

# Introducción al Análisis Exploratorio de Datos (EDA) en R

## Módulo 3

2024-03-16

- 1 Recordatorio: sintaxis de ggplot2
- 2 Introducción al EDA
- 3 Análisis exploratorio univariante
- 4 Análisis exploratorio multivariante
- 5 Detección univariante de valores atípicos
- 6 Recursos alternativos
- 7 Bibliografía de consulta



## Recordatorio: sintaxis de ggplot2

# Recordatorio: sintaxis de ggplot2

El paquete **ggplot2** proporciona un sistema coherente para visualizar datos y crear gráficos. La versatilidad de **ggplot2** radica en el uso de la Gramática de Gráficos (*Grammar of Graphics*).

```
ggplot(dataset, aes()) + geometría + faceta + opciones
```

donde:

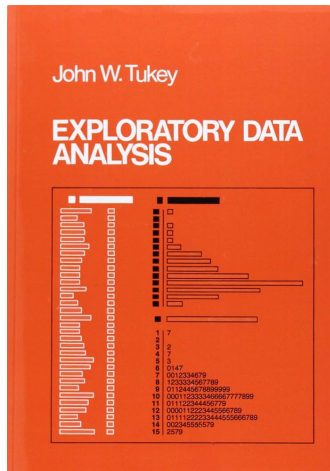
- 1 *dataset* es un data frame
- 2 Las características del mapa **aes()** describe los ejes ( $x, y$ ), el color exterior (**color** o **colour**), el color interior (**fill**), la forma de los puntos (**shape**), el tipo de línea (**linetype**) y el tamaño (**size**)
- 3 Los objetos geométricos (**geometría**) determinan el tipo de gráfico:
  - Puntos (*geom\_point*)
  - Líneas (*geom\_lines*)
  - Histogramas (*geom\_histogram*)
  - Boxplot (*geom\_boxplot*)
- 4 La **faceta** permite dividir un gráfico en múltiples gráficos de acuerdo con grupos

# Introducción al EDA

## ¿Qué es el EDA?

El Análisis Exploratorio de Datos (EDA, por sus siglas en inglés) proporciona una estrategia robusta para ampliar el entendimiento sobre los datos. El principio general es el siguiente:

“Es importante comprender lo que podemos hacer antes de aprender a medir lo bien que parece que lo hemos hecho” (Tukey, 1977).



**Figure 1:** Tukey (1977). Exploratory Data Analysis

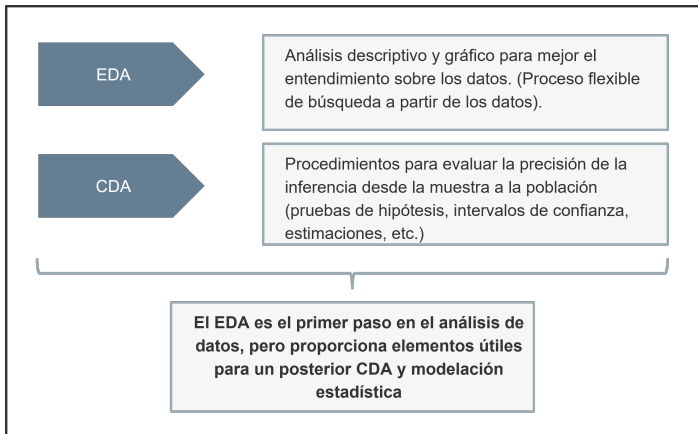
# Utilidad

Mediante métodos gráficos y descriptivos, el EDA permite ([Carranza, 2021](#)):

- Revelar la estructura de los datos
- Determinar las variables relevantes
- Determinar valores atípicos
- Proponer hipótesis
- Proponer estrategias para la modelación



# EDA y CDA



**Figure 2:** Análisis Exploratorio de Datos (EDA) y Análisis Confirmatorio de Datos (CDA)

## Análisis exploratorio univariante

## Base de datos

La base de datos usada es extraída de los microdatos de la **Gran Encuesta Integrada de Hogares (GEIH)** para diciembre de 2023. El análisis considera las siguientes 13 ciudades y áreas metropolitanas:

- Medellín A.M.
- Barranquilla A.M.
- Bogotá
- Cartagena
- Manizales A.M.
- Montería
- Villavicencio
- Pasto
- Cucuta A.M.
- Pereira A.M.
- Pereira A.M.
- Bucaramanga A.M.
- Ibagué
- Cali A.M.

La información es extraída de dos módulos de la GEIH:

- **Ocupados** (horas trabajadas , ingreso laboral, actividad económica, etc.)
- **Características generales, seguridad social en salud y educación** (edad, sexo, nivel de educación, etc.)

## Base de datos (cont.)

Para importar la base de datos (.xlsx),

```
library(readxl)
dataset <- readxl::read_excel("Datos/Formatos/geih_dataset.xlsx")
cont_ds <- dataset %>% dplyr::select(ingreso, edad, horas_semana, t_actual)
cont_ds <- cont_ds %>% dplyr::filter(t_actual > 0)
```

La siguiente tabla muestra un resumen de la base de datos:

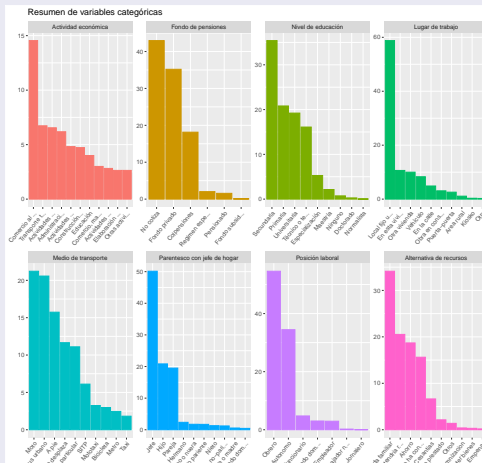
Variable	Clase	Descripción
area	Factor	Área metropolitana
dpto	Factor	Departamento
sexo	Factor	Sexo al nacer
parent	Factor	Parentesco con el jefe o jefa del hogar
edad	Númerica	Años cumplidos
edu	Factor	Mayor nivel educativo alcanzado
ingreso	Númerica	Ingreso laboral
horas_semana	Númerica	Horas trabajadas normalmente a la semana
cotiza	Factor	¿Cotiza a un fondo de pensiones?
lugar	Factor	Lugar principal de trabajo
meses	Númerica	¿Cuántos meses trabajó en los últimos 12 meses?
rama_4	Cadena	Rama de actividad CIU REV 4 (4 dígitos)
rama_2	Cadena	Rama de actividad CIU REV 4 (4 dígitos)
posic	Factor	Posición laboral
fondo	Factor	¿A cuál fondo cotiza?
cambiar	Factor	¿Desea cambiar su trabajo?
estable	Factor	¿Considera que su empleo es estable?
t_actual	Númerica	¿Cuánto tiempo lleva en su empleo actual?
t_viaje	Númerica	Tiempo de desplazamiento al trabajo
mas_h	Factor	¿Quiere trabajar más horas?
medio	Factor	Medio de transporte
sintrab	Factor	¿Si no tuviera trabajo, de dónde obtendría sus recursos?
n_comp	Factor	¿Cuántas personas tiene la empresa donde trabaja?
srl	Factor	¿Afiliación a ARL?
caja	Factor	¿Afiliación a caja de compensación familiar?
actividad	Factor	Actividad económica recodificada
cotiza_fondo	Factor	Fondo de pensiones recodificado
factor_exp	Númerica	Factor de expansión

# Datos cualitativos:

## Resumen de datos cualitativos

Considérese las siguientes variables cualitativas previamente recodificadas:

- La función `forcats::fct_lump_n()` es usada para agregar las categoría en “otros”.
- Las gráficas muestran las 10 categorías más frecuentes.
- La función `ggplot2::facet_wrap()` es usada para obtener los gráficos múltiples



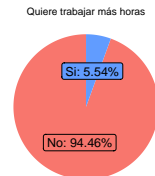
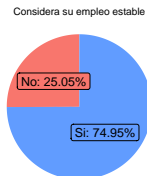
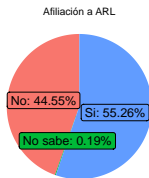
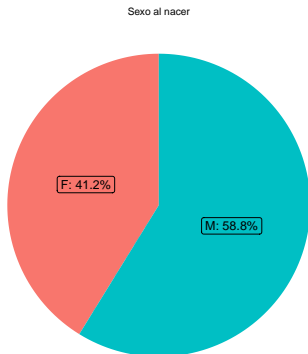
## Resumen de datos cualitativos (cont.):

Idéntica información puede ser representada mediante el siguiente resumen:

	N	Proporción (%)
<b>Actividad económica</b>		
Actividades de atención de la salud humana	53	0.60
Actividades de los hogares individuales como empleadores de personal doméstico	31	0.35
Actividades de servicios de comidas y bebidas	72	0.82
Administración pública y defensa; planes de seguridad social de afiliación obligatoria	68	0.77
Comercio al por menor (incluso el comercio al por menor de combustibles), excepto el de vehículos automotores y motocicletas	160	1.82
Comercio, mantenimiento y reparación de vehículos automotores y motocicletas, sus partes, piezas y accesorios	33	0.38
Construcción de edificios	52	0.59
Educación	44	0.50
Elaboración de productos alimenticios	29	0.33
Other	455	5.17
Otras actividades de servicios personales	29	0.33
Transporte terrestre; transporte por tuberías	74	0.84
<b>Fondo de pensiones</b>		
Colpensiones	227	2.58
Fondo privado	360	4.09
No cotiza	464	5.27
Pensionado	23	0.26
Regimen especial	26	0.30
<b>Educación</b>		
Doctorado	1	0.01
Especialización	62	0.70
Maestría	24	0.27
Ninguno	9	0.10
Primaria	224	2.55
Secundaria	391	4.44
Técnico o tecnológico	177	2.01
Universitaria	212	2.41

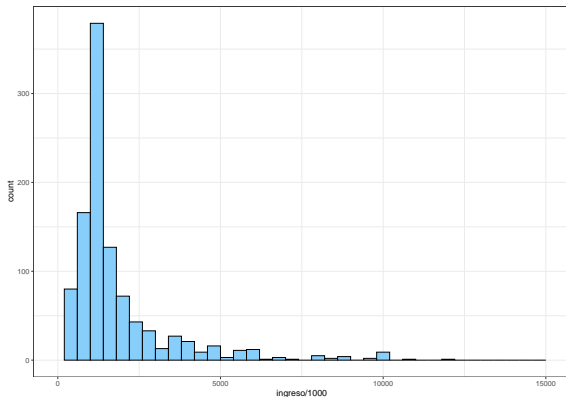
## Resumen de datos cualitativos (cont.):

Un panel de **gráficos circulares** es útil para presentar un resumen sobre las variables categóricas con un número menor de niveles. Considérese las siguientes variables



# Datos cuantitativos

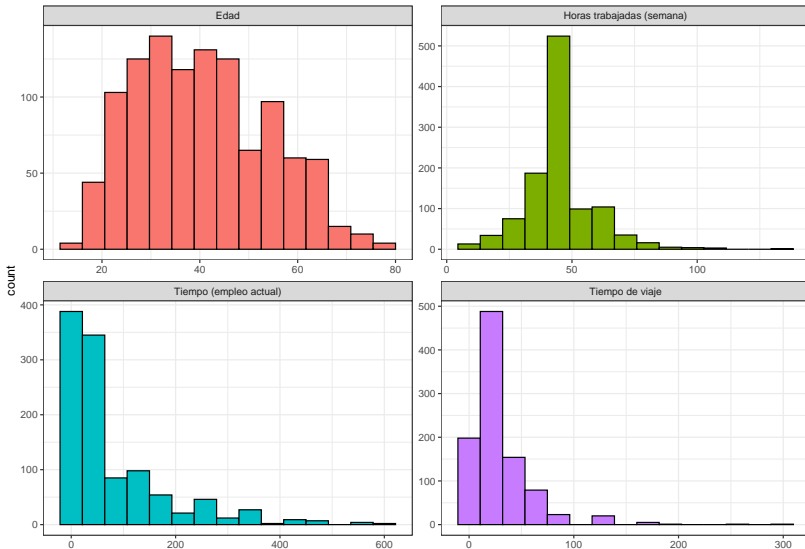
El **Histograma** es una representación gráfica de los datos que muestra la frecuencia de los casos (valores) en categorías de datos (véase la tabla inferior).



