

Tarea 7 de Modelos Probabilistas Aplicados

Transformaciones mediante Escalera de poderes de Tukey

5271

20 de octubre de 2020

1. Introducción

En este trabajo se realiza un breve acercamiento al tema Transformaciones mediante Escalera de poderes de Tukey y regresión lineal. Además presenta un algoritmo desarrollado en lenguaje R para extraer la función que genera la variable dependiente así como sus coeficientes. La experimentación se realiza en el programa R versión 4.0.2 [1] en el entorno de desarrollo Rstudio [2]

2. Transformaciones mediante Escalera de poderes de Tukey

Para entender mejor como funciona las transformaciones mediante Escalera de poderes de Tukey creamos un *data.frame* con datos bivariados (x_1, y) de manera que x_1 se distribuye uniformemente e $y = ax_1 + b$. Lo anterior se realiza con el código 2 en R.

```
1 x_1 = runif(100,10,120)
2 a = 8
3 b = 15
4 y = a*(x_1) + b
5 datas = as.data.frame(y)
6 datas = cbind(datas, x_1)
```

Tarea7n.R

Como primer paso se trazan los datos en un diagrama de dispersión. En la figura 1 de la página 2 se puede observar como se relaciona las variables x_1 e y . Como se aprecia en dicha figura la relación entre x_1 e y es una dependencia lineal y si aplicamos una regresión lineal obtendremos exactamente los coeficientes y formulación con la cual se crea la variable dependiente y . A continuación se metra el resultado de aplicar la regresión lineal. De donde se pueden extraer los coeficientes a y b de la ecuación, quedando $y = 15x_1 + 8$.

Call:

```
lm(formula = datas$y ~ datas$x_1)
```

Residuals:

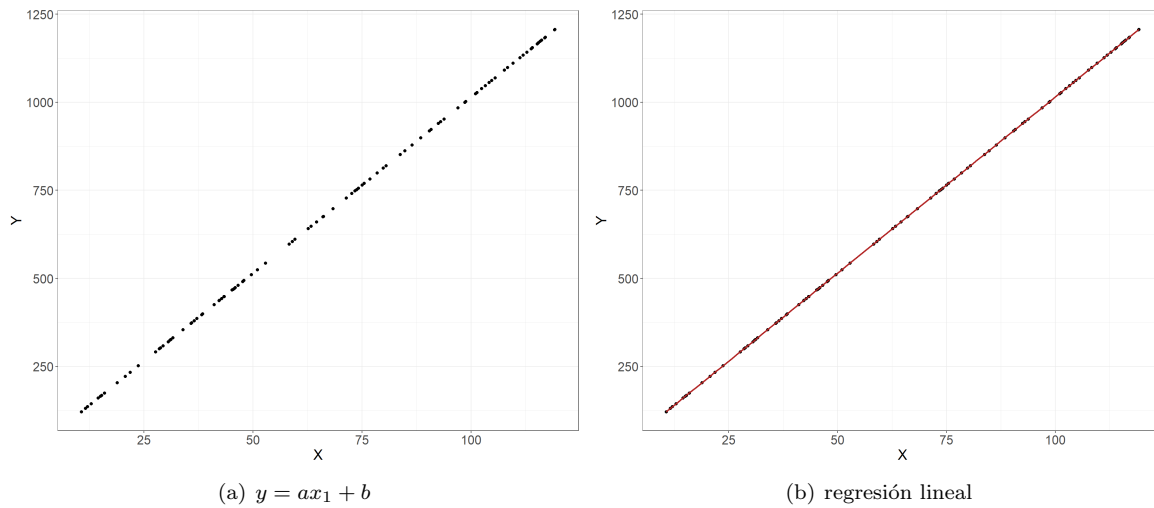


Figura 1: Diagramas de dispersión de los datos bivariados (x_1, y)

	Min	1Q	Median	3Q	Max
	-8.102e-13	-3.110e-15	8.040e-15	1.811e-14	9.951e-14

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.500e+01	2.016e-14	7.442e+14	<2e-16 ***
datas\$x_1	8.000e+00	2.744e-16	2.915e+16	<2e-16 ***

