

# Vacunación a la población mayor de 65 años

Alicia Perdices Guerra

3 de mayo, 2021

## Contents

### 1.PROCESAMIENTO DE LOS DATOS.

- En primer lugar leemos el fichero:

```
vacuna<-read.csv("C:/temp/Vacunacion_+65.csv",sep= ",")
```

- Realicemos una breve inspección de los datos

```
str(vacuna)
```

```
## 'data.frame':   380 obs. of  5 variables:
## $ TIME          : int   2010 2010 2010 2010 2010 2010 2010 2010 2010 2010 ...
## $ GEO           : Factor w/ 38 levels "Austria","Belgium",...: 9 10 2 3 6 7 13 8 17 14 ...
## $ UNIT          : Factor w/ 1 level "Percentage": 1 1 1 1 1 1 1 1 1 1 ...
## $ Value         : Factor w/ 223 levels ":", "0.90", "1.00",...: 141 157 1 1 47 131 173 3 195 1 ...
## $ Flag.and.Footnotes: Factor w/ 7 levels "","b","bd","be",...: 6 6 1 1 2 1 1 1 1 1 ...
```

```
colnames(vacuna) #Nombre de las variables
```

```
## [1] "TIME"          "GEO"           "UNIT"
## [4] "Value"         "Flag.and.Footnotes"
```

```
nrow(vacuna) #Número de registros
```

```
## [1] 380
```

```
ncol(vacuna) #Número de variables
```

```
## [1] 5
```

\*Observamos las siguientes variables:

- **TIME**: variable cuantitativa. Indica el año en el que se ha realizado la medida, en este caso el valor de la variable "Value". Se ha cargado bien como número entero.
- **GEO**: variable cualitativa. Indica el país o región en el que se ha realizado la medida. Se ha cargado bien como factor.
- **UNIT**: variable cualitativa. Indica la medida de la variable valor. Se ha cargado bien como factor. Porcentaje.
- **Value**: Variable cuantitativa. Indica el porcentaje de población mayor de 65 años a la que se le ha vacunado. Se ha cargado mal como factor.
- **Flag.and.footnotes**. Notas sobre etiquetas. Eliminamos esta columna.

\*Años de las mediciones:

```
unique(vacuna$TIME)
```

```
## [1] 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019
```

\*Países:

```
unique(vacuna$GEO)
```

```
## [1] European Union - 27 countries (from 2020)
## [2] European Union - 28 countries (2013-2020)
## [3] Belgium
## [4] Bulgaria
## [5] Czechia
## [6] Denmark
## [7] Germany (until 1990 former territory of the FRG)
## [8] Estonia
## [9] Ireland
## [10] Greece
## [11] Spain
## [12] France
## [13] Croatia
## [14] Italy
## [15] Cyprus
## [16] Latvia
## [17] Lithuania
## [18] Luxembourg
## [19] Hungary
## [20] Malta
## [21] Netherlands
## [22] Austria
## [23] Poland
## [24] Portugal
## [25] Romania
## [26] Slovenia
## [27] Slovakia
## [28] Finland
## [29] Sweden
## [30] Iceland
## [31] Liechtenstein
## [32] Norway
## [33] Switzerland
## [34] United Kingdom
## [35] Montenegro
## [36] North Macedonia
## [37] Serbia
## [38] Turkey
## 38 Levels: Austria Belgium Bulgaria Croatia Cyprus Czechia Denmark ... United Kingdom
```

\*Unidad de las mediciones:

```
unique(vacuna$UNIT)
```

```
## [1] Percentage
## Levels: Percentage
```

- Eliminamos la columna Fal.and.footnotes

```
vacuna<-vacuna[, -5]
```

- Tendríamos que resolver las posibles inconsistencias en relación al formato del valor numérico de la variable **Value** y convertirla a valor numérico.

```
vacuna$Value<-as.character(vacuna$Value)
vacuna$Value<-as.numeric (gsub(',', '.', vacuna$Value) )
```

```
## Warning: NAs introducidos por coerción
```

- Comprobamos que valores tenemos en la columna **Value**:

```
tail(table(vacuna$Value, useNA = "ifany"))
```

```
##  
## 72.8 73.32 73.46 73.5 74 <NA>  
## 2 1 1 1 1 119
```

- Observamos que tenemos **119 valores perdidos**. Guardamos en la variable **idx** los índices de los registros con valores **NA** de la variable **Value**.

