

08.- Data Split_CU_04_19_vacunacion_completo_v_01

June 8, 2023

#

CU04_Optimización de vacunas

Citizenlab Data Science Methodology > II - Data Processing Domain *** > # 08.- Data Split

Data Split is the process of selecting the appropriate division of the data set into train, test and validation set.

0.1 Tasks

Train and Test Rate Evaluation

Split Train and Test Datasets

0.2 Consideraciones casos CitizenLab programados en R

- Puede que algunas de las tareas de este proceso se realicen en los notebooks de los procesos de deploy al estar relacionadas con otras de esos procesos. En esos casos, en este notebook se referencia al notebook del proceso correspondiente
- Puede que el proceso no aplique a los ficheros del caso de uso, y se indique “No aplica” de forma generalizada.
- Si en el nombre de archivo del notebook no aparece ningún sufijo, el notebook se refiere al caso globalmente

0.3 File

- Input File: CU_04_07_20_vacunacion_gripe_completo.csv
- Output File: CU_04_08_20_vacunacion_gripe_train_and_test.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.

```
[1]: 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

```
[2]: library(readr)
      library(dplyr)
      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_04_07_20_vacunacion_gripe_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_04_07_20_vacunacion_gripe_completo.csv

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

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

Rows: 21736 Columns: 48

Column specification

Delimiter: ","

chr (3): GEOCODIGO, DESBDT, nombre_zona

dbl (45): ano, semana, n_vacunas, n_citas, tmed, prec, velmedia, presMax, be...

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: 21,736

Columns: 48

\$ GEOCODIGO <chr> "001", "001", "001", "001", "001",
"001", "001", "00..."

\$ DESBDT <chr> "Abrantes", "Abrantes", "Abrantes",
"Abrantes", "Abr..."

\$ ano <dbl> 2021, 2021, 2021, 2021, 2021, 2021,
2021, 2021, 2021...

\$ semana <dbl> 36, 37, 38, 39, 40, 41, 42, 43, 44,
45, 46, 47, 48, ...

\$ n_vacunas <dbl> 0, 0, 0, 0, 0, 328, 344, 353, 371,
341, 389, 349, 40...

\$ n_citas <dbl> 0, 0, 0, 0, 0, 305, 327, 328, 341,
334, 373, 332, 37...

\$ tmed <dbl> 23.822231, 20.160014, 18.551058,
18.815387, 17.49447...

\$ prec <dbl> -0.0063023484, 3.5537258042,
3.9769052178, 1.4806472...

\$ velmedia <dbl> 3.573728, 2.494744, 3.316148,
2.384262, 1.850839, 1...

\$ presMax <dbl> 940.9841, 940.7610, 943.1540,
944.6697, 944.2973, 94...

\$ benzene <dbl> 0.1713567, 0.1573829, 0.1858059,

0.1486437, 0.142803...
 \$ co <dbl> 0.1680325, 0.2138607, 0.2034376,
 0.2399882, 0.269345...
 \$ no <dbl> 4.098371, 6.515572, 5.477654,
 9.593391, 18.860535, 1...
 \$ no2 <dbl> 20.09480, 27.42594, 20.74836,
 37.08524, 40.19475, 44...
 \$ nox <dbl> 26.48135, 37.45944, 25.61128,
 52.43745, 74.04903, 75...
 \$ o3 <dbl> 50.03434, 42.41281, 56.29918,
 46.79483, 41.06600, 44...
 \$ pm10 <dbl> 17.447652, 17.658399, 12.844436,
 16.395896, 14.90938...
 \$ pm2.5 <dbl> 3.008675, 10.083070, 7.218588,
 9.426029, 8.131753, 1...
 \$ so2 <dbl> 6.861545, 6.589638, 4.364304,
 3.123598, 1.291137, 1...
 \$ campana <dbl> 2021, 2021, 2021, 2021, 2021, 2021,
 2021, 2021, 2021...
 \$ scampana <dbl> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11,
 12, 13, 14, 15, 1...
 \$ capacidad_zona <dbl> 7477, 7477, 7477, 7477, 7477, 7477,
 7477, 7477, 7477...
 \$ prop_riesgo <dbl> 0.1953542, 0.1953542, 0.1953542,
 0.1953542, 0.195354...
 \$ tasa_riesgo <dbl> 0.006867731, 0.006867731,
 0.006867731, 0.006867731, ...
 \$ tasa_mayores <dbl> 0.03324259, 0.03324259, 0.03324259,
 0.03324259, 0.03...
 \$ poblacion_mayores <dbl> 0.1781619, 0.1781619, 0.1781619,
 0.1781619, 0.178161...
 \$ nombre_zona <chr> "Abrantes", "Abrantes", "Abrantes",
 "Abrantes", "Abr...
 \$ nsec <dbl> 23, 23, 23, 23, 23, 23, 23, 23, 23,
 23, 23, 23, 23, ...
 \$ t3_1 <dbl> 42.31249, 42.31249, 42.31249,
 42.31249, 42.31249, 42...
 \$ t1_1 <dbl> 29872, 29872, 29872, 29872, 29872,
 29872, 29872, 298...
 \$ t2_1 <dbl> 0.5345094, 0.5345094, 0.5345094,
 0.5345094, 0.534509...
 \$ t2_2 <dbl> 0.4654906, 0.4654906, 0.4654906,
 0.4654906, 0.465490...
 \$ t4_1 <dbl> 0.1561935, 0.1561935, 0.1561935,
 0.1561935, 0.156193...
 \$ t4_2 <dbl> 0.6656459, 0.6656459, 0.6656459,
 0.6656459, 0.665645...
 \$ t4_3 <dbl> 0.1781619, 0.1781619, 0.1781619,

```

0.1781619, 0.178161...
$ t5_1          <dbl> 0.2183993, 0.2183993, 0.2183993,
0.2183993, 0.218399...
$ t6_1          <dbl> 0.3450855, 0.3450855, 0.3450855,
0.3450855, 0.345085...
$ t7_1          <dbl> 0.04090796, 0.04090796, 0.04090796,
0.04090796, 0.04...
$ t8_1          <dbl> 0.02880326, 0.02880326, 0.02880326,
0.02880326, 0.02...
$ t9_1          <dbl> 0.2685716, 0.2685716, 0.2685716,
0.2685716, 0.268571...
$ t10_1         <dbl> 0.1665607, 0.1665607, 0.1665607,
0.1665607, 0.166560...
$ t11_1         <dbl> 0.4723181, 0.4723181, 0.4723181,
0.4723181, 0.472318...
$ t12_1         <dbl> 0.5637381, 0.5637381, 0.5637381,
0.5637381, 0.563738...
$ area          <dbl> 1571619, 1571619, 1571619, 1571619,
1571619, 1571619...
$ densidad_hab_km <dbl> 19007.15, 19007.15, 19007.15,
19007.15, 19007.15, 19...
$ tuits_gripe   <dbl> 97, 79, 112, 143, 112, 130, 254,
190, 198, 160, 206,...
$ interes_gripe <dbl> 13, 15, 19, 29, 38, 65, 100, 94,
63, 70, 64, 64, 57,...
$ Target        <dbl> 0, 0, 0, 0, 0, 328, 344, 353, 371,
341, 389, 349, 40...

```

Muestra de los primeros datos:

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

	GEOCODIGO <chr>	DESBDT <chr>	ano <dbl>	semana <dbl>	n_vacunas <dbl>	n_citas <dbl>	tmed <dbl>	prec <dbl>
A spec_tbl_df: 5 × 48	001	Abrantes	2021	36	0	0	23.82223	-0.006
	001	Abrantes	2021	37	0	0	20.16001	3.5537
	001	Abrantes	2021	38	0	0	18.55106	3.9769
	001	Abrantes	2021	39	0	0	18.81539	1.4806
	001	Abrantes	2021	40	0	0	17.49447	-0.033

0.6 Split Rate Evaluation

An soft evaluation is performed in order to estimate Split Rate.

0.6.1 Classification Rate Evaluation

```
[9]: # Implementar código solo si aplica
```

Métricas para la Soft-Evaluation

```
[10]: # Implementar código solo si aplica
```

Modelos para la Soft-Evaluation

```
[11]: # Implementar código solo si aplica
```

Operation

```
[12]: # Implementar código solo si aplica
```

0.6.2 Regression Rate Evaluation

Parámetros para la Soft-Evaluation

```
[13]: # Implementar código solo si aplica
```

Métricas para la Soft-Evaluation

```
[14]: # Implementar código solo si aplica
```

Modelos para la Soft-Evaluation

```
[15]: # Implementar código solo si aplica
```

Operation

```
[16]: # Implementar código solo si aplica
```

Clustering Rate Evaluation You do not use training and valing in unsupervised learning. There is no objective function in unsupervised learning to val the performance of the algorithm.

0.7 Split Train and Test datasets

Estimada el porcentaje de división (Rate of Split) procedemos a realiar la división correspondiente del fichero entre Train y Test.

Parámetros para el Split

```
[17]: # Set the seed for reproducibility
set.seed(123)

# Split train and test datasets
train_ratio <- 0.8 # Adjust this value to set the proportion of data for
  ↪ training
```

Operation

```
[18]: # Generate the train dataset
train_data <- data %>%
  slice_sample(prop = train_ratio, replace = FALSE) %>%
```

```
mutate(is_train = TRUE)

# Generate the test dataset
test_data <- data %>%
  anti_join(train_data)%>%
  mutate(is_train = FALSE)
```

```
Joining with `by = join_by(GEOCODIGO, DESBDT, ano, semana, n_vacunas,
n_citas, tmed, prec, velmedia, presMax, benzene, co, no, no2,
nox, o3, pm10, pm2.5, so2, campana, scampana, capacidad_zona, prop_riesgo,
tasa_riesgo, tasa_mayores, poblacion_mayores, nombre_zona,
nsec, t3_1, t1_1, t2_1, t2_2, t4_1, t4_2, t4_3, t5_1, t6_1, t7_1, t8_1, t9_1,
t10_1, t11_1, t12_1, area, densidad_hab_km, tuits_gripe,
interes_gripe, Target)`
```

0.8 Data Save

- Solo si se han hecho cambios
- No aplica

Identificamos los datos a guardar

```
[ ]:
```

```
[19]: data_to_save <- rbind(train_data, test_data)
data_to_save |> slice_head(n = 5)
```

	GEOCODIGO <chr>	DESBDT <chr>	ano <dbl>	semana <dbl>	n_vacunas <dbl>	n_citas <dbl>	tmed <dbl>	prec <dbl>
	259	V Centenario	2022	34	0	0	27.278748	0.16995
A tibble: 5 × 49	260	Valdeacederas	2022	8	0	0	9.577289	1.26491
	041	Canillejas	2022	9	0	0	8.536554	3.12288
	025	Barajas	2022	49	292	280	9.065363	7.31388
	046	Castelló	2022	24	0	0	29.905728	0.01366

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

0.8.1 Proceso 08

```
[20]: caso <- "CU_04"
      proceso <- '_08'
      tarea <- "_20"
      archivo <- ""
      proper <- "_vacunacion_gripe_train_and_test"
      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

```
[21]: # 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/.

```
[22]: 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_04_08_20_vacunacion_gripe_train_and_test.csv'

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

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

TRUE

0.9 REPORT

A continuación se realizará un informe de las acciones realizadas

0.10 Main Actions Carried Out

- No aplica el proceso al caso

0.11 Main Conclusions

- En este caso no se hace división de datos

0.12 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

[24]: `# incluir código`