809
```

El mejor modelo es:

Fertility ~ Examination + Education + Infant

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico**
- 5.Calcular el Error de test
- 6.Conclusión

4. DIAGNÓSTICO

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

4.1.Linealidad, normalidad y homocedasticidad

- **Linealidad.**

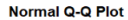
- **Normalidad.**

```
> #Normalidad  
> #Vamos a usar el test de Shapiro-wilk  
> shapiro.test(resid(modelo_mejor))
```

Shapiro-Wilk normality test

```
data:  resid(modelo_mejor)  
W = 0.95708, p-value = 0.08238
```

- ## 6. Conclusión

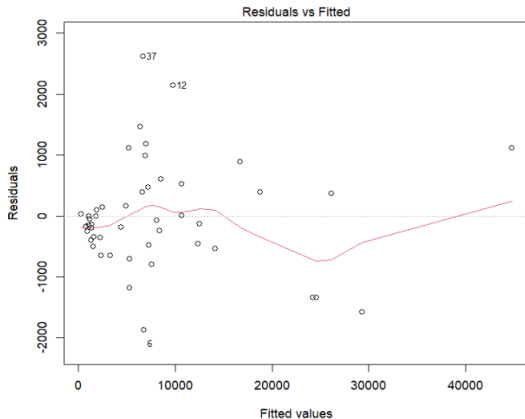


- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

● Homocedastidad.

studentized Breusch-Pagan test

data: modelo_mejor
BP = 3.3069, df = 3, p-value = 0.3467



- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

4.2.Autocorrelación

Como el p-valor es menor que 0.05 existe autocorrelación entre los errores.

```
> durbinWatsonTest(bestmodel)
lag Autocorrelation D-W Statistic p-value
1      0.04244051      1.90737    0.702
Alternative hypothesis: rho != 0
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico**
- 5.Calcular el Error de test
- 6.Conclusión

4.3 ESTUDIO DE OUTLIERS, INFLUYENTES Y LEVERAGE

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

4.3.1.Outliers

Aplicamos el método de Bonferroni con un nivel de significación del 0.05 para detectar outliers:

```
> outlierTest(bestmodel)
No studentized residuals with Bonferroni p < 0.05
Largest |rstudent|:
      rstudent unadjusted p-value Bonferroni p
sierra 3.261806          0.0022017      0.10348
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

4.3.1.Outliers

Obtenemos la observación Sierre (la número 37) como outlier.

```
> bestmodel_sin_37<-lm(Fertility~Examination+Education+Infant, data=conjunto_swiss_sin37)
> bestmodel_sin_37
```

call:

```
lm(formula = Fertility ~ Examination + Education + Infant, data = conjunto_swiss_sin37)
```

Coefficients:

(Intercept)	Examination	Education	Infant
89.8321	-0.5522	-0.3667	4.1261

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

4.3.1.Outliers

Como el p-valor es menor que 0.05 existe autocorrelación entre los errores.

```
> durbinwatsonTest(bestmodel_sin_37)
lag Autocorrelation D-W statistic p-value
1      -0.1036541      2.200708    0.486
Alternative hypothesis: rho != 0
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

4.3.1.Outliers

Como el p-valor es menor que 0.05 existe normalidad.

```
> shapiro.test(resid(bestmodel_sin_37))
```

```
shapiro-wilk normality test
```

```
data: resid(bestmodel_sin_37)
```

```
W = 0.96724, p-value = 0.2182
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

4.3.1.Outliers

Perdemos la homocedasticidad pues el p-valor es menor que 0.05

```
> bptest(bestmodel_sin_37)
```

```
studentized Breusch-Pagan test
```

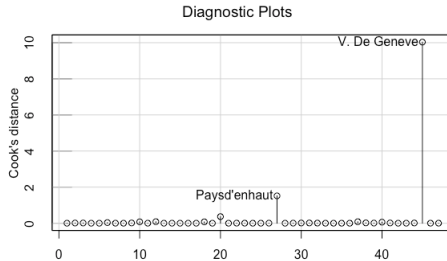
```
data: bestmodel_sin_37
```

```
BP = 7.9765, df = 3, p-value = 0.0465
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico**
- 5.Calcular el Error de test
- 6.Conclusión

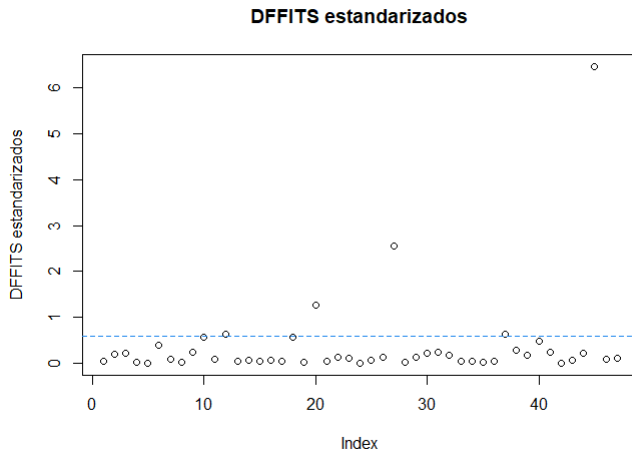
4.3.2.Observaciones influyentes

Para ver las observaciones influyentes hemos utilizado tres métodos distintos: la distancia de Cook, los DFFITS y los DFBETAS. Estos son los resultados que hemos obtenido:



- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico**
- 5.Calcular el Error de test
- 6.Conclusión

4.3.2.Observaciones influyentes



- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

Estudio del modelo sin la observación V. de Geneve.

Shapiro-Wilk normality test

```
data: resid(possible_modelo_sin45)
W = 0.94078, p-value = 0.02094
```

```
> bptest(possible_modelo_sin45)

studentized Breusch-Pagan test

data: possible_modelo_sin45
BP = 1.6853, df = 3, p-value = 0.6402
```

```
> durbinwatsonTest(possible_modelo_sin45)
lag Autocorrelation D-W Statistic p-value
1      0.06527015      1.863147    0.584
Alternative hypothesis: rho != 0
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

Estudio del modelo sin la observación Paysd'enhaut.

```
> shapiro.test(resid(possible_modelo_sin27))
```

Shapiro-Wilk normality test

```
data: resid(possible_modelo_sin27)
W = 0.9159, p-value = 0.002724
```

```
> bptest(possible_modelo_sin27)
```

studentized Breusch-Pagan test

```
data: possible_modelo_sin27
BP = 2.0947, df = 3, p-value = 0.553
```

```
> durbinWatsonTest(possible_modelo_sin27)
lag Autocorrelation D-W Statistic p-value
1      0.1193181      1.74736      0.346
Alternative hypothesis: rho != 0
.
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

Estudio del modelo sin la observación Lavaux.

```
> shapiro.test(resid(possible_modelo_sin20))
```

```
Shapiro-Wilk normality test
```

```
data: resid(possible_modelo_sin20)  
W = 0.94157, p-value = 0.02241
```

```
> bptest(possible_modelo_sin20)
```

```
studentized Breusch-Pagan test
```

```
data: posible_modelo_sin20  
BP = 1.9598, df = 3, p-value = 0.5808
```

```
> durbinWatsonTest(possible_modelo_sin20)  
lag Autocorrelation D-W Statistic p-value  
1 0.02323095 1.947479 0.784  
Alternative hypothesis: rho != 0
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

Estudio del modelo sin la observación Porrentruy.

```
> shapiro.test(resid(possible_modelo_sin6))
```

Shapiro-Wilk normality test

```
data: resid(possible_modelo_sin6)
W = 0.94645, p-value = 0.03429
```

```
> bptest(possible_modelo_sin6)
```

studentized Breusch-Pagan test

```
data: possible_modelo_sin6
BP = 3.4924, df = 3, p-value = 0.3217
```

```
> durbinWatsonTest(possible_modelo_sin6)
lag Autocorrelation D-W Statistic p-value
1 0.02629031 1.936485 0.738
Alternative hypothesis: rho != 0
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

Estudio del modelo sin la observación Sarine.

```
> shapiro.test(resid(posible_modelo_sin10))
```

