

Cuidados Domiciliarios

Alicia Perdices Guerra

3 de mayo, 2021

Contents

1.PROCESAMIENTO DE LOS DATOS.

- En primer lugar leemos el fichero:

```
cuidados<-read.csv("C:/temp/CuidadosDomiciliarios.csv",sep= ",")
```

- Realicemos una breve inspección de los datos

```
str(cuidados)
```

```
## 'data.frame':   396 obs. of  8 variables:
## $ ISCED11      : Factor w/ 4 levels "All ISCED 2011 levels ",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ GEO          : Factor w/ 33 levels "Austria","Belgium",...: 9 9 9 10 10 10 2 2 2 3 ...
## $ UNIT         : Factor w/ 1 level "Percentage": 1 1 1 1 1 1 1 1 1 1 ...
## $ TIME         : int   2014 2014 2014 2014 2014 2014 2014 2014 2014 2014 ...
## $ SEX          : Factor w/ 3 levels "Females","Males",...: 3 2 1 3 2 1 3 2 1 3 ...
## $ AGE          : Factor w/ 1 level "Total": 1 1 1 1 1 1 1 1 1 1 ...
## $ Value        : Factor w/ 96 levels ":", "0.1", "0.2",...: 54 44 63 54 43 63 96 83 29 44 ...
## $ Flag.and.Footnotes: Factor w/ 2 levels "","u": 1 1 1 1 1 1 1 1 1 1 ...
```

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

```
## [1] "ISCED11"      "GEO"          "UNIT"
## [4] "TIME"         "SEX"          "AGE"
## [7] "Value"        "Flag.and.Footnotes"
```

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

```
## [1] 396
```

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

```
## [1] 8
```

*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.
- **SEX**: Variable cualitativa. Indica el sexo de la población estudiada, Males, Females o Total.
- **AGE**: Variable cualitativa. En Cómputo total. No hay distinción en edades.
- **ISCED11**: Variable cualitativa. Estándar en estadísticas de educación en el que se hacen las mediciones. Se ha cargado bien como factor.
- **Value**: Variable cuantitativa. Indica el porcentaje de población a la que se le han aplicado cuidados domiciliarios. Se ha cargado mal como factor.
- **Flag.and.footnotes**. Notas sobre etiquetas. Eliminamos esta columna.

*Años de las mediciones:

```
unique(cuidados$TIME)
```

```
## [1] 2014
```

*Países:

```
unique(cuidados$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] Norway
## [32] United Kingdom
## [33] Turkey
## 33 Levels: Austria Belgium Bulgaria Croatia Cyprus Czechia Denmark ... United Kingdom
```

*Unidad de las mediciones:

```
unique(cuidados$UNIT)
```

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

*Estándar de las mediciones.

```
unique(cuidados$ISCED11)
```

```
## [1] All ISCED 2011 levels
## [2] Less than primary, primary and lower secondary education (levels 0-2)
```

```
## [3] Upper secondary and post-secondary non-tertiary education (levels 3 and 4)
## [4] Tertiary education (levels 5-8)
## 4 Levels: All ISCED 2011 levels ...
```

- Sexo de la población estudiada.

```
unique(cuidados$SEX)
```

```
## [1] Total    Males    Females
## Levels: Females Males Total
```

- Eliminamos la columna Fal.and.footnotes y AGE ya que no nos aporta información relevante.

```
cuidados<-cuidados[,-8]
cuidados<-cuidados[,-6]
```

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

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

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

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

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

```
##
##    12 12.9    14 14.1 15.9 <NA>
##     1    1    1    2    1    12
```

- Observamos que tenemos **12 valores perdidos**. Guardamos en la variable **idx** los índices de los registros con valores **NA** de la variable **Value**.

