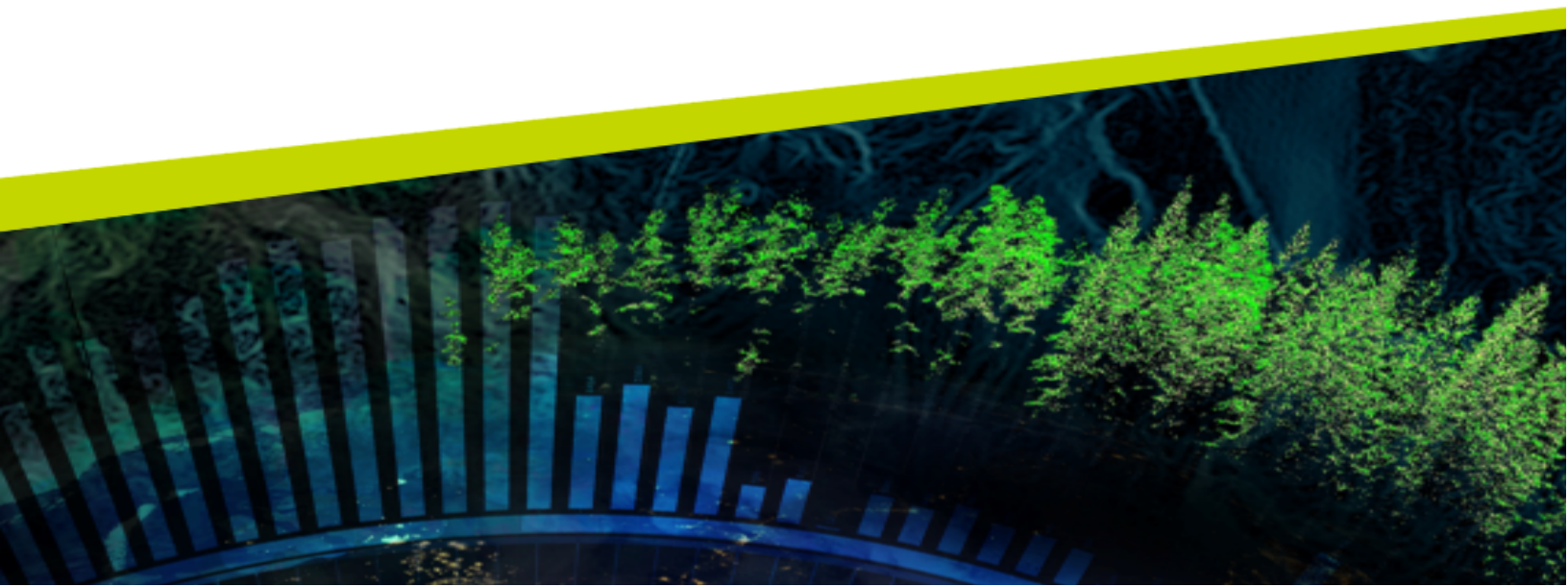




# INFERENCIA Y MODELOS ESTADÍSTICOS

Jacqueline Köhler C. y José Luis Jara V.



## CAPÍTULO 11. MÉTODOS CLÁSICOS PARA ENFRENTAR VARIABLES NUMÉRICAS PROBLEMÁTICAS

Como ya sabemos, muchos procedimientos estadísticos requieren que los datos cumplan con ciertas propiedades o condiciones, lo que no siempre ocurre. En este capítulo abordaremos algunos métodos para enfrentar estos problemas cuando se intenta inferir sobre las medias poblacionales de variables aleatorias numéricas.

Para ello nos basaremos principalmente en Lane (s.f., pp. 1, 16), Lowry (1999, caps. 11a, 12a, 14a, 15a), Glen (2021b) y Lærd Statistics (2020).

### 11.1 TRANSFORMACIÓN DE DATOS

Un fenómeno que ocurre a menudo en los estudios, además del incumplimiento de ciertas condiciones, es la necesidad de convertir los datos de una escala a otra diferente. Para hacer tales transformaciones, debemos aplicar una determinada función a una variable aleatoria  $X$ , lo que nos entrega como resultado una nueva variable aleatoria  $Y$ .

Existen diversos métodos que podemos usar para transformar datos, dependiendo de la forma que tengan los datos originales y la que deseemos obtener como resultado. En esta sección conoceremos, entonces, algunas transformaciones de uso frecuente en estadística.

#### 11.1.1 Transformación lineal

Las **transformaciones lineales** son las más sencillas, y para hacerlas nos basta con aplicar una función lineal, es decir, de la forma presentada en la ecuación 11.1, donde  $m$  y  $n$  son constantes.

$$y_i = m \cdot x_i + n \quad (11.1)$$

La física nos ofrece muchos escenarios en que es necesario aplicar este tipo de transformaciones, pues es el tipo de operación que realizamos cuando convertimos de una unidad a otra. A modo de ejemplo, consideremos la conversión de grados Celsius a grados Fahrenheit:  $F = 1,8 \cdot C + 32$ .

En R, podemos hacer este tipo de transformaciones de forma muy sencilla mediante operaciones aritméticas que se aplican vectorialmente, como muestra el script 11.1. En él se transforma un vector con 4 temperaturas en grados Celsius a grados Fahrenheit. El resultado se muestra en la figura 11.1.

Script 11.1: transformación lineal para convertir grados Celsius a grados Fahrenheit.

```
1 # Crear un vector con cuatro observaciones en grados Celsius.
2 Celsius <- c(-8, 0, 29.8, 100)
3
4 # Aplicar transformación lineal para convertir a grados Fahrenheit.
5 Fahrenheit <- 1.8 * Celsius + 32
6
7 # Mostrar los resultados.
8 cat("Temperaturas en grados Celsius\n")
9 print(Celsius)
10 cat("\nTemperaturas en grados Fahrenheit\n")
11 print(Fahrenheit)
```

```

Temperaturas en grados Celsius
[1] -8.0  0.0 29.8 100.0

Temperaturas en grados Fahrenheit
[1] 17.60 32.00 85.64 212.00

```

Figura 11.1: resultado de la transformación lineal del script 11.1.

