

Fundamentos de Análisis de Datos

Unidad 5

Análisis Exploratorio de Datos
(EDA)

Referencias

- Maldonado, S. (2022). *Analytics y Big Data: ciencia de los Datos aplicada al mundo de los negocios*: (1 ed.). RIL editores.
<https://elibro.net/es/ereader/bibliotecautpl/225562?page=8> (apartado 2.1, pág 41)
- Diego, I. M. D. & Fernández Isabel, A. (2020). *Ciencia de datos para la ciberseguridad*: (1 ed.). RA-MA Editorial.
<https://elibro.net/es/ereader/bibliotecautpl/222714?page=129> (apartado Análisis exploratorio de datos, pág 114)

Agenda

- Introducción al EDA
 - Conceptos
 - Objetivos
 - Pasos
 - Herramientas
- Comprensión del conjunto de datos
 - Detección y tratamiento de valores faltantes
 - Identificación de tipos de variables
- Análisis univariado de variables categóricas
- Análisis univariado de variables numéricas
 - Medidas de tendencia central y de dispersión
 - Gráficas de distribución, dispersión, y correlación
 - Detección y tratamiento de outliers

Introducción al Análisis Exploratorio de Datos (EDA)

El Análisis Exploratorio de Datos (AED), conocido como Exploratory Data Analysis (EDA) en inglés, es un enfoque crucial en el análisis de datos que busca comprender y resumir las características principales de un conjunto de datos

El EDA es una etapa crítica en cualquier proyecto de análisis de datos o de ciencia de datos. Proporciona una comprensión profunda de los datos, lo que es esencial para tomar decisiones informadas sobre el procesamiento y modelado de los datos

Introducción al EDA

OBJETIVOS

- **Comprender la estructura del conjunto de datos:** Identificar las variables y sus tipos (numéricas, categóricas, fechas, etc.).
- **Resumir las características principales:** Usar estadísticas descriptivas y visualizaciones.
- **Detectar anomalías y valores faltantes:** Identificar y tratar outliers y valores ausentes.
- **Descubrir patrones y relaciones:** Identificar correlaciones y tendencias entre variables.
- **Preparar los datos para el modelado:** Realizar transformaciones y limpieza de datos.

Introducción al EDA

PASOS

1. Carga y visualización inicial de los datos

Importar datos desde diferentes fuentes / Visualizar una muestra de los datos para familiarizarse con ellos.

2. Detección y corrección de valores faltantes

Identificar la presencia de valores faltantes. / Decidir cómo manejarlos (eliminación, imputación, etc.)

3. Estadísticas descriptivas

Calcular medidas de tendencia central / Calcular medidas de dispersión / Obtener resúmenes estadísticos generales.

4. Detección de valores atípicos

Usando el Rango Intercuartil (IQR) / Usando Diagramas de Caja (Boxplots) / Usando el Test de Grubbs

5. Visualización de datos

Histogramas y diagramas de densidad / Diagramas de cajas (boxplots) / Diagramas de dispersión (scatter plots) / Gráficos de barras / Matrices de correlación

6. Análisis de correlación

Calcular y visualizar matrices de correlación. / Identificar relaciones lineales y no lineales entre variables.

7. Filtrado y limpieza de datos

Filtrar datos irrelevantes o fuera de rango. / Corregir errores tipográficos e inconsistencias

8. Transformación y enriquecimiento de datos

Aplicar transformaciones logarítmicas, de raíz cuadrada, estandarización y normalización. / Crear variables derivadas si es necesario.

Introducción al EDA

HERRAMIENTAS

- **R y Python:** Lenguajes de programación muy utilizados para análisis de datos y muy efectivos para realizar EDA.
- **Bibliotecas de R:** dplyr, tidyr, ggplot2, corrplot, data.table.
- **Bibliotecas de Python:** pandas, numpy, matplotlib, seaborn, scipy.

Introducción al EDA

RESUMEN

El **EDA** es el primer paso clave en un análisis de datos. Su objetivo es **comprender qué hay en el conjunto de datos, detectar errores o valores atípicos, explorar relaciones entre variables y preparar los datos para el modelado posterior**. No se trata solo de ver los datos, sino de analizarlos críticamente.



Comprensión del conjunto de datos

Alcance

- Aquí se identifican los tipos de variables (numéricas, categóricas), sus nombres, si hay valores faltantes, etc.

```
5 # Cargar datos
6 df <- read.csv("titanic.csv")
7
8 dim(df)           # Cantidad de observaciones y variables
9 colnames(df)      # Nombres de las variables
10 head(df)         # Ver las primeras filas
11 str(df)          # Ver la estructura de los datos
12 summary(df)      # Estadísticas fundamentales de cada variable
13 class(df$Age)     # Consulta la clase lógica de una variable
14 unique(df$Pclass) # Obtener dominio de valores de una variables
15
16 # Estadísticas valores nulos (NA)
17 colSums(is.na(df)) # Cantidad valores NA
18 colMeans(is.na(df)) # Media valores NA
19
20 # Estadísticas valores faltantes (nulos y vacíos)
21 sapply(df, function(x) sum(is.na(x) | !nzchar(x))) # Cantidad valores faltantes
22 sapply(df, function(x) mean(is.na(x) | !nzchar(x))) # Media valores faltantes
23
24 # Inspeccionar observaciones con valores faltantes
25 df[is.na(df$Embarked) | !nzchar(df$Embarked),]
```

Comprensión del conjunto de datos

Tratamiento de valores faltantes

- Buscar variables que contengan valores NULOS o cadenas vacías
- Determinar la proporción y cantidad de valores faltantes
- Inspeccionar los valores faltantes
- Determinar el tratamiento a aplicar para dichos valores:
 - Eliminación de filas
 - Eliminación de columnas
 - Imputación de valores nulos
 - Relleno adelante o atrás

Comprensión del conjunto de datos

Tratamiento de valores faltantes

```
28 library(dplyr)
29
30 # Opción 1: Eliminación de filas
31 df1 <- df %>% filter(!is.na(Embarked) & nzchar(Embarked))
32
33 # Opción 2: Eliminación de columnas
34 df2 <- df %>% select(-Cabin)
35
36 # Opción 3: Imputación de valores nulos (media, mediana, u otro)
37 # Ej: asignar la mediana del resto de observaciones
38 mediana_Age <- median(df$Age, na.rm = TRUE)
39 df3 <- df %>% mutate(Age = ifelse(is.na(Age), mediana_Age, Age))
40
41 # Opción 4: Imputación con valores múltiples
42 # Ej: asignar la media de la edad según la clase a la que pertenece
43 df4 <- df %>%
44   mutate (Age = case_when(
45     is.na(Age) & Pclass == 1 ~ mean(Age[Pclass == 1], na.rm = TRUE),
46     is.na(Age) & Pclass == 2 ~ mean(Age[Pclass == 2], na.rm = TRUE),
47     is.na(Age) & Pclass == 3 ~ mean(Age[Pclass == 3], na.rm = TRUE),
48     TRUE ~ Age
49   ))
```

Comprensión del conjunto de datos

Identificación de tipos de variables

- Como parte del EDA es muy importante identificar los tipos de variables del dataset, en el contexto estadístico. Que pueden ser:
 - **Numéricas/cuantitativas:** discretas o continuas
 - **Categóricas:** nominales u ordinales
 - **Bimodales:** Que pueden ser tratadas como numéricas o categóricas dependiendo del tipo de análisis a realizar.
 - **Variables de texto:** no relevantes desde el contexto estadístico
- En función del tipo de variables sabremos como tratarlas en el contexto de EDA, y para otros tipos de análisis.

Comprensión del conjunto de datos

Identificación de tipos de variables

```
58 names(df)
59 sapply(df, function(x) mode(x))
60 str(df)
61 unique(df$Parch)
62 unique(df$Ticket)
63 unique(df$Cabin)
64 unique(df$Embarked)
65
66 # TIPOS DE VARIABLES EN DATASET TITANIC
67
68 # Numericas/cuantitativas:
69 #                               Age, Fare
70 # Categóricas:
71 #                               Survived, Pclass, Sex, Embarked
72 # Bimodales:
73 #                               SibSp, Parch
74 # Variables texto:
75 #                               PassengerId, Name, Ticket, Cabin
```

Análisis univariado V. Categóricas

Sobrevista

- Se analiza una sola variable categórica para ver cómo se distribuyen sus categorías
- **Importante:** En R, las variables categóricas deben convertirse en **factores** antes de usarlas en gráficos o análisis estadísticos:
 - Los factores permiten que R entienda que la variable representa categorías, no texto libre.
 - Mejora la interpretación en modelos estadísticos y visualizaciones.
 - Permite controles precisos sobre el orden y agrupamiento de los niveles.
 - Si una variable no se va a usar en modelos de análisis y visualización, no es necesario factorizar

Análisis univariado V. Categóricas

¿Cómo factorizar variables categóricas?

```

84 str(df)
85
86
87
88
89
90 # Convertir variables categóricas a factores
91 df$Sex <- as.factor(df$Sex)
92 df$Pclass <- as.factor(df$Pclass)
93
94 str(df)
95 sapply(df, function(x) mode(x))
96
97 # o
98
99 df <- df %>%
100   mutate(Sex = as.factor(Sex),
101          Pclass = as.factor(Pclass))
102

```

```
$ PassengerId: int 1 2 3 4 5 6 7 8 9 10 ...
$ Survived : int 0 1 1 1 0 0 0 0 1 1 ...
$ Pclass : int 3 1 3 1 3 3 1 3 3 2 ...
$ Name : chr "Braund, Mr. Owen Harris"
"Malenka, Mrs. Jacques Heath (Lily May
$ Sex : chr "male" "female" "female"
$ Age : num 22 38 26 35 35 NA 54 2 27
$ SibSp : int 1 1 0 1 0 0 0 3 0 1 ...
$ Parch : int 0 0 0 0 0 0 0 1 2 0 ...
$ Ticket : chr "A/5 21171" "PC 17599" "S
```

```
$ PassengerId: int 1 2 3 4 5 6 7 8 9 10 ...
$ Survived : int 0 1 1 1 0 0 0 0 1 1 ...
$ Pclass : Factor w/ 3 levels "1","2","3": 3 1 3 1
$ Name : chr "Braund, Mr. Owen Harris" "Cumings,
Laina" "Eutrelle, Mrs. Jacques Heath (Lily May Peel)"
$ Sex : Factor w/ 2 levels "female","male": 2 1
$ Age : num 22 38 26 35 35 NA 54 2 27 14 ...
$ SibSp : int 1 1 0 1 0 0 0 3 0 1 ...
$ Parch : int 0 0 0 0 0 0 0 1 2 0 ...
$ Ticket : chr "A/5 21171" "PC 17599" "STON/O2. 31
```