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

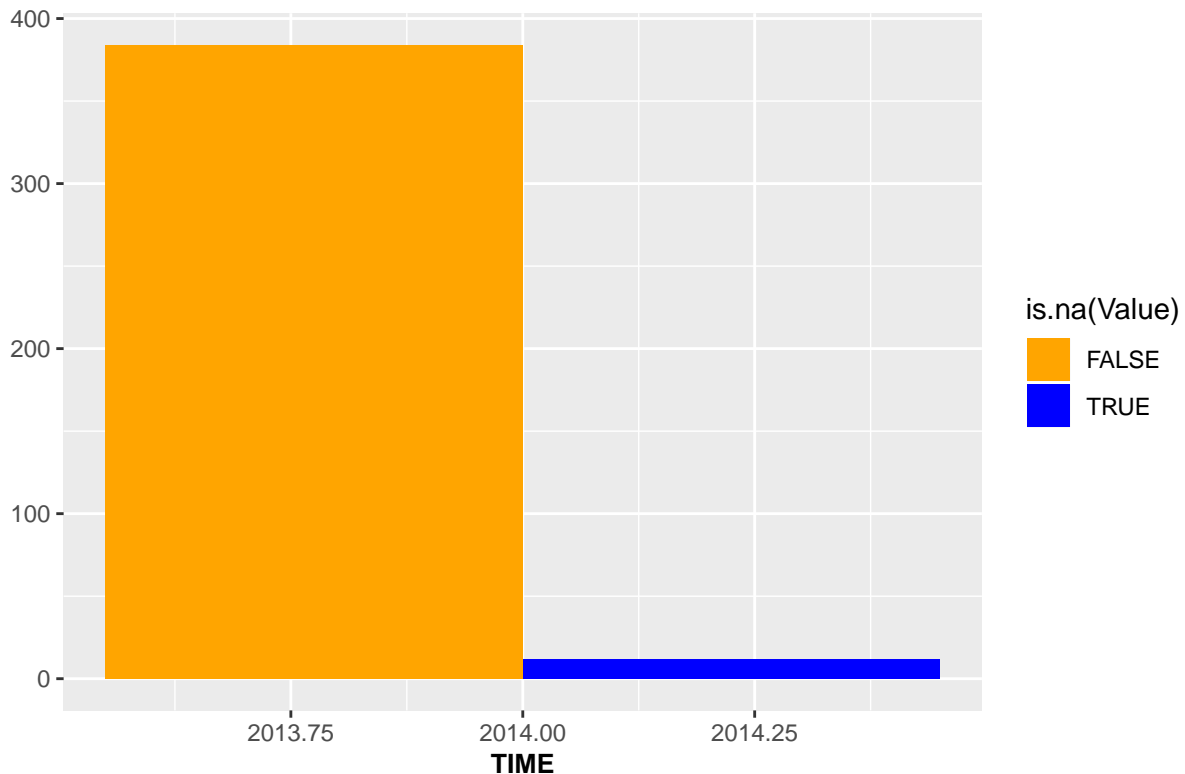
```
## [1] 12
```

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

```
library(ggplot2)
library(scales)
g = ggplot(cuidados, 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(cuidados, variable=c("Value"),k=3)
```

