

Universidad del Valle de Guatemala

Data Science

Julio Avila

Elisa Samayoa

Laboratorio 10

Para este laboratorio, no se especificó que variable debía predecirse, tampoco que base de datos de las dos usar o que modelos debían utilizarse específicamente, por lo que se decidió predecir el nivel de pobreza tanto para hogares como para personas según distintas características encontradas en las bases de datos.

Carga y Limpieza de Datos

Para la carga y limpieza de datos primero se descargaron de canvas, los datos estaban en formato .sav; por lo que tuvo que investigarse como trabajar ese tipo de datos.

Además, se notó que habían demasiados valores faltantes en la base de datos, por lo que fue necesario hacer un proceso de limpieza de los datos para eliminar las columnas que tenían una gran cantidad de datos faltantes y en las que no eran muchos los datos faltantes se colocó la media ya que si se eliminaban todas las filas o columnas con datos faltantes, las bases de datos quedaban demasiado pequeñas.

```
import pyreadstat

# Cargar datos desde un archivo SPSS
hogar, meta = pyreadstat.read_sav("ENCOVI_Hogar.sav")

# Acceder a Los datos y La metadata
print(hogar.head()) # Muestra Las primeras filas de datos
#print(hogar)       # Muestra La metadata

# Cargar datos desde un archivo SPSS
personas, meta = pyreadstat.read_sav("ENCOVI_Personas.sav")

# Acceder a Los datos y La metadata
print(personas.head()) # Muestra Las primeras filas de datos
#print(personas)
```

```

missing_values_hogar = hogar.isnull().sum()
missing_values_personas = personas.isnull().sum()

print("Missing values in 'hogar' data frame:")
print(missing_values_hogar)

print("\nMissing values in 'personas' data frame:")
print(missing_values_personas)

```

Missing values in 'hogar' data frame:

```

REGION      0
DEPTO       0
AREA        0
UPM         0
NUMHOG      0

```

...

```

P01H15      2641
P01H16      2641
DIA_ENC     0
MES_ENC     1
A_ENC       1

```

Length: 163, dtype: int64

Missing values in 'personas' data frame:

```

REGION      0
DEPTO       0
AREA        0
UPM         0
NUMHOG      0

```

...

```

P11B02B     54247
P11B03A     8999
P11B03B     54612
P11B04A     9001
P11B04B     54549

```

Length: 465, dtype: int64

Aquí se puede observar la carga de datos y donde se encontró la gran cantidad de valores faltantes, el resto del proceso de limpieza y pre procesamiento puede observarse en el notebook ya que fue un proceso bastante extenso.

Después de la limpieza las bases de datos quedaron así:

per_sonho3_cleaned_1_and_2_green_imputed = per_sonho3_cleaned_1_and_2_121118(per_sonho3_cleaned_1_and_2_121118)

Final cleaned 'hogar' data frame with mean imputation:

	REGION	DEPTO	AREA	UPM	NUMHOG	FACTOR	FACTOR3	POBREZA	THOGAR	\
0	1.0	1.0	1.0	1.0	1.0	525.0	1575.0	3.0	3.0	
1	1.0	1.0	1.0	1.0	2.0	525.0	1575.0	3.0	3.0	
2	1.0	1.0	1.0	1.0	3.0	525.0	9450.0	2.0	18.0	
3	1.0	1.0	1.0	1.0	4.0	525.0	3150.0	2.0	6.0	
4	1.0	1.0	1.0	1.0	5.0	525.0	1575.0	3.0	3.0	
...	
11531	8.0	17.0	2.0	833.0	11532.0	305.0	915.0	3.0	3.0	
11532	8.0	17.0	2.0	833.0	11533.0	305.0	3355.0	1.0	11.0	
11533	8.0	17.0	2.0	833.0	11534.0	305.0	1220.0	2.0	4.0	
11534	8.0	17.0	2.0	833.0	11535.0	305.0	915.0	3.0	3.0	
11535	8.0	17.0	2.0	833.0	11536.0	305.0	1830.0	1.0	6.0	

	PPB01	...	P01H10	P01H11	P01H12	P01H13	P01H14	P01H15	\
0	2.0	...	1.582125	1.592917	2.006633	1.839236	1.879708	2.047218	
1	2.0	...	1.582125	1.592917	2.006633	1.839236	1.879708	2.047218	
2	2.0	...	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	
3	2.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
4	2.0	...	1.582125	1.592917	2.006633	1.839236	1.879708	2.047218	
...	
11531	2.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
11532	2.0	...	1.000000	1.000000	2.000000	1.000000	1.000000	2.000000	
11533	2.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
11534	2.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
11535	2.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	

	P01H16	DIA_ENC	MES_ENC	A_ENC
0	2.063631	22.0	8.0	2014.0
1	2.063631	27.0	8.0	2014.0
2	1.000000	22.0	8.0	2014.0
3	2.000000	24.0	8.0	2014.0
4	2.063631	22.0	8.0	2014.0
...
11531	2.000000	5.0	12.0	2014.0
11532	2.000000	5.0	12.0	2014.0

11533 2.000000 5.0 12.0 2014.0

	REGION	DEPTO	AREA	UPM	NUMHOG	FACTOR	POBREZA	THOGAR	ID	\
0	1.0	1.0	1.0	1.0	1.0	525.0	3.0	3.0	1.0	
1	1.0	1.0	1.0	1.0	1.0	525.0	3.0	3.0	2.0	
2	1.0	1.0	1.0	1.0	1.0	525.0	3.0	3.0	3.0	
3	1.0	1.0	1.0	1.0	2.0	525.0	3.0	3.0	1.0	
4	1.0	1.0	1.0	1.0	2.0	525.0	3.0	3.0	2.0	
...	
54817	8.0	17.0	2.0	833.0	11536.0	305.0	1.0	6.0	2.0	
54818	8.0	17.0	2.0	833.0	11536.0	305.0	1.0	6.0	3.0	
54819	8.0	17.0	2.0	833.0	11536.0	305.0	1.0	6.0	4.0	
54820	8.0	17.0	2.0	833.0	11536.0	305.0	1.0	6.0	5.0	
54821	8.0	17.0	2.0	833.0	11536.0	305.0	1.0	6.0	6.0	

	PPA02	...	P11A05A	P11A06A	P11A07A	P11A08A	P11A09A	P11A10A	\
0	1.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
1	2.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
2	2.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
3	2.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
4	1.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
...	
54817	2.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
54818	1.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
54819	2.0	...	2.000000	2.000000	2.000000	2.000000	2.000000	2.000000	
54820	2.0	...	1.945115	1.970561	1.998276	1.992515	1.997359	1.99952	
54821	2.0	...	1.945115	1.970561	1.998276	1.992515	1.997359	1.99952	

	P11B01A	P11B02A	P11B03A	P11B04A
0	2.000000	2.000000	2.000000	2.000000
1	2.000000	2.000000	2.000000	2.000000
2	2.000000	2.000000	2.000000	2.000000
3	2.000000	2.000000	2.000000	2.000000
4	2.000000	2.000000	2.000000	2.000000
...
54817	2.000000	2.000000	2.000000	2.000000
54818	2.000000	2.000000	2.000000	2.000000
54819	2.000000	2.000000	2.000000	2.000000
54820	1.965345	1.987452	1.995417	1.994042
54821	1.965345	1.987452	1.995417	1.994042

[54822 rows x 100 columns]

Tras este proceso se obtuvieron las variables con una correlación con la variable POBREZA mayor o igual al valor absoluto de 0.1 para ser las variables con las que se trabajaría en el análisis estadístico y posteriormente elegir cuales se incluirían en los modelos.

```
# Calcular correlación en la base de datos 'hogar'
correlation_hogar = hogar_cleaned_final_mean_imputed.corr()['POBREZA']

# Filtrar variables con correlación mayor o igual a 0.5 en 'hogar'
high_corr_variables_hogar = correlation_hogar[abs(correlation_hogar) >= 0.1].index.tolist()

# Imprimir lista de variables con correlación mayor o igual a 0.1 en 'hogar'
print("Variables en 'hogar' con correlación mayor o igual a 0.1 con 'POBREZA':")
print(high_corr_variables_hogar)

# Calcular correlación en la base de datos 'personas'
correlation_personas = personas_cleaned_final_mean_imputed.corr()['POBREZA']

# Filtrar variables con correlación mayor o igual a 0.5 en 'personas'
high_corr_variables_personas = correlation_personas[abs(correlation_personas) >= 0.1].index.tolist()

# Imprimir lista de variables con correlación mayor o igual a 0.1 en 'personas'
print("\nVariables en 'personas' con correlación mayor o igual a 0.1 con 'POBREZA':")
print(high_corr_variables_personas)
```

C:\Users\USUARIO\AppData\Local\Temp\ipykernel_1344\2691544929.py:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

```
correlation_hogar = hogar_cleaned_final_mean_imputed.corr()['POBREZA']
```

Variables en 'hogar' con correlación mayor o igual a 0.1 con 'POBREZA':

```
['REGION', 'DEPTO', 'AREA', 'UPM', 'NUMHOG', 'FACTOR3', 'POBREZA', 'THOGAR', 'PPB04', 'PPD04A', 'PPD06', 'PPD07', 'PPD08', 'P01A02', 'P01A04', 'P01A05A', 'P01A05B', 'P01A05C', 'P01A05D', 'P01A05E', 'P01A05F', 'P01A06', 'P01B02', 'P01D01', 'P01D02', 'P01D04', 'P01D08', 'P01D09', 'P01D17', 'P01D19A', 'P01D19B', 'P01D19C', 'P01D19D', 'P01D19E', 'P01D19F', 'P01E01_1', 'P01E01_2', 'P01E01_3', 'P01E01_4', 'P01E01_5', 'P01E01_6', 'P01E02_6', 'P01E03_6', 'P01E01_8', 'P01H09']
```

C:\Users\USUARIO\AppData\Local\Temp\ipykernel_1344\2691544929.py:12: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

```
correlation_personas = personas_cleaned_final_mean_imputed.corr()['POBREZA']
```

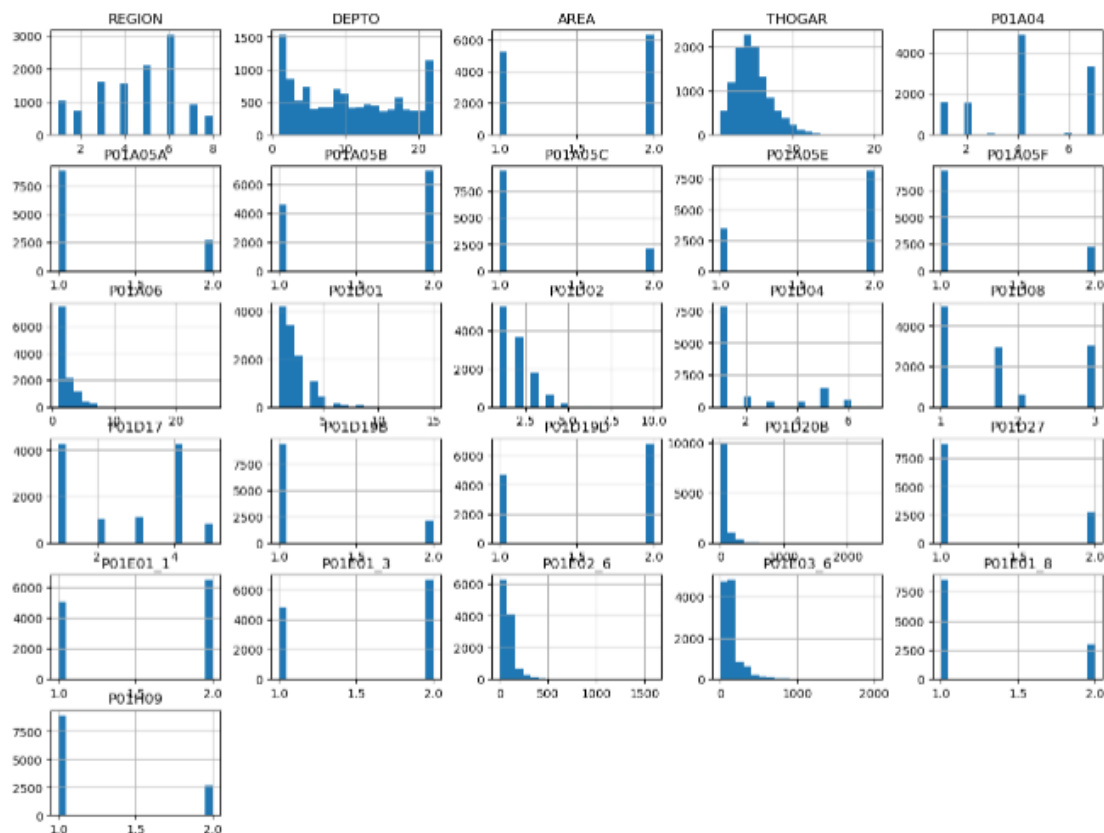
Variables en 'personas' con correlación mayor o igual a 0.1 con 'POBREZA':

```
['REGION', 'DEPTO', 'AREA', 'UPM', 'NUMHOG', 'POBREZA', 'THOGAR', 'ID', 'PPA03', 'PPA04C', 'PPA06', 'PPA08', 'P04A01A', 'P04A04A', 'P04A07A', 'P04A08A', 'P04A11A', 'P06B01', 'P06B03', 'P06B02', 'P06B05', 'P06B06', 'P06B07A', 'P06B08A', 'P06B09A', 'P06B10A', 'P06B11A', 'P06B12A', 'P06B13A', 'P06B14A', 'P06B15A', 'P06B16A', 'P06B17A', 'P06B18A', 'P06B19A', 'P06B20A', 'P06B21A', 'P06B22A', 'P06B23A', 'P06B24A', 'P06B25A', 'P06B26A', 'P06B27A', 'P06B28A', 'P06B29A', 'P06B30A', 'P06B31A', 'P06B32A', 'P06B33A', 'P06B34A', 'P06B35A', 'P06B36A', 'P06B37A', 'P06B38A', 'P06B39A', 'P06B40A', 'P06B41A', 'P06B42A', 'P06B43A', 'P06B44A', 'P06B45A', 'P06B46A', 'P06B47A', 'P06B48A', 'P06B49A', 'P06B50A', 'P06B51A', 'P06B52A', 'P06B53A', 'P06B54A', 'P06B55A', 'P06B56A', 'P06B57A', 'P06B58A', 'P06B59A', 'P06B60A', 'P06B61A', 'P06B62A', 'P06B63A', 'P06B64A', 'P06B65A', 'P06B66A', 'P06B67A', 'P06B68A', 'P06B69A', 'P06B70A', 'P06B71A', 'P06B72A', 'P06B73A', 'P06B74A', 'P06B75A', 'P06B76A', 'P06B77A', 'P06B78A', 'P06B79A', 'P06B80A', 'P06B81A', 'P06B82A', 'P06B83A', 'P06B84A', 'P06B85A', 'P06B86A', 'P06B87A', 'P06B88A', 'P06B89A', 'P06B90A', 'P06B91A', 'P06B92A', 'P06B93A', 'P06B94A', 'P06B95A', 'P06B96A', 'P06B97A', 'P06B98A', 'P06B99A', 'P07A01A', 'P07A02A', 'P07A03A', 'P07A04A', 'P07A05A', 'P07A06A', 'P07A07A', 'P07A08A', 'P07A09A', 'P07A10A', 'P07A11A', 'P07A12A', 'P07A13A', 'P07A14A', 'P07A15A', 'P07A16A', 'P07A17A', 'P07A18A', 'P07A19A', 'P07A20A', 'P07A21A', 'P07A22A', 'P07A23A', 'P07A24A', 'P07A25A', 'P07A26A', 'P07A27A', 'P07A28A', 'P07A29A', 'P07A30A', 'P07A31A', 'P07A32A', 'P07A33A', 'P07A34A', 'P07A35A', 'P07A36A', 'P07A37A', 'P07A38A', 'P07A39A', 'P07A40A', 'P07A41A', 'P07A42A', 'P07A43A', 'P07A44A', 'P07A45A', 'P07A46A', 'P07A47A', 'P07A48A', 'P07A49A', 'P07A50A', 'P07A51A', 'P07A52A', 'P07A53A', 'P07A54A', 'P07A55A', 'P07A56A', 'P07A57A', 'P07A58A', 'P07A59A', 'P07A60A', 'P07A61A', 'P07A62A', 'P07A63A', 'P07A64A', 'P07A65A', 'P07A66A', 'P07A67A', 'P07A68A', 'P07A69A', 'P07A70A', 'P07A71A', 'P07A72A', 'P07A73A', 'P07A74A', 'P07A75A', 'P07A76A', 'P07A77A', 'P07A78A', 'P07A79A', 'P07A80A', 'P07A81A', 'P07A82A', 'P07A83A', 'P07A84A', 'P07A85A', 'P07A86A', 'P07A87A', 'P07A88A', 'P07A89A', 'P07A90A', 'P07A91A', 'P07A92A', 'P07A93A', 'P07A94A', 'P07A95A', 'P07A96A', 'P07A97A', 'P07A98A', 'P07A99A', 'P08A01A', 'P08A02A', 'P08A03A', 'P08A04A', 'P08A05A', 'P08A06A', 'P08A07A', 'P08A08A', 'P08A09A', 'P08A10A', 'P08A11A', 'P08A12A', 'P08A13A', 'P08A14A', 'P08A15A', 'P08A16A', 'P08A17A', 'P08A18A', 'P08A19A', 'P08A20A', 'P08A21A', 'P08A22A', 'P08A23A', 'P08A24A', 'P08A25A', 'P08A26A', 'P08A27A', 'P08A28A', 'P08A29A', 'P08A30A', 'P08A31A', 'P08A32A', 'P08A33A', 'P08A34A', 'P08A35A', 'P08A36A', 'P08A37A', 'P08A38A', 'P08A39A', 'P08A40A', 'P08A41A', 'P08A42A', 'P08A43A', 'P08A44A', 'P08A45A', 'P08A46A', 'P08A47A', 'P08A48A', 'P08A49A', 'P08A50A', 'P08A51A', 'P08A52A', 'P08A53A', 'P08A54A', 'P08A55A', 'P08A56A', 'P08A57A', 'P08A58A', 'P08A59A', 'P08A60A', 'P08A61A', 'P08A62A', 'P08A63A', 'P08A64A', 'P08A65A', 'P08A66A', 'P08A67A', 'P08A68A', 'P08A69A', 'P08A70A', 'P08A71A', 'P08A72A', 'P08A73A', 'P08A74A', 'P08A75A', 'P08A76A', 'P08A77A', 'P08A78A', 'P08A79A', 'P08A80A', 'P08A81A', 'P08A82A', 'P08A83A', 'P08A84A', 