### 11.1.2 Transformación logarítmica

La transformación logarítmica nos puede servir cuando tenemos distribuciones muy asimétricas, pues ayuda a reducir la desviación y así facilita el cumplimiento de la condición de normalidad requerida por muchas de las pruebas estadísticas que ya conocemos. Para ver este efecto de manera más clara, usaremos un conjunto de datos que registra el peso corporal (en kilogramos) y el peso del cerebro (en gramos) de diversos animales, algunos de ellos extintos (Rousseeuw & Leroy, 1987, p. 57). En R, esta transformación puede hacerse gracias a la función `log(x, base)`, aunque debemos tener cuidado con posibles valores iguales a 0. El script 11.2 aplica esta transformación al peso corporal y al peso cerebral de los animales (líneas 22–23). La figura 11.2 (script 11.2, líneas 28–42) muestra gráficamente el resultado de esta transformación para el peso cerebral de los animales.

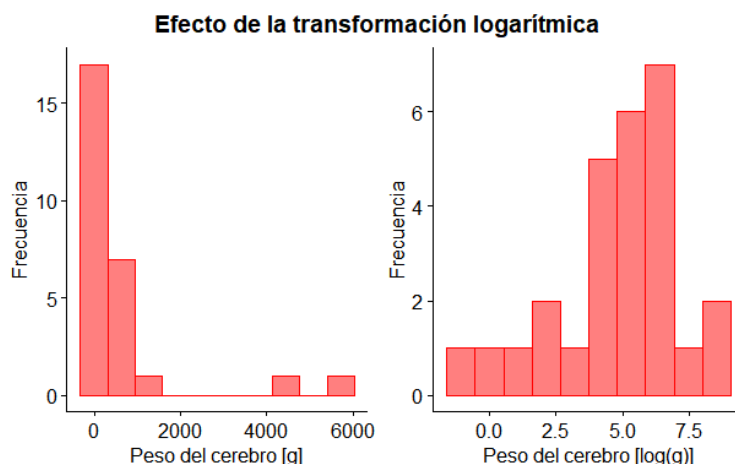


Figura 11.2: histogramas del peso cerebral antes y después de la transformación logarítmica.

Muchas veces la transformación logarítmica hace que nos sea más fácil interpretar los datos, evidenciando patrones más claros para la relación entre variables. En la figura 11.3 (script 11.2, líneas 46–60), a la derecha, se evidencia una fuerte relación entre el peso corporal y el peso del cerebro después de transformar ambas variables, relación que no podemos percibir con los datos originales (izquierda).

Script 11.2: transformación logarítmica.

```

1 library(ggpubr)
2
3 # Cargar datos
4 animal <- c("Mountain beaver", "Cow", "Grey wolf", "Goat", "Guinea pig",
5            "Dipliodocus", "Asian elephant", "Donkey", "Horse",
6            "Potar monkey", "Cat", "Giraffe", "Gorilla", "Human",
7            "African elephant", "Triceratops", "Rhesus monkey", "Kangaroo",
8            "Golden hamster", "Mouse", "Rabbit", "Sheep", "Jaguar",
9            "Chimpanzee", "Brachiosaurus", "Mole", "Pig")
10
11 body_weight <- c(1.35, 465, 36.33, 27.66, 1.04, 11700, 2547, 187.1, 521, 10,
12                3.3, 529, 207, 62, 6654, 9400, 6.8, 35, 0.12, 0.023, 2.5,

```

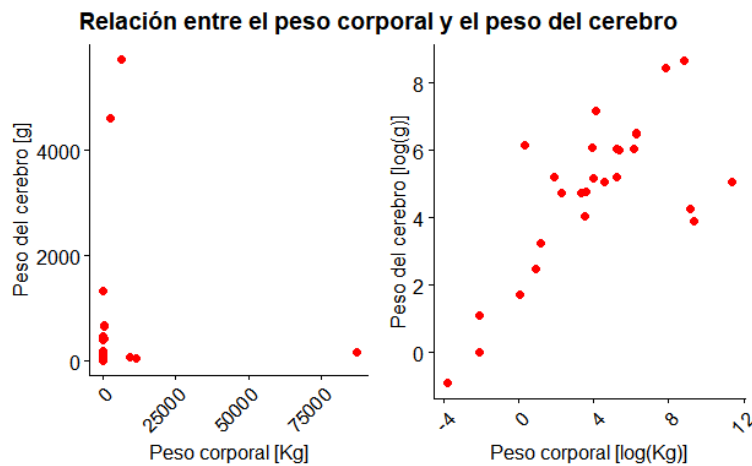


Figura 11.3: gráficos de dispersión para el peso corporal y el peso del cerebro antes y después de las transformaciones logarítmicas.

```

13         55.5, 100, 52.16, 87000, 0.122, 192)
14
15 brain_weight <- c(465, 423, 119.5, 115, 5.5, 50, 4603, 419, 655, 115, 25.6,
16                  680, 406, 1320, 5712, 70, 179, 56, 1, 0.4, 12.1, 175, 157,
17                  440, 154.5, 3, 180)
18
19 datos <- data.frame(animal, body_weight, brain_weight)
20
21 # Aplicar transformación logarítmica
22 log_cuerpo <- log(body_weight)
23 log_cerebro <- log(brain_weight)
24 datos <- data.frame(datos, log_cuerpo, log_cerebro)
25
26 # Histogramas para el peso cerebral antes y después de la transformación
27 # logarítmica.
28 g3 <- ggghistogram(datos, x = "brain_weight", bins = 10,
29                    xlab = "Peso del cerebro [g]", ylab = "Frecuencia",
30                    color = "red", fill = "red")
31
32 g4 <- ggghistogram(datos, x = "log_cerebro", bins = 10,
33                    xlab = "Peso del cerebro [log(g)]", ylab = "Frecuencia",
34                    color = "red", fill = "red")
35
36 # Crear una única figura con ambos histogramas.
37 histograma <- ggarrange(g3, g4, ncol = 2, nrow = 1)
38 titulo <- text_grob("Efecto de la transformación logarítmica",
39                    face = "bold", size = 14)
40 histograma <- annotate_figure(histograma, top = titulo)
41
42 print(histograma)
43
44 # Gráficos de dispersión para la relación entre peso corporal y peso del
45 # cerebro, antes y después de aplicar la transformación logarítmica.
46 g1 <- ggscatter(datos, x = "body_weight", y = "brain_weight",
47                 color = "red", xlab = "Peso corporal [Kg]",
48                 ylab = "Peso del cerebro [g]") + rotate_x_text(45)
49
50 g2 <- ggscatter(datos, x = "log_cuerpo", y = "log_cerebro",
51                 color = "red", xlab = "Peso corporal [log(Kg)]",

```

```

52     ylab = "Peso del cerebro [log(g)]") + rotate_x_text(45)
53
54 # Crear una única figura con los gráficos de dispersión.
55 dispersion <- ggarrange(g1, g2, ncol = 2, nrow = 1)
56 texto <- "Relación entre el peso corporal y el peso del cerebro"
57 titulo <- text_grob(texto, face = "bold", size = 14)
58 dispersion <- annotate_figure(dispersion, top = titulo)
59
60 print(dispersion)

```

Eso sí, tenemos que ser muy cuidadosos al usar esta transformación, porque cuando comparamos medias de datos tras una transformación logarítmica, en realidad estamos comparando **medias geométricas**! Recordemos que la media geométrica se calcula de acuerdo a la ecuación 11.2 y suele ocuparse para representar tasas de crecimiento o de interés (Glen, 2021a).

$$\left(\prod_{i=1}^n x_i\right)^{\frac{1}{n}} = \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n} \quad (11.2)$$

Para ilustrar esta idea, supongamos que aplicamos una transformación logarítmica con base 10 al vector  $[1, 10, 100]$ , obteniendo como resultado  $[0, 1, 2]$ .

La media aritmética del vector transformado es:

$$\frac{0 + 1 + 2}{3} = 1$$

Al revertir la transformación para la media, tenemos que:

$$10^1 = 10$$

Lo que es distinto a la media aritmética de los datos originales:

$$\frac{1 + 10 + 100}{3} = 37$$

A su vez, la media geométrica del vector original es:

$$\sqrt[3]{1 \cdot 10 \cdot 100} = 10$$

Así, si dos variables a las que se ha aplicado la transformación logarítmica tienen igual media, entonces **las medias geométricas de las variables originales son iguales**.

### 11.1.3 Escalera de potencias de Tukey

Más general que la transformación logarítmica, la **escalera de potencias de Tukey** nos ayuda a cambiar la forma de una distribución asimétrica para que se asemeje a la normal. Este método consiste en explorar relaciones de la forma que muestra la ecuación 11.3, donde  $\lambda$  puede tomar cualquier valor real y se escoge de modo que la distribución de los datos transformados sea lo más cercana a la normal posible. También es útil al explorar la relación entre dos variables, en cuyo caso se busca obtener un gráfico de dispersión en que los puntos se asemejen a una recta.

$$y = x^\lambda \quad (11.3)$$

Formalmente, la transformación de Tukey se define según la ecuación 11.4, aunque (por la falta de computadores) suelen usarse únicamente aquellas que se muestran en la tabla 11.1. Fijémonos en que si  $\lambda = 1$ , no

se realiza transformación alguna, y que para el caso de  $\lambda = 0$ , se tiene que  $x^0 = 1$ , por lo que se reemplaza en este caso por la transformación logarítmica.

$$\tilde{x}_\lambda = \begin{cases} x^\lambda & \lambda > 0 \\ \log(x) & \lambda = 0 \\ -(x^\lambda) & \lambda < 0 \end{cases} \quad (11.4)$$

$\lambda$	-2	-1	$-\frac{1}{2}$	0	$\frac{1}{2}$	1	2
$\tilde{x}$	$-\frac{1}{x^2}$	$-\frac{1}{x}$	$-\frac{1}{\sqrt{x}}$	$\log(x)$	$\sqrt{x}$	$x$	$x^2$

Tabla 11.1: escalera de transformaciones de Tukey.

Usemos ahora la población total de Estados Unidos entre los años 1610 y 1850 (United States Census Bureau, 2004, 2021) como ejemplo para entender mejor esta transformación. La figura 11.4 (líneas 18–30 del script 11.3) muestra, a la izquierda, el histograma para la población (en millones de habitantes) y, a la derecha, un gráfico de dispersión para la población por año. Podemos ver claramente que la distribución del número de habitantes presenta una fuerte asimetría positiva (o hacia la derecha) y que la población parece aumentar de manera exponencial a partir de mediados del siglo XVIII.

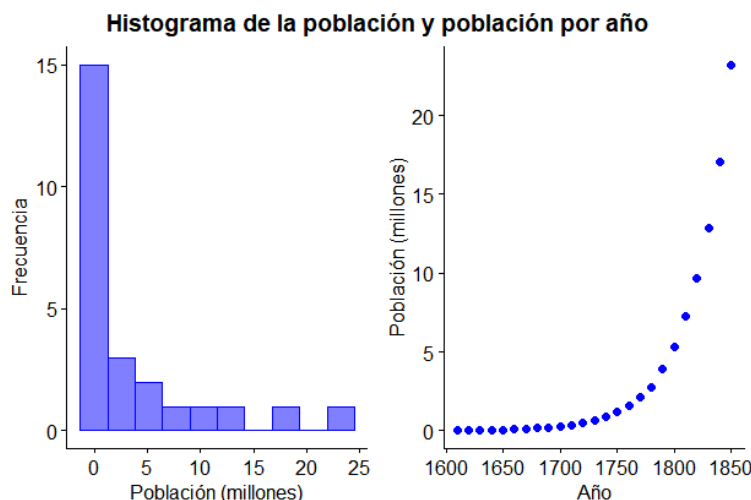


Figura 11.4: histograma de la población histórica de Estados Unidos y gráfico de dispersión de la población por año.