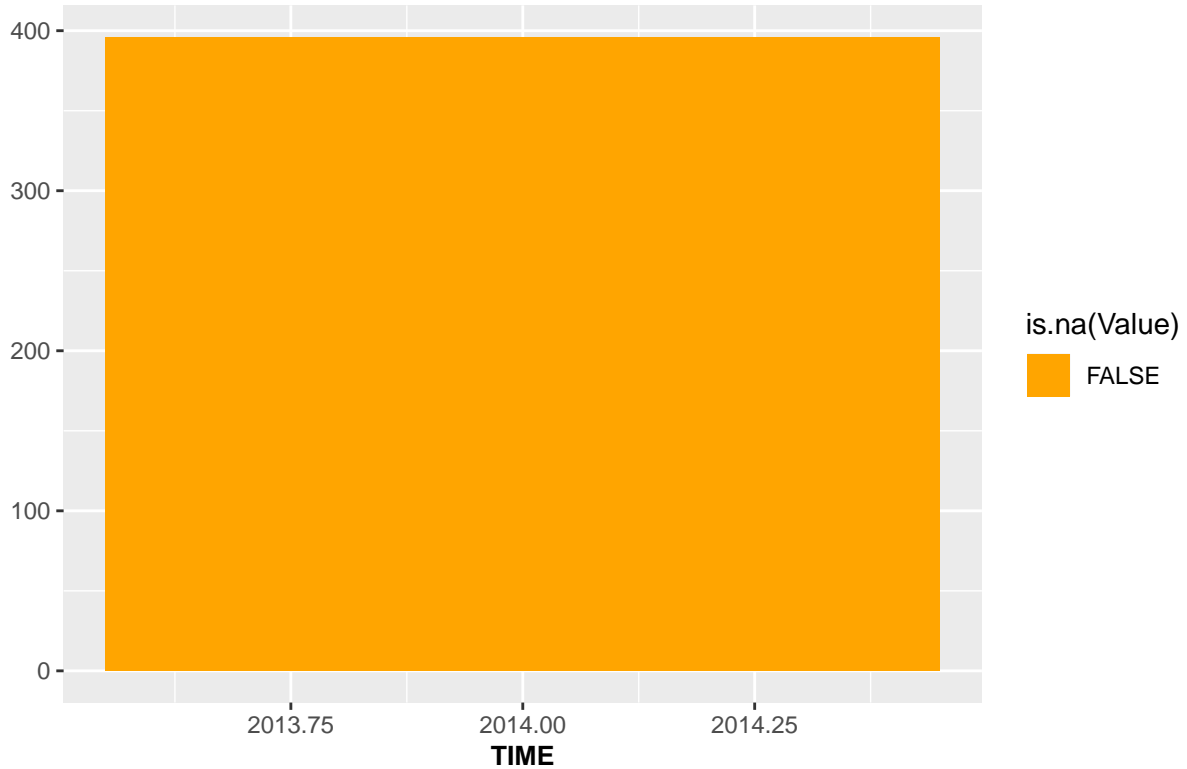
```
cuidados<-output
```

- Comprobamos que no tenemos valores nulos después de la imputación

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

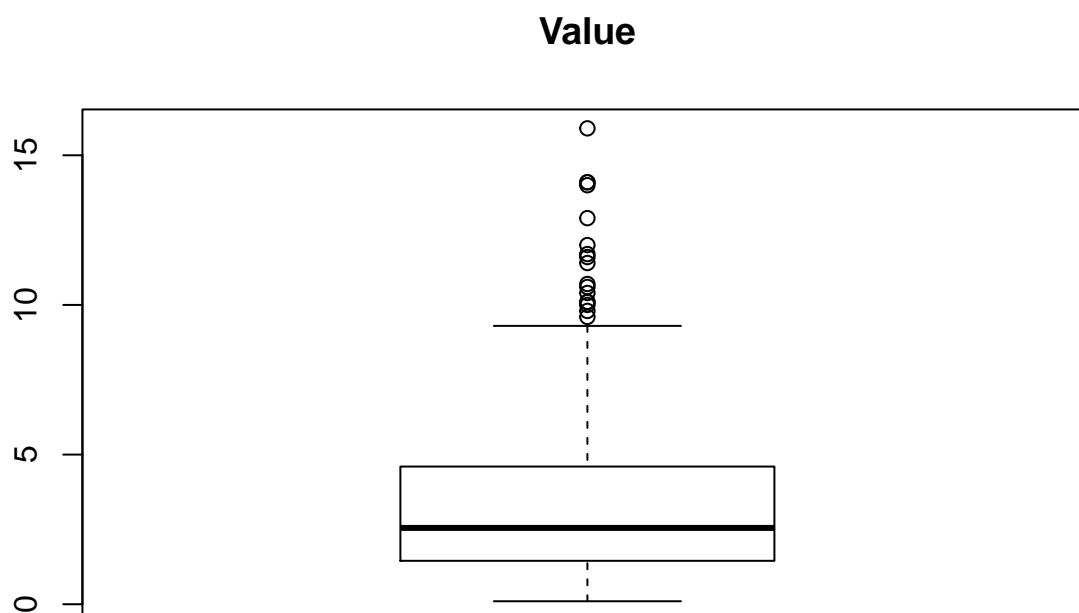
```
g+geom_bar(position="dodge") + scale_fill_manual(values = alpha(c("orange", "blue"), 1)) +  
theme(axis.title.x = element_text(face="bold", size=10))
```

Valores Nulos



- Con el siguiente gráfico, observaremos que la variable **Value** tiene outliers o valores extremos

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



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

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

```
##
## 2014
## 396
```

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

```
##
## Austria
## 12
## Belgium
## 12
## Bulgaria
## 12
## Croatia
## 12
## Cyprus
## 12
## Czechia
## 12
## Denmark
## 12
## Estonia
## 12
## European Union - 27 countries (from 2020)
```

```

## 12
## European Union - 28 countries (2013-2020)
## 12
## Finland
## 12
## France
## 12
## Germany (until 1990 former territory of the FRG)
## 12
## Greece
## 12
## Hungary
## 12
## Iceland
## 12
## Ireland
## 12
## Italy
## 12
## Latvia
## 12
## Lithuania
## 12
## Luxembourg
## 12
## Malta
## 12
## Netherlands
## 12
## Norway
## 12
## Poland
## 12
## Portugal
## 12
## Romania
## 12
## Slovakia
## 12
## Slovenia
## 12
## Spain
## 12
## Sweden
## 12
## Turkey
## 12
## United Kingdom
## 12

```

```

table(cuidados$UNIT, useNA = "ifany")

```

```

##
## Percentage
## 396

```

```
table(cuidados$SEX, useNA = "ifany")
```

```
##
## Females    Males    Total
##      132      132     132
```

```
table(cuidados$ISCED11, useNA = "ifany")
```

```
##
##                                     All ISCED 2011 levels
##                                     99
##      Less than primary, primary and lower secondary education (levels 0-2)
##                                     99
##                                     Tertiary education (levels 5-8)
##                                     99
## Upper secondary and post-secondary non-tertiary education (levels 3 and 4)
##                                     99
```

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(cuidados)
```

```
## 'data.frame':    396 obs. of  7 variables:
## $ ISCED11 : Factor w/ 4 levels "All ISCED 2011 levels ",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ GEO      : Factor w/ 33 levels "Austria","Belgium",...: 9 9 9 10 10 10 2 2 2 3 ...
## $ UNIT     : Factor w/ 1 level "Percentage": 1 1 1 1 1 1 1 1 1 1 ...
## $ TIME     : int  2014 2014 2014 2014 2014 2014 2014 2014 2014 2014 ...
## $ SEX      : Factor w/ 3 levels "Females","Males",...: 3 2 1 3 2 1 3 2 1 3 ...
## $ Value    : num  4 3 5 4 2.9 5 9.8 7.5 12 3 ...
## $ Value_imp: logi  FALSE FALSE FALSE FALSE FALSE FALSE ...
```

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

```
write.csv(cuidados, file="CuidadosDomiciliarios_clean.csv", row.names = FALSE)
```