[0,1000]	(1000,2000]	(2000,3000]	(3000,4000]	(4000,5000]	(5000,6000]	(6000,7000]	(7000,8000]
246	556	98	57	29	25	5	6

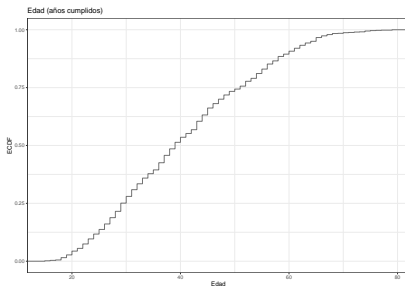
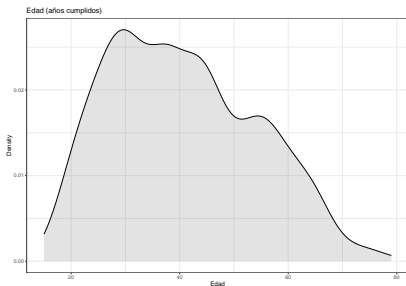


# Histograma



## PDF y ECDF

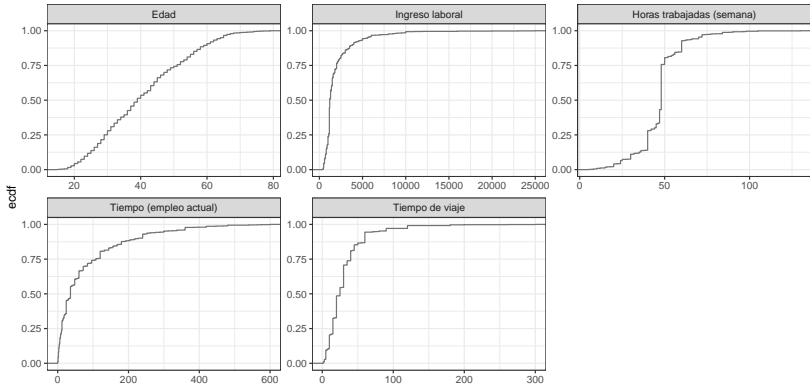
La **función de densidad empírica (PDF)**  $f(x)$  y la **función de distribución acumulada empírica (ECDF)**  $F(x)$  son obtenidas mediante las funciones `density()` y `ecdf()` del paquete Stats.



La función de densidad  $f(x)$  satisface que  $\int_a^b f(x)dx = P[a \leq X \leq b]$ , donde  $P[a \leq X \leq b]$  significa la probabilidad de que  $X$  se encuentre en el intervalo  $a$  a  $b$ . Por definición,  $F(x) = P(X \leq x)$ , es decir, expresa la probabilidad de que  $X$  no sea mayor que el valor de  $x$ .

# Función de Distribución Acumulada Empírica

Un ejercicio análogo es implementado para las variables cuantitativas restantes.



## Diagrama de caja

Un **diagrama de caja** es un resumen gráfico de los datos con base en la los cuartiles Q1 y Q3, la mediana y el rango intercuartílico (RIC). El siguiente diagrama muestra su elaboración:

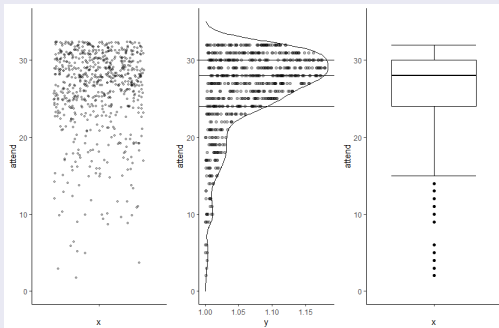


Figure 3: Construcción de un diagrama de caja

# Un resumen de datos cuantitativos usando StatDA()

La librería StatDA() proporciona una utilidad para representar la distribución y los principales elementos descriptivos de las variables continuas.

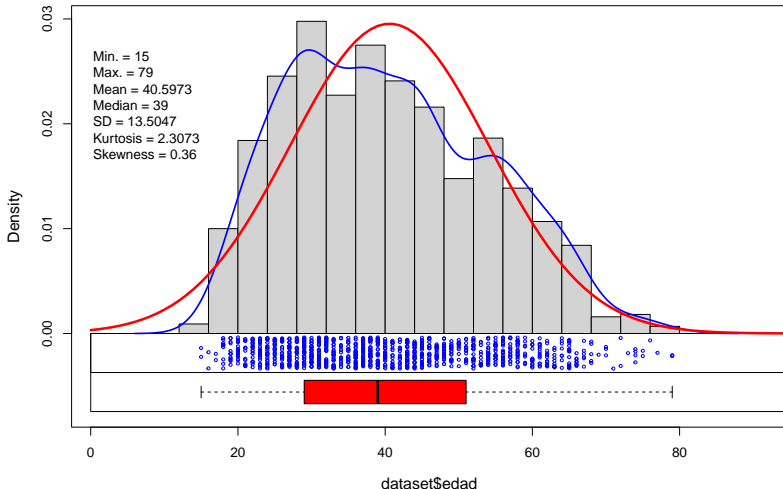
```
library(StatDA)
library(moments)

me = mean(dataset$edad)
sd = sd(dataset$edad)

StatDA::edaplot(dataset$edad, scatter=TRUE, H.freq=FALSE, box=TRUE,
  H.breaks=seq(0,100, by = 4),H.col="lightgray", H.border=TRUE, H.labels=FALSE,
  S.pch=1, S.col="blue", S.cex=0.5, D.lwd=2, D.lty=1, D.plot=FALSE,
  P.xlim=c(1, 91), P.cex.lab =1.2, P.log=FALSE, P.main="Histograma,
  función de densidad, gráfico de dispersión y diagrama de caja de la edad",
  P.xlab="Edad (años)", P.plot=TRUE,
  P.ylab="Densidad",
  B.pch=1,B.cex=0.5, B.col="red")
lines(density(dataset$edad), lwd=2, col='blue')
curve(dnorm(x, mean=me, sd=sd), from=0, to=100, add=T,
  col='red', lwd=3)
leg.txt <- c(paste0("Min. = ", round(min(dataset$edad),4)),
  paste0("Max. = ", round(max(dataset$edad),4)),
  paste0("Mean = ", round(mean(dataset$edad),4)),
  paste0("Mediana = ", round(median(dataset$edad),4)),
  paste0("SD = ", round(sd(dataset$edad),4)),
  paste0("Kurtosis = ", round(kurtosis(dataset$edad),4)),
  paste0("Skewness = ", round(skewness(dataset$edad),4)))
legend(x=-3, y=0.028, bty="n", leg.txt)
```