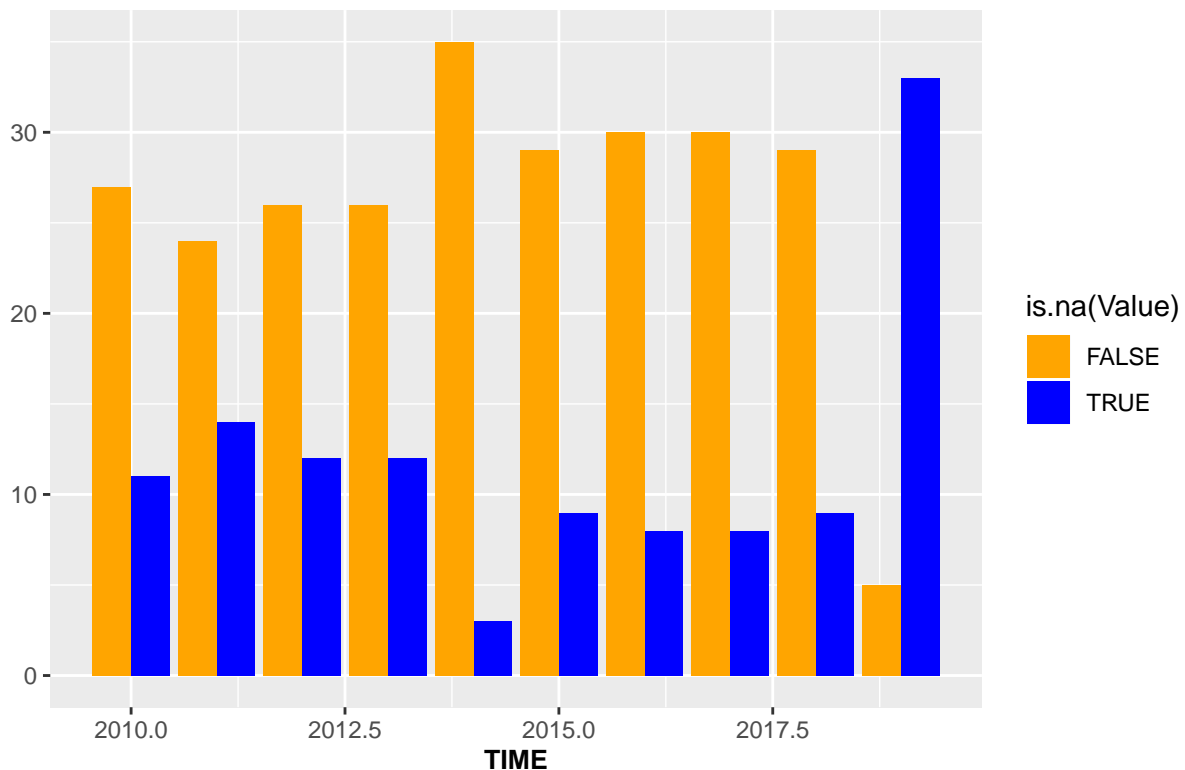
```
idx<-which(is.na(vacuna$Value))  
length(idx)
```

```
## [1] 119
```

- Grafiquemos la información que contiene la variable **Value**.

```
library(ggplot2)  
library(scales)  
g = ggplot(vacuna, aes(TIME, fill=is.na(Value)) ) +  
labs(title = "Valores Nulos")+ylab("") +  
theme(plot.title = element_text(size = rel(2), colour = "blue"))  
  
g+geom_bar(position="dodge") + scale_fill_manual(values = alpha(c("orange", "blue"), 1)) +  
theme(axis.title.x = element_text(face="bold", size=10))
```

