

# Limpieza y Análisis de Datos

Lucas Gómez Torres y Joan Amengual Mesquida

13 de enero, 2023

## Índice General

<b>1 Descripción del dataset. ¿Por qué es importante y qué pregunta/problema pretende responder?</b>	<b>2</b>
<b>2 Integración y selección de los datos de interés a analizar.</b>	<b>5</b>
<b>3 Limpieza de los datos.</b>	<b>6</b>
3.1 ¿Los datos contienen ceros o elementos vacíos? Gestiona cada uno de estos casos. . . . .	6
3.2 Identifica y gestiona los valores extremos. . . . .	6
<b>4 Análisis de los datos.</b>	<b>6</b>
4.1 Selección de los grupos de datos que se quieren analizar/comparar (p.ej., si se van a comparar grupos de datos, ¿cuáles son estos grupos y qué tipo de análisis se van a aplicar?). . . . .	6
4.2 Comprobación de la normalidad y homogeneidad de la varianza. . . . .	6
4.3 Aplicación de pruebas estadísticas para comparar los grupos de datos. En función de los datos y el objetivo del estudio, aplicar pruebas de contraste de hipótesis, correlaciones, regresiones, etc. Aplicar al menos tres métodos de análisis diferentes. . . . .	6
<b>5 Representación de los resultados a partir de tablas y gráficas.</b>	<b>6</b>
<b>6 Resolución del problema. A partir de los resultados obtenidos, ¿cuáles son las conclusiones? ¿Los resultados permiten responder al problema?</b>	<b>6</b>
<b>7 Código.</b>	<b>6</b>
<b>8 Vídeo.</b>	<b>6</b>

# 1 Descripción del dataset. ¿Por qué es importante y qué pregunta/problema pretende responder?

Actualmente cada vez sufren más personas ataques al corazón originados por diferentes factores como pueden ser el exceso de colesterol, el nivel de azúcar en la sangre, el consumo de tabaco, la presión arterial, la obesidad, la edad o la falta de ejercicio, entre muchos otros más, que pueden dar lugar a un daño permanente en el corazón como la insuficiencia cardíaca o a la muerte.

Por ello, los ataques al corazón son un problema muy grave que hay que intentar prevenir, analizando las diferentes variables que pueden influir a la hora de que una persona sufra un ataque al corazón o no, pudiendo responder a preguntas como por ejemplo:

- ¿Los hombres son más probables a sufrir un ataque?
- ¿El nivel de azúcar en sangre es determinante para que una persona pueda padecer un ataque ?
- ¿Las personas mayores tienen más probabilidad de sufrir un ataque?
- ¿Qué factor es el más influye en un ataque ?

El dataset tiene 303 observaciones (casos de personas que han sufrido o no un ataque al corazón) y 14 atributos. De estos 14 atributos, 13 son variables independientes y 1 la variable dependiente (nuestra variable objetivo que servirá para predecir si ese paciente tendrá un ataque o no). A continuación, se describen todos los atributos del dataset:

- **age**: Variable de tipo numérica. Determina la edad de la persona.
- **sex**: Variable de tipo numérica. Refleja el género de la persona ( $1 = \text{masculino}$ ,  $0 = \text{femenino}$ ).
- **cp**: Variable de tipo numérica. Identifica el tipo de dolor en el pecho ( $0 = \text{angina típica}$ ,  $1 = \text{angina atípica}$ ,  $2 = \text{dolor no anginoso}$ ,  $3 = \text{asintomático}$ ).
- **trtbps**: Variable de tipo numérica. Indica la presión arterial en reposo en mg/dl.
- **chol**: Variable de tipo numérica. Hace referencia al nivel de colesterol en mg/dl.
- **fbs**: Variable de tipo numérica. Indica si el nivel de azúcar en sangre en ayunas es mayor a 120 mg/dl ( $1 = \text{verdadero}$ ,  $0 = \text{falso}$ ).
- **restecg**: Variable de tipo numérica. Muestra los resultados electrocardiográficos en reposo ( $0 = \text{normal}$ ,  $1 = \text{anomalía de onda ST-T (inversiones de onda T y/o elevación o depresión ST de } > 0,05 \text{ mV)}$ ,  $2 = \text{hipertrofia ventricular izquierda probable o definida por los criterios de Estes}$ ).
- **thalachh**: Variable de tipo numérica. Determina la frecuencia cardíaca máxima alcanzada.
- **exng::**: Variable de tipo numérica. Indica si la angina ha sido inducida por el ejercicio ( $1 = \text{sí}$ ,  $0 = \text{no}$ ).
- **oldpeak**: Variable de tipo numérica. Señala la depresión ST inducida por el ejercicio en relación con el descanso.
- **slp**: Variable de tipo numérica. Muestra la pendiente del segmento ST de ejercicio máximo ( $0 = \text{inclinación hacia abajo}$ ,  $1 = \text{plano}$ ,  $2 = \text{inclinación hacia arriba}$ ).
- **caa**: Variable de tipo numérica. Indica el número de buques principales ( $0, 1, 2, 3$ ).
- **thall**: Variable de tipo numérica. Señala si el ratio de un trastorno sanguíneo llamado talasemia ( $0 = \text{no tiene}$ ,  $1 = \text{defecto fijo (sin flujo sanguíneo en alguna parte del corazón)}$ ,  $2 = \text{flujo sanguíneo normal}$ ,  $3 = \text{defecto reversible (se observa un flujo sanguíneo, pero no es normal)}$ ).

