

09.2.- Data Cleansing-Missing_05_servicios_completo_v_01

June 16, 2023

#

CUxx_Nombre del caso de uso

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: xxxxxxxxxxxx
- Output File: No aplica

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.

```
[19]: 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

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

0.4.3 Paths

```
[21]: 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

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

OPCION B: Especificar el nombre de archivo

```
[23]: iFile <- "CU_34_091_05_servicios_completo.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_34_091_05_servicios_completo.csv

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

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

Rows: 274792 Columns: 19

Column specification

Delimiter: ","

chr (5): Servicio, CMUN, CDIS, CSEC, NSEC

dbl (12): Futbol, nservicios, capacidad, tmed, prec, velmedia,
presMax, t1_...

lgl (1): is_train

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:

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

Rows: 274,792

Columns: 19

\$ Fecha <date> 2022-01-12, 2022-01-31, 2022-01-28,

```

2022-01-06, 2022...
$ Servicio      <chr> "Delivery", "Taxi", "Taxi",
"Delivery", "Delivery", "...
$ CMUN          <chr> "079", "079", "903", "079", "007",
"022", "079", "079...
$ CDIS          <chr> "14", "01", "01", "04", "04", "01",
"16", "01", "16",...
$ CSEC          <chr> "050", "048", "006", "080", "012",
"004", "041", "033...
$ Futbol        <dbl> 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0,
0, 0, 0, 0, 0, 0,...
$ nservicios    <dbl> 58, 5, 0, 14, 60, 50, 13, 4, 3, 9,
68, 12, 0, 1, 14, ...
$ capacidad     <dbl> 80, 69, 56, 80, 70, 56, 80, 69, 69,
69, 80, 80, 69, 6...
$ tmed          <dbl> 7.366319, 8.823406, 7.854915,
4.226603, 4.982656, 7.2...
$ prec          <dbl> -0.009468616, 0.000000000,
0.000000000, 0.010181896, ...
$ velmedia      <dbl> 1.5999961, 1.5114967, 2.2536168,
1.0279945, 1.0387037...
$ presMax       <dbl> 954.7939, 948.9795, 940.2553,
945.1884, 948.9570, 943...
$ t1_1          <dbl> 1094, 1251, 2232, 746, 1080, 2256,
692, 1270, 2229, 8...
$ t3_1          <dbl> 45.4360, 41.6091, 44.2016, 47.1729,
48.5361, 43.2877,...
$ NSEC          <chr> "Madrid - 14.050", "Madrid -
01.048", "Tres Cantos - ...
$ area          <dbl> 38753.96, 15289.89, 124539.78,
89206.78, 24473.30, 34...
$ elevation     <dbl> 658, 635, 719, 710, 693, 710, 702,
635, 690, 710, 690...
$ densidad_hab_km2 <dbl> 28229.3737, 81818.7738, 17921.9842,
8362.5928, 44129....
$ is_train      <lgl> TRUE, TRUE, TRUE, TRUE, TRUE, TRUE,
TRUE, TRUE, TRUE,...

```

Muestra de los primeros datos:

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

	Fecha	Servicio	CMUN	CDIS	CSEC	Futbol	nservicios	capacidad	tmed
	<date>	<chr>	<chr>	<chr>	<chr>	<dbl>	<dbl>	<dbl>	<d
A spec_tbl_df: 5 × 19	2022-01-12	Delivery	079	14	050	1	58	80	7.3
	2022-01-31	Taxi	079	01	048	0	5	69	8.8
	2022-01-28	Taxi	903	01	006	0	0	56	7.8
	2022-01-06	Delivery	079	04	080	0	14	80	4.2
	2022-01-21	Delivery	007	04	012	1	60	70	4.9

0.6 Missing Values

0.6.1 Missing Values Identification

Missing Values Per Sample

```
[27]: missing_values <- is.na(data)
      missing_values_per_sample <- rowSums(missing_values)
```

```
[28]: data = data[missing_values_per_sample == 0, ]
```

Missing Values Per Feature

```
[29]: # No missings
```

Zero Missing Values

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

```
[31]: # Frequency of values by column
```

Select column to replace

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

Operation

```
[33]: # Replace zero missing values by nan
```

```
[ ]:
```

Other Missing Values Select column to replace

```
[34]: # Select column to replace and missing value
```

Operation

```
[35]: # Replace other missing values by nan
```

Null/NaN Missing Values

```
[36]: # Intuitivamente: miramos n° datos en todas las columnas
      # los null no los cuenta --> debe hacer el mismo n° por columna
```

```
[37]: # Podemos mirar directamente info donde viene
```

```
[38]: # Contamos los nulos de forma explícita
```

```
[39]: # summarize the number of rows with missing values for each column
```

```
[ ]:
```

0.6.2 Delete Missing Values

Deleting Rows with Missing Values in Target Column

```
[ ]:
```

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

```
[40]: # Eliminamos las filas con valores nulos
```

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

```
[41]: # Selecciono las columnas con algún valor missing:
```

Deleting Features using Rate Missing Values

```
[42]: # Number of data
```

```
[43]: # Number of missing data
```

```
[44]: # Rate (%) of missing data
```

```
[ ]:
```

Select column to delete

```
[45]: # Select column to delete
```

Operation

```
[46]: # Deleting Features selected
```

0.6.3 Basic Imputation

Imputation by Previous Row Value

```
[47]: # Sustituimos valores null por otro valor: VALOR FILA ANTERIOR
```

Imputation by Next Row Value

```
[48]: # Sustituimos valores null por otro valor: VALOR FILA SIGUIENTE
```

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.

Selection of Imputation Strategy

```
[49]: # 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.  
#
```

```
[50]: # 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

```
[51]: # Select constant value
```

Operation

```
[52]: # Constant imputation
```

Mean Imputation

```
[53]: # Sustituimos valores null por otro valor: MEDIA  
# Miro la media
```

```
[54]: # Sustituyo
```

Median Imputation

```
[55]: # Sustituimos valores null por otro valor: MEDIA  
# Miro la media
```

```
[56]: # Sustituyo
```

Most Frequent Imputation

```
[ ]:
```

```
[57]: # Sustituyo
```

Interpolation Imputation

```
[58]: # Sustituimos valores null por otro valor: INTERPOLANDO  
# Métodos de interpolación  
# 'linear', 'time', 'index', 'values', 'nearest', 'zero', 'slinear',  
# 'quadratic', 'cubic', 'barycentric', 'krogh', 'polynomial', 'spline'  
# 'piecewise_polynomial', 'pchip'
```

0.6.5 Prediction Imputation (KNN Imputation)

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

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

```
[59]: # Numbers of neighbors to evaluate
```

Operation

```
[ ]:
```

Applying KNN Imputation

```
[ ]:
```

Select numbers of neighbors to evaluate

```
[60]: # Number of neighbors
```

Operation

```
[ ]:
```

```
[ ]:
```

```
[61]: # Generating de new Data dataframe
```

0.6.6 Iterative Imputation

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.

```
[62]: # compare iterative imputation strategies for the horse colic dataset
```

Applying Iterative Imputation Select strategie

```
[63]: # Selecting strategie  
# strategies = ['ascending', 'descending', 'roman', 'arabic', 'random']
```

Operation

```
[ ]:
```

```
[ ]:
```

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

0.7 Data Save

- Solo si se han hecho cambios
- No aplica

Identificamos los datos a guardar

```
[65]: 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_34"  
proceso <- '_092'  
tarea <- "_05"  
archivo <- ""  
proper <- "_servicios_completo"  
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_to_save, path_out)

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

File saved as:

'Data/Output/CU_34_092_05_servicios_completo.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

```
[ ]:
```