Shapiro-Wilk normality test

data: resid(posible_modelo_sin10)
W = 0.95504, p-value = 0.0735

```
> bptest(posible_modelo_sin10)
```

studentized Breusch-Pagan test

data: posible_modelo_sin10
BP = 2.1482, df = 3, p-value = 0.5422

```
> durbinWatsonTest(posible_modelo_sin10)
```

lag	Autocorrelation	D-W Statistic	p-value
1	0.04657698	1.899036	0.74

Alternative hypothesis: rho != 0
~ |

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

Estudio del modelo sin la observación La Chauxdfnd.

```
> shapiro.test(resid(possible_modelo_sin40))
```

```
Shapiro-Wilk normality test
```

```
data: resid(possible_modelo_sin40)  
W = 0.95488, p-value = 0.07243
```

```
> bptest(possible_modelo_sin40)
```

```
studentized Breusch-Pagan test
```

```
data: possible_modelo_sin40  
BP = 3.0309, df = 3, p-value = 0.3869
```

```
> durbinWatsonTest(possible_modelo_sin40)  
lag Autocorrelation D-W Statistic p-value  
1 0.05081376 1.890029 0.688  
Alternative hypothesis: rho != 0  
> |
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico**
- 5.Calcular el Error de test
- 6.Conclusión

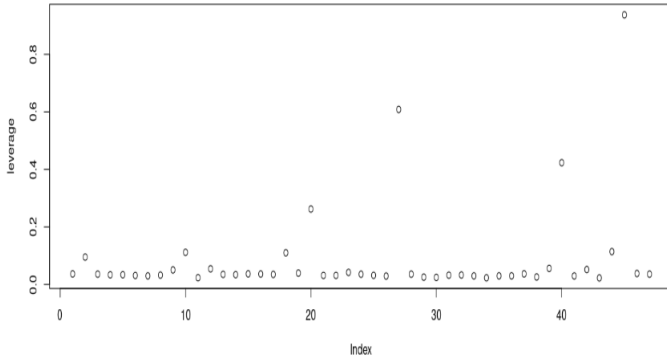
4.4.Colinealidad

Como el valor del determinante es muy grande, descartamos la opción de que la matriz sea casi singular. Concluimos que no existe multicolinealidad.

```
> X<-cbind(rep(1,length(Fertility)),Examination,Education,Infant)
> det(t(X)%*%X)
[1] 6.218467e+26
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico**
- 5.Calcular el Error de test
- 6.Conclusión

4.4.1.Leverage



- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico**
- 5.Calcular el Error de test
- 6.Conclusión

5. CALCULAR EL ERROR DE TEST

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

5.Calcular el Error de test

Separeremos los datos entre 70 % entrenamiento y 30 % test y realizaremos la validación cruzada.

```
> model.exh
Subset selection object
Call: regsubsets.formula(Fertility ~ ., data = conjunto_swiss[train,
  1:6], method = "exhaustive")
5 Variables (and intercept)
              Forced in Forced out
Agriculture      FALSE      FALSE
Examination      FALSE      FALSE
Education        FALSE      FALSE
Catholic         FALSE      FALSE
Infant.Mortality FALSE      FALSE
1 subsets of each size up to 5
Selection Algorithm: exhaustive
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

```
> summary(model.exh) #todos los modelos posibles para los `predictores
Subset selection object
Call: regsubsets.formula(Fertility ~ ., data = conjunto_swiss[train,
  1:6], method = "exhaustive")
5 Variables (and intercept)
      Forced in Forced out
Agriculture      FALSE      FALSE
Examination      FALSE      FALSE
Education         FALSE      FALSE
Catholic          FALSE      FALSE
Infant.Mortality  FALSE      FALSE
1 subsets of each size up to 5
Selection Algorithm: exhaustive
      Agriculture Examination Education Catholic Infant.Mortality
1 ( 1 ) " " " " " " " " " "
2 ( 1 ) "*" " " " " " " " "
3 ( 1 ) "*" " " "*" " " " "
4 ( 1 ) "*" "*" "*" " " " "
5 ( 1 ) "*" "*" "*" "*" " "