Análisis univariado V. Categóricas

Herramientas

- **Tablas de frecuencia:** muestran el conteo de cada categoría.
- **Gráficos de barras:** ayudan a visualizar la distribución de variables categóricas
- **Gráficos de barras apiladas:** para ver la relación entre dos variables categóricas.
- **Diagramas de caja (combinado con una variable numérica):** permite comparar distribuciones de una variable numérica a través de las categorías de una variable categórica
- **Tablas de contingencia:** muestra la relación entre 2 variables categóricas.
- **Gráficos de Mosaico:** útiles para visualizar tablas de contingencia

Análisis univariado V. Categóricas

Ejemplo

```
82 # TABLAS DE FRECUENCIA
```

```
83 table(df$Pclass)
```

1	2	3
216	184	491

```
85  
86 table(df$Sex)
```

female	male
314	577

```
87  
88  
89  
90 prop.table(table(df$Pclass))
```

1	2	3
0.2424242	0.2065095	0.5510662

```
91  
92 # GRAFICO DE BARRAS
```

```
93  
94 install.packages("ggplot2")
```

```
95 library(ggplot2)
```

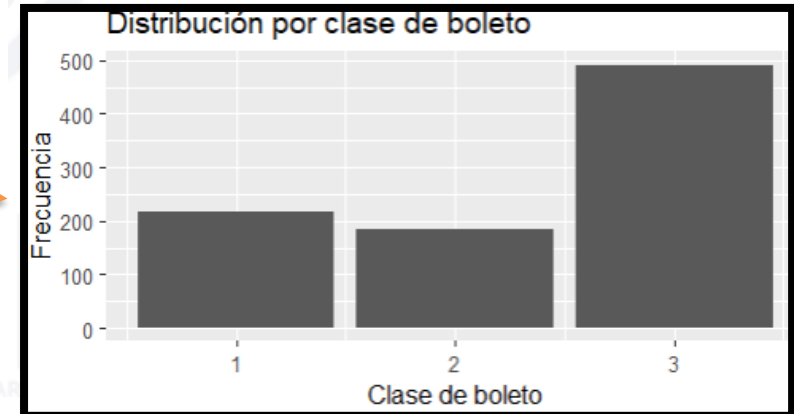
```
96  
97  
98 ggplot(df, aes(x = Pclass)) +
```

```
99   geom_bar() +
```

```
100   labs(title = "Distribución por clase de boleto",
```

```
101         x = "Clase de boleto",
```

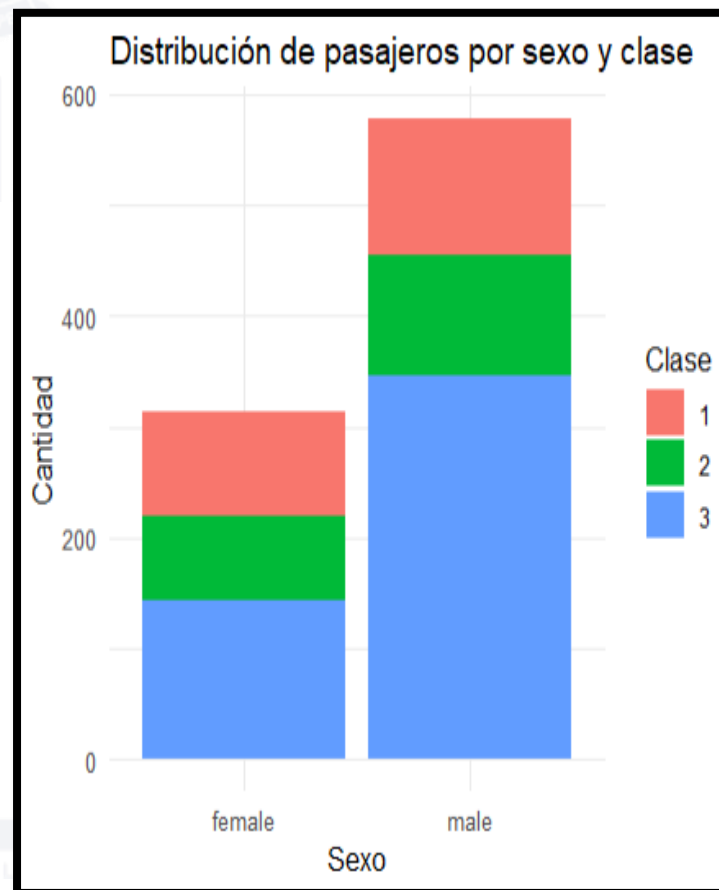
```
102         y = "Frecuencia")
```



Análisis univariado V. Categóricas

Ejemplo

```
124 # GRÁFICO DE BARRAS APILADAS
125
126 ggplot(df, aes(x = Sex, fill = Pclass)) +
127   geom_bar() +
128   labs(
129     title = "Distribución de pasajeros por sexo y clase",
130     x = "Sexo",
131     y = "Cantidad",
132     fill = "Clase"
133   ) +
134   theme_minimal()
```