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

Residual standard error: 8.501e-14 on 98 degrees of freedom

Multiple R-squared: 1, Adjusted R-squared: 1

F-statistic: 8.498e+32 on 1 and 98 DF, p-value: < 2.2e-16

Warning message:

In summary.lm(modelo_lineal) :

essentially perfect fit: summary may be unreliable

El ejemplo anterior es sencillo pues con solo una regresión lineal se pueden obtener los coeficientes y la función que crea a y . En cambio si creamos representamos los datos bivariados creados con $y = a * \log(x_1) + b$, se puede observar en la figura 2 de la página 3 que la relación de dependencia no es lineal, por lo que una simple regresión no arrojará los resultados deseados, por lo cual antes de aplicar la regresión se debe hacer una transformación a una de las variable, es decir tratar de linealizar la relación de dependencia.

No hay ninguna restricción sobre los valores de λ que podamos considerar. Obviamente, elegir $\lambda = 1$ deja los datos sin cambios. Los valores negativos de λ también son razonables. Lo que da lugar a la ecuación 2. El cuadro 1 de la página 3 se muestran ejemplos de la escalera de transformaciones de Tukey.

Tukey sugiere explorar relaciones como:

$$y = ax_1^\lambda + b \quad (1)$$

Cuadro 1: Escalera de transformaciones de Tukey

λ	-2	-1	-1/2	0	1/2	1	2
y	$\frac{-1}{x^2}$	$\frac{-1}{x}$	$\frac{-1}{\sqrt{x}}$	$\log x$	\sqrt{x}	x	x^2

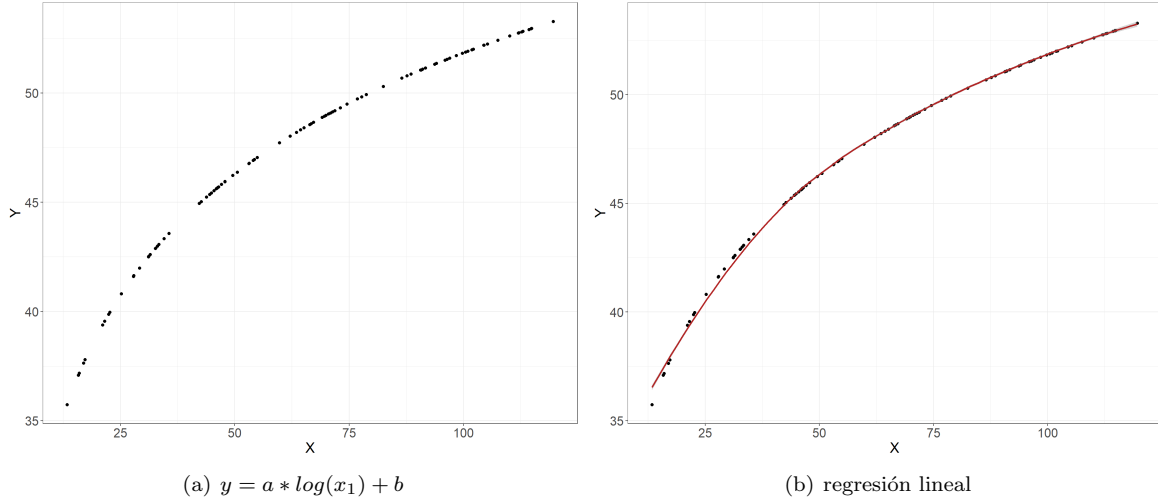


Figura 2: Diagramas de dispersión de los datos bivariados (x_1, y) con dependencia no lineal

donde λ es una parámetro que elegido con el fin de linealizar la relación.

$$y = \begin{cases} x^\lambda & \text{si } \lambda > 0 \\ \log x & \text{si } \lambda = 0 \\ -(x^\lambda) & \text{si } \lambda < 0 \end{cases} \quad (2)$$

Se Aplica transformaciones mediante la escalera de Tukey a los datos, en forma que se calcula el logaritmo de la variable x_1 , y se vuelve a representar los datos y se realiza la regresión lineal. En la figura 3 de la página 4 se puede observar la transformación de los datos.

3. Algoritmo propuesto

Según lo visto en la sección anterior podemos mediante las transformaciones de la Escalera de Tukey, podemos acercar la relación entre los datos bivariados a una dependencia lineal. A partir de esta conclusión se creo una función en R que la pasarle como parámetros un *data.frame* y la cantidad de variables de cuales depende y (máximo tres variables), devuelve un *data.frame* con los valores de λ encontrados para la linealización, los coeficientes de la función y el error estándar de la estimación de los coeficientes. En el código ?? se muestra dicha función.

```

1 ajuste1 = function(datos,v){
2   rsq=numeric()
3   landax1=numeric()
4   landax2=numeric()
5   landax3=numeric()
6   coefi_a=numeric()

```

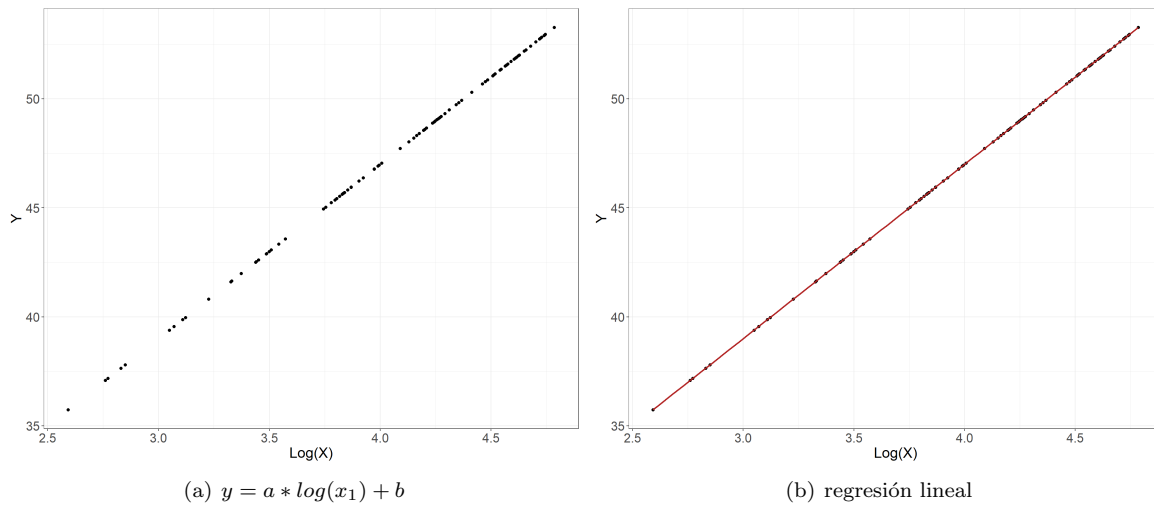


Figura 3: Diagramas de dispersión de los datos bivariados (x_1, y) con dependencia no lineal

```

7  coefi_b=numeric()
8  coefi_c=numeric()
9  error_a=numeric()
10 error_c=numeric()
11 error_b=numeric()
12 corr = numeric()
13 diferencia = numeric()
14 s=seq(-2,2,0.25)
15 if(v==1){
16   for (i in s) {
17     x1=numeric()
18     if(i<0){
19       x1=-1/((datos$x_1)**(i*-1))
20     } else if(i>0){
21       x1=(datos$x_1)**i
22     } else {
23       x1=log(datos$x_1)}
24
25     modelo_lineal = lm(datos$y ~ x1)
26     rsq = c(rsq,summary(modelo_lineal)$r.squared)
27     diferencia = c(diferencia,(1-summary(modelo_lineal)$r.squared))
28     coefi_a=c(coefi_a,modelo_lineal$coefficient[1])
29     error_a=c(error_a,round(summary(modelo_lineal)$coefficient[1,2],4))
30     coefi_b=c(coefi_b,modelo_lineal$coefficient[2])
31     error_b=c(error_b,round(summary(modelo_lineal)$coefficient[2,2],4))
32     landax1=c(landax1,i)
33   }
34   corre= as.data.frame(rsq)
35   corrl= cbind(corre,landax1,diferencia,coefi_a,error_a,coefi_b,error_b)
36   resultado=filter(corrl,diferencia==min(corrl$diferencia))
37
38 } else if(v==2){
39   for (i in s) {
40     x1=numeric()
41     if(i<0){
42       x1=-1/((datos$x_1)**(i*-1))
43     } else if(i>0){
44       x1=(datos$x_1)**i
45     } else {
46