La figura 11.5 (líneas 46–75 del script 11.3) muestra los gráficos de dispersión para la población por año tras aplicar la transformación de Tukey con diferentes valores de  $\lambda$  a la población. En ella podemos observar que la curva cambia gradualmente de convexa a cóncava a medida que aumenta el valor de  $\lambda$ , lo que deja entrever que, para obtener un resultado que sea lo más cercano posible a una línea recta, tenemos que resolver un problema de optimización que minimice los valores residuales tras ajustar una recta a los puntos transformados. De las transformaciones presentadas en la figura 11.5, la más cercana a una recta se obtiene para  $\lambda = 0$ .

Como podemos suponer, reducir la asimetría nos ayuda a cumplir el requisito de normalidad que imponen muchas pruebas estadísticas, permitiéndonos así lograr resultados teóricamente más precisos. Sin embargo, una vez más tenemos que ser cuidadosos y tener en cuenta la transformación realizada al momento de interpretar los resultados. Si bien tenemos certeza que si se encuentran diferencias significativas en la variable transformada, estas diferencias **también existen en la variable original**, los estadísticos y los intervalos de confianza **no son los mismos** que arrojarían las pruebas con los datos originales!

En R, la función `transformTukey(x, start, end, int, plotit, verbose, quiet, statistic, returnLambda)` incluida en el paquete `rcompanion` permite aplicar la escalera fácilmente, donde:

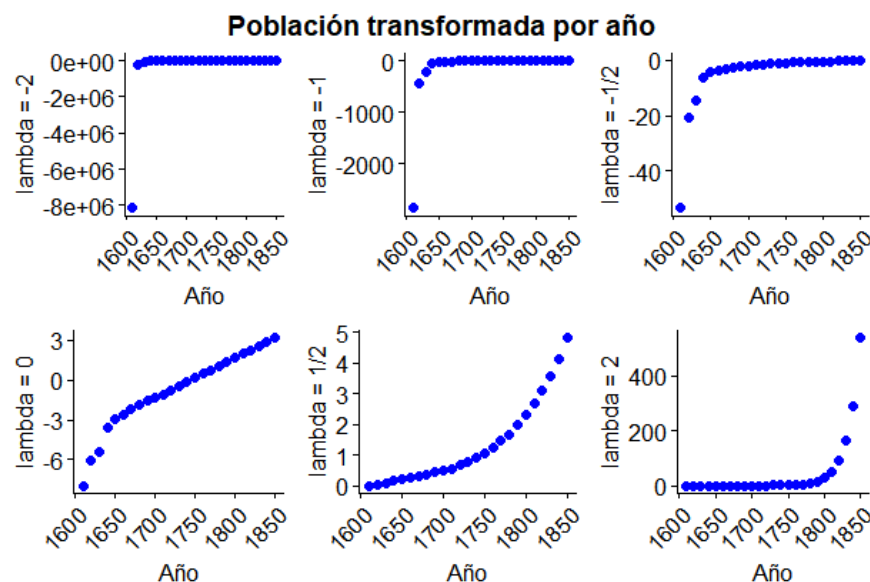


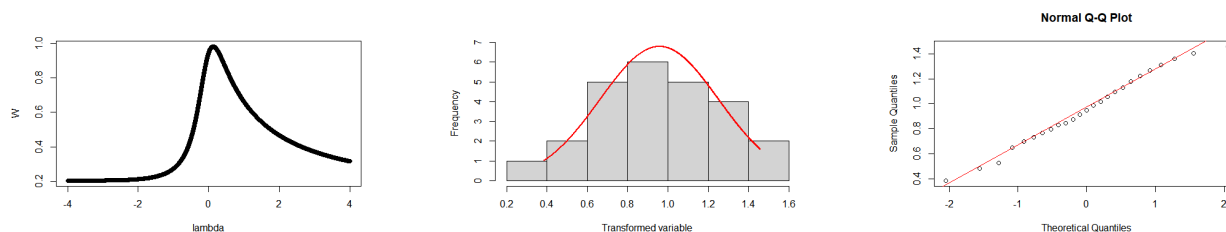
Figura 11.5: población de Estados Unidos por año tras aplicar la transformación de Tukey con distintos valores de  $\lambda$ .

- `x`: vector de valores a transformar.
- `start`: valor inicial de  $\lambda$  a evaluar.
- `end`: valor final de  $\lambda$  a evaluar.
- `int`: intervalo entre los valores de  $\lambda$  a evaluar.
- `plotit`: si toma valor `TRUE`, entrega los siguientes gráficos:
  - Estadístico de la prueba de normalidad versus  $\lambda$ .
  - Histograma de los valores transformados.
  - Gráfico Q-Q de los valores transformados.
- `verbose`: si toma valor `TRUE`, muestra información adicional sobre la prueba de normalidad con respecto a  $\lambda$ .
- `quiet`: si toma valor `TRUE`, no muestra información alguna por pantalla.
- `statistic`: si toma valor 1, usa la prueba de normalidad de Shapiro-Wilk. Con valor 2, usa la prueba de Anderson-Darling.
- `returnLambda`: si toma valor `TRUE`, devuelve el valor de  $\lambda$ . Si toma valor `FALSE`, devuelve los datos transformados.

Tras llamar a la función `transformTukey()` con los datos de la población de Estados Unidos (script 11.3, líneas 111–112), obtenemos que el valor óptimo de lambda es  $\lambda^* = 0,12$  y los gráficos de la figura 11.6. La figura 11.6a nos muestra gráficamente que el valor óptimo  $\lambda^*$  corresponde al que maximiza el estadístico entregado por la prueba de normalidad. A su vez, en la figura 11.6b vemos que la distribución obtenida aplicando  $\lambda^*$  se asemeja mucho más a una distribución normal, lo que se ve confirmado por el gráfico Q-Q de los mismos datos de la figura 11.6c.

Script 11.3: transformación de Tukey para la población total de Estados Unidos.

```
1 library(ggpubr)
2 library(rcompanion)
3
4 # Cargar datos
5 Year <- c(1610, 1620, 1630, 1640, 1650, 1660, 1670, 1680, 1690, 1700, 1710,
6           1720, 1730, 1740, 1750, 1760, 1770, 1780, 1790, 1800, 1810, 1820,
7           1830, 1840, 1850)
8
9 Population <- c(0.00035, 0.002302, 0.004646, 0.026634, 0.050368, 0.075058,
10                0.111935, 0.151507, 0.210372, 0.250888, 0.331711, 0.466185,
```



(a) estadístico W de la prueba de Shapiro-Wilk por cada valor de  $\lambda$ .

