

# 18.- Feature Construction\_04\_19\_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 > III - Feature Engineering Domain \*\*\* > # 18.- Feature Construction

Feature Construction is the process related to create new features from your existing ones to improve model performance.

## 0.1 Tasks

Feature Construction - Create Interaction Features - Create derived variables - Combine Sparse Classes - Explore Binning for Feature Construction

## 0.2 Consideraciones casos CitizenLab programados en R

- Algunas 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
- Otras tareas típicas de este proceso se realizan en los notebooks del dominio IV al ser más eficiente realizarlas en el propio pipeline de modelización.
- 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\_08\_03\_turismo\_receptor.csv
- Sampled Input File: CU\_45\_07\_03\_turismo\_receptor.csv
- Output File: No aplica

### 0.3.1 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.

```
[57]: 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 Settings

### 0.4.1 Libraries to use

```
[58]: library(readr)  
library(dplyr)  
library(tidyr)  
library(forcats)  
library(lubridate)
```

### 0.4.2 Paths

```
[59]: 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. Ucomment the line if using this option

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

OPCION B: Especificar el nombre de archivo

```
[62]: iFile <- "CU_45_08_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\_08\_03\_turismo\_receptor.csv

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

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

```
Rows: 50294 Columns: 9  
Column specification
```

```
Delimiter: ","
```

```
chr (5): mes, pais_orig_cod, pais_orig, mun_dest, CMUN
```

```
dbl (3): mun_dest_cod, turistas, Target
```

```
lgl (1): is_train
```

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.

Estructura de los datos:

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

```
Rows: 50,294
```

```
Columns: 9
```

```
$ mes      <chr> "2019-08", "2021-07", "2021-07",  
"2022-01", "2019-08", "...
```

```
$ pais_orig_cod <chr> "110", "010", "010", "000", "128",  
"000", "011", "126", ...
```

```
$ pais_orig   <chr> "Francia", "Total Europa", "Total  
Europa", "Total", "Rum...
```

```
$ mun_dest_cod <dbl> 28161, 28176, 28132, 28141, 28130,  
28126, 28075, 28005, ...
```

```
$ mun_dest    <chr> "Valdemoro", "Villanueva de la Cañada",  
"San Martín de l...
```

```
$ turistas    <dbl> 466, 1375, 465, 54, 135, 30, 285, 768,  
31, 1646, 116, 36...
```

```
$ CMUN        <chr> "161", "176", "132", "141", "130",  
"126", "075", "005", ...
```

```
$ Target      <dbl> 466, 1375, 465, 54, 135, 30, 285, 768,  
31, 1646, 116, 36...
```

```
$ is_train    <lgl> TRUE, TRUE, TRUE, TRUE, TRUE, TRUE,  
TRUE, TRUE, TRUE, TR...
```

Muestra de los primeros datos:

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

	mes	pais_orig_cod	pais_orig	mun_dest_cod	mun_dest	
	<chr>	<chr>	<chr>	<dbl>	<chr>	<
A spec_tbl_df: 5 × 9	2019-08	110	Francia	28161	Valdemoro	4
	2021-07	010	Total Europa	28176	Villanueva de la Cañada	1
	2021-07	010	Total Europa	28132	San Martín de la Vega	4
	2022-01	000	Total	28141	Sevilla la Nueva	5
	2019-08	128	Rumania	28130	San Fernando de Henares	1

## 0.6 Creating Interaction Features

Ver notebooks del proceso 05 Data Collectio

## 0.7 Creating derived variables

Ver notebooks del proceso 05 Data Collectio

## 0.8 Combining Sparse Classes

Ver notebooks del proceso 05 Data Collectio

## 0.9 Binning for Feature Construction

Ver notebooks del proceso 05 Data Collectio