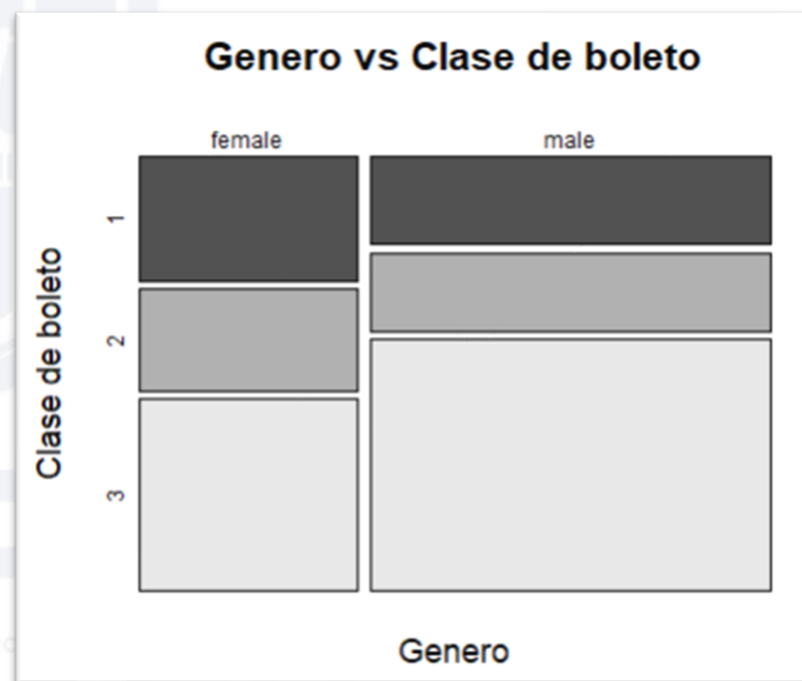


Análisis univariado V. Categóricas

Ejemplo

```
105 # TABLAS DE CONTINGENCIA
106
107 table(df$Sex, df$Pclass)
108
109
110
111 # GRAFICO DE MOSAICO
112
113 tabla_contingencia <- table(df$Sex, df$Pclass)
114
115 mosaicplot(tabla_contingencia,
116             color = TRUE,
117             main = "Genero vs Clase de boleto",
118             xlab = "Genero",
119             ylab = "Clase de boleto")
120
```

	1	2	3
female	94	76	144
male	122	108	347



Análisis univariado V. Numéricas

Sobrevista y herramientas

- Se estudia una variable numérica para ver su comportamiento
- Para EDA de variables numéricas, se puede usar:
 - **Medidas de tendencia central y dispersión:** media, mediana, moda, desviación típica, varianza, rango
 - **Histogramas:** que muestra la distribución de una variable numérica
 - **Diagrama de dispersión (Scatter Plot):** para mostrar la relación entre dos variables numéricas.
 - **Matriz de Correlación:** muestra las relaciones entre todas las variables numéricas del conjunto de datos.
 - **Diagramas de caja (bloxplots):** para visualizar la dispersión, la mediana, y detectar outliers.
 - **Análisis de outliers:** obtenidos mediante consultas directas o utilizando boxplots

Análisis univariado V. Numéricas

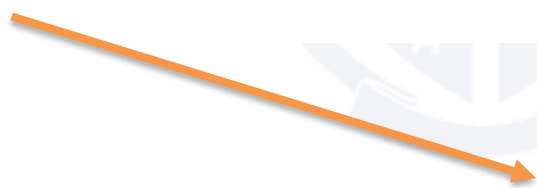
Medidas de tendencia central y de dispersión

- **Media:** Promedio aritmético de un conjunto de valores
- **Mediana:** valor central de un conjunto de datos ordenados
- **Moda:** valor que aparece con mayor frecuencia en un conjunto de datos.
- **Desviación típica:** mide la dispersión de un conjunto de datos en relación con su media. Indica cuánto se desvían, en promedio, los valores respecto a la media
- **Varianza:** media de los cuadrados de las desviaciones respecto a la media
- **Rango:** la diferencia entre el valor máximo y el valor mínimo de un conjunto de datos.
- **Rango intercuartil (IQR):** la diferencia entre el tercer cuartil (Q3) y el primer cuartil (Q1). Mide la dispersión del 50% central de los datos

Análisis univariado V. Numéricas

Ejemplo

```
157 # MEDIDAS DE TENDENCIA CENTRAL Y DISPERSIÓN
158
159 mean(df$Age, na.rm = TRUE)           # Media
160
161 median(df$Age, na.rm = TRUE)         # Mediana
162
163 sd(df$Age, na.rm = TRUE)             # Desviación estándar
164
165 var(df$Age, na.rm = TRUE)            # Varianza
166
167 max(df$Age, na.rm = TRUE) - min(df$Age, na.rm = TRUE) # Rango
168
169 diff(range(df$Age, na.rm = TRUE))    # Rango
170
171 summary(df$Age)                      # Resumen estadístico
```