(b) histograma de los datos transformados con  $\lambda^*$ .

(c) gráfico Q-Q de los datos transformados con  $\lambda^*$ .

Figura 11.6: gráficos entregados por `transformTukey()` para los datos de población de Estados Unidos.

```

11      0.629445, 0.905563, 1.17076, 1.593625, 2.148076, 2.780369,
12      3.929214, 5.308483, 7.239881, 9.638453, 12.86602, 17.069453,
13      23.191876)
14
15 datos <- data.frame(Year, Population)
16
17 # Gráfico de dispersión e histograma.
18 g1 <- gghistogram(datos, x = "Population", bins = 10,
19                  xlab = "Población (millones)", ylab = "Frecuencia",
20                  color = "blue", fill = "blue")
21
22 g2 <- ggscatter(datos, x = "Year", y = "Population", color = "blue",
23                xlab = "Año", ylab = "Población (millones)")
24
25 # Histograma de la población y población por año
26 original <- ggarrange(g1, g2, ncol = 2, nrow = 1)
27 texto <- "Histograma de la población y población por año"
28 titulo <- text_grob(texto, face = "bold", size = 14)
29 original <- annotate_figure(original, top = titulo)
30 print(original)
31
32 # Transformaciones de la población
33 lambda_menos_dos <- -1 / (datos$Population ** 2)
34 lambda_menos_uno <- -1 / datos$Population
35 lambda_menos_un_medio <- -1 / sqrt(datos$Population)
36 lambda_cero <- log(datos$Population)
37 lambda_un_medio <- sqrt(datos$Population)
38 lambda_dos <- datos$Population ** 2
39
40 transformaciones <- data.frame(datos, lambda_menos_dos, lambda_menos_uno,
41                               lambda_menos_un_medio, lambda_cero,
42                               lambda_un_medio, lambda_dos)
43
44 # Gráficos de dispersión para la transformación de Tukey de la población y el
45 # año, usando distintos valores de lambda.
46 gt1 <- ggscatter(transformaciones, x = "Year", y = "lambda_menos_dos",
47                  color = "blue", xlab = "Año",
48                  ylab = "lambda = -2") + rotate_x_text(45)
49
50 gt2 <- ggscatter(transformaciones, x = "Year", y = "lambda_menos_uno",
51                  color = "blue", xlab = "Año",
52                  ylab = "lambda = -1") + rotate_x_text(45)
53
54 gt3 <- ggscatter(transformaciones, x = "Year", y = "lambda_menos_un_medio",
55                  color = "blue", xlab = "Año",

```



```

56         ylab = "lambda = -1/2") + rotate_x_text(45)
57
58 gt4 <- ggscatter(transformaciones, x = "Year", y = "lambda_cero",
59                 color = "blue", xlab = "Año",
60                 ylab = "lambda = 0") + rotate_x_text(45)
61
62 gt5 <- ggscatter(transformaciones, x = "Year", y = "lambda_un_medio",
63                 color = "blue", xlab = "Año",
64                 ylab = "lambda = 1/2") + rotate_x_text(45)
65
66 gt6 <- ggscatter(transformaciones, x = "Year", y = "lambda_dos",
67                 color = "blue", xlab = "Año",
68                 ylab = "lambda = 2") + rotate_x_text(45)
69
70 # Crear una única figura con todos los gráficos de dispersión.
71 dispersion <- ggarrange(gt1, gt2, gt3, gt4, gt5, gt6, ncol = 3, nrow = 2)
72 texto <- "Población transformada por año"
73 titulo <- text_grob(texto, face = "bold", size = 14)
74 dispersion <- annotate_figure(dispersion, top = titulo)
75 print(dispersion)
76
77 # Histogramas para la transformación de Tukey de la población y el año,
78 # usando distintos valores de lambda.
79 h1 <- gghistogram(transformaciones, bins = 10, x = "lambda_menos_dos",
80                 color = "blue", fill = "blue",
81                 xlab = "lambda = -2") + rotate_x_text(45)
82
83 h2 <- gghistogram(transformaciones, bins = 10, x = "lambda_menos_uno",
84                 color = "blue", fill = "blue",
85                 xlab = "lambda = -1") + rotate_x_text(45)
86
87 h3 <- gghistogram(transformaciones, bins = 10, x = "lambda_menos_un_medio",
88                 color = "blue", fill = "blue",
89                 xlab = "lambda = -1/2") + rotate_x_text(45)
90
91 h4 <- gghistogram(transformaciones, bins = 10, x = "lambda_cero",
92                 color = "blue", fill = "blue",
93                 xlab = "lambda = 0") + rotate_x_text(45)
94
95 h5 <- gghistogram(transformaciones, bins = 10, x = "lambda_un_medio",
96                 color = "blue", fill = "blue",
97                 xlab = "lambda = 1/2") + rotate_x_text(45)
98
99 h6 <- gghistogram(transformaciones, bins = 10, x = "lambda_dos",
100                 color = "blue", fill = "blue",
101                 xlab = "lambda = 2") + rotate_x_text(45)
102
103 # Crear una única figura con todos los gráficos de dispersión.
104 histograma <- ggarrange(h1, h2, h3, h4, h5, h6, ncol = 3, nrow = 2)
105 texto <- "Histogramas de la población transformada"
106 titulo <- text_grob(texto, face = "bold", size = 14)
107 histograma <- annotate_figure(histograma, top = titulo)
108 print(histograma)
109
110 # Buscar la mejor transformación de Tukey usando una función de R.
111 transformacion <- transformTukey(datos$Population, start = -4, end = 4,
112                                 int = 0.001, returnLambda = TRUE)

```

### 11.1.4 Transformaciones Box-Cox

La **transformación Box-Cox** es una versión escalada de la transformación de Tukey, dada por la ecuación 11.5:

$$x'_\lambda = \frac{x^\lambda - 1}{\lambda} \quad (11.5)$$

Si bien a primera vista parece muy diferente a la ecuación 11.4, podemos reescribirla como:

$$x'_\lambda = \frac{x^\lambda - 1}{\lambda} \approx \frac{(1 + \lambda \log(x) + \frac{1}{2}\lambda^2 \log(x)^2 + \dots) - 1}{\lambda}$$

De donde:

$$\lim_{\lambda \rightarrow 0} \frac{(1 + \lambda \log(x) + \frac{1}{2}\lambda^2 \log(x)^2 + \dots) - 1}{\lambda} = \frac{0}{0}$$

Tras aplicar la regla de l'Hôpital, obtenemos finalmente que:

$$\lim_{\lambda \rightarrow 0} x'_\lambda = \log(x)$$

Por lo que, al igual que en la escalera de potencias de Tukey, pero de forma más natural, empleamos la transformación logarítmica para  $\lambda = 0$ .

