

# Aprendizaje Supervisado

Abraham Zamudio

VI Programa de Especialización en Machine Learning con Python

2020



**CTIC-UNI**  
BUSINESS SCHOOL

# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

En el ejemplo que vamos a utilizar, vamos a imaginarnos que una organización sin fines de lucro soporta su operación mediante la organización periódica de una campaña para recaudar fondos por correo. Esta organización ha creado una base de datos con más de 40 mil personas que por lo menos una vez en el pasado ha sido donante. La campaña de recaudación de fondos se realiza mediante el envío a una lista de correo (o un subconjunto de ella) de un regalo simbólico y la solicitud de una donación. Una vez que se planifica la campaña, el costo total de la misma se conoce de forma automática:

$$[\text{número de potenciales donantes a contactar}] \times ([\text{costo de regalo}] + [\text{costo de correo}])$$

# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

Sin embargo, el resultado de la recaudación de fondos depende tanto del número de donantes que responde a la campaña, como del importe medio de dinero que es donado.

La idea es que, utilizando las técnicas de Machine Learning sobre la base de datos de esta organización, podamos ayudarla a maximizar los beneficios de la campaña de recaudación, esto es, lograr el máximo importe posible de dinero recaudado, minimizando lo más que se pueda el costo total de la campaña. Debemos tener en cuenta que un miembro de la organización le enviará el correo a un potencial donante, siempre que el rendimiento esperado del pedido excede el costo del correo con la solicitud de donación.

# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

Para nuestro ejemplo, el costo por donante de la campaña va a ser igual al **[costo de regalo] + [costo de correo]**, y esto va a ser igual a \$ 0.75 por correo enviado. Los ingresos netos de la campaña se calculan como la suma (**importe de donación real - \$ 0.75**) sobre todos los donantes a los que se ha enviado el correo. Nuestro objetivo es ayudar a esta organización sin fines de lucro a seleccionar de su lista de correo los donantes a los que debe abordar a los efectos de maximizar los beneficios de la campaña de recaudación.

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

### Los Datos

El dataset que vamos a utilizar, consiste en la base de datos de la organización sin fines de lucro con la lista de correo de los donantes de sus campañas anteriores.

► Data1marzo

El contenido debe ser :

1. Entrenamiento.csv
2. Validacion.csv
3. DescripcionData.txt

# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

### Los Datos

Algunos otros datos a tener en cuenta, son los siguientes:

1. El dataset de aprendizaje contiene 47720 registros y 481 columnas. La primera fila / cabecera del mismo contiene los nombres de cada campo.
2. El dataset de validación contiene 47692 registros y 479 columnas. Al igual que en el caso anterior, la primera fila contiene los nombres de cada campo.
3. Los registros del dataset de validación son idénticos a los registros del dataset de aprendizaje, excepto que los valores para nuestros campos objetivo que necesitamos para el aprendizaje, no existen(es decir, las columnas DONOR\_FLAG y DONOR\_AMOUNT no están incluidas en el dataset de validación).

# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

### Los Datos

4. Los espacios en blanco en los campos de tipo texto y los puntos en los campos de tipo numérico corresponden a valores faltantes o perdidos.
5. Cada registro tiene un identificador único de registro o índice (campo `IDX`). Para cada registro, hay dos variables objetivo (campos `DONOR_FLAG` y `DONOR_AMOUNT`). `DONOR_FLAG` es una variable binaria que indica si ese registro fue donante o no; mientras que `DONOR_AMOUNT` contiene el importe de la donación para los casos que fueron donantes.

# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

### Los Datos

6. Algunos de los valores en el dataset pueden contener errores de formato o de ingreso. Por lo que se deberían corregir o limpiar.
7. Una descripción detallada del significado de cada columna del dataset, la pueden encontrar en el archivo `DescripcionData.txt`.



# Preprocesamiento de Datos en Python

## Análisis exploratorio y preprocesamiento

El primer paso que deberíamos emprender, es realizar un pequeño análisis exploratorio de nuestro dataset; es decir, valernos de algunas herramientas de la estadística, junto con algunas visualizaciones para entender un poco más los datos de los que disponemos.

```
1  # Importando las librerías que vamos a utilizar
2  import pandas as pd
3  import numpy as np
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6  from sklearn.preprocessing import LabelEncoder
7
8  # graficos incrustados
9  %matplotlib inline
10
11 # parametros esteticos de seaborn
12 sns.set_palette("deep", desat=.6)
13 sns.set_context(rc={"figure.figsize": (8, 4)})
```