- **output:** Variable de tipo numérica. Indica si el paciente sufre un ataque al corazón o no (0 = No, 1 = Sí). Se trata de la variable objetivo o dependiente que pretendemos predecir.

```
# Se carga el dataset
heart_data <- read.csv("heart.csv", header = TRUE)

# Modificamos los nombres de las variables para que sean más intuitivos
colnames(heart_data) <- c("age", "sex", "chest_pain_type", "resting_blood_pressure", "cholesterol",
  "fasting_blood_sugar", "rest_ecg_type", "max_heart_rate_achieved",
  "exercise_induced_angina", "st_depression", "st_slope_type", "num_major_vessels",
  "thalassemia_type", "heart_attack")

# Mostramos los tipos de datos de las variables tal y como las interpreta R
sapply(heart_data, class)
```

```
##           age           sex      chest_pain_type
##      "integer"      "integer"      "integer"
## resting_blood_pressure      cholesterol      fasting_blood_sugar
##      "integer"      "integer"      "integer"
##      rest_ecg_type max_heart_rate_achieved exercise_induced_angina
##      "integer"      "integer"      "integer"
##      st_depression      st_slope_type      num_major_vessels
##      "numeric"      "integer"      "integer"
##      thalassemia_type      heart_attack
##      "integer"      "integer"
```

```
# Mostramos un resumen de los datos
summary(heart_data)
```

```
##      age           sex      chest_pain_type resting_blood_pressure
## Min.   :29.00   Min.   :0.0000   Min.   :0.000   Min.   : 94.0
## 1st Qu.:47.50   1st Qu.:0.0000   1st Qu.:0.000   1st Qu.:120.0
## Median :55.00   Median :1.0000   Median :1.000   Median :130.0
## Mean   :54.37   Mean   :0.6832   Mean   :0.967   Mean   :131.6
## 3rd Qu.:61.00   3rd Qu.:1.0000   3rd Qu.:2.000   3rd Qu.:140.0
## Max.   :77.00   Max.   :1.0000   Max.   :3.000   Max.   :200.0
## cholesterol      fasting_blood_sugar rest_ecg_type      max_heart_rate_achieved
## Min.   :126.0     Min.   :0.0000     Min.   :0.0000     Min.   : 71.0
## 1st Qu.:211.0     1st Qu.:0.0000     1st Qu.:0.0000     1st Qu.:133.5
## Median :240.0     Median :0.0000     Median :1.0000     Median :153.0
## Mean   :246.3     Mean   :0.1485     Mean   :0.5281     Mean   :149.6
## 3rd Qu.:274.5     3rd Qu.:0.0000     3rd Qu.:1.0000     3rd Qu.:166.0
## Max.   :564.0     Max.   :1.0000     Max.   :2.0000     Max.   :202.0
## exercise_induced_angina st_depression st_slope_type      num_major_vessels
## Min.   :0.0000     Min.   :0.00     Min.   :0.000     Min.   :0.0000
## 1st Qu.:0.0000     1st Qu.:0.00     1st Qu.:1.000     1st Qu.:0.0000
## Median :0.0000     Median :0.80     Median :1.000     Median :0.0000
## Mean   :0.3267     Mean   :1.04     Mean   :1.399     Mean   :0.7294
## 3rd Qu.:1.0000     3rd Qu.:1.60     3rd Qu.:2.000     3rd Qu.:1.0000
## Max.   :1.0000     Max.   :6.20     Max.   :2.000     Max.   :4.0000
## thalassemia_type heart_attack
## Min.   :0.000     Min.   :0.0000
## 1st Qu.:2.000     1st Qu.:0.0000
## Median :2.000     Median :1.0000
## Mean   :2.314     Mean   :0.5446
```

```
## 3rd Qu.:3.000    3rd Qu.:1.0000
## Max.    :3.000    Max.    :1.0000
```

```
# Se muestran las 6 primeras observaciones de los datos
head(heart_data)
```

```
##   age sex chest_pain_type resting_blood_pressure cholesterol
## 1  63   1             3             145             233
## 2  37   1             2             130             250
## 3  41   0             1             130             204
## 4  56   1             1             120             236
## 5  57   0             0             120             354
## 6  57   1             0             140             192
##   fasting_blood_sugar rest_ecg_type max_heart_rate_achieved
## 1                   1             0                   150
## 2                   0             1                   187
## 3                   0             0                   172
## 4                   0             1                   178
## 5                   0             1                   163
## 6                   0             1                   148
##   exercise_induced_angina st_depression st_slope_type num_major_vessels
## 1                      0           2.3             0                0
## 2                      0           3.5             0                0
## 3                      0           1.4             2                0
## 4                      0           0.8             2                0
## 5                      1           0.6             2                0
## 6                      0           0.4             1                0
##   thalassemia_type heart_attack
## 1                 1           1
## 2                 2           1
## 3                 2           1
## 4                 2           1
## 5                 2           1
## 6                 1           1
```

```
dim(heart_data)
```

```
## [1] 303 14
```

```
# Se carga el dataset
```

```
O2_saturation <- read.csv("o2Saturation.csv", header = TRUE)
```

```
# Mostramos los tipos de datos de las variables tal y como las interpreta R
sapply(O2_saturation,class)
```

```
##      X98.6
## "numeric"
```

```
# Mostramos un resumen de los datos
```

```
summary(O2_saturation)
```

```
##      X98.6
## Min.   :96.50
## 1st Qu.:97.60
## Median :98.60
## Mean   :98.24
## 3rd Qu.:98.60
```

```
## Max.      :99.60
# Se muestran las 6 primeras observaciones de los datos
head(O2_saturation)

##      X98.6
## 1  98.6
## 2  98.6
## 3  98.6
## 4  98.1
## 5  97.5
## 6  97.5
```