La figura 11.7 (creada con el script 11.4, líneas 39–64) muestra los gráficos de dispersión para la población total de Estados Unidos por año tras aplicar la transformación de Box-Cox con diferentes valores de  $\lambda$ . Podemos ver que el resultado se parece al que obtuvimos con la transformación de Tukey (figura 11.5).

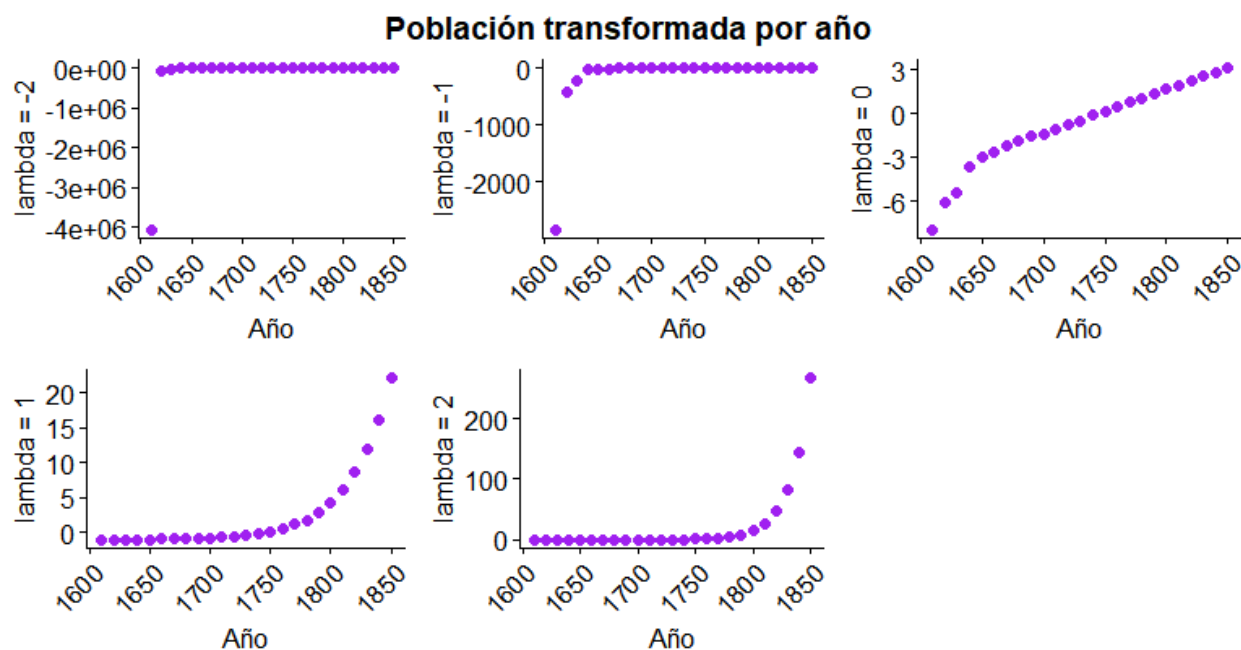


Figura 11.7: población de Estados Unidos por año tras aplicar la transformación de Box-Cox con distintos valores de  $\lambda$ .

Una característica interesante de esta transformación es que, para cualquier valor de  $\lambda$ ,  $x'_\lambda = 0$  cuando  $x = 1$ . Podemos observar esto claramente en la figura 11.8, que compara transformaciones Box-Cox con distintos valores de  $\lambda$  con  $\log(x)$ .

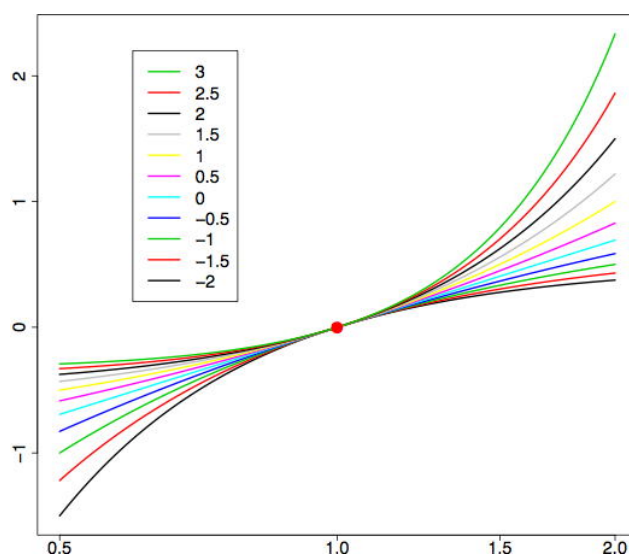


Figura 11.8: ejemplos de la transformación Box-Cox versus  $\log(x)$ . Fuente: (Lane, s.f., p. 16).

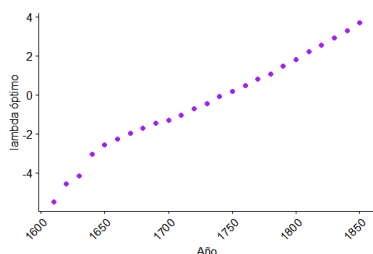
El paquete `DescTools` de R incluye varias funciones que permiten efectuar la transformación Box-Cox (Carchedi y col., s.f.). Destacan entre ellas:

- `BoxCoxLambda(x, lower, upper)`: devuelve el valor óptimo de  $\lambda$  para la transformación Box-Cox del vector `x`.
- `BoxCox(x, lambda)`: devuelve un vector correspondiente a la transformación Box-Cox de `x` con parámetro `lambda`.
- `BoxCoxInv(x, lambda)`: revierte la transformación Box-Cox del vector `x` con parámetro `lambda`.