'P08A85A', 'P08A86A', 'P08A87A', 'P08A88A', 'P08A89A', 'P08A90A', 'P08A91A', 'P08A92A', 'P08A93A', 'P08A94A', 'P08A95A', 'P08A96A', 'P08A97A', 'P08A98A', 'P08A99A', 'P09A01A', 'P09A02A', 'P09A03A', 'P09A04A', 'P09A05A', 'P09A06A', 'P09A07A', 'P09A08A', 'P09A09A', 'P09A10A', 'P09A11A', 'P09A12A', 'P09A13A', 'P09A14A', 'P09A15A', 'P09A16A', 'P09A17A', 'P09A18A', 'P09A19A', 'P09A20A', 'P09A21A', 'P09A22A', 'P09A23A', 'P09A24A', 'P09A25A', 'P09A26A', 'P09A27A', 'P09A28A', 'P09A29A', 'P09A30A', 'P09A31A', 'P09A32A', 'P09A33A', 'P09A34A', 'P09A35A', 'P09A36A', 'P09A37A', 'P09A38A', 'P09A39A', 'P09A40A', 'P09A41A', 'P09A42A', 'P09A43A', 'P09A44A', 'P09A45A', 'P09A46A', 'P09A47A', 'P09A48A', 'P09A49A', 'P09A50A', 'P09A51A', 'P09A52A', 'P09A53A', 'P09A54A', 'P09A55A', 'P09A56A', 'P09A57A', 'P09A58A', 'P09A59A', 'P09A60A', 'P09A61A', 'P09A62A', 'P09A63A', 'P09A64A', 'P09A65A', 'P09A66A', 'P09A67A', 'P09A68A', 'P09A69A', 'P09A70A', 'P09A71A', 'P09A72A', 'P09A73A', 'P09A74A', 'P09A75A', 'P09A76A', 'P09A77A', 'P09A78A', 'P09A79A', 'P09A80A', 'P09A81A', 'P09A82A', 'P09A83A', 'P09A84A', 'P09A85A', 'P09A86A', 'P09A87A', 'P09A88A', 'P09A89A', 'P09A90A', 'P09A91A', 'P09A92A', 'P09A93A', 'P09A94A', 'P09A95A', 'P09A96A', 'P09A97A', 'P09A98A', 'P09A99A', 'P10A01A', 'P10A02A', 'P10A03A', 'P10A04A', 'P10A05A', 'P10A06A', 'P10A07A', 'P10A08A', 'P10A09A', 'P10A10A', 'P10A11A', 'P10A12A', 'P10A13A', 'P10A14A', 'P10A15A', 'P10A16A', 'P10A17A', 'P10A18A', 'P10A19A', 'P10A20A', 'P10A21A', 'P10A22A', 'P10A23A', 'P10A24A', 'P10A25A', 'P10A26A', 'P10A27A', 'P10A28A', 'P10A29A', 'P10A30A', 'P10A31A', 'P10A32A', 'P10A33A', 'P10A34A', 'P10A35A', 'P10A36A', 'P10A37A', 'P10A38A', 'P10A39A', 'P10A40A', 'P10A41A', 'P10A42A', 'P10A43A', 'P10A44A', 'P10A45A', 'P10A46A', 'P10A47A', 'P10A48A', 'P10A49A', 'P10A50A', 'P10A51A', 'P10A52A', 'P10A53A', 'P10A54A', 'P10A55A', 'P10A56A', 'P10A57A', 'P10A58A', 'P10A59A', 'P10A60A', 'P10A61A', 'P10A62A', 'P10A63A', 'P10A64A', 'P10A65A', 'P10A66A', 'P10A67A', 'P10A68A', 'P10A69A', 'P10A70A', 'P10A71A', 'P10A72A', 'P10A73A', 'P10A74A', 'P10A75A', 'P10A76A', 'P10A77A', 'P10A78A', 'P10A79A', 'P10A80A', 'P10A81A', 'P10A82A', 'P10A83A', 'P10A84A', 'P10A85A', 'P10A86A', 'P10A87A', 'P10A88A', 'P10A89A', 'P10A90A', 'P10A91A', 'P10A92A', 'P10A93A', 'P10A94A', 'P10A95A', 'P10A96A', 'P10A97A', 'P10A98A', 'P10A99A', 'P11A01A', 'P11A02A', 'P11A03A', 'P11A04A', 'P11A05A', 'P11A06A', 'P11A07A', 'P11A08A', 'P11A09A', 'P11A10A', 'P11A11A', 'P11A12A', 'P11A13A', 'P11A14A', 'P11A15A', 'P11A16A', 'P11A17A', 'P11A18A', 'P11A19A', 'P11A20A', 'P11A21A', 'P11A22A', 'P11A23A', 'P11A24A', 'P11A25A', 'P11A26A', 'P11A27A', 'P11A28A', 'P11A29A', 'P11A30A', 'P11A31A', 'P11A32A', 'P11A33A', 'P11A34A', 'P11A35A', 'P11A36A', 'P11A37A', 'P11A38A', 'P11A39A', 'P11A40A', 'P11A41A', 'P11A42A', 'P11A43A', 'P11A44A', 'P11A45A', 'P11A46A', 'P11A47A', 'P11A48A', 'P11A49A', 'P11A50A', 'P11A51A', 'P11A52A', 'P11A53A', 'P11A54A', 'P11A55A', 'P11A56A', 'P11A57A', 'P11A58A', 'P11A59A', 'P11A60A', 'P11A61A', 'P11A62A', 'P11A63A', 'P11A64A', 'P11A65A', 'P11A66A', 'P11A67A', 'P11A68A', 'P11A69A', 'P11A70A', 'P11A71A', 'P11A72A', 'P11A73A', 'P11A74A', 'P11A75A', 'P11A76A', 'P11A77A', 'P11A78A', 'P11A79A', 'P11A80A', 'P11A81A', 'P11A82A', 'P11A83A', 'P11A84A', 'P11A85A', 'P11A86A', 'P11A87A', 'P11A88A', 'P11A89A', 'P11A90A', 'P11A91A', 'P11A92A', 'P11A93A', 'P11A94A', 'P11A95A', 'P11A96A', 'P11A97A', 'P11A98A', 'P11A99A', 'P12A01A', 'P12A02A', 'P12A03A', 'P12A04A', 'P12A05A', 'P12A06A', 'P12A07A', 'P12A08A', 'P12A09A', 'P12A10A', 'P12A11A', 'P12A12A', 'P12A13A', 'P12A14A', 'P12A15A', 'P12A16A', 'P12A17A', 'P12A18A', 'P12A19A', 'P12A20A', 'P12A21A', 'P12A22A', 'P12A23A', 'P12A24A', 'P12A25A', 'P12A26A', 'P12A27A', 'P12A28A', 'P12A29A', 'P12A30A', 'P12A31A', 'P12A32A', 'P12A33A', 'P12A34A', 'P12A35A', 'P12A36A', 'P12A37A', 'P12A38A', 'P12A39A', 'P12A40A', 'P12A41A', 'P12A42A', 'P12A43A', 'P12A44A', 'P12A45A', 'P12A46A', 'P12A47A', 'P12A48A', 'P12A49A', 'P12A50A', 'P12A51A', 'P12A52A', 'P12A53A', 'P12A54A', 'P12A55A', 'P12A56A', 'P12A57A', 'P12A58A', 'P12A59A', 'P12A60A', 'P12A61A', 'P12A62A', 'P12A63A', 'P12A64A', 'P12A65A', 'P12A66A', 'P12A67A', 'P12A68A', 'P12A69A', 'P12A70A', 'P12A71A', 'P12A72A', 'P12A73A', 'P12A74A', 'P12A75A', 'P12A76A', 'P12A77A', 'P12A78A', 'P12A79A', 'P12A80A', 'P12A81A', 'P12A82A', 'P12A83A', 'P12A84A', 'P12A85A', 'P12A86A', 'P12A87A', 'P12A88A', 'P12A89A', 'P12A90A', 'P12A91A', 'P12A92A', 'P12A93A', 'P12A94A', 'P12A95A', 'P12A96A', 'P12A97A', 'P12A98A', 'P12A99A', 'P13A01A', 'P13A02A', 'P13A03A', 'P13A04A', 'P13A05A', 'P13A06A', 'P13A07A', 'P13A08A', 'P13A09A', 'P13A10A', 'P13A11A', 'P13A12A', 'P13A13A', 'P13A14A', 'P13A15A', 'P13A16A', 'P13A17A', 'P13A18A', 'P13A19A', 'P13A20A', 'P13A21A', 'P13A22A', 'P13A23A', 'P13A24A', 'P13A25A', 'P13A26A', 'P13A27A', 'P13A28A', 'P13A29A', 'P13A30A', 'P13A31A', 'P13A32A', 'P13A33A', 'P13A34A', 'P13A35A', 'P13A36A', 'P13A37A', 'P13A38A', 'P13A39A', 'P13A40A', 'P13A41A', 'P13A42A', 'P13A43A', 'P13A44A', 'P13A45A', 'P13A46A', 'P13A47A', 'P13A48A', 'P13A49A', 'P13A50A', 'P13A51A', 'P13A52A', 'P13A53A', 'P13A54A', 'P13A55A', 'P13A56A', 'P13A57A', 'P13A58A', 'P13A59A', 'P13A60A', 'P13A61A', 'P13A62A', 'P13A63A', 'P13A64A', 'P13A65A', 'P13A66A', 'P13A67A', 'P13A68A', 'P13A69A', 'P13A70A', 'P13A71A', 'P13A72A', 'P13A73A', 'P13A74A', 'P13A75A', 'P13A76A', 'P13A77A', 'P13A78A', 'P13A79A', 'P13A80A', 'P13A81A', 'P13A82A', 'P13A83A', 'P13A84A', 'P13A85A', 'P13A86A', 'P13A87A', 'P13A88A', 'P13A89A', 'P13A90A', 'P13A91A', 'P13A92A', 'P13A93A', 'P13A94A', 'P13A95A', 'P13A96A', 'P13A97A', 'P13A98A', 'P13A99A', 'P14A01A', 'P14A02A', 'P14A03A', 'P14A04A', 'P14A05A', 'P14A06A', 'P14A07A', 'P14A08A', 'P14A09A', 'P14A10A', 'P14A11A', 'P14A12A', 'P14A13A', 'P14A14A', 'P14A15A', 'P14A16A', 'P14A17A', 'P14A18A', 'P14A19A', 'P14A20A', 'P14A21A', 'P14A22A', 'P14A23A', 'P14A24A', 'P14A25A', 'P14A26A', 'P14A27A', 'P14A28A', 'P14A29A', 'P14A30A', 'P14A31A', 'P14A32A', 'P14A33A', 'P14A34A', 'P14A35A', 'P14A36A', 'P14A37A', 'P14A38A', 'P14A39A', 'P14A40A', 'P14A41A', 'P14A42A', 'P14A43A', 'P14A44A', 'P14A45A', 'P14A46A', 'P14A47A', 'P14A48A', 'P14A49A', 'P14A50A', 'P14A51A', 'P14A52A', 'P14A53A', 'P14A54A', 'P14A55A', 'P14A56A', 'P14A57A', 'P14A58A', 'P14A59A', 'P14A60A', 'P14A61A', 'P14A62A', 'P14A63A', 'P14A64A', 'P14A65A', 'P14A66A', 'P14A67A', 'P14A68A', 'P14A69A', 'P14A70A', 'P14A71A', 'P14A72A', 'P14A73A', 'P14A74A', 'P14A75A', 'P14A76A', 'P14A77A', 'P14A78A', 'P14A79A', 'P14A80A', 'P14A81A', 'P14A82A', 'P14A83A', 'P14A84A', 'P14A85A', 'P14A86A', 'P14A87A', 'P14A88A', 'P14A89A', 'P14A90A', 'P14A91A', 'P14A92A', 'P14A93A', 'P14A94A', 'P14A95A', 'P14A96A', 'P14A97A', 'P14A98A', 'P14A99A', 'P15A01A', 'P15A02A', 'P15A03A', 'P15A04A', '
```

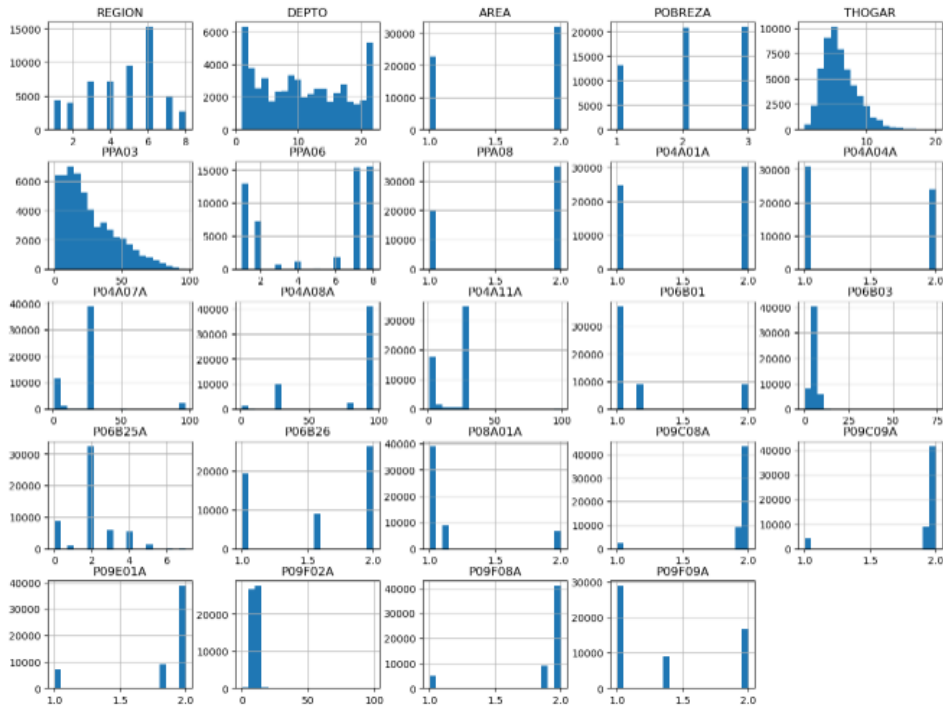
Análisis Univariable

Se buscó la distribución de las variables que cumplían con ese parámetro de correlación para ir viendo el comportamiento de las variables y así poder ir descartando cuales no iban a servir en los modelos.

Hogares:



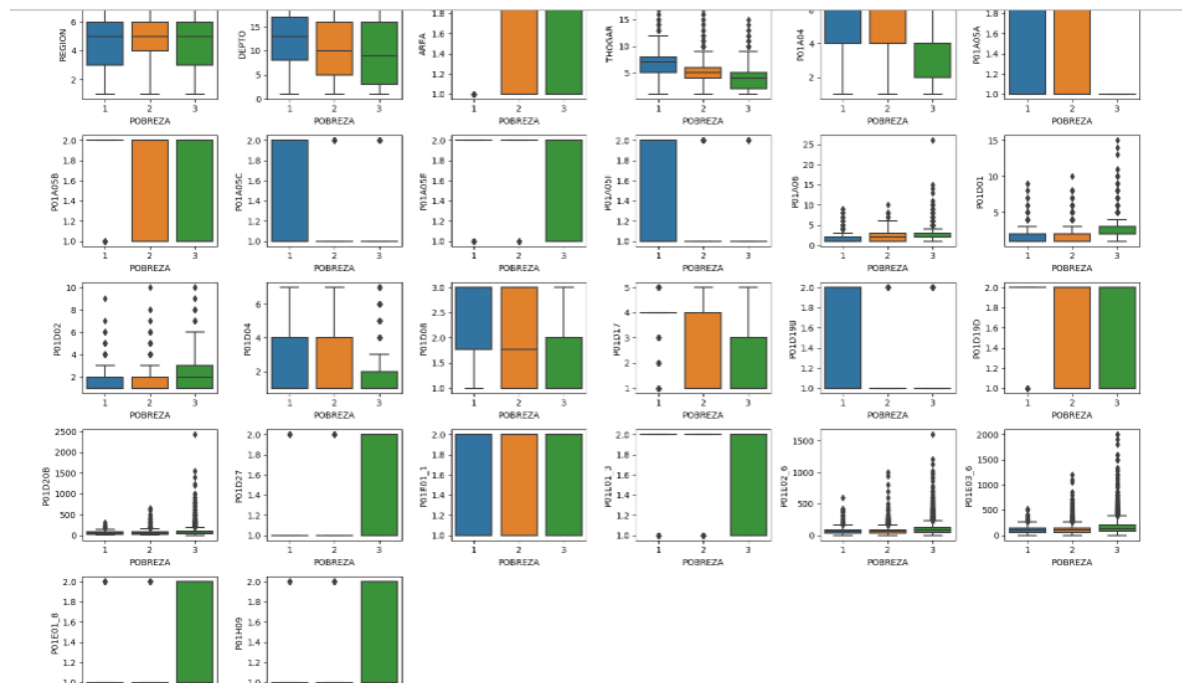
Personas:



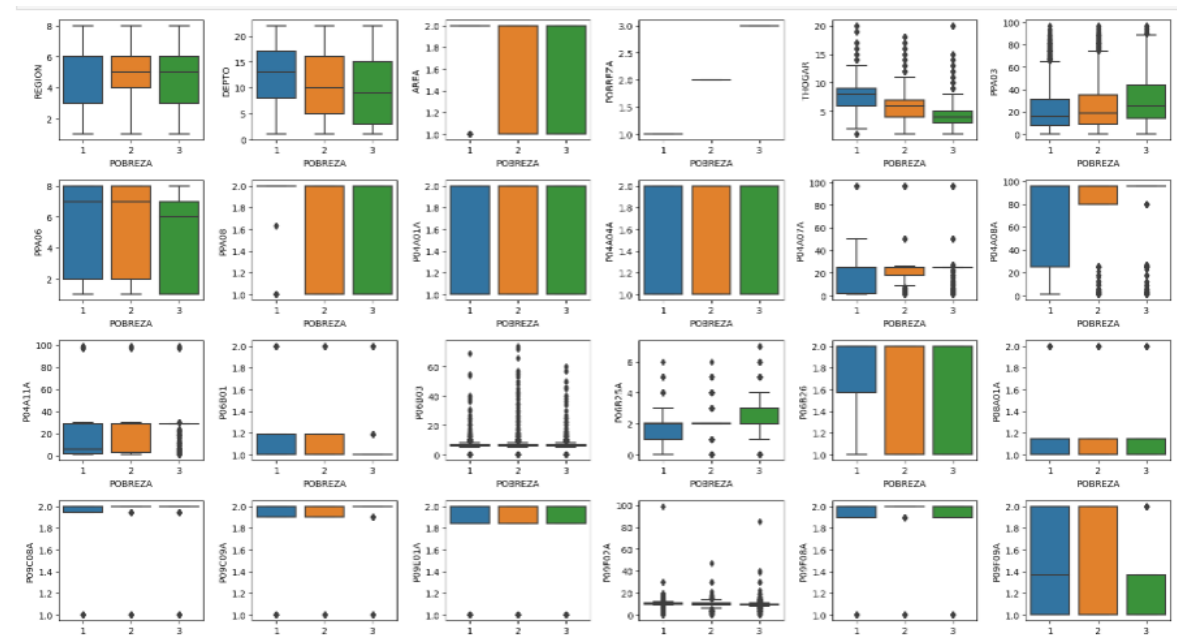
Análisis Bivariable

Posteriormente se analizó mediante boxplots el comportamiento de las variables respecto a la variable a predecir (POBREZA) y así seleccionar finalmente cuales serían las que se usarían en el modelo

Hogares:



Personas:



Preparación de Datos para Modelado:

Así se decidió que las variables a utilizar en los modelos serían:

Después de hacer el análisis estadístico de las variables con una correlación considerable, elegimos las variables a usar en los modelos.

En el modelo de hogar, serían:

- Region
- Depto
- Area
- THOGAR
- P01A04
- P01A05A
- P01A05B
- P01A06
- P01D01
- P01D02
- P01D04
- P01D08
- P01D17
- P01D19D
- P01D20B

- P01E01_1
- P01E02_6
- P01E03_6

y en el modelo de personas serían:

- REGION
- DEPTO
- AREA
- THOGAR
- PPA03
- PPA06
- PPA08
- P04A01A
- P04A04A
- P04A07A
- P04A08A
- P04A11A
- P06B01
- P06B03
- P06B25A
- P06B26
- P08A01A
- P09C08A
- P09C09A
- P09E01A
- P09F02A
- P09F08A
- P09F09A

Además se hizo la división de los datos en train y test

Modelos de hogares:


```

import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.neural_network import MLPRegressor
from sklearn.metrics import mean_squared_error
from math import sqrt

# Seleccionar Las variables predictoras
selected_features = ['REGION', 'DEPTO', 'THOGAR', 'P01A04', 'P01A05A', 'P01A05B', 'P01A06',
                    'P01D01', 'P01D02', 'P01D04', 'P01D08', 'P01D17', 'P01D19D', 'P01D20B',
                    'P01E02_6', 'P01E03_6']

# Seleccionar La variable objetivo
target_variable = 'POBREZA'

# Filtrar el DataFrame con Las columnas seleccionadas
df_hogar_selected = hogar_cleaned_final_mean_imputed[selected_features + [target_variable]].dropna()

# Dividir el conjunto de datos en entrenamiento y prueba
X_train, X_test, y_train, y_test = train_test_split(df_hogar_selected[selected_features],
                                                    df_hogar_selected[target_variable],
                                                    test_size=0.2, random_state=42)

```

Modelos de personas:

```

69]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import mean_squared_error
from math import sqrt
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout
from tensorflow.keras.optimizers import Adam

70]: # Seleccionar Las variables predictoras
selected_features_personas = ['REGION', 'DEPTO', 'AREA', 'THOGAR', 'PPA03', 'PPA06', 'PPA08', 'P04A01A', 'P04A04A',
                             'P04A07A', 'P04A08A', 'P04A11A', 'P06B01', 'P06B03', 'P06B25A', 'P06B26', 'P08A01A',
                             'P09C08A', 'P09C09A', 'P09E01A', 'P09F02A', 'P09F08A', 'P09F09A']

# Seleccionar La variable objetivo
target_variable_personas = 'POBREZA'

# Filtrar el DataFrame con Las columnas seleccionadas y eliminar filas con valores nulos
df_personas_selected = personas_cleaned_final_mean_imputed[selected_features_personas + [target_variable_personas]].dropna()

# Dividir el conjunto de datos en entrenamiento y prueba
X_train_personas, X_test_personas, y_train_personas, y_test_personas = train_test_split(
    df_personas_selected[selected_features_personas], df_personas_selected[target_variable_personas],
    test_size=0.2, random_state=42
)

# Normalizar Los datos
scaler_personas = StandardScaler()
X_train_normalized_personas = scaler_personas.fit_transform(X_train_personas)
X_test_normalized_personas = scaler_personas.transform(X_test_personas)

```

Construcción del Modelo

Se realizaron para ambas bases de datos modelos de regresión, random forest y una red neuronal

Modelos hogares:

```
# Modelo de Regresión Lineal
regression_model = LinearRegression()
regression_model.fit(X_train, y_train)
regression_predictions = regression_model.predict(X_test)
regression_rmse = sqrt(mean_squared_error(y_test, regression_predictions))
print(f"Regresión Lineal RMSE: {regression_rmse}")
```

Regresión Lineal RMSE: 0.514009442541057

```
# Modelo Random Forest
random_forest_model = RandomForestRegressor()
random_forest_model.fit(X_train, y_train)
rf_predictions = random_forest_model.predict(X_test)
rf_rmse = sqrt(mean_squared_error(y_test, rf_predictions))
print(f"Random Forest RMSE: {rf_rmse}")
```

```
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
from math import sqrt
from tensorflow.keras.optimizers import Adam

# Dividir el conjunto de datos en entrenamiento y prueba
X_train, X_test, y_train, y_test = train_test_split(hogar_cleaned_final_mean_imputed[['REGION', 'DEPTO', 'AREA', 'THOGAR',
'P01A04', 'P01A05A', 'P01A05B', 'P01A05C',
'P01A05E', 'P01A05F', 'P01A06', 'P01D01', 'P01D02', 'P01D04', 'P01D08',
'P01D17', 'P01D19D', 'P01D20B', 'P01E01_1', 'P01E02_6', 'P01E03_6']], hogar_cleaned_final_mean_imputed['POBREZA'], test_size=0.2, random_state=42)

# Normalizar los datos
X_train_normalized = (X_train - X_train.mean()) / X_train.std()
X_test_normalized = (X_test - X_train.mean()) / X_train.std()

# Construir el modelo de red neuronal
model = Sequential()
model.add(Dense(64, input_dim=X_train_normalized.shape[1], activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(32, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(1, activation='linear'))

model.compile(optimizer=Adam(learning_rate=0.001), loss='mean_squared_error')

# Entrenar el modelo
model.fit(X_train_normalized, y_train, epochs=50, batch_size=32, validation_split=0.2)

# Evaluar el modelo en el conjunto de prueba
nn_predictions = model.predict(X_test_normalized)
nn_rmse = sqrt(mean_squared_error(y_test, nn_predictions))
print(f"Red Neuronal RMSE en conjunto de prueba: {nn_rmse}")
```

Modelos personas

```
# Modelo de Regresión Lineal
regression_model_personas = LinearRegression()
regression_model_personas.fit(X_train_normalized_personas, y_train_personas)
regression_predictions_personas = regression_model_personas.predict(X_test_normalized_personas)
regression_rmse_personas = sqrt(mean_squared_error(y_test_personas, regression_predictions_personas))
print(f"Regresión Lineal RMSE para personas: {regression_rmse_personas}")
```

Regresión Lineal RMSE para personas: 0.5964444889804388

```
# Modelo Random Forest
random_forest_model_personas = RandomForestRegressor()
random_forest_model_personas.fit(X_train_normalized_personas, y_train_personas)
rf_predictions_personas = random_forest_model_personas.predict(X_test_normalized_personas)
rf_rmse_personas = sqrt(mean_squared_error(y_test_personas, rf_predictions_personas))
print(f"Random Forest RMSE para personas: {rf_rmse_personas}")
```

Random Forest RMSE para personas: 0.5448930320541682

```

# Modelo de Red Neuronal
model_personas = Sequential()
model_personas.add(Dense(64, input_dim=X_train_normalized_personas.shape[1], activation='relu'))
model_personas.add(Dropout(0.5))
model_personas.add(Dense(32, activation='relu'))
model_personas.add(Dropout(0.5))
model_personas.add(Dense(1, activation='linear'))

model_personas.compile(optimizer=Adam(learning_rate=0.001), loss='mean_squared_error')

# Entrenar el modelo
model_personas.fit(X_train_normalized_personas, y_train_personas, epochs=50, batch_size=32, validation_split=0.2)

# Evaluar el modelo en el conjunto de prueba
nn_predictions_personas = model_personas.predict(X_test_normalized_personas)
nn_rmse_personas = sqrt(mean_squared_error(y_test_personas, nn_predictions_personas))
print(f"Red Neuronal RMSE para personas: {nn_rmse_personas}")

```

Evaluación del Modelo:

Se hizo la comparación entre los 3 modelos de cada una de las bases de datos, llegando a concluirse que en ambos casos el mejor modelo para predecir POBREZA es mediante random forests ya que obtuvieron los valores de rmse mas bajos.

Hogares:

```

# Mostrar resultados
print(f"Regresión Lineal RMSE: {regression_rmse}")
print(f"Random Forest RMSE: {rf_rmse}")
print(f"Red Neuronal RMSE: {nn_rmse}")

Regresión Lineal RMSE: 0.514009442541057
Random Forest RMSE: 0.4950994511744218
Red Neuronal RMSE: 0.5050387118355069

```

Personas:

```

# Imprimir resultados
print(f"Regresión Lineal RMSE para personas: {regression_rmse_personas}")
print(f"Random Forest RMSE para personas: {rf_rmse_personas}")
print(f"Red Neuronal RMSE para personas: {nn_rmse_personas}")

Regresión Lineal RMSE para personas: 0.5964444889804388
Random Forest RMSE para personas: 0.5448930320541682
Red Neuronal RMSE para personas: 0.5775876367919405

```

Debe recalcar que se probó con diferentes parámetros e hiperparámetros hasta llegar a estos que fueron los mejores valores de rmse, también se intentaron otros tipos de redes y regresiones, sin embargo, no se obtuvieron resultados tan buenos como los seleccionados. A continuación, se pueden ver resultados de algunos modelos que no llegaron a considerarse como los más eficientes para resolver este problema:

```
# Crear y entrenar el modelo de árbol de decisión
model_tree = DecisionTreeClassifier(random_state=42)
model_tree.fit(X_train, y_train)

# Hacer predicciones en el conjunto de prueba
predictions_tree = model_tree.predict(X_test)

# Evaluar el modelo
accuracy_tree = accuracy_score(y_test, predictions_tree)
print(f'Test accuracy (Decision Tree): {accuracy_tree:.4f}\n')

# Mostrar el informe de clasificación
print('Classification Report (Decision Tree):\n')
print(classification_report(y_test, predictions_tree))
```

Test accuracy (Decision Tree): 0.6438

Classification Report (Decision Tree):

	precision	recall	f1-score	support
0	0.51	0.55	0.53	364
1	0.56	0.53	0.54	875
2	0.76	0.77	0.76	1069
accuracy			0.64	2308
macro avg	0.61	0.62	0.61	2308
weighted avg	0.64	0.64	0.64	2308

Epoch 20/20

462/462 [=====] - 2s 3ms/step - loss: 0.6835 - accuracy: 0.6855 - val_loss: 0.7095 - val_accuracy: 0.6777
 73/73 [=====] - 0s 2ms/step - loss: 0.6836 - accuracy: 0.6902

Test accuracy: 0.6902079582214355

Reflexión:

Este laboratorio fue bastante retador debido a que se tuvo mucha libertad para seleccionar modelos, variables a predecir y análisis a realizar. Permite poner en práctica un rango más amplio de conocimientos adquiridos a lo largo del curso.

Otra de las principales dificultades para realizar este laboratorio fue la gran cantidad de procesamiento que los datos requerían ya que mostraban una cantidad grandísima de datos faltantes.

Este es el laboratorio donde se ha tenido la mayor libertad en cuanto a las decisiones a tomar y, por lo mismo, da la oportunidad de experimentar con muchas cosas y encontrar así los modelos más eficientes posibles. Esto permite también conocer un poco más como funcionan, posiblemente escenarios reales en un ambiente laboral ya que es muy común que se soliciten trabajos o resultados sin una gran explicación o guía de como hacerlo y como científicos de datos es nuestro trabajo encontrar la mejor forma de resolver dichos problemas.