> val.errors
[1] 50899007 20279332 2991010 3047738 3051125
> coef(model.exh, which.min(val.errors))
      (Intercept)      Agriculture      Education Infant.Mortality
      -39.7102481       0.2150309      -0.4950186       3.3301860
```

El mejor modelo es: $Fertility \sim Agriculture + Education + Infant$

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

6. CONCLUSIÓN

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

6.1.Nuevas observaciones

Creamos un nuevo `data.frame()` con 3 nuevas provincias suizas

```
> newdataframe
  Fertility Agriculture Education Examination Catholic Infant
1    13.20      40.31      23      20.11     50.31    62.31
2    50.30      52.54       2      12.19     36.54    21.00
3    18.93      39.99       5      20.51     27.89    24.70
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

6.1.Nuevas observaciones

El valor estimado de la fertilidad y sus intervalos de confianza al 95 %:

	Yhat	S_yhat	IntervConfInf	IntervConfSup
1	37673.40	177.1037	37233.36	38113.45
2	12901.66	178.3231	12458.59	13344.74
3	13973.17	178.1431	13530.53	14415.80

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

6.2.Conclusión final

Mejor modelo,

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i, \quad i = 1, \dots, 46$$

siendo $\beta_0 = 145,57623$, $\beta_1 = -0,63120$, $\beta_2 = -0,33337$ y
 $\beta_3 = 4,18837$

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

6.2.Conclusión final

Su tabla ANOVA es:

```
> anova(bestmodel_sin_10)
Analysis of Variance Table
```

Response: Fertility

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Examination	1	1326184684	1326184684	1675.57	< 2.2e-16 ***
Education	1	608156072	608156072	768.38	< 2.2e-16 ***
Infant	1	1563443859	1563443859	1975.34	< 2.2e-16 ***
Residuals	42	33242274	791483		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

6.2.Conclusión final

Su valor R^2 es 0.9899132. En el summary observamos:

```
> summary(bestmodel_sin_10)

Call:
lm(formula = Fertility ~ Examination + Education + Infant, data = conjunto_swiss_sin10)

Residuals:
    Min       1Q   Median       3Q      Max
-1905.3  -464.5  -124.2   376.7  2575.7

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 145.57623   179.19163    0.812  0.421140
Examination  -0.63120    0.14011   -4.505  5.22e-05 ***
Education    -0.33337    0.08314   -4.010  0.000244 ***
Infant        4.18837    0.09424   44.445 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 889.7 on 42 degrees of freedom
Multiple R-squared:  0.9906,    Adjusted R-squared:  0.9899
F-statistic: 1473 on 3 and 42 DF,  p-value: < 2.2e-16
```

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

6.2 Conclusión final

También hemos observado que un modelo igual de bueno es el que relaciona las variables 'Agriculture', 'Education' e 'Infant'.

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i, \quad i = 1, \dots, 47$$

siendo $\beta_0 = 849,82333$, $\beta_1 = 0,84578$, $\beta_2 = 1,60544$ y $\beta_3 = -0,46890$

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

6.2.Conclusión final

Su tabla ANOVA es:

Analysis of Variance Table

Response: Fertility

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Agriculture	1	2490455546	2490455546	311.369	< 2.2e-16 ***
Examination	1	897597291	897597291	112.222	1.457e-13 ***
Education	1	24947226	24947226	3.119	0.08448 .
Residuals	43	343931429	7998405		

signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

- 1.Introducción
2. Estudio y evaluación del modelo completo
- 3.Selección del mejor modelo
- 4.Diagnóstico
- 5.Calcular el Error de test
- 6.Conclusión

6.2.Conclusión final

Su valor R^2 es 0.9947118. En el summary observamos:

```
call:
lm(formula = Fertility ~ Agriculture + Examination + Education,
    data = conjunto_swiss)

Residuals:
    Min       1Q   Median       3Q      Max
-10623.2  -1071.9   -614.9   1121.7   6186.5

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  849.82333   557.81880    1.523  0.1350
Agriculture    0.84578     0.06492   13.029 < 2e-16 ***
Examination    1.60544     0.35824    4.482 5.42e-05 ***
Education     -0.46890     0.26550   -1.766  0.0845 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2828 on 43 degrees of freedom
Multiple R-squared:  0.9085,    Adjusted R-squared:  0.9021
F-statistic: 142.2 on 3 and 43 DF,  p-value: < 2.2e-16
```