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

```
1  # importando el dataset a un Dataframe de Pandas  
2  ONG_data = pd.read_csv('data/Entrenamiento.csv', header=0)  
3  
4  # Examinando las primeras 5 filas y 5 columnas del dataset  
5  ONG_data.iloc[:5, :5]  
6  
7  # Contamos la cantidad de registros  
8  ONG_data['DONOR_AMOUNT'].count()
```

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

Comenzemos a explorar un poco más en detalle este dataset!. En primer lugar, lo que deberíamos hacer es controlar si existen valores faltantes o nulos; esto lo podemos realizar utilizando el método `isnull()` del siguiente modo:

► `pandas.isnull`

```
1  # Controlando valores nulos
2  ONG_data.isnull().any().any()
```

Como podemos ver, el método nos devuelve el valor "True", lo que indica que existen valores nulos en nuestro dataset. Estos valores pueden tener una influencia significativa en nuestro modelo predictivo, por lo que siempre es una decisión importante determinar la forma en que los vamos a manejar.

# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

Las alternativas que tenemos son:

1. Dejarlos como están, lo que a la larga nos va a traer bastantes dolores de cabeza ya que en general los algoritmos no los suelen procesar correctamente y provocan errores.
2. Eliminarlos, lo que es una alternativa viable aunque, dependiendo la cantidad de valores nulos, puede afectar significativamente el resultado final de nuestro modelo predictivo.
3. Inferir su valor. En este caso, lo que podemos hacer es tratar de inferir el valor faltante y reemplazarlo por el valor inferido. Esta suele ser generalmente la mejor alternativa a seguir.

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

Vamos a inferir los valores faltantes utilizando la media aritmética para los datos cuantitativos y la moda para los datos categóricos. Como vamos a utilizar dos métodos distintos para reemplazar a los valores faltantes, dependiendo de si son numéricos o categóricos, el primer paso que debemos realizar es tratar de identificar que columnas de nuestro dataset corresponde a cada tipo de datos; para realizar esto vamos a utilizar el atributo `dtypes` del Dataframe de Pandas.

```
1  # Agrupando columnas por tipo de datos
2  tipos = ONG_data.columns.to_series()
3  tipos = tipos.groupby(ONG_data.dtypes).groups
4
5  # Armando lista de columnas categoricas
6  ctext = tipos[np.dtype('object')]
7  len(ctext)  # cantidad de columnas con datos categoricos.
```

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

```
1  # Armando lista de columnas numericas  
2  columnas = ONG_data.columns # lista de todas las columnas  
3  cnum = list(set(columnas) - set(ctext))  
4  len(cnum)
```

Ahora ya logramos separar a las 481 columnas que tiene nuestro dataset. 68 columnas contienen datos categóricos y 413 contienen datos cuantitativos.

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

```
1  # Completando valores faltantes : datos cuantitativos
2  for c in cnum:
3      mean = ONG_data[c].mean()
4      ONG_data[c] = ONG_data[c].fillna(mean)
5
6  # Completando valores faltantes : datos categoricos
7  for c in ctext:
8      mode = ONG_data[c].mode()[0]
9      ONG_data[c] = ONG_data[c].fillna(mode)
10
11 # Verificando que no hayan valores faltantes
12 ONG_data.isnull().any().any()
```

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

Con este criterio ya realizado procedemos a guardar nuestro dataframe en disco duro.

```
1  # Guardando el dataset preprocesado
2  # Con las transformaciones ya realizadas
3  ONG_data.to_csv("data/Entrenamiento_procesado.csv",
4                  index=False)
```

Perfecto! Ahora tenemos un dataset limpio de valores faltantes.



# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

Ya estamos listos para comenzar a explorar los datos, comencemos por determinar el porcentaje de personas que alguna vez fue donante de la ONG y están incluidos en la base de datos con la que estamos trabajando.

```
1  # Calculando el porcentaje de donantes sobre
2  # toda la base de datos
3  percent_donantes = (ONG_data[ONG_data.DONOR_AMOUNT > 0]
4                      ['DONOR_AMOUNT'].count() *
5                      1.0 / ONG_data['DONOR_AMOUNT']
6                      .count()) * 100.0
7  print("El porcentaje de donantes de la base de datos es
8        {0:.2f}%".format(percent_donantes))
```

Vemos que el resultado es: El porcentaje de donantes de la base de datos es 5.08%

# Preprocesamiento de Datos en Python

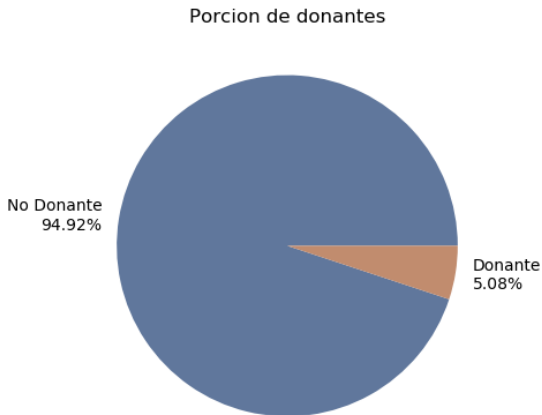
## Preprocesamiento y exploración

### Proporción de Donantes

```
1  # Grafico de totas del porcentaje de donantes
2  # Agrupando por DONOR_FLAG
3  donantes = ONG_data.groupby('DONOR_FLAG').IDX.count()
4  # Creando las leyendas del grafico.
5  labels = [ 'Donante\n' + str(round(x * 1.0/donantes.sum()*
6                                100.0, 2)) + '%'
7             for x in donantes ]
8  labels[0] = 'No ' + labels[0]
9  plt.pie(donantes, labels=labels)
10 plt.title('Porcion de donantes')
11 plt.savefig("data/Pie_Donantes.png")
```

# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*



# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

```
1  # Creando subset con solo los donates  
2  ONG_donantes = ONG_data[ONG_data.DONOR_AMOUNT > 0]  
3  
4  # cantidad de donantes  
5  len(ONG_donantes)
```

Aquí podemos ver que el porcentaje de personas que fueron donantes en el pasado es realmente muy bajo, solo un 5% del total de la base de datos (**2423** personas). Este es un dato importante a tener en cuenta ya que al existir tanta diferencia entre las clases a clasificar, esto puede afectar considerablemente a nuestro algoritmo de aprendizaje.

# Preprocesamiento de Datos en Python

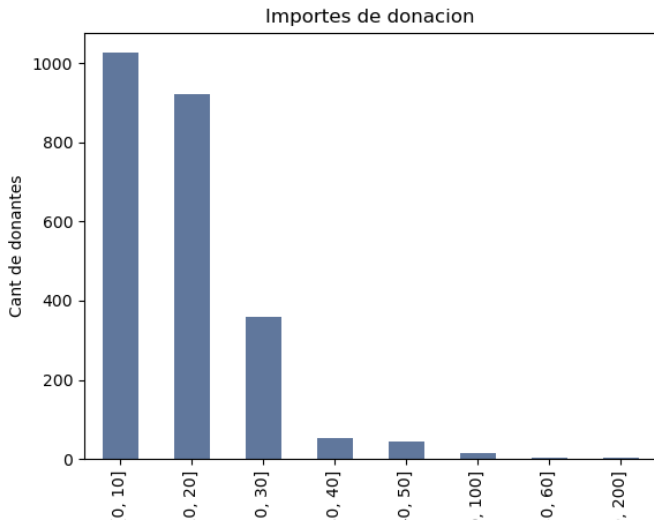
## Preprocesamiento y exploración

Exploremos también un poco más en detalle a este grupo pequeño de personas que fueron donantes; veamos por ejemplo como se dividen de acuerdo a la cantidad de dinero donado.

```
1  # Analizando el importe de donaciones
2  # Creando un segmentos de importes
3  imp_segm = pd.cut(ONG_donantes['DONOR_AMOUNT'],
4                    [0, 10, 20, 30, 40, 50, 60, 100, 200])
5  # Creando el grafico de barras desde pandas
6  plot = pd.value_counts(imp_segm).plot(kind='bar',
7                                       title='Importes[donacion]')
8  plot.set_ylabel('Cant de donantes')
9  plot.set_xlabel('Rango de importes')
10 plt.savefig("data/Histograma_Donantes.png")
```

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración



# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

```
1 # Agrupacion por segmento segun importe donado.  
2 pd.value_counts(imp_seg)
```

(0, 10]	1026
(10, 20]	921
(20, 30]	358
(30, 40]	53
(40, 50]	43
(60, 100]	15
(50, 60]	4
(100, 200]	3

Name: DONOR\_AMOUNT, dtype: int64

# Preprocesamiento de Datos en Python

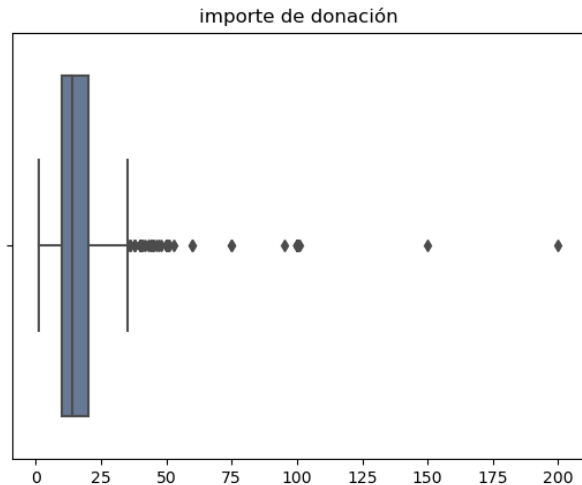
## Preprocesamiento y exploración

```
1  # importe de donacion promedio
2  ONG_donantes['DONOR_AMOUNT'].mean()
3
4
5  # Grafico de cajas del importe de donacion
6  sns.boxplot(list(ONG_donantes['DONOR_AMOUNT']))
7  plt.title('importe de donacion')
8  plt.savefig("data/Boxplot_Donacion.png")
```



# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración



# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

Este análisis nos muestra que la mayor cantidad de donaciones caen en un rango de importes entre 0 y 30, siendo la donación promedio 15.60. También podemos ver que donaciones que superen un importe de 50 son casos realmente poco frecuentes, por lo que constituyen valores atípicos y sería prudente eliminar estos casos al entrenar nuestro modelo para que no distorsionen los resultados.

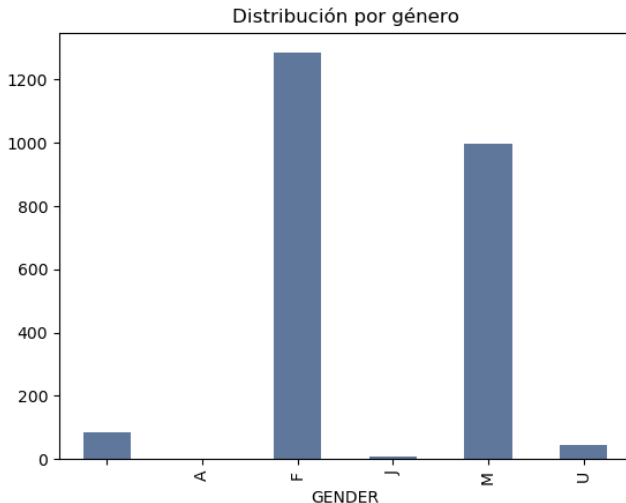
# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

```
1 # Grafico del genero de los donantes
2 ONG_donantes.groupby('GENDER').size().plot(kind='bar')
3 plt.title('Distribucion por genero')
4 plt.savefig("data/GrafBarras_Genero.png")
```

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración



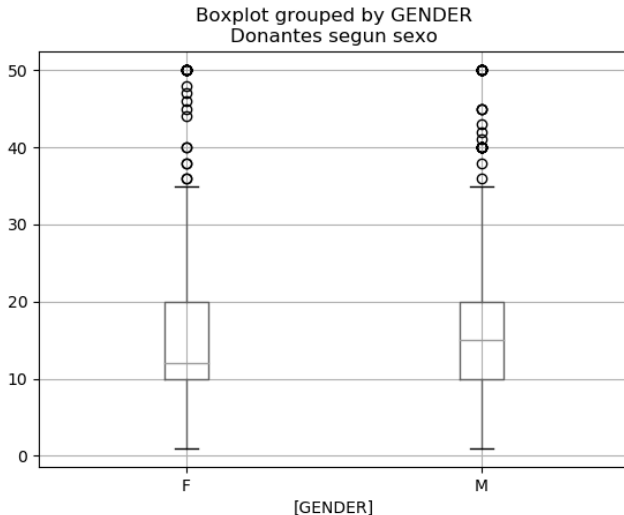
# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

```
1  # Donaciones segun el genero
2  ONG_donantes[(ONG_donantes.DONOR_AMOUNT <= 50)
3                & (ONG_donantes.GENDER.isin(['F', 'M']))
4                ][['DONOR_AMOUNT', 'GENDER']]
5                .boxplot(by='GENDER')
6  plt.title('Donantes segun sexo')
7  plt.savefig("data/Boxplot_genero_FM-png")
```

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración



# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

```
1  # Media de impote donado por mujeres
2  ONG_donantes[ONG_donantes.GENDER == 'F'][['DONOR_AMOUNT']]
3      .mean()
4
5
6  # Media de impote donado por hombres
7  ONG_donantes[ONG_donantes.GENDER == 'M'][['DONOR_AMOUNT']]
8      .mean()
```

Aquí vemos que las mujeres suelen estar más propensas a donar, aunque donan un importe promedio menor (14.61) al que donan los hombres (16.82).

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

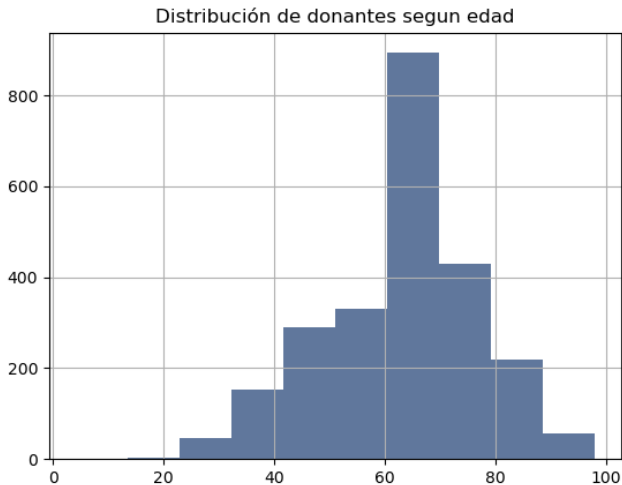
Veamos ahora como se comportan las donaciones respecto a la

```
1  # Distribucion de la edad de los donantes
2  ONG_donantes['AGE'].hist()
3      .set_title('Distribucion de
4                  donantes segun edad')
5  plt.savefig("data/hist_edad-png")
```



# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración



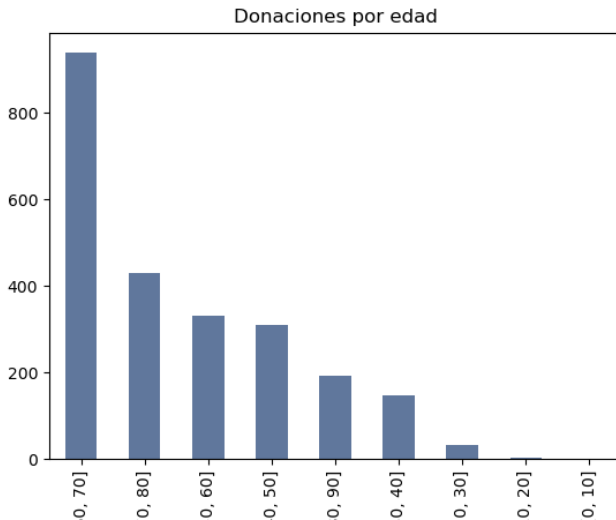
# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

```
1  # Agrupando la edad por rango de a 10
2  AGE2 = pd.cut(ONG_donantes['AGE'], range(0, 100, 10))
3  ONG_donantes['AGE2'] = AGE2
4
5  # Grafico de barras de donaciones por edad
6  pd.value_counts(AGE2).plot(kind='bar',
7                               title='Donaciones por edad')
8  plt.savefig("data/barras_edad.png")
```

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración



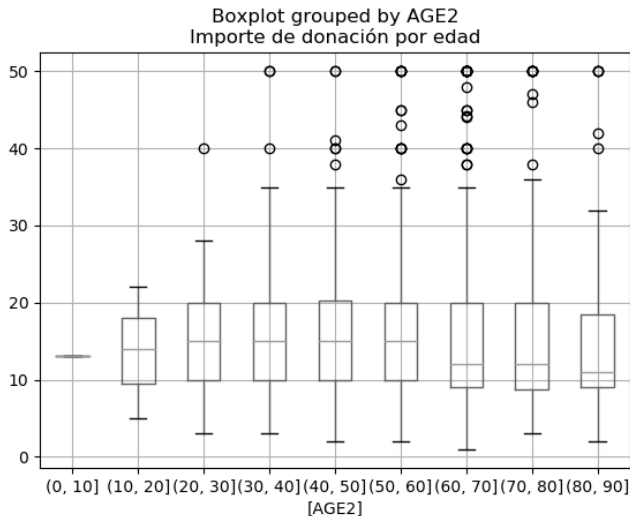
# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración

```
1  # Importes de donacion por grango de edad
2  ONG_donantes[ONG_donantes.DONOR_AMOUNT <= 50][['DONOR_AMOUNT',
3          'AGE2']].boxplot(by='AGE2')
4  plt.title('Importe de donacion por edad')
5  plt.savefig("data/ImporteEdad.png")
```

# Preprocesamiento de Datos en Python

## Preprocesamiento y exploración



# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

En este último análisis podemos ver que la mayor cantidad de los donantes son personas de entre 60 y 70 años, aunque la media de importe donado más alta la tienen las personas que van desde los 30 a los 60 años.

# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

1

# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

1



# *Preprocesamiento de Datos en Python*

## *Preprocesamiento y exploración*

1