## Resumen de datos cuantitativos (cont.)

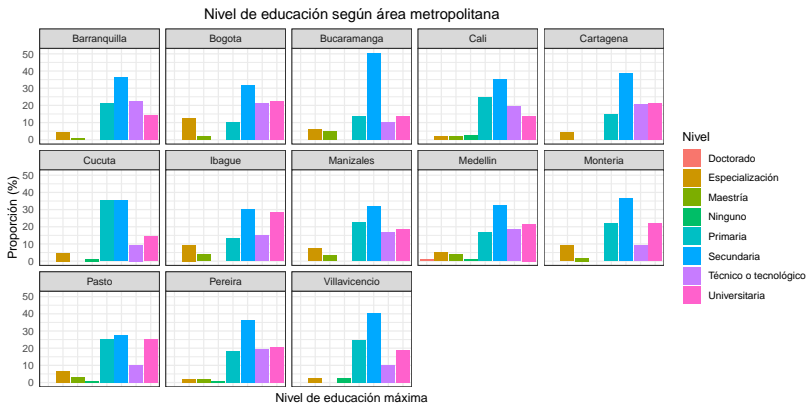
Histograma,  
función de densidad, gráfico de dispersión y diagrama de caja de la edad



## Análisis exploratorio multivariante

# Independencia de variables categóricas

La siguiente gráfica representa la asociación entre el nivel educativo máximo alcanzado por el trabajador (**educación**) y el área metropolitana en que reside.





# Tabla de contingencia

La **tabla de contingencia** permite resumir la información reportada en la figura anterior:

	Doctorado	Especialización	Maestría	Ninguno	Primaria	Secundaria	Técnico o tecnológico	Universitaria	Sum
Barranquilla	0	4	1	0	19	33	20	13	90
Bogotá	0	11	2	0	9	28	19	20	89
Bucaramanga	0	5	4	0	11	40	8	11	79
Cali	0	2	2	3	27	38	21	15	108
Cartagena	0	4	0	0	14	36	19	20	93
Cúcuta	0	4	0	1	30	30	8	12	85
Ibagué	0	5	2	0	7	16	8	15	53
Manizales	0	7	3	0	22	31	16	18	97
Medellín	1	5	4	1	17	33	19	22	102
Montería	0	5	1	0	12	20	5	12	55
Pasto	0	6	3	1	22	24	9	22	87
Pereira	0	2	2	1	17	34	18	19	93
Villavicencio	0	2	0	2	17	28	7	13	69
Sum	1	62	24	9	224	391	177	212	1100

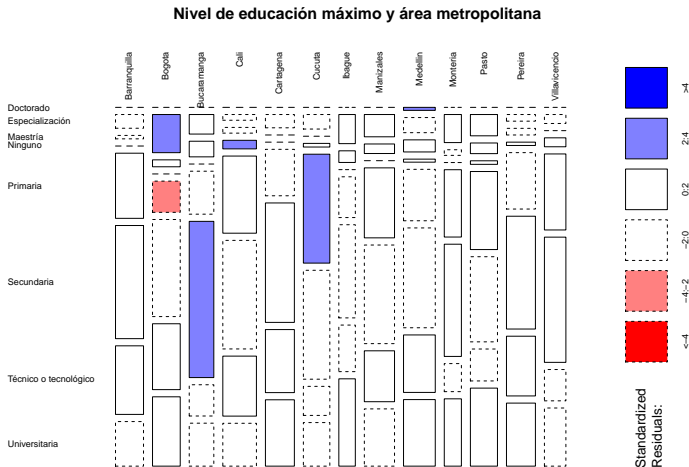
La **prueba  $\chi^2$  de independencia** permite determinar si dos variables asociadas a **una muestra** son independientes o no.

```
##
## Pearson's Chi-squared test
##
## data: dataset$area and dataset$edu
## X-squared = 109.12, df = 84, p-value = 0.03415
```

$H_0$  : Las dos variables son independientes

## Tabla de contingencia (cont.)

Los **gráficos de mosaico** proporcionan una visualización sobre las tablas de contingencia:



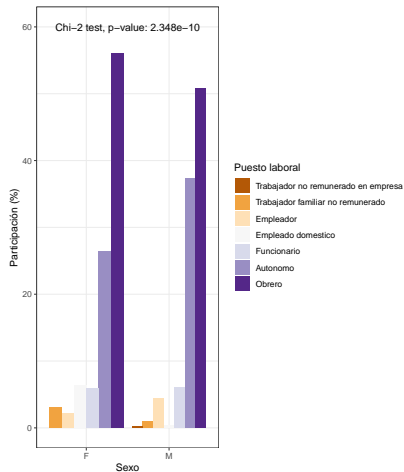
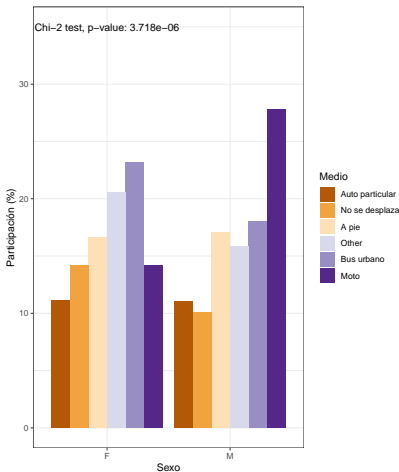
## Resumen: independencia de variables categóricas

Usando tablas de contingencia, un análisis análogo es implementado para las variables categóricas restantes:

V1	actividad	arl	caja	cotiza_fondo	edu	estable	medio	parent	posic	sexo
actividad	<0.001	0.013	0.046	<0.001	0.035	0.001	0.004	0.017	<0.001	<0.001
arl	0.013	<0.001	<0.001	<0.001	0.002	<0.001	<0.001	0.005	0.001	0.186
caja	0.039	<0.001	<0.001	<0.001	0.003	<0.001	0.008	0.452	<0.001	<0.001
cotiza_fondo	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	0.048	<0.001	0.004
edu	0.036	0.002	0.004	<0.001	<0.001	<0.001	0.004	0.096	0.001	0.006
estable	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	0.434	<0.001	0.58
medio	0.006	<0.001	0.013	<0.001	0.003	<0.001	<0.001	0.001	0.004	<0.001
parent	0.02	0.006	0.428	0.052	0.101	0.429	0.002	<0.001	0.004	<0.001
posic	<0.001	0.002	<0.001	<0.001	0.002	<0.001	0.003	0.003	<0.001	<0.001
sexo	<0.001	0.18	0.001	0.006	0.002	0.573	<0.001	<0.001	<0.001	<0.001

## Resumen: independencia de variables categóricas

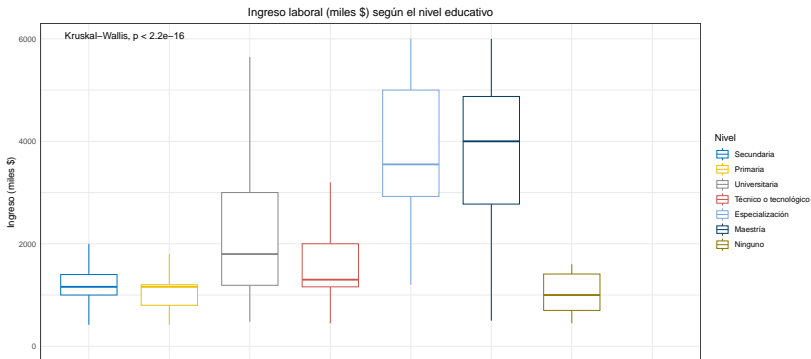
Una forma alternativa de presentar los resultados sobre la independencia de variables categóricas a través de un gráfico:



## Diferencias en variables continuas

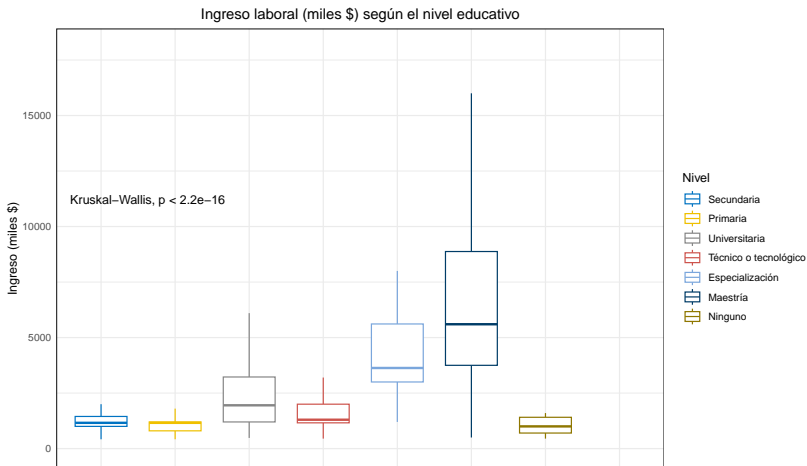
La **prueba de Kruskal-Wallis** es un prueba no-paramétrica cuyo objetivo es determinar si todas las  $k$  poblaciones son idénticas o si al menos una de las poblaciones proporciona observaciones distintas a las poblaciones restantes. Así:

$H_0$  : No hay diferencia entre las  $k$  poblaciones



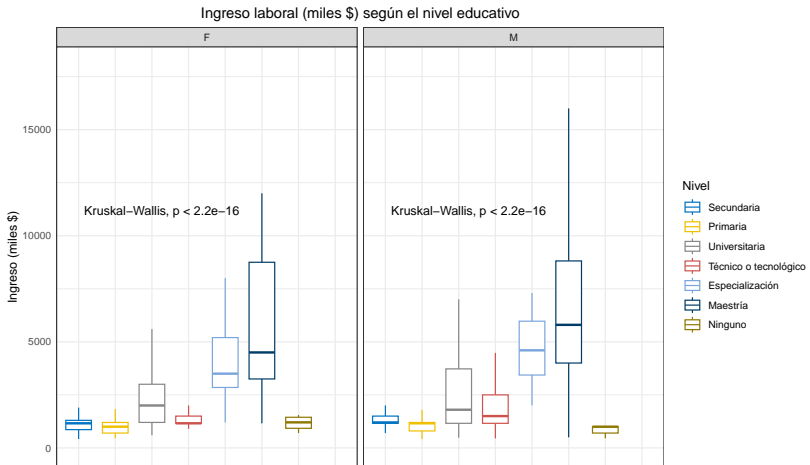
## Diferencias en variables continuas (cont.)

Nótese que la **prueba de Wilcoxon** puede ser usada con el mismo propósito cuando  $k = 2$ , es decir, para considerar pares de poblaciones.