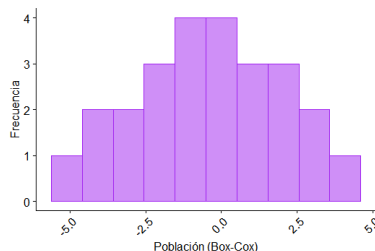
Donde:

- `x`: vector numérico.
- `lower`: límite inferior para los posibles valores de  $\lambda$ .
- `upper`: límite superior para los posibles valores de  $\lambda$ .
- `lambda`: parámetro de la transformación.

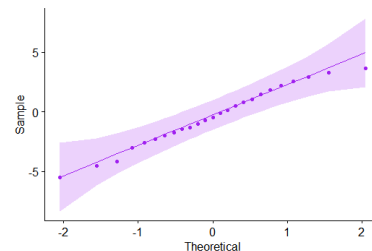
En las líneas 67–70 del script 11.4 se determina el valor óptimo del parámetro  $\lambda$ , obteniéndose como resultado  $\lambda = 0,09942656$ , para luego efectuar la transformación Box-Cox correspondiente de la población de Estados Unidos. La figura 11.9 (script 11.4, líneas 84–87) muestra gráficamente el resultado de la transformación. En el gráfico 11.9a podemos ver que la relación entre la población transformada y el año se asemeja a una recta, mientras que el histograma de la figura 11.9b se parece bastante a una distribución normal, hecho que vemos confirmado por el gráfico Q-Q de la figura 11.9c.



(a) población de Estados Unidos por año tras aplicar la transformación de Box-Cox con  $\lambda$  óptimo.



(b) histograma de la población transformada con el valor óptimo de  $\lambda$ .



(c) gráfico Q-Q de la población transformada con el valor óptimo de  $\lambda$ .

Figura 11.9: gráficos de población de Estados Unidos por año usando la transformación de Box-Cox.

Script 11.4: transformación de Box-Cox para la población total de Estados Unidos.

```

1 library(ggpubr)
2 library(DescTools)
3
4 # Cargar datos
5 Year <- c(1610, 1620, 1630, 1640, 1650, 1660, 1670, 1680, 1690, 1700, 1710,
6           1720, 1730, 1740, 1750, 1760, 1770, 1780, 1790, 1800, 1810, 1820,
7           1830, 1840, 1850)
8
9 Population <- c(0.00035, 0.002302, 0.004646, 0.026634, 0.050368, 0.075058,
10                0.111935, 0.151507, 0.210372, 0.250888, 0.331711, 0.466185,
11                0.629445, 0.905563, 1.17076, 1.593625, 2.148076, 2.780369,
12                3.929214, 5.308483, 7.239881, 9.638453, 12.86602, 17.069453,
13                23.191876)
14
15 datos <- data.frame(Year, Population)
16
17 # Transformación de Box-cox
18 box_cox <- function(x, lambda) {
19   if(lambda == 0) {
20     return(log(x))
21   }
22
23   resultado <- (x ** lambda -1) / lambda
24   return(resultado)
25 }
26
27 # Transformaciones de la población
28 lambda_menos_dos <- box_cox(datos$Population, -2)
29 lambda_menos_uno <- box_cox(datos$Population, -1)
30 lambda_cero <- box_cox(datos$Population, 0)
31 lambda_uno <- box_cox(datos$Population, 1)
32 lambda_dos <- box_cox(datos$Population, 2)
33
34 transformaciones <- data.frame(datos, lambda_menos_dos, lambda_menos_uno,
35                                lambda_cero, lambda_uno, lambda_dos)
36
37 # Gráficos de dispersión para la transformación de Box-Cox de la población y
38 # el año, usando distintos valores de lambda.
39 gt1 <- ggscatter(transformaciones, x = "Year", y = "lambda_menos_dos",
40                  color = "purple", xlab = "Año",
41                  ylab = "lambda = -2") + rotate_x_text(45)
42
43 gt2 <- ggscatter(transformaciones, x = "Year", y = "lambda_menos_uno",
44                  color = "purple", xlab = "Año",
45                  ylab = "lambda = -1") + rotate_x_text(45)
46
47 gt3 <- ggscatter(transformaciones, x = "Year", y = "lambda_cero",
48                  color = "purple", xlab = "Año",
49                  ylab = "lambda = 0") + rotate_x_text(45)
50
51 gt4 <- ggscatter(transformaciones, x = "Year", y = "lambda_uno",
52                  color = "purple", xlab = "Año",
53                  ylab = "lambda = 1") + rotate_x_text(45)
54
55 gt5 <- ggscatter(transformaciones, x = "Year", y = "lambda_dos",
56                  color = "purple", xlab = "Año",
57                  ylab = "lambda = 2") + rotate_x_text(45)
58
59 # Crear una única figura con todos los gráficos de dispersión.

```

```

60 dispersion <- ggarrange(gt1, gt2, gt3, gt4, gt5, ncol = 3, nrow = 2)
61 texto <- "Población transformada por año"
62 titulo <- text_grob(texto, face = "bold", size = 14)
63 dispersion <- annotate_figure(dispersion, top = titulo)
64 print(dispersion)
65
66 # Buscar la mejor transformación Box-Cox usando funciones de R.
67 lambda <- BoxCoxLambda(datos$Population, lower = -4, upper = 4)
68 cat("Lambda óptimo:", lambda)
69 transformacion <- BoxCox(datos$Population, lambda)
70 datos <- data.frame(datos, transformacion)
71
72 # Graficar los datos transformados.
73 g1 <- ggqqplot(transformacion, color = "purple")
74 print(g1)
75
76 g2 <- gghistogram(datos, bins = 10, x = "transformacion", color = "purple",
77                 fill = "purple", xlab = "Población (Box-Cox)",
78                 ylab = "Frecuencia") + rotate_x_text(45)
79
80 print(g2)
81
82 # Gráfico de dispersión para la transformación de Box-Cox de la población y
83 # el año, usando lambda óptimo.
84 g3 <- ggscatter(datos, x = "Year", y = "transformacion", color = "purple",
85                 xlab = "Año", ylab = "lambda óptimo") + rotate_x_text(45)
86
87 print(g3)