```

```

47     x1=log(datos$x_1)}
48   for (j in s) {
49     x2=numeric()
50     if(j<0){
51       x2=-1/((datos$x_2)**(j*-1))
52     }else if(j>0){
53       x2=(datos$x_2)**j
54     }else{
55       x2=log(datos$x_2)}
56
57     modelo_lineal = lm(datos$y ~ x1+x2)
58     rsq = c(rsq,summary(modelo_lineal)$r.squared)
59     diferencia = c(diferencia,(1-summary(modelo_lineal)$r.squared))
60     coefi_a=c(coefi_a,modelo_lineal$coefficient[2])
61     error_a=c(error_a,round(summary(modelo_lineal)$coefficient[2,2],4))
62     coefi_b=c(coefi_b,modelo_lineal$coefficient[3])
63     error_b=c(error_b,round(summary(modelo_lineal)$coefficient[3,2],4))
64     landax1=c(landax1,i)
65     landax2=c(landax2,j)
66
67   }
68 }
69 corre= as.data.frame(rsq)
70 corrl= cbind(corre,landax1,landax2,diferencia,coefi_a,error_a,coefi_b,error_b)
71 resultado=filter(corrl,diferencia==min(corrl$diferencia))
72 }else if(v==3){
73
74   for (i in s) {
75     x1=numeric()
76     if(i<0){
77       x1=-1/((datos$x_1)**(i*-1))
78     }else if(i>0){
79       x1=(datos$x_1)**i
80     }else{
81       x1=log(datos$x_1)}
82   for (j in s) {
83     x2=numeric()
84     if(j<0){
85       x2=-1/((datos$x_2)**(j*-1))
86     }else if(j>0){
87       x2=(datos$x_2)**j
88     }else{
89       x2=log(datos$x_2)}
90   for (k in s) {
91     x3=numeric()
92     if(k<0){
93       x3=-1/((datos$x_3)**(k*-1))
94     }else if(k>0){
95       x3=(datos$x_3)**k
96     }else{
97       x3=log(datos$x_3)}
98
99     modelo_lineal = lm(datos$y ~ x1+x2+x3)
100    rsq = c(rsq,summary(modelo_lineal)$r.squared)
101    diferencia = c(diferencia,(1-summary(modelo_lineal)$r.squared))
102    coefi_a=c(coefi_a,modelo_lineal$coefficient[2])
103    error_a=c(error_a,round(summary(modelo_lineal)$coefficient[2,2],4))
104    coefi_b=c(coefi_b,modelo_lineal$coefficient[3])
105    error_b=c(error_b,round(summary(modelo_lineal)$coefficient[3,2],4))
106    coefi_c=c(coefi_c,modelo_lineal$coefficient[4])
107    error_c=c(error_c,round(summary(modelo_lineal)$coefficient[4,2],4))
108    landax1=c(landax1,i)
109    landax2=c(landax2,j)
110    landax3=c(landax3,k)
111  }

```

```

112     }
113   }
114   corre= as.data.frame(rsq)
115   corrl= cbind(corre,landax1,landax2,landax3,diferencia,coefi_a,error_a,coefi_b,
116   error_b,coefi_c,error_c)
117   resultado=filter(corrl,diferencia == min(corrl$diferencia))
118 }
119
120 return(resultado)
121 }

```

Tarea7n.R

Resultados El código general se encuentra disponible en el repositorio. <https://github.com/Albertomnoa/Tareas>

Referencias

- [1] R Core Team. R: Un lenguaje y un entorno para la informática estadística, 2020.
- [2] RStudio Team. Rstudio: Entorno de desarrollo integrado para R, 2020.