## 2 Integración y selección de los datos de interés a analizar.

Puede ser el resultado de adicionar diferentes datasets o una subselección útil de los datos originales, en base al objetivo que se quiera conseguir.

```
sum(rowSums(is.na(heart_data))>0)

## [1] 0

colMeans(is.na(heart_data))

##           age           sex      chest_pain_type
##           0           0           0
## resting_blood_pressure      cholesterol      fasting_blood_sugar
##           0           0           0
##      rest_ecg_type max_heart_rate_achieved exercise_induced_angina
##           0           0           0
##      st_depression      st_slope_type      num_major_vessels
##           0           0           0
##      thalassemia_type      heart_attack
##           0           0
```

### **3 Limpieza de los datos.**

**3.1** ¿Los datos contienen ceros o elementos vacíos? Gestiona cada uno de estos casos.

**3.2** Identifica y gestiona los valores extremos.

### **4 Análisis de los datos.**

**4.1** Selección de los grupos de datos que se quieren analizar/comparar (p.ej., si se van a comparar grupos de datos, ¿cuáles son estos grupos y qué tipo de análisis se van a aplicar?).

**4.2** Comprobación de la normalidad y homogeneidad de la varianza.

**4.3** Aplicación de pruebas estadísticas para comparar los grupos de datos. En función de los datos y el objetivo del estudio, aplicar pruebas de contraste de hipótesis, correlaciones, regresiones, etc. Aplicar al menos tres métodos de análisis diferentes.

### **5 Representación de los resultados a partir de tablas y gráficas.**

Este apartado se puede responder a lo largo de la práctica, sin necesidad de concentrar todas las representaciones en este punto de la práctica.

### **6 Resolución del problema. A partir de los resultados obtenidos, ¿cuáles son las conclusiones? ¿Los resultados permiten responder al problema?**

### **7 Código.**

Hay que adjuntar el código, preferiblemente en R, con el que se ha realizado la limpieza, análisis y representación de los datos. Si lo preferís, también podéis trabajar en Python.

### **8 Vídeo.**

Realizar un breve vídeo explicativo de la práctica (máximo 10 minutos), donde ambos integrantes del equipo expliquen con sus propias palabras el desarrollo de la práctica, basándose en las preguntas del enunciado para justificar y explicar el código desarrollado. Este vídeo se deberá entregar a través de un enlace al Google Drive de la UOC, junto con enlace al repositorio Git entregafo.