```

### 11.1.5 Ejercicios propuestos (sección 11.1)

- 11.1 En la década del 1920 se hicieron los primeros estudios sobre la relación entre la velocidad de un automóvil con la distancia que necesita para detenerse. Los datos de estas pruebas se pueden encontrar el conjunto `cars` del paquete `datasets`. Con ellos, responde la siguiente pregunta: en promedio, ¿se requieren más de 40 pies para detener un vehículo que viaja a más de 10 millas por hora?
- 11.2 El conjunto `airquality` del paquete `datasets` contiene mediciones diarias de la calidad del aire en la ciudad de New York, EE.UU., registradas de mayo a septiembre de 1973. Verifica si las mediciones de ozono [partes por billón] son, en promedio, iguales en mayo y junio. Usa la escalera de potencias de Tukey si los datos presentan problemas para su análisis.
- 11.3 Para el conjunto `airquality` del paquete `datasets`, verifica si las mediciones de ozono [partes por billón] son, en promedio, iguales en agosto y septiembre. Usa la transformación Box-Cox si los datos presentan problemas para su análisis.
- 11.4 El conjunto `WorldPhones` del paquete `datasets` contiene la cantidad de teléfonos (en miles) instalados para diferentes regiones del planeta en algunos años entre 1951 y 1961. Verifica si, en promedio, la cantidad de teléfonos instalados en aquella época eran iguales en Centroamérica, Sudamérica y África. Usa la escalera de potencias de Tukey si los datos presentan problemas para su análisis.
- 11.5 Para el conjunto `WorldPhones` del paquete `datasets`, verifica si, en promedio, los datos sugieren que la cantidad de teléfonos instalados en Asia, Oceanía y África era similares por aquellos años. Usa la transformación Box-Cox si los datos presentan problemas para su análisis.
- 11.6 Da un ejemplo de una pregunta de investigación sobre las asignaturas comunes en ingeniería que necesite aplicar una prueba T de Student para dos muestras independientes, pero que probablemente, por la naturaleza de las variables estudiadas, requiera una transformación de los datos. Identifica bien las variables involucradas y enuncia las hipótesis a docimar. (*Ayuda: variables que presentan distribuciones*

*exponenciales o logarítmicas, entre otras, suelen mostrar mucha asimetría. Ejemplos se encuentran en variables relacionadas con tiempos, precios, ingresos, etc.)*

- 11.7 Da un ejemplo de una pregunta de investigación sobre los conciertos realizados en Santiago que necesite aplicar una prueba T de Student para dos muestras apareadas, pero que probablemente, por la naturaleza de las variables estudiadas, requiera una transformación de los datos. Identifica bien las variables involucradas y enuncia las hipótesis a docimar. (Ver la nota de ayuda en 11.6)
- 11.8 Da un ejemplo de una pregunta de investigación sobre el estado de la salud mental de estudiantes universitarios que requiera utilizar una prueba ANOVA para muestras independientes, pero que probablemente, por la naturaleza de las variables estudiadas, requiera una transformación de los datos. Identifica bien las variables involucradas y enuncia las hipótesis a docimar. (Ver la nota de ayuda en 11.6)
- 11.9 Da un ejemplo de una pregunta de investigación sobre el estado de la salud mental de estudiantes universitarios que requiera utilizar una prueba ANOVA para medidas repetidas, pero que probablemente, por la naturaleza de las variables estudiadas, requiera una transformación de los datos. Identifica bien las variables involucradas y enuncia las hipótesis a docimar. (Ver la nota de ayuda en 11.6)
- 11.10 Justificando tus suposiciones, inventa muestras plausibles de los datos que se podrían encontrar en el estudio propuesto en la pregunta 11.6 que tengan entre 15 y 18 observaciones cada una.
- 11.11 Justificando tus suposiciones, inventa muestras plausibles de los datos que se podrían encontrar en el estudio propuesto en la pregunta 11.7 que tengan 22 pares observaciones cada una, y que cumplan las condiciones requeridas por la prueba.
- 11.12 Justificando tus suposiciones, inventa muestras plausibles de los datos que se podrían encontrar en el estudio propuesto en la pregunta 11.8 que tengan entre 21 y 28 observaciones cada una.
- 11.13 Justificando tus suposiciones, inventa muestras plausibles de los datos que se podrían encontrar en el estudio propuesto en la pregunta 11.9 con 20 casos.
- 11.14 Usando las muestras inventadas en 11.10, aplica la escalera de potencias de Tukey y verifica que los datos transformados cumplen las condiciones requeridas por la prueba de hipótesis propuesta.
- 11.15 Usando las muestras inventados en 11.11, aplica la transformación de Box-Cox y verifica que los datos transformados cumplen las condiciones requeridas por la prueba de hipótesis propuesta.
- 11.16 Usando las muestras inventadas en 11.12, aplica la escalera de potencias de Tukey y verifica que los datos transformados cumplen las condiciones requeridas por la prueba de hipótesis propuesta.
- 11.17 Usando las muestras inventados en 11.13, aplica la transformación de Box-Cox y verifica que los datos transformados cumplen las condiciones requeridas por la prueba de hipótesis propuesta.
- 11.18 Usando los datos transformados en 11.14, realiza la prueba propuesta y entrega una conclusión a la pregunta planteada.
- 11.19 Usando los datos transformados en 11.15, realiza la prueba propuesta y entrega una conclusión a la pregunta planteada.
- 11.20 Usando los datos transformados en 11.16, realiza la prueba propuesta, usando la prueba HSD de Tukey para el análisis post-hoc de ser necesario, y entrega una conclusión a la pregunta planteada.
- 11.21 Usando los datos transformados en 11.17, realiza la prueba propuesta, usando el método de Benjamini y Hochberg para el análisis post-hoc de ser necesario, y entrega una conclusión a la pregunta planteada.