## Diferencias en variables continuas (cont.)

Un ejercicio análogo es realizado después de incluir otra variable categórica:



## Resumen: diferencias en variables continuas según variables categóricas

El análisis puede ser extendido a todas las variables (continuas y categóricas) incluidas en la base de datos:

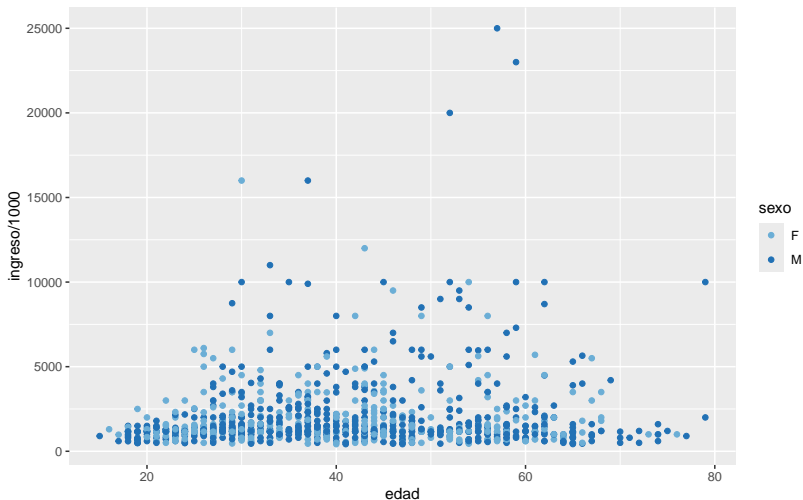
	area	sexo	parent	edu	lugar	medio	fondo	actividad
edad	18.1 (12)	2.9 (1)	250.5 (10)	113.1 (7)	77.7 (9)	74.5 (5)	80.4 (4)	133.2 (67)
horas_semana	25.5 (12)	13.2 (1)	14.7 (10)	33.4 (7)	43.8 (9)	14.2 (5)	12.4 (4)	133.4 (67)
ingreso	28.3 (12)	1.5 (1)	27 (10)	305.7 (7)	130.1 (9)	171.1 (5)	241.4 (4)	300.6 (67)
t_actual	15.9 (12)	3 (1)	73 (10)	22.4 (7)	41.7 (9)	51.6 (5)	57.3 (4)	141.9 (67)
t_viaje	91.6 (12)	5 (1)	3.4 (9)	6.1 (7)	65.9 (9)	273.6 (4)	45 (4)	135 (67)

	area	sexo	parent	edu	lugar	medio	fondo	actividad
edad	18.1 (12)	2.9 (1)	250.5 (10)	113.1 (7)	77.7 (9)	74.5 (5)	80.4 (4)	133.2 (67)
horas_semana	25.5 (12)	13.2 (1)	14.7 (10)	33.4 (7)	43.8 (9)	14.2 (5)	12.4 (4)	133.4 (67)
ingreso	28.3 (12)	1.5 (1)	27 (10)	305.7 (7)	130.1 (9)	171.1 (5)	241.4 (4)	300.6 (67)
t_actual	15.9 (12)	3 (1)	73 (10)	22.4 (7)	41.7 (9)	51.6 (5)	57.3 (4)	141.9 (67)
t_viaje	91.6 (12)	5 (1)	3.4 (9)	6.1 (7)	65.9 (9)	273.6 (4)	45 (4)	135 (67)

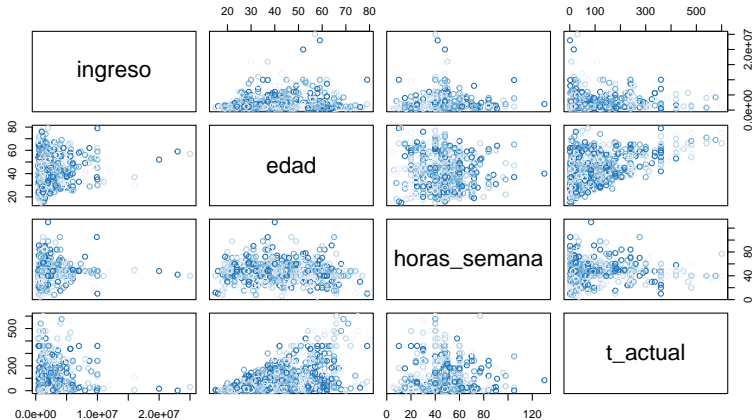


# Covariación en variables continuas

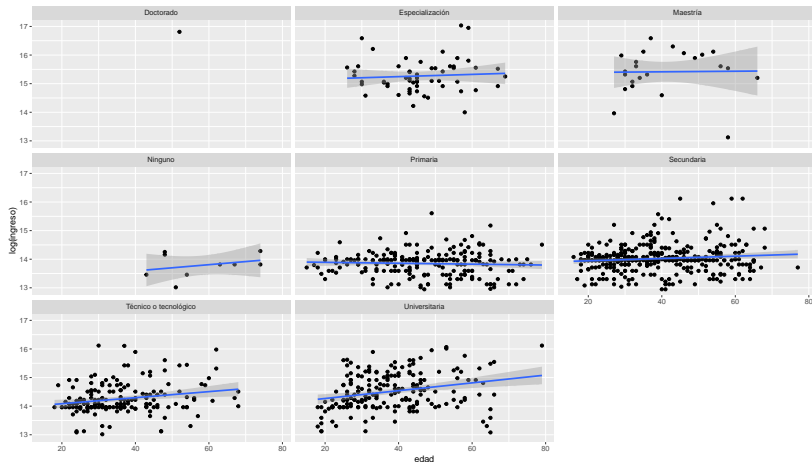


# Covariación con múltiples variables continuas

Diagrama de dispersión sobre las variables continuas



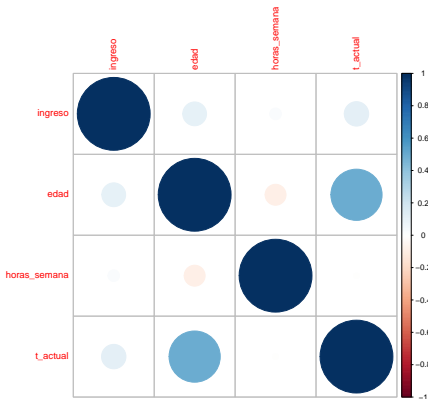
# Ajuste con múltiples variables continuas



# Correlación lineal simple

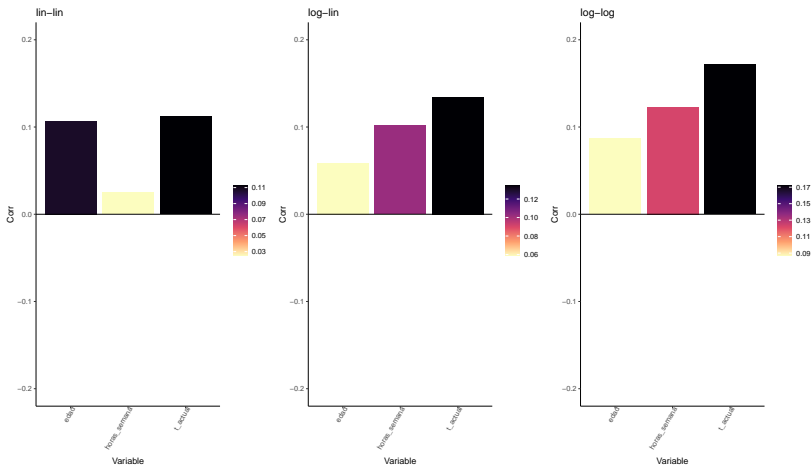
El **coeficiente de correlación simple** es una medida sobre la intensidad de la relación lineal entre dos variables cuantitativas.

```
##           ingreso  edad horas_semana t_actual
## ingreso      1.00  0.11      0.03      0.11
## edad         0.11  1.00     -0.08      0.49
## horas_semana 0.03 -0.08      1.00     -0.01
## t_actual      0.11  0.49     -0.01      1.00
```



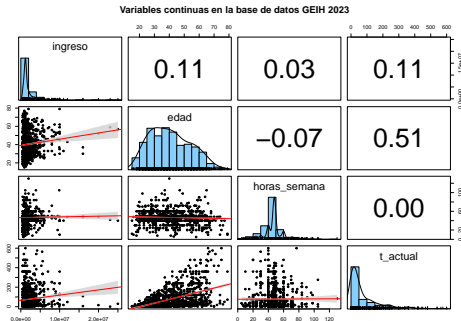
## Correlación sobre variables transformadas

El **coeficiente de correlación** es calculado, adicionalmente, considerando sobre la variable **ingreso** y las variables continuas restantes.



## Resumen sobre variables continuas: psych()

La función `pairs.panels()` de la librería `psych()` proporciona el siguiente resumen: **diagramas de dispersión y ajuste lineal** (bajo la diagonal principal); el **coeficiente de correlación de Pearson** (sobre la diagonal principal) y los **histogramas** para cada variable continua (en la diagonal principal).



## Detección univariante de valores atípicos

## Valores atípicos

La perspectiva univariante selecciona como **valores atípicos** u **outliers** aquellas observaciones que caen fuera de los rangos de la distribución. Un valor atípico se puede producir por alguna de las siguientes cuatro causas ([Aldás y Uriel, 2017](#)):

- **Errores en los datos:** errores en la recolección o introducción de los datos.
- **Errores voluntarios:** errores intencionados en la respuesta del entrevistado.
- **Errores de muestreo:** errores que son el resultado de introducir en la muestra a individuos pertenecientes a una población distinta a la **población objetivo**.
- **Outliers legítimos:** caso de la población objetivo que, por la variabilidad de las muestras, difiere del resto de casos.



# Detección univariada de valores atípicos

Considérese las siguientes alternativas para la detección univariante de *outliers*:

## Criterio intercuartílico

$$x^* \in [q_{0.25} - 1.5/IQR, q_{0.75} + 1.5/IQR]$$

## Criterio de valores estandarizados (Hair et al., 2014)

- Para muestras pequeñas ( $n < 80$ ),  $x^*$  tiene valores estándar de 2.5 o superiores.
- Para muestras mayores ( $n \geq 80$ ),  $x^*$  tiene valores estándar de 3-4 o superiores.

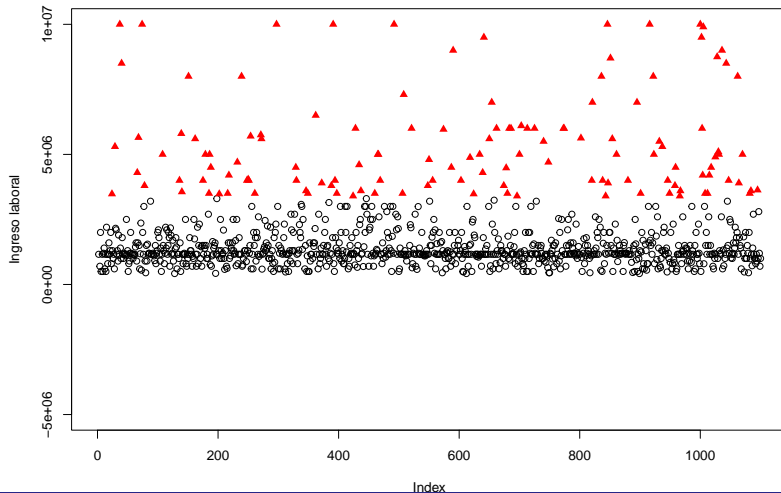
## Test de Grubbs

El **Test de Grubbs** supone la normalidad de la distribución (Grubbs, 1969; Stefansky, 1971). La hipótesis nula (no hay *outliers*) se rechaza si

$$G > \frac{n-1}{n} \sqrt{\frac{t_{(\alpha/2n, n-2)}^2}{n-2 + t_{(\alpha/2n, n-2)}^2}}$$

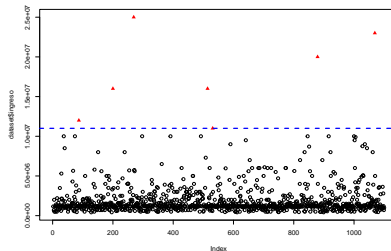
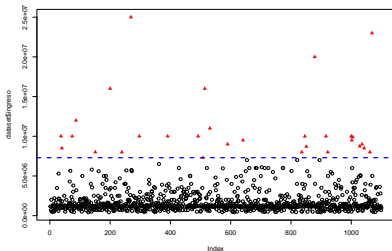
donde  $G = \max |x_i - \bar{x}| / \sigma$

# Criterio intercuartílico



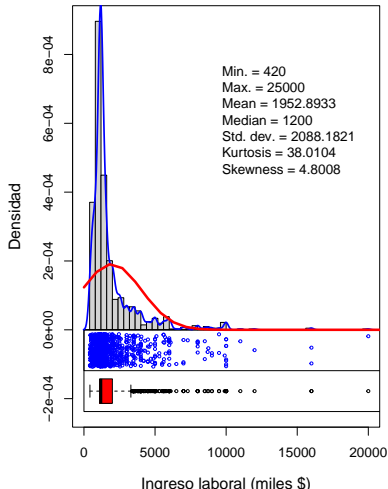
# Criterio de valores estandarizados

Considérese la detección de valores atípicos mediante los siguientes valores estandarizados:  $x^*$  es un valor atípico cuando  $z^* \geq 2.5$  (**figura A**); y  $x^*$  es un valor atípico cuando  $z^* \geq 4$  (**figura B**).

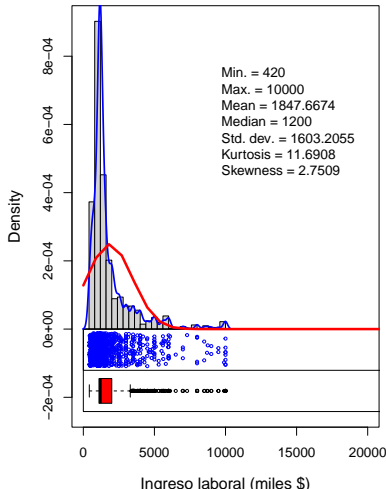


# Resumen de la variable continua con y sin outliers

## Resumen con outliers



## Resumen sin outliers



# Alternativas

Algunas alternativas para el tratamiento de valores atípicos son las siguientes:

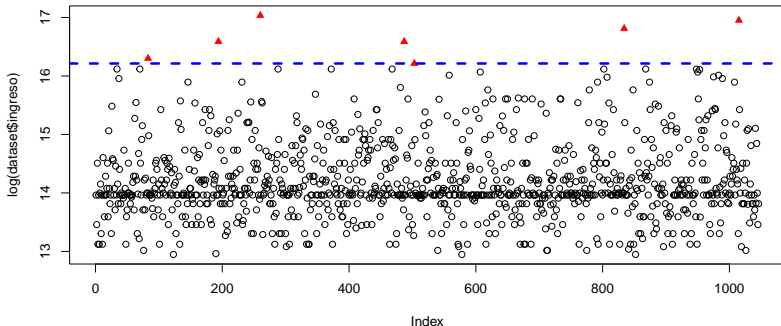
- 1 Eliminación de los valores atípicos para garantizar estimaciones correctas sobre la mayoría de la población ([Judd et al., 2009](#)).
- 2 Suavizar la influencia de los valores atípicos mediante el uso de transformaciones (raíces o logaritmos) para reducir su rango ([Hamilton, 1992](#)).
- 3 Análisis estadístico **robusto**.

## Desventajas

- Pérdida de información
- No todas las transformaciones conservan el sentido teórico de la escala original.

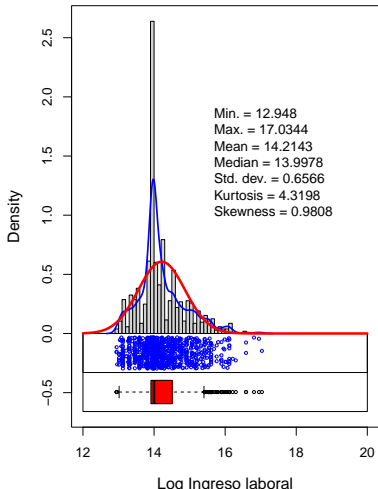
## Transformaciones para reducir su rango

En lo sucesivo, consideramos el efecto de una transformación logarítmica. Nótese que la transformación logra **reducir el rango** y suavizar, en consecuencia, la influencia de los valores atípicos. (La línea azul indica el umbral a partir del cual un valor es considerado atípico según el **criterio de valores estandarizados**).

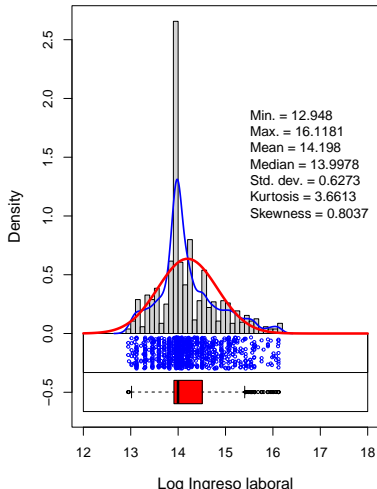


# Transformaciones

## Resumen con outliers



## Resumen sin outliers



## Recursos alternativos



## Recursos alternativos

- La librería `swirl` proporciona un tutorial sobre elementos básicos en R

```
install.packages("swirl")  
library (swirl)  
swirl()
```

- Data wrangling with `dplyr` and `tidyr` (Cheat Sheet): [Recurso 1.2](#)
- Visualización de datos usando `ggplot2` (Guía Rápida): [Recurso 1.3](#)
- Factors with `forcats` (Cheat Sheet): [Recurso 1.4](#)

## Bibliografía de consulta

## Bibliografía de consulta

- Wickham, H. (2016) GGplot2. Elegant Graphics for Data Analysis. Springer
- Golemund, G. (2014). Hands-On Programming with R. O'Reilly Media: Sebastopol, CA.
- Schutt, R. & O'Neil, C. (2014). Doing Data Science. O'Reilly Media: Sebastopol, CA.
- Wickham & Golemund, G. (2016). R for Data Science: Import, Tidy, Transform, Visualize, and Model Data. O'Reilly Media: Sebastopol, CA.
- Aldás J. & Uriel, E. (2017). Análisis multivariante aplicado con R. Madrid: Ediciones Paraninfo