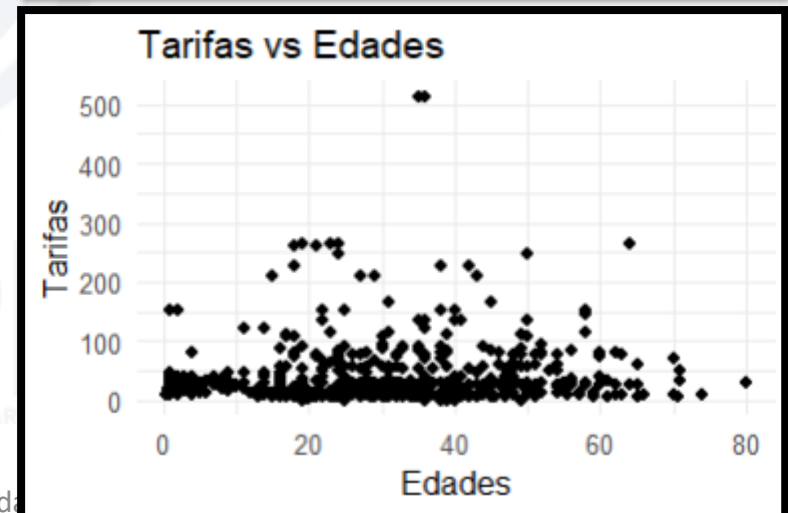
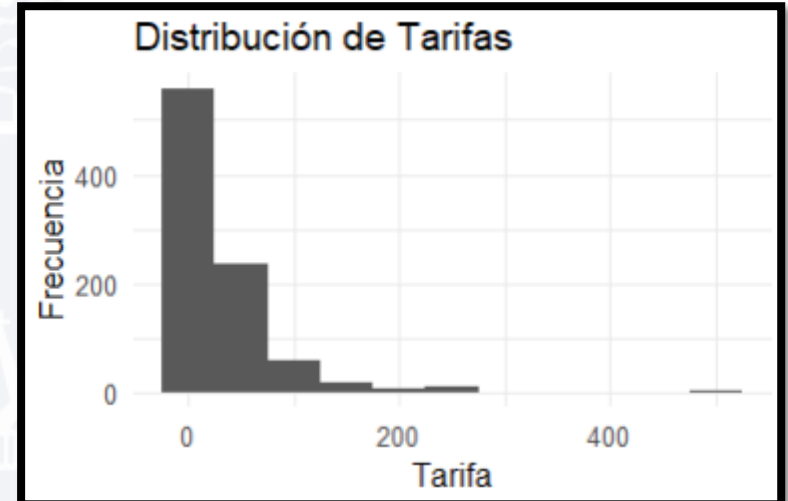


Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
0.42	20.12	28.00	29.70	38.00	80.00	177

Análisis univariado V. Numéricas

Ejemplo

```
174 # HISTOGRAMAS
175
176 library(ggplot2)
177
178 ggplot(df, aes(x=Fare)) +
179   geom_histogram(binwidth = 50) +
180   labs(title = "Distribución de Tarifas",
181         x = "Tarifa",
182         y = "Frecuencia") +
183   theme_minimal()
184
185 # DIAGRAMAS DE DISPERSIÓN
186
187 ggplot(df, aes(x = Age, y = Fare)) +
188   geom_point() +
189   labs(title = "Tarifas vs Edades",
190         x = "Edades",
191         y = "Tarifas") +
192   theme_minimal()
193
```

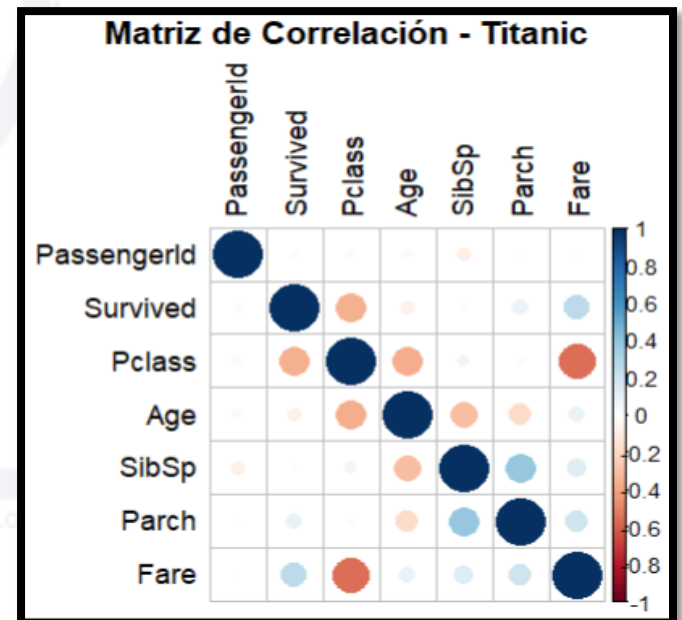


Análisis univariado V. Numéricas

Ejemplo

```
196 # MATRIZ DE CORRELACIÓN
197
198 install.packages("corrplot")
199 library(corrplot)
200 library(dplyr)
201
202 # Seleccionar solo variables numéricas
203 df_num <- df %>%
204   select(where(is.numeric))
205
206 # Ver estructura de variables numéricas seleccionadas
207 str(df_num)
208
209 # Calcular la matriz de correlación (omitendo NA)
210 cor_matrix <- cor(df_num, use = "complete.obs")
211
212 # Mostrar la matriz en consola
213 print(cor_matrix)
214
215 # Graficar la matriz de correlación
216 corrplot(cor_matrix, method = "circle",
217           tl.col = "black",
218           title = "Matriz de Correlación - Titanic",
219           mar = c(0,0,2,0))
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
PassengerId	1.00000000	0.02934016	-0.03534911	0.03684720	-0.08239772	-0.01161741	0.00959178
Survived	0.02934016	1.00000000	-0.35965268	-0.07722109	-0.01735836	0.09331701	0.26818862
Pclass	-0.03534911	-0.35965268	1.00000000	-0.36922602	0.06724737	0.02568307	-0.55418247
Age	0.03684720	-0.07722109	-0.36922602	1.00000000	0.06724737	0.02568307	-0.55418247
SibSp	-0.08239772	-0.01735836	0.06724737	0.06724737	1.00000000	0.02568307	-0.55418247
Parch	-0.01161741	0.09331701	0.02568307	0.02568307	0.02568307	1.00000000	-0.55418247
Fare	0.00959178	0.26818862	-0.55418247	-0.55418247	-0.55418247	-0.55418247	1.00000000



Análisis univariado V. Numéricas

Análisis de outliers (valores atípicos)

- Es muy importante identificar y tratar los outliers adecuadamente para evitar que influyan negativamente en los resultados de un análisis.
- El método comúnmente usado para detectar outliers está basado en el Rango Inter cuartil (IQR), tomando como base los límites de dicho IQR (Q1, Q3).
- Para detectar valores atípicos, buscamos aquellos muy alejados del "centro" del conjunto de datos. Usamos el IQR para definir qué tan lejos es "demasiado lejos".
- Los valores que se encuentran por debajo de $Q1 - 1.5 \times IQR$ o por encima de $Q3 + 1.5 \times IQR$ son considerados outliers.

Análisis univariado V. Numéricas

Detección de outliers

