

## 09.2.- Data Cleansing-Missing\_45\_06\_turismo\_origen\_completo\_v\_01

June 11, 2023

#

CU45\_Planificación y promoción del destino en base a los patrones en origen de los turistas

Citizenlab Data Science Methodology > II - Data Processing Domain \*\*\* > # 09.2.- Data Cleansing  
- Missing

Data Cleaning refers to identifying and correcting (or removing) errors in the dataset that may negatively impact a predictive model, replacing, modifying, or deleting the dirty or coarse data.



## 0.1 Tasks

Basic operations	Text data analysis
	Delete Needless/Irrelevant/Private Columns Inconsistent Data. Expected values Zeroes Columns with a Single Value Columns with Very Few Values Columns with Low Variance Duplicates (rows/samples) & (columns/features) Data
Missing Values	Missing Values Identification
	Missing Values Per Sample Missing Values Per Feature Zero Missing Values Other Missing Values Null/NaN Missing Values Delete Missing Values Deleting Rows with Missing Values in Target Column Deleting Rows with Missing Values Deleting Features with some Missing Values Deleting Features using Rate Missing Values
	Basic Imputation
	Imputation by Previous Row Value Imputation by Next Row Value
	Statistical Imputation
	Selection of Imputation Strategy Constant Imputation Mean Imputation Median Imputation Most Frequent Imputation Interpolation Imputation
	Prediction Imputation (KNN Imputation )
	Evaluating k-hyperparameter in KNN Imputation Applying KNN Imputation
	Iterative Imputation
	Evaluating Different Imputation Order Applying Iterative Imputation
Outliers	Outliers - Univariate
	Visualizing Outliers Distribution Box Plots Isolation Forest Outliers Identification Grubbs' Test Z-Score Standard Deviation Method Interquartile Range Method Tukey's method Internally studentized residuals AKA z-score method Median Absolute Deviation method
	Outliers - MultiVariate
	Visualizing Outliers ScatterPlots Outliers Identification Mahalanobis Distance Robust Mahalanobis Distance DBSCAN Clustering PyOD Library
	Automatic Detection and Removal of Outliers
	Compare Algorithms LocalOutlierFactor IsolationForest Minimum Covariance Determinant

## 0.2 Consideraciones casos CitizenLab programados en R

- La mayoría de las tareas de este proceso se han realizado en los notebooks del proceso 05 Data Collection porque eran necesarias para las tareas ETL. En esos casos, en este notebook se referencia al notebook del proceso 05 correspondiente
- Por tanto en los notebooks de este proceso de manera general se incluyen las comprobaciones necesarias, y comentarios si procede
- Las tareas del proceso se van a aplicar solo a los archivos que forman parte del despliegue, ya que hay muchos archivos intermedios que no procede pasar por este proceso
- El nombre de archivo del notebook hace referencia al nombre de archivo del proceso 05 al que se aplica este proceso, por eso pueden no ser correlativa la numeración
- Las comprobaciones se van a realizar teniendo en cuenta que el lenguaje utilizado en el despliegue de este caso es R

## 0.3 File

- Input File: CU\_45\_06\_03\_turismo\_receptor.csv
- Sampled Input File: CU\_45\_07\_03\_turismo\_receptor.csv
- Output File: CU\_45\_09.2\_03\_turismo\_receptor.csv

## 0.4 Settings

### 0.4.1 Encoding

Con la siguiente expresión se evitan problemas con el encoding al ejecutar el notebook. Es posible que deba ser eliminada o adaptada a la máquina en la que se ejecute el código.

```
[40]: Sys.setlocale(category = "LC_ALL", locale = "es_ES.UTF-8")
```

```
Warning message in Sys.setlocale(category = "LC_ALL", locale = "es_ES.UTF-8"):  
"OS reports request to set locale to "es_ES.UTF-8" cannot be honored"  
"
```

### 0.4.2 Libraries to use

```
[41]: library(readr)  
library(dplyr)  
library(tidyr)  
library(stringr)
```

### 0.4.3 Paths

```
[42]: iPath <- "Data/Input/"  
oPath <- "Data/Output/"
```

## 0.5 Data Load

OPCION A: Seleccionar fichero en ventana para mayor comodidad

Data load using the {tcltk} package. Uncomment the line if using this option

```
[43]: # file_data <- tcltk::tk_choose.files(multi = FALSE)
```

OPCION B: Especificar el nombre de archivo

```
[44]: iFile <- "CU_45_06_03_turismo_receptor.csv"
file_data <- paste0(iPath, iFile)

if(file.exists(file_data)){
  cat("Se leerán datos del archivo: ", file_data)
} else{
  warning("Cuidado: el archivo no existe.")
}
```

Se leerán datos del archivo: Data/Input/CU\_45\_06\_03\_turismo\_receptor.csv

**Data file to dataframe** Usar la función adecuada según el formato de entrada (xlsx, csv, json, ...)

```
[45]: data <- read_csv(file_data)
```

Rows: 50294 Columns: 8  
Column specification

Delimiter: ","

chr (5): mes, pais\_orig\_cod, pais\_orig, mun\_dest, CMUN

dbl (3): mun\_dest\_cod, turistas, Target

Use `spec()` to retrieve the full column specification for this data.

Specify the column types or set `show\_col\_types = FALSE` to quiet this message.

Visualizo los datos.

Estructura de los datos:

```
[46]: data |> glimpse()
```

Rows: 50,294

Columns: 8

\$ mes <chr> "2019-07", "2019-07", "2019-07",  
"2019-07", "2019-07", "...

\$ pais\_orig\_cod <chr> "000", "010", "011", "030", "110",  
"121", "123", "126", ...

\$ pais\_orig <chr> "Total", "Total Europa", "Total Unión  
Europea", "Total A...

```
$ mun_dest_cod <dbl> 28002, 28002, 28002, 28002, 28002,
28002, 28002, 28002, ...
$ mun_dest <chr> "Ajalvir", "Ajalvir", "Ajalvir",
"Ajalvir", "Ajalvir", "...
$ turistas <dbl> 338, 290, 268, 37, 56, 54, 37, 40, 157,
116, 109, 8461, ...
$ CMUN <chr> "002", "002", "002", "002", "002",
"002", "002", "002", ...
$ Target <dbl> 338, 290, 268, 37, 56, 54, 37, 40, 157,
116, 109, 8461, ...
```

Muestra de los primeros datos:

```
[47]: data |> slice_head(n = 5)
```

	mes <chr>	pais_orig_cod <chr>	pais_orig <chr>	mun_dest_cod <dbl>	mun_dest <chr>	turistas <dbl>
	2019-07	000	Total	28002	Ajalvir	338
A spec_tbl_df: 5 × 8	2019-07	010	Total Europa	28002	Ajalvir	290
	2019-07	011	Total Unión Europea	28002	Ajalvir	268
	2019-07	030	Total América	28002	Ajalvir	37
	2019-07	110	Francia	28002	Ajalvir	56

## 0.6 Missing Values

### 0.6.1 Missing Values Identification

#### Missing Values Per Sample

```
[48]: # Create a dataframe with ID and number of missing values per sample
missing_per_sample_df <- data.frame(ID = 1:nrow(data), NAs = rowSums(is.
  ↪na(data)))

# Filter rows with more than one missing value
filtered_missing_per_sample_df <-
  ↪missing_per_sample_df[missing_per_sample_df$NAs > 1, ]

# Print the filtered dataframe
print(filtered_missing_per_sample_df)
```

```
[1] ID NAs
<0 rows> (or 0-length row.names)
```

#### Missing Values Per Feature

```
[49]: # Calculate the number of missing values per feature
missing_per_feature <- colSums(is.na(data))

# Print the number of missing values per feature
print(missing_per_feature)
```

	mes	pais_orig_cod	pais_orig	mun_dest_cod	mun_dest
	0	0	0	0	0
turistas		CMUN	Target		
	0	0	0		

## Zero Missing Values

```
[50]: # Detecting columns with minimum value of zero (0).

# Calculate the frequency of zeros in each column
zero_counts <- sapply(data, function(x) sum(x == 0, na.rm = TRUE))

# Calculate the proportion of zeros in each column
zero_proportions <- zero_counts / nrow(data)

# Set a threshold for the proportion of zeros
zero_threshold <- 0.9 # Adjust as needed

# Identify columns with a high proportion of zeros
columns_with_high_zeros <- names(zero_proportions[zero_proportions >=
  ↪zero_threshold])

# Print the columns with a high proportion of zeros
print(columns_with_high_zeros)
```

character(0)

Select column to replace

```
[51]: # Select column to replace
```

Operation

```
[52]: # Replace zero missing values by nan
# Replace columns with a high proportion of zeros with NaN
data[, columns_with_high_zeros] <- NA
```

```
[ ]:
```

## 0.6.2 Delete Missing Values

Esta acción no aplica ya que no se han encontrado ningún valor faltante en los datos.