## Valores Nulos



- En caso de detectar algún valor anómalo (en nuestro caso los NAs) en las variables tendríamos que

realizar una imputación de esos valores o bien sustituyéndolos por la media o usando el algoritmo KNN ( k-Nearest Neighbour) con los 3 vecinos más cercanos usando la distancia que consideremos, en este caso usaremos Gower(Mediana), por ser una medida más robusta frente a extremos.

```
library(VIM)

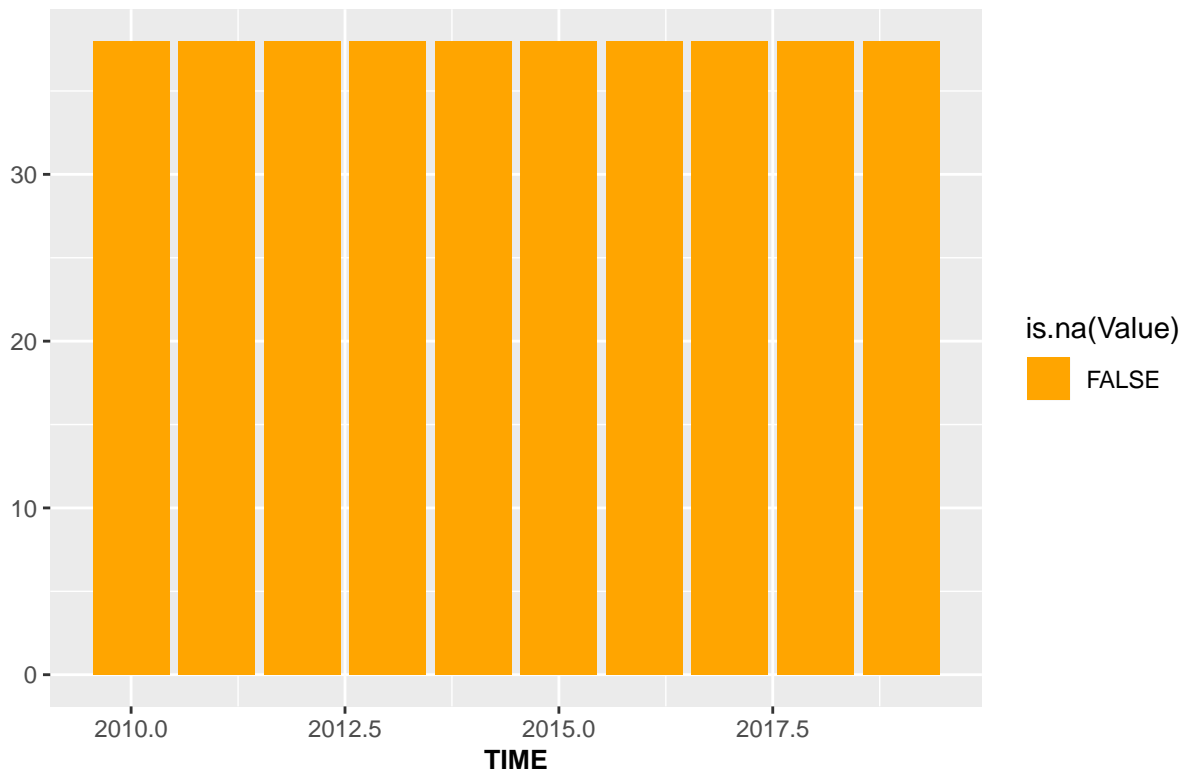
## Loading required package: colorspace
## Loading required package: grid
## VIM is ready to use.
## Suggestions and bug-reports can be submitted at: https://github.com/statistikat/VIM/issues
##
## Attaching package: 'VIM'
## The following object is masked from 'package:datasets':
##
##      sleep
output<-kNN(vacuna, variable=c("Value"),k=3)
vacuna<-output
```

- Comprobamos que no tenemos valores nulos después de la imputación

```
g = ggplot(vacuna, aes(TIME, fill=is.na(Value)) ) +
labs(title = "Valores Nulos")+ylab("") +
theme(plot.title = element_text(size = rel(2), colour = "blue"))

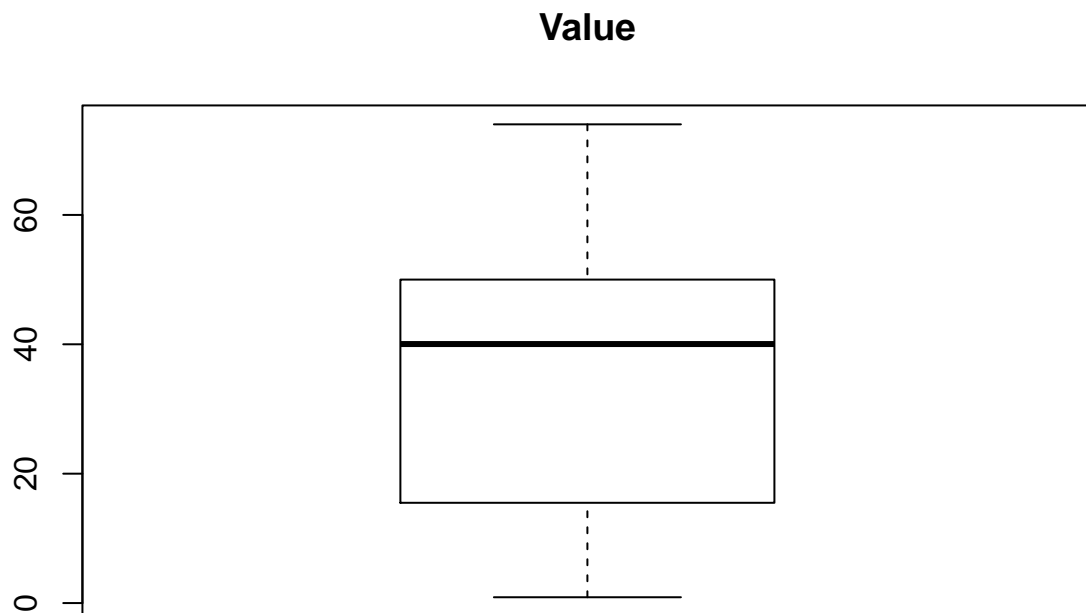
g+geom_bar(position="dodge") + scale_fill_manual(values = alpha(c("orange", "blue"), 1)) +
theme(axis.title.x = element_text(face="bold", size=10))
```

## Valores Nulos



- Con el siguiente gráfico, observaremos que la variable **Value** no tiene outliers o valores extremos

```
boxplot(vacuna$Value, main="Value")
```



- Por otro lado, revisamos para el resto de columnas si tenemos valores NA.(desconocidos o perdidos)

```
table(vacuna$TIME, useNA = "ifany")
```

```
##
## 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019
##   38   38   38   38   38   38   38   38   38   38
```

```
table(vacuna$GEO, useNA = "ifany")
```

```
##
##                               Austria
##                               10
##                               Belgium
##                               10
##                               Bulgaria
##                               10
##                               Croatia
##                               10
##                               Cyprus
##                               10
##                               Czechia
##                               10
##                               Denmark
##                               10
##                               Estonia
##                               10
## European Union - 27 countries (from 2020)
```

##	10
##	European Union - 28 countries (2013-2020)
##	10
##	Finland
##	10
##	France
##	10
##	Germany (until 1990 former territory of the FRG)
##	10
##	Greece
##	10
##	Hungary
##	10
##	Iceland
##	10
##	Ireland
##	10
##	Italy
##	10
##	Latvia
##	10
##	Liechtenstein
##	10
##	Lithuania
##	10
##	Luxembourg
##	10
##	Malta
##	10
##	Montenegro
##	10
##	Netherlands
##	10
##	North Macedonia
##	10
##	Norway
##	10
##	Poland
##	10
##	Portugal
##	10
##	Romania
##	10
##	Serbia
##	10
##	Slovakia
##	10
##	Slovenia
##	10
##	Spain
##	10
##	Sweden
##	10
##	Switzerland

```
##              10
##              Turkey
##              10
##              United Kingdom
##              10
```

```
table(vacuna$UNIT, useNA = "ifany")
```

```
##
## Percentage
##      380
```

Observamos que no existen ahora valores perdidos después de la imputación. La suma de las cantidades de cada variable, suman el total.

La estructura de los datos quedaría:

```
str(vacuna)
```

```
## 'data.frame':   380 obs. of  5 variables:
## $ TIME      : int  2010 2010 2010 2010 2010 2010 2010 2010 2010 2010 ...
## $ GEO       : Factor w/ 38 levels "Austria","Belgium",...: 9 10 2 3 6 7 13 8 17 14 ...
## $ UNIT      : Factor w/ 1 level "Percentage": 1 1 1 1 1 1 1 1 1 1 ...
## $ Value     : num  49 51.8 58 49 17.4 ...
## $ Value_imp: logi  FALSE FALSE TRUE TRUE FALSE FALSE ...
```

- Finalmente, creamos un fichero con toda la información corregida.

```
write.csv(vacuna, file="Vacunacion_+65_clean.csv", row.names = FALSE)
```