```
226 # Mostrar cuartiles
227 quantile(df$Fare, na.rm = TRUE)
228
229 # Calcular IQR
230 Q1 <- quantile(df$Fare,0.25, na.rm = TRUE)
231 Q3 <- quantile(df$Fare,0.75, na.rm = TRUE)
232 IQR <- Q3 - Q1
233
234 # Calcular límites
235 limite_inf <- Q1 - 1.5 * IQR
236 limite_sup <- Q3 + 1.5 * IQR
237
238 # Detectar outliers usando los límites inferior y superior
239 outliers <- df$Fare[df$Fare < limite_inf | df$Fare > limite_sup]
240 print(outliers)
241
242 # Cantidad absoluta y relativa de outliers
243 sum(df$Fare < limite_inf | df$Fare > limite_sup)
244 mean(df$Fare < limite_inf | df$Fare > limite_sup)
245
246 # Observaciones con valores atípicos
247 df_out <- df %>% filter(Fare < limite_inf | Fare > limite_sup)
248
249 # Visualización con boxplot
250 boxplot(df$Fare, main = "Boxplot de Fare",
251         ylab = "Precio del pasaje")
252 boxplot.stats(df$Fare)$out
```

Análisis univariado V. Numéricas

Detección de outliers

```
226 # Mostrar cuartiles
227 quantile(df$Fare, na.rm = TRUE)
228
229 # Calcular IQR
230 Q1 <- quantile(df$Fare,0.25, na.rm = TRUE)
231 Q3 <- quantile(df$Fare,0.75, na.rm = TRUE)
232 IQR <- Q3 - Q1
233
234 # Calcular límites
235 limite_inf <- Q1 - 1.5 * IQR
236 limite_sup <- Q3 + 1.5 * IQR
237
238 # Detectar outliers usando los límites inferior y superior
239 outliers <- df$Fare[df$Fare < limite_inf | df$Fare > limite_sup]
240 print(outliers)
241
242 # Cantidad absoluta y relativa de outliers
243 sum(df$Fare < limite_inf | df$Fare > limite_sup)
244 mean(df$Fare < limite_inf | df$Fare > limite_sup)
245
246 # Observaciones con valores atípicos
247 df_out <- df %>% filter(Fare < limite_inf | Fare > limite_sup)
248
249 # Visualización con boxplot
250 boxplot(df$Fare, main = "Boxplot de Fare",
251         ylab = "Precio del pasaje")
252 boxplot.stats(df$Fare)$out
```

0%	25%	50%	75%	100%
0.0000	7.9104	14.4542	31.0000	512.3292

23.09

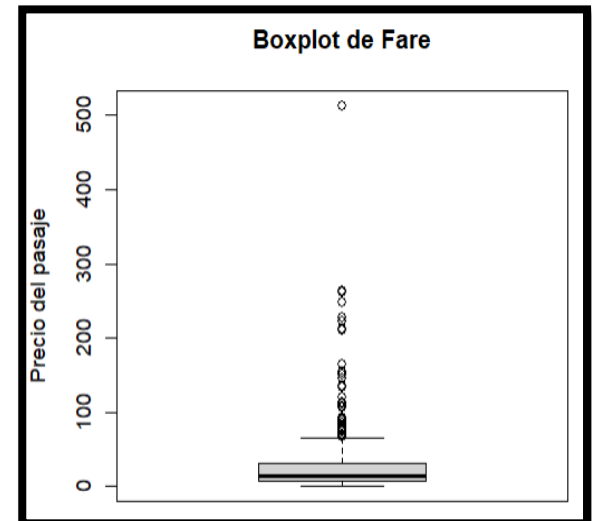
-26.72

65.63

[1]	71.2833	263.0000	146.
[16]	69.5500	69.5500	146.
[31]	77.9583	78.8500	91.
[46]	133.6500	66.6000	134.
[61]	263.0000	81.8583	89.
[76]	227.5250	79.6500	110.
[91]	76.7292	211.3375	110.
[106]	79.2000	69.5500	120.

116

0.1302

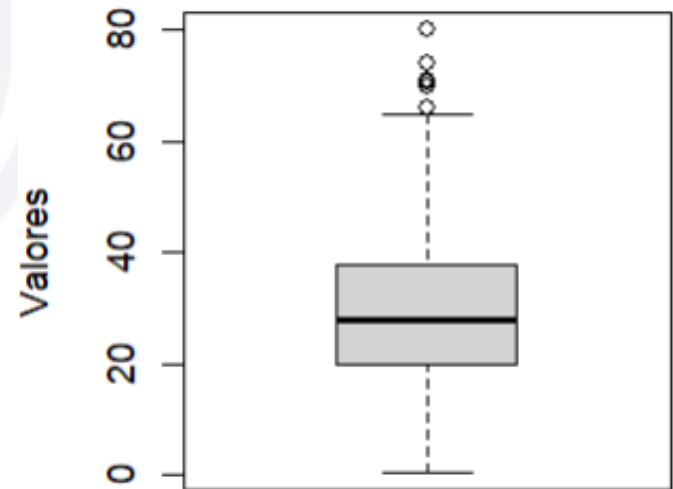


Análisis univariado V. Numéricas

Detección de outliers

- Se puede visualizar y analizar outliers usando diagramas de caja (boxplots)
- Un boxplot muestra la distribución de una variable y destaca los valores atípicos.
- La caja es el IQR
- La línea dentro de la caja es la mediana
- Los bigotes se extienden a los límites superior e inferior de valores normales.
- Los puntos fuera de los bigotes son considerados outliers

Boxplot de Datos de Edades



Análisis univariado V. Numéricas

Tratamiento de outliers

- Los valores atípicos puede incidir significativamente en los resultados del análisis. valores
- No hay una regla fija del % de outliers que son significativos para un análisis. Pero en general en distribuciones normales, si más de 5% de lo datos son atípicos, pueden incidir significativamente en los resultados.
- Opciones para tratar valores atípicos:
 - **Imputación de outliers** (con la media, mediana, limite superior, limite inferior)
 - **Transformación de datos** (logarítmo natural, raíz cuadrada)
 - **Eliminación**
- Si los outliers corresponden a errores en los datos, se los debe eliminar.
- Si lo que se busca es analizar comportamiento atípicos (Ej lavado de dinero), entonces los outliers son relevantes para el análisis

Análisis univariado V. Numéricas

Tratamiento de outliers