### Deleting Rows with Missing Values in Target Column

```
[53]: # data <- data[!is.na(data$Target), ]
```

**Deleting Rows with Missing Values** Only in case of high data size

```
[54]: # Eliminamos las filas con valores nulos
# data <- data[complete.cases(data), ]
```

**Deleting Features with some Missing Values** Only with many features and for non-relevant features

```
[55]: # # Selecciono las columnas con algún valor missing:
# # Identify columns with any missing values
# cols_with_na <- sapply(data, function(x) any(is.na(x)))

# # Keep only those columns that do not have missing values
# data <- data[, !cols_with_na]
```

### 0.6.3 Basic Imputation

Esta acción no aplica ya que no se han encontrado ningún valor faltante en los datos.

#### Imputation by Previous Row Value

```
[56]: # Sustituimos valores null por otro valor: VALOR FILA ANTERIOR
# # data <- data %>% fill(everything(), .direction = "down")
```

#### Imputation by Next Row Value

```
[116]: # # Sustituimos valores null por otro valor: VALOR FILA SIGUIENTE
# data <- data %>% fill(everything(), .direction = "up")
```

### 0.6.4 Statistical Imputation

A popular approach for data imputation is to calculate a statistical value for each column (such as a mean) and replace all missing values for that column with the statistic.

Esta acción no aplica ya que no se han encontrado ningún valor faltante en los datos.

#### Selection of Imputation Strategy

```
[ ]: # The mean accuracy of each approach can then be compared.
#
# Specific results may vary given the stochastic nature of
# the learning algorithm, the evaluation procedure, or
# differences in numerical precision. Consider running the
# example a few times and compare the average performance.
#
```

```
[ ]: # Plot model performance for comparison
# box and whisker plot is created for each set of results,
# allowing the distribution of results to be compared.
```



### Constant Imputation Select constant value

```
[23]: ## Select constant value
      # constant = 0
```

#### Operation

```
[25]: ## Constant imputation
      # data <- replace_na(data, list(everything() <- constant))
```

### Mean Imputation

```
[26]: ## Identify numeric columns
      # numeric_cols <- sapply(data, is.numeric)

      ## Apply mean imputation only on numeric columns
      # data_imputed <- data
      # for(column in names(data)[numeric_cols]) {
      #   data_imputed[is.na(data_imputed[,column]), column] <- mean(data[,column],
      # ↪na.rm = TRUE)
      # }
```

### Median Imputation

```
[27]: ## Identify numeric columns
      # numeric_cols <- sapply(data, is.numeric)

      ## Apply median imputation only on numeric columns
      # data_imputed <- data
      # for(column in names(data)[numeric_cols]) {
      #   data_imputed[is.na(data_imputed[,column]), column] <-
      # ↪median(data[,column], na.rm = TRUE)
      # }
```

### Most Frequent Imputation

```
[28]: ## Create a function to calculate the mode
      # getmode <- function(v) {
      #   unquv <- unique(v)
      #   unquv[which.max(tabulate(match(v, unquv)))]
      # }

      ## Apply most frequent imputation to each column
      # data_imputed <- data
      # for(column in names(data)) {
      #   data_imputed[is.na(data_imputed[,column]), column] <-
      # ↪getmode(data[,column])
      # }
```

## Interpolation Imputation

```
[29]: ## If not already installed, install the zoo package
      # if(!require(zoo)) install.packages('zoo')

      ## Load the zoo package
      # library(zoo)

      ## Choose an interpolation method
      # interpolation_method <- "linear" # "linear" or "spline"

      ## Apply interpolation imputation
      # data_imputed <- data
      # if(interpolation_method == "linear") {
      #   data_imputed <- na.approx(data_imputed, na.rm = FALSE)
      # } else if(interpolation_method == "spline") {
      #   data_imputed <- na.spline(data_imputed, na.rm = FALSE)
      # } else {
      #   stop("Invalid interpolation method")
      # }
```

### 0.6.5 Prediction Imputation (KNN Imputation )

An approach to missing data imputation is to use a model to predict the missing values.

Esta acción no aplica ya que no se han encontrado ningún valor faltante en los datos.

**Evaluating k-hyperparameter in KNN Imputation** Select numbers of neighbors to evaluate

```
[30]: ## Numbers of neighbors to evaluate
      # k_values <- c(3, 5, 7, 9, 11)
```

Operation

```
[31]: ## If not already installed, install the VIM package
      # if(!require(VIM)) install.packages('VIM')

      ## Load the VIM package
      # library(VIM)

      ## Cross-validation
      # cv_errors <- sapply(k_values, function(k) {
      #   # Apply KNN imputation
      #   data_imputed <- kNN(data, k = k)

      #   # Calculate and return the mean squared error
      #   mean((data - data_imputed)^2, na.rm = TRUE)
      # })
```

```
# # Print the cross-validation errors
# print(cv_errors)
```

## Applying KNN Imputation

[ ]:

Select numbers of neighbors to evaluate

```
[33]: # Number of neighbors
# optimal_k <- k_values[which.min(cv_errors)]
```

Operation

```
[34]: # data_imputed <- kNN(data, k = optimal_k)
```

### 0.6.6 Iterative Imputation

Esta acción no aplica ya que no se han encontrado ningún valor faltante en los datos.

**Evaluating Different Imputation Order** We can experiment with different imputation order strategies, such as descending, right-to-left (Arabic), left-to-right (Roman), and random.

```
[35]: # # If not already installed, install the mice package
# if(!require(mice)) install.packages('mice')

# # Load the mice package
# library(mice)

# # Determine the percentage of missing values in each column
# missing_values <- sapply(data, function(x) sum(is.na(x))/length(x))

# # Order the variables based on the percentage of missing values
# ascending_order <- order(missing_values)
# descending_order <- order(missing_values, decreasing = TRUE)
# random_order <- sample(length(missing_values))

# # Print the imputation orders
# print(ascending_order)
# print(descending_order)
# print(random_order)
```

**Applying Iterative Imputation** Select strategie

```
[36]: # # Choose an imputation order strategy
# imputation_order <- ascending_order # or descending_order, or random_order
```

Operation

```
[38]: # # Perform iterative imputation
# mice_output <- mice(data[, imputation_order], m = 5, maxit = 50, method = 'pmm', seed = 500)

# # Get the imputed data
# data_imputed <- complete(mice_output, 1)
```

```
[ ]:
```

```
[39]: # Generating the new Data dataframe
```

## 0.7 Data Save

- Solo si se han hecho cambios
- No aplica

Identificamos los datos a guardar

```
[61]: data_to_save <- data
```

Estructura de nombre de archivos:

- Código del caso de uso, por ejemplo “CU\_04”
- Número del proceso que lo genera, por ejemplo “\_06”.
- Resto del nombre del archivo de entrada
- Extensión del archivo

Ejemplo: "CU\_04\_06\_01\_01\_zonasgeo.json, primer fichero que se genera en la tarea 01 del proceso 05 (Data Collection) para el caso de uso 04 (vacunas) y que se ha transformado en el proceso 06

Importante mantener los guiones bajos antes de proceso, tarea, archivo y nombre

### 0.7.1 Proceso 09.2

```
[66]: caso <- "CU_45"
proceso <- '_09.2'
tarea <- "_03"
archivo <- ""
proper <- "_turismo_receptor"
extension <- ".csv"
```

OPCION A: Uso del paquete “tcltk” para mayor comodidad

- Buscar carpeta, escribir nombre de archivo SIN extensión (se especifica en el código)
- Especificar sufijo2 si es necesario
- Cambiar datos por datos\_xx si es necesario

```
[67]: # file_save <- paste0(caso, proceso, tarea, tcltk::tkgetSaveFile(), proper,
      ↪extension)
      # path_out <- paste0(oPath, file_save)
      # write_csv(data_to_save_XXXXX, path_out)

      # cat('File saved as: ')
      # path_out
```

OPCION B: Especificar el nombre de archivo

- Los ficheros de salida del proceso van siempre a Data/Output/.

```
[68]: file_save <- paste0(caso, proceso, tarea, archivo, proper, extension)
      path_out <- paste0(oPath, file_save)
      write_csv(data, path_out)

      cat('File saved as: ')
      path_out
```

File saved as:

'Data/Output/CU\_45\_09.2\_03\_turismo\_receptor.csv'

**Copia del fichero a Input** Si el archivo se va a usar en otros notebooks, copiar a la carpeta Input

```
[69]: path_in <- paste0(iPath, file_save)
      file.copy(path_out, path_in, overwrite = TRUE)
```

TRUE

## 0.8 REPORT

A continuación se realizará un informe de las acciones realizadas

### 0.9 Main Actions Carried Out

- Si eran necesarias se han realizado en el proceso 05 por cuestiones de eficiencia

### 0.10 Main Conclusions

- Los datos están limpios para el despliegue

### 0.11 CODE TO DEPLOY (PILOT)

A continuación se incluirá el código que deba ser llevado a despliegue para producción, dado que se entiende efectúa operaciones necesarias sobre los datos en la ejecución del prototipo

Description

- No hay nada que desplegar en el piloto, ya que estos datos son estáticos o en todo caso cambian con muy poca frecuencia, altamente improbable durante el proyecto.

CODE

[ ]: