

## 09.1.- Data Cleansing-Basic\_CU\_18\_20\_infra\_meteo\_v\_01

June 13, 2023

#

CU18\_Infraestructuras\_eventos

Citizenlab Data Science Methodology > II - Data Processing Domain \*\*\* > # 09.1.- Data Cleansing  
- Basic

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\_18\_06\_20\_diario\_infra
- Output File: CU\_18\_09.1\_20\_diario\_infra

## 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.

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

```
'LC_CTYPE=es_ES.UTF-8;LC_NUMERIC=C;LC_TIME=es_ES.UTF-8;LC_COLLATE=es_ES.UTF-8;LC_MONETARY=es_ES.UTF-8;LC_MESSAGES=en_US.UTF-8;LC_PAPER=es_ES.UTF-8;LC_NAME=C;LC_ADDRESS=C;LC_TELEPHONE=C;LC_MEASUREMENT=es_ES.UTF-8;LC_IDENTIFICATION=C'
```

### 0.4.2 Libraries to use

```
[2]: library(readr)
library(dplyr)
# library(sf)
library(tidyr)
library(stringr)
```

Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

filter, lag

The following objects are masked from 'package:base':

```
intersect, setdiff, setequal, union
```

### 0.4.3 Paths

```
[3]: 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

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

OPCION B: Especificar el nombre de archivo

```
[5]: iFile <- "CU_18_06_20_diario_infra.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\_18\_06\_20\_diario\_infra.csv

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

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

```
Rows: 415370 Columns: 10
  Column specification
```

```
Delimiter: ","
```

```
dbl (9): id_inf, capacidad, demanda, evento_infra, evento_zona, tmed,
prec,...
```

```
date (1): fecha
```

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:

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

```
Rows: 415,370
Columns: 10
$ id_inf      <dbl> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12,
13, 14, 15, 16, 17...
$ fecha       <date> 2019-01-01, 2019-01-01, 2019-01-01,
2019-01-01, 2019-01-...
$ capacidad   <dbl> 993, 996, 1036, 1020, 992, 1026, 1007,
976, 1037, 972, 94...
$ demanda     <dbl> 883, 888, 922, 1134, 1103, 1139, 897,
1086, 1150, 861, 83...
$ evento_infra <dbl> 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0,
0, 0, 1, 0, 1, 1, ...
$ evento_zona <dbl> 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0,
1, 1, 1, 1, 0, 1, ...
$ tmed        <dbl> 6.953211, 6.196420, 6.483569, 5.875797,
6.212680, 5.87854...
$ prec        <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, ...
$ velmedia    <dbl> 0.6433886, 0.3417523, 0.4132169,
0.1820178, 0.2118110, 0.1...
$ presMax     <dbl> 952.9357, 950.2191, 950.8051, 951.6768,
953.5118, 952.168...
```

Muestra de los primeros datos:

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

	id_inf	fecha	capacidad	demanda	evento_infra	evento_zona	tmed	p
	<dbl>	<date>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
	1	2019-01-01	993	883	1	1	6.953211	0
A spec_tbl_df: 5 × 10	2	2019-01-01	996	888	0	0	6.196420	0
	3	2019-01-01	1036	922	0	0	6.483569	0
	4	2019-01-01	1020	1134	1	0	5.875797	0
	5	2019-01-01	992	1103	1	1	6.212680	0

## 0.6 Text data analysis

Select columns

```
[9]: # Select column
```

Operation

```
[10]: # Analizar datos de texto y verificar su corrección
      # e.g. faltas ortografía, etc
```

```
[11]: # pasar a mayúsculas todas las columnas de texto
```

## 0.7 Delete Columns Needless/Irrelevant/Private

Select columns

```
[12]: # Select columns
```

Operation

```
[13]: # Eliminamos columnas que consideramos irrelevantes o innecesarias
```

```
[14]: # Eliminamos columnas que consideramos irrelevantes o innecesarias
```

## 0.8 Inconsistent Data

```
[ ]:
```

Select columns and value

```
[15]: # Select column and value
      # e.g. age > 100
```

Operation

```
[16]: # Inconsistent data is unique to each data set and
      # must be searched manually
```

## 0.9 Expected values

```
[17]: # Check for expected value
```

## 0.10 Zeros

```
[18]: # Check for zeroes in data
      zero_counts <- colSums(data == 0, na.rm = TRUE)

      # Print variables with zero counts
      for (variable in names(zero_counts[zero_counts > 0])) {
        num_zeroes <- zero_counts[variable]
        print(paste("Variable:", variable, "- Number of zeroes:", num_zeroes))
      }
```

```
[1] "Variable: evento_infra - Number of zeroes: 207662"
[1] "Variable: evento_zona - Number of zeroes: 207855"
[1] "Variable: prec - Number of zeroes: 218496"
```

## 0.11 Single Value

```
[19]: # Obtener el número de valores diferentes en cada columna
unique_counts <- summarise(data, across(everything(), ~length(unique(.))))

# Imprimir los resultados
print(unique_counts)

# A tibble: 1 × 10
  id_inf fecha capacidad demanda evento_infra evento_zona tmed prec velmedia
  <int> <int>
<int> <int>
<int> <int>
<int> <int> <int>
1 1138 365 267 485 2 2
388782 178287
404343
# 1 more variable: presMax <int>
```

```
[20]: # Columns with a single unique value
#
# NOTE: WE MUST EXCLUDE TARGET COLUMN TO DELETE
```

## 0.12 Very Few Values

Select rate

```
[21]: # Select rate
rate <- 0.1
```

Operation

```
[22]: # Show features with over rate rows being the same value
features_with_same_values <- names(data)[apply(data, 2, function(x) {
  length(unique(x)) <= (1 - rate) * nrow(data)}
)]
print(features_with_same_values)
```

```
[1] "id_inf" "fecha" "capacidad" "demanda" "evento_infra"
[6] "evento_zona" "tmed" "prec" "presMax"
```

```
[23]: # Summarize the number of unique values in each column
# followed by the percentage of unique values for each
# variable as a percentage of the total number of rows
# in the dataset.
```



```

# Summarize the number of unique values in each column
unique_counts <- sapply(data, function(x) length(unique(x)))

# Calculate the percentage of unique values for each variable
percentage_unique <- unique_counts / nrow(data) * 100

# Create a data frame with the results
summary_data <- data.frame(Variable = names(unique_counts), Unique_Count = unique_counts, Percentage_Unique = percentage_unique)

# Print the summary data
print(summary_data)

```

	Variable	Unique_Count	Percentage_Unique
id_inf	id_inf	1138	2.739726e-01
fecha	fecha	365	8.787346e-02
capacidad	capacidad	267	6.428004e-02
demanda	demanda	485	1.167634e-01
evento_infra	evento_infra	2	4.814984e-04
evento_zona	evento_zona	2	4.814984e-04
tmed	tmed	388782	9.359896e+01
prec	prec	178287	4.292245e+01
velmedia	velmedia	404343	9.734526e+01
presMax	presMax	413505	9.955100e+01

Select percent of the number of rows

```

[24]: # Select percent of the number of rows
percentage_threshold <- 0.05

```

Operation

```

[25]: # Summarize columns that have unique values that are less than
# "percentage_threshold" percent of the number of rows.
columns_to_delete <- names(data)[sapply(data, function(x) length(unique(x)) <=
  percentage_threshold * nrow(data))]
columns_to_delete <- setdiff(columns_to_delete, "target_column") # Excluir la
  columna objetivo
columns_to_delete

```

1. 'id\_inf' 2. 'fecha' 3. 'capacidad' 4. 'demanda' 5. 'evento\_infra' 6. 'evento\_zona'

Select percent of the number of rows

```

[26]: # Select percent of the number of rows
# print(colnames(data))

```

Operation

```
[27]: # Delete columns that have unique values that are less than
# "percentage_threshold" percent of the number of rows.
#
# NOTE: WE MUST EXCLUDE TARGET COLUMN TO DELETE
# data <- data[, -which(names(data) %in% columns_to_delete)]
```

```
[ ]:
```

## 0.13 Low Variance

### A) Calculating variances

```
[28]: # Calculate variance for all numeric variables
variance <- summarise(data, across(where(is.numeric), var))

# Print the variance
print(variance)
```

# A tibble: 1 × 9

id_inf	capacidad	demanda	evento_infra	evento_zona	tmed	prec	velmedia	presMax
<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	1.08e5	1002.	13101.					
0.250	0.250	NA	NA	NA	591.			

### B) Automatic calculation and representation of variances

Define thresholds to check

```
[29]: # define thresholds to check
threshold <- 0.05 # Example thresholds, adjust as needed
```

Operation

```
[30]: data_without_date <- data[, !(names(data) %in% "fecha")]
variances <- apply(data_without_date, 2, function(x) var(x, na.rm = TRUE))
```

```
[31]: variances
```

id\_inf	107920.509817776	capacidad	1001.90218346458	demanda	13100.9513797354
evento\_infra	0.250000598808376	evento\_zona	0.250000434369146	tmed	59.1737234836182
prec	16.5613867030549	velmedia	2.66823547053618	presMax	590.751700729737

### C) Delete variables with low variance

Select column

```
[32]: # Select column  
low_variance_columns <- names(variances)[variances < threshold]
```

```
[33]: low_variance_columns
```

Operation

```
[34]: # drop columns with low variance  
data <- data[, !(names(data) %in% low_variance_columns)]
```

## 0.14 Duplicates

Entendido como ERROR -> Eliminar duplicados

### 0.14.1 Identificación de Duplicates

```
[35]: # Buscamos filas duplicadas  
duplicated_rows <- data[duplicated(data),]
```

### 0.14.2 Eliminación de primera fila duplicada

```
[36]: # Eliminamos filas duplicadas  
# Deja la primera y elimina el resto duplicadas  
data_no_duplicates <- data[!duplicated(data),]
```

### 0.14.3 Eliminación de todas las filas duplicadas

```
[37]: # Eliminamos filas duplicadas  
# Elimina todas las filas que están duplicadas (no deja una)  
  
# data_no_duplicates_all <- data %>% distinct()
```

### 0.14.4 Eliminación de filas duplicadas en columnas concretas

Select columns

```
[38]: # Select columns  
# columns <- c("columna1", "columna2")
```

Operation

```
[39]: # Eliminamos filas duplicadas  
# Elimina las filas que están duplicadas en columnas seleccionadas  
# data_no_duplicates_columns <- data %>% distinct(across(all_of(columns)))
```

## 0.15 Data Save

- Solo si se han hecho cambios
- No aplica

Identificamos los datos a guardar

```
[40]: 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

### 0.15.1 Proceso 09.1

```
[41]: caso <- "CU_18"  
proceso <- '_09.1'  
tarea <- "_20"  
archivo <- ""  
proper <- "_diario_infra"  
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

```
[42]: # 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/.

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

```
write_csv(data_to_save, path_out)

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

File saved as:

'Data/Output/CU\_18\_09.1\_20\_diario\_infra.csv'

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

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

TRUE

## 0.16 REPORT

A continuación se realizará un informe de las acciones realizadas

## 0.17 Main Actions Carried Out

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

## 0.18 Main Conclusions

- Los datos están limpios para el despliegue

## 0.19 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

```
[ ]:
```