```
259 # SIN TRATAMIENTO
260
261 sd(df$Fare)
262 quantile(df$Fare)
263 boxplot(df$Fare)
264
265 # OPCIÓN 1: Imputar atípicos con mediana
266
267 df_op1 <- df %>%
268   mutate(Fare = ifelse(Fare < limite_inf | Fare > limite_sup,
269                       median(Fare), Fare))
270 sd(df_op1$Fare)
271 quantile(df_op1$Fare)
272 boxplot(df_op1$Fare)
273
274 # OPCIÓN 2: Imputar atípicos con límites
275
276 df_op2 <- df %>%
277   mutate(Fare = ifelse(Fare < limite_inf, 0,
278                       ifelse(Fare > limite_sup,
279                             limite_sup, Fare)))
280 sd(df_op2$Fare)
281 quantile(df_op2$Fare)
282 boxplot(df_op2$Fare)
```



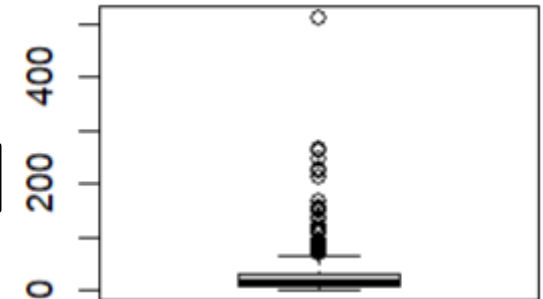
Análisis univariado V. Numéricas

Tratamiento de outliers

```
259 # SIN TRATAMIENTO
```

```
260  
261 sd(df$Fare)  
262 quantile(df$Fare)  
263 boxplot(df$Fare)
```

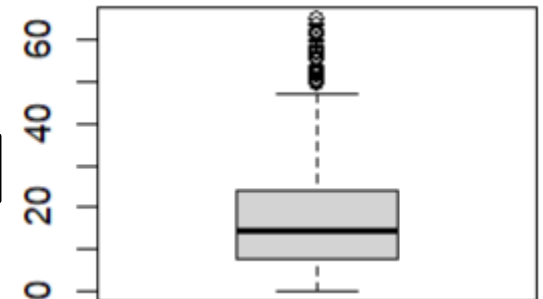
sd = 49.69343



```
265 # OPCIÓN 1: Imputar atípicos con mediana
```

```
266  
267 df_op1 <- df %>%  
268   mutate(Fare = ifelse(Fare < limite_inf | Fare > limite_sup,  
269                       median(Fare), Fare))  
270  
271 sd(df_op1$Fare)  
272 quantile(df_op1$Fare)  
273 boxplot(df_op1$Fare)
```

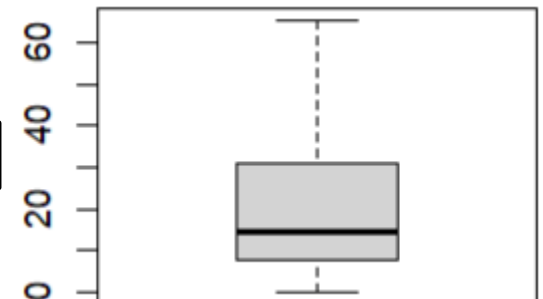
sd = 12.71302



```
274 # OPCIÓN 2: Imputar atípicos con limites
```

```
275  
276 df_op2 <- df %>%  
277   mutate(Fare = ifelse(Fare < limite_inf, 0,  
278                       ifelse(Fare > limite_sup,  
279                             limite_sup, Fare)))  
280  
281 sd(df_op2$Fare)  
282 quantile(df_op2$Fare)  
283 boxplot(df_op2$Fare)
```

sd = 20.48162



Análisis univariado V. Numéricas

Tratamiento de outliers

```
285 # SIN TRATAMIENTO
286
287 sd(df$Fare)
288 quantile(df$Fare)
289 boxplot(df$Fare)
290
291 # OPCIÓN 3: Transforma datos por Raiz cuadrada
292
293 df_op3 <- df %>%
294   mutate(Fare = sqrt(Fare))
295 sd(df_op3$Fare)
296 quantile(df_op3$Fare)
297 boxplot(df_op3$Fare)
298
299
300 # OPCIÓN 4: Transformar datos por logaritmo
301
302 df_op4 <- df %>%
303   mutate(Fare = log(Fare + 1))
304 sd(df_op4$Fare)
305 quantile(df_op4$Fare)
306 boxplot(df_op4$Fare)
```



Análisis univariado V. Numéricas

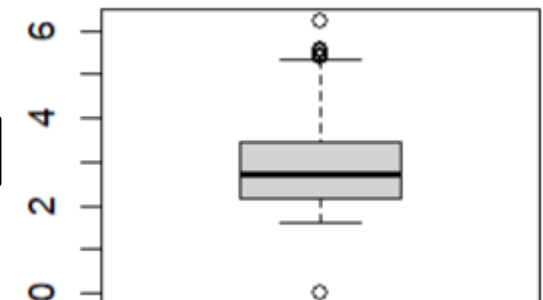
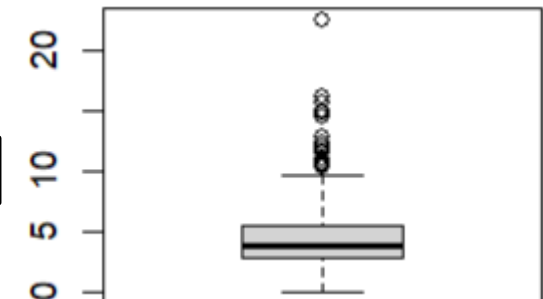
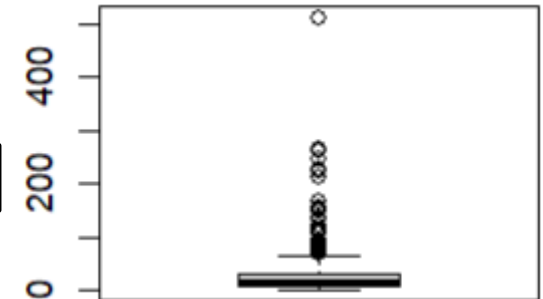
Tratamiento de outliers

```
285 # SIN TRATAMIENTO
286
287 sd(df$Fare)
288 quantile(df$Fare)
289 boxplot(df$Fare)
290
291 # OPCIÓN 3: Transforma datos por Raiz cuadrada
292
293 df_op3 <- df %>%
294   mutate(Fare = sqrt(Fare))
295 sd(df_op3$Fare)
296 quantile(df_op3$Fare)
297 boxplot(df_op3$Fare)
298
299
300 # OPCIÓN 4: Transformar datos por logaritmo
301
302 df_op4 <- df %>%
303   mutate(Fare = log(Fare + 1))
304 sd(df_op4$Fare)
305 quantile(df_op4$Fare)
306 boxplot(df_op4$Fare)
```

sd = 49.69343

sd = 2.946119

sd = 0.969048



Análisis univariado V. Numéricas

Tratamiento de outliers

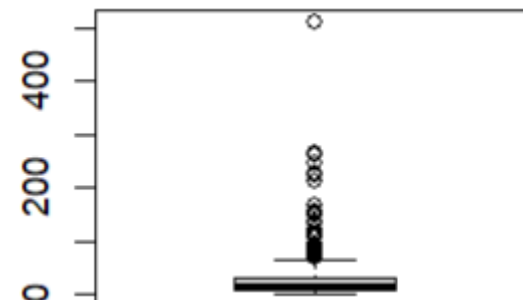
```
309 # SIN TRATAMIENTO
```

```
310  
311 sd(df$Fare)  
312 quantile(df$Fare)  
313 boxplot(df$Fare)
```

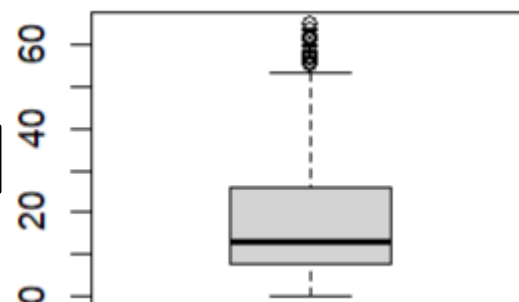
```
314  
315 # OPCIÓN 5: Eliminar atípicos
```

```
316  
317 df_op5 <- df %>%  
318   filter(df$Fare >= limite_inf & df$Fare <= limite_sup)  
319 sd(df_op5$Fare)  
320 quantile(df_op5$Fare)  
321 boxplot(df_op5$Fare)
```

sd = 49.69343



sd = 13.57809



Análisis univariado V. Numéricas

Tratamiento de outliers

- **Imputación con media o mediana:** Cuando outlier es el resultado de error o ruido, no de un valor real importante.
 - **Mediana:** si la distribución está muy sesgada
 - **Media:** si la distribución es simétrica o no hay sesgo fuerte
- **Imputación con límites:** Queremos reducir el impacto de los outliers sin eliminarlos ni cambiarlos totalmente. El valor extremo es válido, pero demasiado influyente. Queremos mantener cierta dispersión sin alterar la estructura de los datos.
- **Transformación con raíz cuadrada:** La variable tiene sesgo moderado. Los valores extremos son válidos, pero necesitamos reducir la variabilidad. No funciona con negativos.
- **Transformación con logaritmo natural:** La variable tiene alto sesgo (muy dispersa, con valores extremos grandes). El rango es amplio y necesitamos comprimir la escala. No funciona con valores menores o iguales a 0.