

# 10.- Imbalanced Analysis\_CU\_53\_02\_spi\_v\_01

June 13, 2023

#

CU53\_impacto de las políticas de inversión en sanidad, infraestructuras y promoción turística en el SPI

Citizenlab Data Science Methodology > II - Data Processing Domain \*\*\* > # 10.- Imbalanced Analysis

Data Balancing is the process to obtain an adequate data balance if is required, in order to have the adequate amount of data that reflects the intrinsic structure of the problem to be solved.

## 0.1 Tasks

Imbalanced Analysis

Evaluate Imbalanced Classification Models

Select appropriate metrics

Data Balancing

- Undersampling the Majority Class
- Oversampling the Minority Class
- Mix under-oversampling
- Evaluate a model with random oversampling and undersampling

Cost-Sensitive Algorithms

## 0.2 File

- Input File: CU\_53\_09.2\_02\_spi
- Output File: No aplica

### 0.2.1 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'
```

## 0.3 Settings

### 0.3.1 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

### 0.3.2 Paths

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

## 0.4 Data Load

OPCION A: Seleccionar fichero en ventana para mayor comodidad

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

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

OPCION B: Especificar el nombre de archivo

```
[4]: iFile <- "CU_53_09.2_02_spi.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\_53\_09.2\_02\_spi.csv

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

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

Rows: 2028 Columns: 18  
Column specification

Delimiter: ","

dbl (17): rank\_score\_spi, score\_spi, score\_bhn, score\_fow, score\_opp,  
score\_...

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.

Visualizo los datos.

Estructura de los datos:

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

Rows: 2,028

Columns: 18

\$ rank\_score\_spi <dbl> 80, 97, 46, 84, 99, 150, 74, 105, 36,  
143, 154, 69, 168...

\$ score\_spi <dbl> 67.59, 60.10, 73.96, 62.86, 61.43,  
45.57, 66.56, 59.45,...

\$ score\_bhn <dbl> 79.16, 74.55, 81.88, 79.45, 77.84,  
47.15, 80.41, 66.16,...

\$ score\_fow <dbl> 65.40, 51.25, 70.69, 61.22, 57.63,  
45.21, 62.82, 54.62,...

\$ score\_opp <dbl> 58.22, 54.49, 69.32, 47.92, 48.83,  
44.34, 56.46, 57.56,...

\$ score\_nbmc <dbl> 86.67, 72.88, 86.33, 83.91, 87.72,  
54.66, 92.38, 72.21,...

\$ score\_ws <dbl> 86.44, 83.35, 88.07, 77.71, 78.15,  
47.82, 78.47, 66.32,...

\$ score\_sh <dbl> 87.69, 77.17, 89.59, 85.11, 86.61,  
36.59, 85.21, 75.91,...

\$ score\_ps <dbl> 55.85, 64.81, 63.55, 71.08, 58.87,  
49.53, 65.57, 50.21,...

\$ score\_abk <dbl> 74.20, 47.04, 89.07, 65.15, 55.79,  
50.36, 81.61, 68.71,...

\$ score\_aic <dbl> 74.19, 37.15, 68.14, 51.25, 78.17,  
33.84, 61.95, 56.61,...

```
$ score_hw      <dbl> 53.55, 64.58, 61.41, 62.00, 45.35,
36.99, 61.64, 41.87,...
$ score_eq      <dbl> 59.66, 56.22, 64.13, 66.47, 51.22,
59.66, 46.07, 51.28,...
$ score_pr      <dbl> 81.60, 71.05, 90.28, 61.56, 60.41,
69.20, 70.02, 74.13,...
$ score_pfc     <dbl> 60.29, 64.77, 67.65, 56.51, 58.62,
40.61, 62.49, 59.83,...
$ score_incl    <dbl> 40.24, 56.12, 68.48, 48.70, 35.57,
41.81, 36.89, 55.73,...
$ score_aae     <dbl> 50.73, 26.03, 50.87, 24.90, 40.72,
25.72, 56.45, 40.54,...
$ is_train      <lgl> TRUE, TRUE, TRUE, TRUE, TRUE, TRUE,
TRUE, TRUE, TRUE, T...
```

Muestra de los primeros datos:

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

	rank_score_spi <dbl>	score_spi <dbl>	score_bhn <dbl>	score_fow <dbl>	score_opp <dbl>	score_nbmc <dbl>	score <dbl>
A spec_tbl_df: 5 × 8	80	67.59	79.16	65.40	58.22	86.67	86.44
	97	60.10	74.55	51.25	54.49	72.88	83.35
	46	73.96	81.88	70.69	69.32	86.33	88.07
	84	62.86	79.45	61.22	47.92	83.91	77.71
	99	61.43	77.84	57.63	48.83	87.72	78.15

## 0.5 Imbalanced Analysis

```
[ ]:
```

```
[8]: # Visualizando los datos para comprobar como de balanceados están
```

```
[9]: # summarize class distribution
```

```
[10]: # Generate and plot imbalanced classification dataset
```

## 0.6 Evaluate Imbalanced Classification Models

```
[ ]:
```

## 0.7 Undersampling the Majority Class

```
[11]: # example of undersampling the majority class
```

```
[12]: # Generate new data set
```

```
[ ]:
```

```
[13]: # Visualizando los datos para comprobar como de balanceados están
```

## 0.8 Oversampling the Minority Class

```
[14]: # example of oversampling the minority class
```

```
[15]: # Generate new data set
```

```
[ ]:
```

```
[16]: # Visualizando los datos para comprobar como de balanceados están
```

## 0.9 Combine Data Undersampling and Oversampling with SMOTEENN

```
[17]: # example of both undersampling and oversampling
```

```
[18]: # Generate new data set
```

```
[ ]:
```

```
[19]: # Visualizando los datos para comprobar como de balanceados están
```

## 0.10 Evaluating a model with random oversampling and undersampling

```
[ ]:
```

## 0.11 Cost-Sensitive Algorithms

```
[20]: # example of cost sensitive logistic regression for imbalanced classification
```

## 0.12 Data Save

- Solo si se han hecho cambios
- No aplica

Identificamos los datos a guardar

```
[21]: 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

### 0.12.1 Proceso 10

```
[22]: caso <- "CU_53"
      proceso <- '_10'
      tarea <- "_02"
      archivo <- ""
      proper <- "_spi"
      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

```
[23]: # 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/.

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

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

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

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

## 0.13 REPORT

A continuación se realizará un informe de las acciones realizadas

### 0.14 Main Actions Carried Out

- No aplica el proceso al caso

## 0.15 Main Conclusions

- Al no haber clasificación no procede analizar el desbalanceo.

## 0.16 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

[ ]: