

UNIVERSIDAD COMPLUTENSE DE MADRID
FACULTAD DE ESTUDIOS ESTADÍSTICOS



TAREA EVALUABLE

**MÓDULO: MINERÍA DE DATOS Y MODELIZACIÓN PREDICTIVA -
PREPARACIÓN DE DATOS Y REGRESIONES**

Jorge González Perea 51553561G

Máster en Big Data, Data Science & Inteligencia Artificial

Curso académico 2024-2025

Índice

1. Introducción.	2
2. Depuración de datos.	3
2.1. Corrección de variables categóricas.	5
2.2. Corrección de datos numéricos.	6
2.3. Valores atípicos y <i>missing</i> . Imputaciones.	7
3. Gráficos y presentación de variables.	9
3.1. Coeficientes V de Cramer.	9
3.2. Variable binaria.	9
3.3. Variable continua.	12
4. Modelos de regresión lineal.	14
4.1. Modelo de regresión lineal. Selección de variables clásica.	14
4.2. Modelo de regresión lineal. Selección de variables aleatoria.	18
4.3. Modelo de regresión lineal. Conclusiones.	19
5. Modelos de regresión logística.	20
5.1. Modelo de regresión logística. Selección de variables clásica.	20
5.2. Modelo de regresión logística. Selección de variables aleatoria.	22
5.3. Modelo de regresión logística. Conclusiones.	24
6. Anexo.	25

1. Introducción.

En este proyecto se trabaja con una base de datos que incluye información demográfica sobre las últimas elecciones en España. Se tienen 8117 filas, correspondientes a distintos municipios, y 41 columnas, que contienen los datos de las variables. De las 41 variables, 7 de ellas son de interés:

- **AbstentionPtge**: porcentaje de abstención.
- **Izda_Pct**: porcentaje de votos a partidos de izquierda
- **Dcha_Pct**: porcentaje de votos a partidos de derecha
- **Otros_Pct**: porcentaje de votos a partidos distintos de PP, Ciudadanos, PSOE y Podemos.
- **AbstencionAlta**: variable dicotómica que toma el valor 1 si el porcentaje de abstención es superior al 30 % y, 0, en otro caso.
- **Izquierda**: variable dicotómica que toma el valor 1 si la suma de los votos de izquierdas es superior a la de derechas y otros y, 0, en otro caso.
- **Derecha**: variable dicotómica que toma el valor 1 si la suma de los votos de derecha es superior a la de izquierda y otros y, 0, en otro caso.

Para este estudio se han elegido dos variables objetivo de las siete anteriores. Una de ellas (numérica) se someterá a un modelo de regresión lineal para su correspondiente predicción. Con la otra (dicotómica) se hará lo mismo pero con un modelo de regresión logística. Las variables escogidas son **Izda Pct** e **Izquierda**. Las cinco variables restantes son eliminadas de la base de datos debido a que no se utilizarán como variables explicativas:

Además de estas, también se tienen diversas variables en la base de datos. Algunas son de carácter geográfico (nombre del municipio, código de provincia, CCAA, etc.), otras son datos demográficos (población, porcentaje de nativos de cada CCAA y extranjeros, datos sobre empleo en diferentes sectores, etc.) y otras están relacionadas con los resultados de las elecciones.

El objetivo es depurar esta base de datos y, posteriormente, crear un modelo de regresión lineal y otro de regresión logística que permitan predecir variables continuas y resolver problemas de clasificación respectivamente para las dos variables objetivo elegidas. El resultado esperado de este proceso es un análisis sobre los votos a partidos de izquierda a nivel municipal para poder observar tendencias y poder hacer predicciones.

2. Depuración de datos.

A continuación se encuentra un listado con todas las librerías y funciones importadas.

```
1 import os
2 import pandas as pd
3 import numpy as np
4 import pickle
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 from sklearn.model_selection import train_test_split
8 from itertools import combinations
9 from collections import Counter
10 from FuncionesMineria_original import analizar_variables_categoricas,
    atipicosAmissing, patron_perdidos, ImputacionCuant, ImputacionCuali,
    graficoVcramer, Vcramer, mosaico_targetbinaria, boxplot_targetbinaria,
    hist_targetbinaria, hist_target_categorica, lm, Rsq, validacion_cruzada_lm,
    modelEffectSizes, crear_data_modelo, Transf_Auto, Rsq, lm, lm_forward,
    lm_backward, lm_stepwise, validacion_cruzada_lm, crear_data_modelo
11 os.chdir(r'C:\Users\jorge\Desktop\Mster\M dulo 7 - Miner a de datos y
    modelizaci n predictiva\Parte 1\Tarea')
```

Listing 1: Librerías y funciones.

Para leer los datos desde el archivo **DatosEleccionesEspaña.xlsx** con Python es necesario ejecutar las siguientes instrucciones:

```
1 datos = pd.read_excel('DatosEleccionesEspa a.xlsx')
2 datos = datos.drop(columns = ['AbstentionPtge', 'Dcha_Pct', 'Otros_Pct', '
    AbstencionAlta', 'Derecha'])
3 datos.dtypes
```

Listing 2: Lectura de datos.

Gracias a esta instrucción se genera una variable llamada '**datos**' que tiene formato de **DataFrame** de la librería Pandas. En esta tabla se pueden visualizar algunos de los valores que toman las distintas variables:

Índice	Name	digoProvinc	CCAA	Population	totalCensus	stentionPt	stencionAl	Izda_Pct	Dcha_Pct	Otros_Pct
0	Abadía	10	Extremadura	336	282	20.213	0	60.444	35.555	1.778
1	Abertura	10	Extremadura	429	364	25.275	0	54.779	44.118	0.368
2	Acebo	10	Extremadura	569	569	27.241	0	44.203	53.14	0.966
3	Acehúche	10	Extremadura	822	704	30.114	1	50.813	45.325	0
4	Aceituna	10	Extremadura	623	540	30.185	1	44.562	49.067	0.796
5	Ahigal	10	Extremadura	1421	1263	22.565	0	53.068	45.901	0.409
6	Alagón del Río	10	Extremadura	923	826	30.024	1	59.343	37.89	0.692
7	Albalá	10	Extremadura	730	629	27.027	0	37.037	60.567	0.218
8	Alcántara	10	Extremadura	1571	1320	22.803	0	46.81	50.736	0.785
9	Alcollarín	10	Extremadura	256	266	24.436	0	37.811	50.209	0.498
10	Alcuéscar	10	Extremadura	2759	2290	40.524	1	55.286	41.923	0.734
11	Aldea del Cano	10	Extremadura	667	595	21.681	0	61.159	37.339	0
12	Aldeacentenera	10	Extremadura	646	560	18.214	0	57.861	38.21	0.655

Figura 1: Algunas columnas de la base de datos sin depurar.

Con el comando **datos.dtypes** se imprime por pantalla el nombre de todas las columnas de la tabla 1 junto con el formato de datos que contiene cada una (que tiene que ser único según el funcionamiento de un DataFrame). En la impresión por pantalla se puede ver una serie de variables que vienen en formato flotante (**float64**) o entero (**int64**) y que deberían ser objetos o cadenas de caracteres:

totalEmpresas	float64
Industria	float64
Construccion	float64
ComercTTEHosteleria	float64
Servicios	float64
ActividadPpal	object
inmuebles	float64
Pob2010	float64
SUPERFICIE	float64
Densidad	object
PobChange_pct	float64
PersonasInmueble	float64
Explotaciones	int64

Figura 2: formatos de algunas columnas de la base de datos.

En el siguiente listado se encuentra el código para dividir las variables en función de su tipo (numéricas o categóricas) así como la creación de una tabla que contiene distintas medidas estadísticas relativas a las variables y que pueden ser de ayuda a la hora de llevar a cabo una depuración de los datos. Las variables 'CodigoProvincia' e 'Izquierda' se pueden convertir a categóricas debido al reducido conjunto de sus valores.

```

1 numericasAcategoricas = ['CodigoProvincia', 'Izquierda']
2 for var in numericasAcategoricas:
3     datos[var] = datos[var].astype(str)
4 variables = datos.columns.tolist()
5 numericas = datos.select_dtypes(include=['int', 'int32', 'int64', 'float',
6     'float32', 'float64']).columns.tolist()
7 categoricas = [variable for variable in variables if variable not in numericas]
8 descriptivos_num = datos.describe().T
9 for num in numericas:
10     descriptivos_num.loc[num, "Asimetria"] = datos[num].skew()
11     descriptivos_num.loc[num, "Kurtosis"] = datos[num].kurtosis()
12     descriptivos_num.loc[num, "Rango"] = np.ptp(datos[num].dropna().values)

```

Listing 3: Separación de variables y análisis estadístico básico.

Y la tabla mencionada incluye una recopilación de medidas estadísticas de relevancia para esta base de datos que permiten localizar valores atípicos, no numéricos, etc. como se verá en el apartado 2.2:

Índice	count	mean	std	min	25%	50%	75%	max	Asimetria	Kurtosis	Rango
CodigoProvincia	8117	26.6647	14.8934	1	13	26	41	59	0.08961447	-1.32302	49
Population	8117	5722.34	46204.2	5	166	548	2427	3.14199e+06	46.0407	2820.33	3.14199e+06
TotalCensus	8117	4247.86	34423.4	5	140	447	1843	2.36383e+06	46.5446	2893.45	2.36382e+06
Izda_pct	8117	34.404	16.4843	0	21.891	35.165	46.832	94.117	0.0598817	-0.493538	94.117
Izquierda	8117	0.222866	0.416194	0	0	0	0	1	1.33208	-0.225612	1
Age_0-4_ptge	8117	3.61827	2.05263	0	1.389	2.975	4.533	13.245	0.343639	-0.206608	13.245
Age_under19_ptge	8117	13.5641	6.77745	0	8.334	13.881	19.055	33.696	-0.104763	-0.79225	33.696
Age_19_65_pct	8117	57.3706	6.81864	23.459	53.045	58.655	61.818	100.002	-0.814264	2.15594	76.543
Age_over65_pct	8117	29.0653	11.767	18.052	19.827	27.559	36.911	76.472	0.584708	0.102323	94.524
WomanPopulationPtge	8117	47.3023	4.36235	11.765	45.725	48.485	50	72.683	-1.6711	5.80063	60.918
ForeignersPtge	8117	5.61832	7.3487	-8.96	1.06	3.59	8.18	71.47	2.49826	11.3568	80.43
SameComAutonPtge	8117	81.6335	12.2873	0	75.806	84.493	90.462	127.156	-1.52276	3.47954	127.156
SameComAutonDiffProvPtge	8117	4.33764	6.39494	0	0.676	2.19	5.277	67.308	3.28683	14.5601	67.308

Figura 3: Algunas medidas estadísticas de la base de datos sin depurar.

2.1. Corrección de variables categóricas.

El conjunto de variables categóricas está formado por **Name**, **CodigoProvincia**, **CCAA**, **Izquierda**, **ActividadPpal** y **Densidad**. Mediante el comando `analizar_variables_categoricas` se pueden visualizar todos los posibles valores de cada columna.

```

1  frecuencias = analizar_variables_categoricas(datos)
2  frecuencias
3
4  # Cambio de '?' por nan en la columna 'Densidad'.
5  datos['Densidad'] = datos['Densidad'].replace('?', np.nan)
6  analizar_variables_categoricas(datos)
7
8  # Cambio de cadenas de caracteres 'nan' a NaN.
9  for x in categoricas:
10     datos[x] = datos[x].replace('nan', np.nan)
11     datos[variables].isna().sum()
12
13 # Agrupación de categorías poco representadas:
14 datos['ActividadPpal'] = datos['ActividadPpal'].replace({'Otro': 'Otro', '
ComercTTEHosteleria': 'ComercioyHosteleria', 'Servicios': '
Servicios_Construccion_Industria', 'Construccion': '
Servicios_Construccion_Industria', 'Industria': '
Servicios_Construccion_Industria'})

```

Listing 4: Comprobación de variables categóricas.

Este código proporciona la siguiente salida:

Name:	n	%	Name:	n	%	Columna	Distintos
Mieres	2	0.000246	Mieres	2	0.000246	CodigoProvincia	50
Cieza	2	0.000246	Cieza	2	0.000246	Population	3595
Noya	2	0.000246	Noya	2	0.000246	TotalCensus	3308
Rebollar	2	0.000246	Rebollar	2	0.000246	Izda_Pct	6568
Villaseca	2	0.000246	Villaseca	2	0.000246	Age_0-4_Ptge	3759
...	Age_under10_Ptge	5889
Muras	1	0.000123	Muras	1	0.000123	Age_19_65_pct	6213
Monterroso	1	0.000123	Monterroso	1	0.000123	Age_over65_pct	7077
Monforte de Lemos	1	0.000123	Monforte de Lemos	1	0.000123	WomanPopulationPtge	4523
Mondofredo	1	0.000123	Mondofredo	1	0.000123	ForeignersPtge	2329
Zulada	1	0.000123	Zulada	1	0.000123	SameComAutonPtge	6150
[8100 rows x 2 columns]			[8100 rows x 2 columns]			SameComAutonDiffProvPtge	4207
CCAA:	n	%	CCAA:	n	%	DiffComAutonPtge	5573
CastillaLeón	2248	0.276950	CastillaLeón	2248	0.276950	UnemployLess25_Ptge	2340
Cataluña	947	0.116669	Cataluña	947	0.116669	Unemploy25_40_Ptge	2679
CastillaLaMancha	919	0.113219	CastillaLaMancha	919	0.113219	UnemployMore40_Ptge	3172
Andalucía	773	0.095232	Andalucía	773	0.095232	AgricultureUnemploymentPtge	2523
Aragón	731	0.090058	Aragón	731	0.090058	IndustryUnemploymentPtge	2536
ComValenciana	542	0.066773	ComValenciana	542	0.066773	ConstructionUnemploymentPtge	2593
Extremadura	387	0.047678	Extremadura	387	0.047678	ServiciosUnemploymentPtge	2993
Galicia	314	0.038684	Galicia	314	0.038684	totalEmpresas	1225
Navarra	272	0.033510	Navarra	272	0.033510	Industria	3088
PaisVasco	251	0.030923	PaisVasco	251	0.030923	Construccion	455
Madrid	179	0.022052	Madrid	179	0.022052	ComercTTEHosteleria	801
Rioja	174	0.021436	Rioja	174	0.021436	Servicios	756
Cantabria	162	0.019841	Cantabria	162	0.019841	inmuebles	3086
Canarias	88	0.010841	Canarias	88	0.010841	Pob2010	3623
Asturias	78	0.009569	Asturias	78	0.009569	SUPERFICIE	8109
Baleares	67	0.008254	Baleares	67	0.008254	PobChange_pct	3048
Murcia	45	0.005544	Murcia	45	0.005544	PersonasInmueble	282
[8100 rows x 2 columns]			[8100 rows x 2 columns]			Explotaciones	758
Izquierda:	n	%	Izquierda:	n	%		
0	6380	0.777134	0	6380	0.777134		
1	1809	0.222866	1	1809	0.222866		
ActividadPpal:	n	%	ActividadPpal:	n	%		
Otro	4932	0.607614	Otro	4932	0.607614		
ComercTTEHosteleria	2538	0.312677	ComercTTEHosteleria	2538	0.312677		
Servicios	620	0.076383	Servicios	620	0.076383		
Construccion	14	0.001725	Construccion	14	0.001725		
Industria	13	0.001602	Industria	13	0.001602		
Densidad:	n	%	Densidad:	n	%		
NuyBaja	6416	0.790440	NuyBaja	6416	0.790440		
Baja	1853	0.229728	Baja	1853	0.229728		
Alta	556	0.069283	Alta	556	0.069283		
?	92	0.011334	?	92	0.011334		

Figura 4: Frecuencias de los valores de las variables categóricas.

En el listado 4 se llevan a cabo los siguientes procesos:

- Se comprueba la frecuencia de cada valor en las columnas categóricas.
- Se sustituyen los caracteres '?' de la columna **Densidad** por NaN.
- Se sustituyen posibles valores NaN que estén escritos como la cadena de caracteres 'nan' por el valor NaN real.
- Se agrupan las categorías poco representadas de la columna **ActividadPpal** para evitar errores en la codificación de los datos más adelante.

Los municipios repetidos podrían estar en la misma CCAA, lo cual sería un error evidente. Con ayuda de la función `analizar_variables_categoricas` de la librería `FuncionesMineria.py` se obtienen las frecuencias de cada municipio. Se observa que 17 de ellos se repiten y se comprueba si es un error o no. Ninguno de los municipios está repetido, ya que los 17 que se repiten dos veces (no hay ninguno que se repita 3 veces o más) están en CCAA diferentes.

Nota: en la descripción de esta variable se indica que simboliza la densidad de habitantes por hectárea (d), y que su valor puede comprenderse entre 3 posibles: 'MuyBaja' si $d < 1$ hab/ha, 'Baja' si $1 < d < 5$, y 'Alta' si $d > 5$ hab/ha. Por lo tanto, se podría calcular la densidad como el cociente de los datos de la columna **Population** entre los de la columna **SUPERFICIE** y asignar una de estas cadenas de caracteres en función del valor numérico para cada municipio de forma muy sencilla. Sin embargo, se desconocen las unidades de la variable **SUPERFICIE** y también se desconoce qué censo poblacional se ha utilizado para calcular los valores que sí están disponibles en la columna **Densidad**, luego no es conveniente hacerlo y es preferible tratar esos valores como *missing*.

2.2. Corrección de datos numéricos.

El siguiente paso consiste en corregir los valores atípicos, porcentajes erróneos, etc. Para ello es necesario una lectura manual de los datos con ayuda de, por ejemplo, la tabla de la figura 3.

Los errores detectados y sus soluciones son los siguientes:

- En la columna **Explotaciones** hay valores atípicos que se parecen haber sido acotados en 9999, por lo que también se sustituyen por **NaN**.
- Algunos porcentajes están mal indicados, pues su valor está fuera del intervalo $[0, 100]$; y ninguno representa un decrecimiento numérico ($< 0\%$) o representan algo que pueda exceder el total de la población de forma coherente ($> 100\%$). Para solucionarlo se sustituyen todos los valores fuera del rango por **NaN**.

A continuación se encuentran las instrucciones para corregir estos errores:

```
1  # Missings no declarados variables cuantitativas (-1, 99999)
2  datos['Explotaciones'] = datos['Explotaciones'].replace(99999, np.nan)
3
4  # Valores fuera de rango. 'PobChange_pct' puede contener porcentajes negativos
5
6  c_no_porcentajes = ['Name', 'CodigoProvincia', ..., 'Explotaciones']
7  c_porcentajes = [c for c in variables if c not in c_no_porcentajes]
8  for var in c_porcentajes:
9      if var != 'PobChange_pct':
10         datos[var] = [x if 0 <= x <= 100 else np.nan for x in datos[var]]
```

Listing 5: Corrección de errores en variables numéricas.

2.3. Valores atípicos y *missing*. Imputaciones.

Antes de proceder, conviene definir las variables objetivo y diferenciar todas las columnas:

```
1      datos = datos.set_index(datos['Name']).drop('Name', axis = 1)
2      varObjCont = datos['Izda_Pct']
3      varObjBin = datos['Izquierda']
4      datos_input = datos.drop(['Izda_Pct', 'Izquierda'], axis = 1)
5      variables_input = list(datos_input.columns)
6
7      numericas_input = datos_input.select_dtypes(include = ['int', 'int32', '
int64', 'float', 'float32', 'float64']).columns
8      categoricas_input = [variable for variable in variables_input if variable
not in numericas_input]
9
```

Listing 6: Definición de variables objetivo e input.

Para el tratamiento de los valores atípicos es necesario calcular la proporción de estos en cada variable, de forma que se puedan identificar aquellas columnas con porcentajes destacables (>50%).

```
1      numAtipicos = {x: atipicosAmissing(datos_input[x])[1] / len(datos_input) for x
in numericas_input}
2      for x in numericas_input:
3          datos_input[x] = atipicosAmissing(datos_input[x])[0]
4
5
```

Listing 7: Corrección de errores en variables numéricas.

El último paso antes de las imputaciones consiste en tratar las variables con un porcentaje notable de datos NaN. En concreto, se calcula la proporción de ellos en cada una y si este es superior al 50% se elimina la variable de la BBDD, tal y como se puede observar en la imagen ??.

```
1      # Mapa de calor que muestra la matriz de correlaci n:
2      patron_perdidos(datos_input)
3
4      # Proporción de valores perdidos por cada variable.
5      prop_missingsVars = datos_input.isna().sum()/len(datos_input)
6      prop_missingsVars
7      # Número de valores perdidos por cada observaci n:
8      datos_input['prop_missings'] = datos_input.isna().mean(axis = 1)
9
10     # Número de valores distintos que tiene la nueva variable.
11     len(datos_input['prop_missings'].unique())
12
13     # Se eliminan las observaciones con m s del 50% de datos missing.
14     eliminar = datos_input['prop_missings'] > 0.5
15     datos_input = datos_input[~eliminar]
16     varObjBin = varObjBin[~eliminar]
17     varObjCont = varObjCont[~eliminar]
18
19     # Transformaci n de la nueva variable a categ rica:
20     datos_input["prop_missings"] = datos_input["prop_missings"].astype(str)
21
22     # Se eliminan las variables con m s del 50% de datos missing.
23     eliminar = [prop_missingsVars.index[x] for x in range(len(prop_missingsVars))
if prop_missingsVars[x] > 0.5]
24     datos_input = datos_input.drop(eliminar, axis = 1)
25
26     # Se crean otra vez las listas de variables por si acaso:
27     variables_input = list(datos_input.columns)
```



```

28     categoricas_input = [variable for variable in variables_input if variable not
29     in numericas_input]

```

Listing 8: Corrección de errores en variables numéricas.

Para esta base de datos, no existe ninguna variable con una proporción de valores atípicos superior al 10%, y tan solo hay dos columnas que se acercan (**Population** y **TotalCensus**), por lo que no se elimina ninguna. Además, se observa que la proporción de valores perdidos tiene 9 valores únicos, por lo que se puede convertir en categórica para su codificación.

Por último, los datos *missing* son rellenados con valores aleatorios introducidos mediante código. Basta con ejecutar las siguientes líneas y el depurado de la base de datos estará completo. Esta nueva base de datos limpia se guarda en el archivo **datosTarea.pickle**, que se empleará para el desarrollo de los modelos de regresión lineal y logística.

```

1     # Variables cuantitativas:
2     for x in numericas_input:
3         datos_input[x] = ImputacionCuant(datos_input[x], 'aleatorio')
4     # Variables cualitativas:
5     for x in categoricas_input:
6         datos_input[x] = ImputacionCuali(datos_input[x], 'aleatorio')
7
8     datos_input.isna().sum()
9
10    datosTarea = pd.concat([varObjBin, varObjCont, datos_input], axis = 1)
11    with open('datosTarea.pickle', 'wb') as archivo:
12        pickle.dump(datosTarea, archivo)
13

```

Listing 9: Corrección de errores en variables numéricas.

3. Gráficos y presentación de variables.

Los siguientes apartados incluyen algunos gráficos que permiten visualizar la correlación entre variables, coeficientes relevantes, etc. Para ello, se cargan los datos depurados y se definen las variables objetivo continua y binaria:

```
1 with open('datosTarea.pickle', 'rb') as f:
2     datos = pickle.load(f)
3
4 # VARIABLES OBJETIVO.
5 varObjCont = datos['Izda_Pct']
6 varObjBin = datos['Izquierda']
7 datos_input = datos.drop(['Izda_Pct', 'Izquierda'], axis = 1)
8 variables = list(datos_input.columns)
```

Listing 10: Librerías y lectura de datos.

3.1. Coeficientes V de Cramer.

La V de Cramer indica el grado de asociación entre dos variables categorías. La función **VCramer** de **FuncionesMineria** permite calcular estos coeficientes. Con el uso de otra función, **graficoVCramer** se genera un gráfico de barras que permite comparar el valor de este coeficiente para todas las variables del DataFrame. El código es el siguiente:

```
1 # Gr ficos comparando los V de Cramer.
2 graficoVcramer(datos_input, varObjBin)
3 graficoVcramer(datos_input, varObjCont)
4
5 # C lculo de los coeficientes V de Cramer:
6 VCramer = pd.DataFrame(columns=['Variable', 'Objetivo', 'Vcramer'])
7 for variable in variables:
8     v_cramer = Vcramer(datos_input[variable], varObjCont)
9     VCramer = VCramer.append({'Variable': variable, 'Objetivo': varObjCont.
10 name, 'Vcramer': v_cramer}, ignore_index=True)
11 for variable in variables:
12     v_cramer = Vcramer(datos_input[variable], varObjBin)
13     VCramer = VCramer.append({'Variable': variable, 'Objetivo': varObjBin.name
14 , 'Vcramer': v_cramer}, ignore_index=True)
```

Listing 11: Código para generar un gráfico y calcular las V de Cramer.

Los bucles de esta celda se emplean para guardar los cálculos en un DataFrame en caso de que sean necesarios más adelante, y se pueden ver en la tabla 7 del anexo. Para estos coeficientes, un valor mayor simboliza una mejor asociación entre variables, mientras que un valor más cercano a 0 tiene significado contrario.

3.2. Variable binaria.

El gráfico obtenido para la variable objetivo binaria (**Izquierda**) es el siguiente:

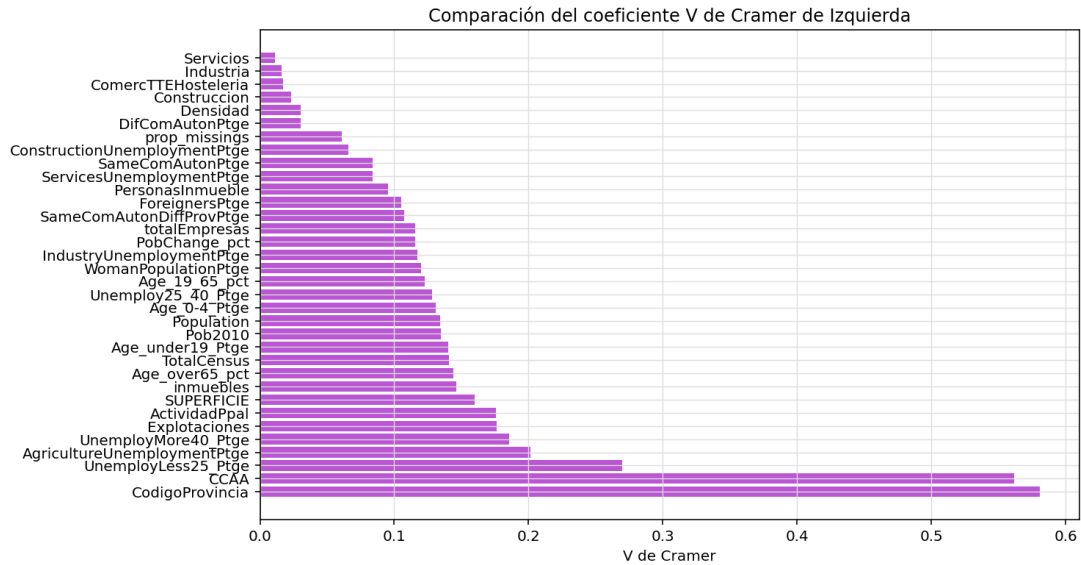


Figura 5: Coeficientes V de Cramer con respecto a la variable binaria **Izquierda**.

Como se puede observar en la imagen 5, las variables categoricas que menos y más se relacionan respectivamente con **Izquierda** son **Densidad** y **CCAA**. Como la variable **Densidad** no provoca mucha variación en los gráficos, se escoge en su lugar **ActividadPpal**, ya que al presentar variación será mejor para predecir la variable objetivo binaria. Los gráficos de mosaico para la relación entre cada una de estas dos variables y la variable objetivo binaria son:

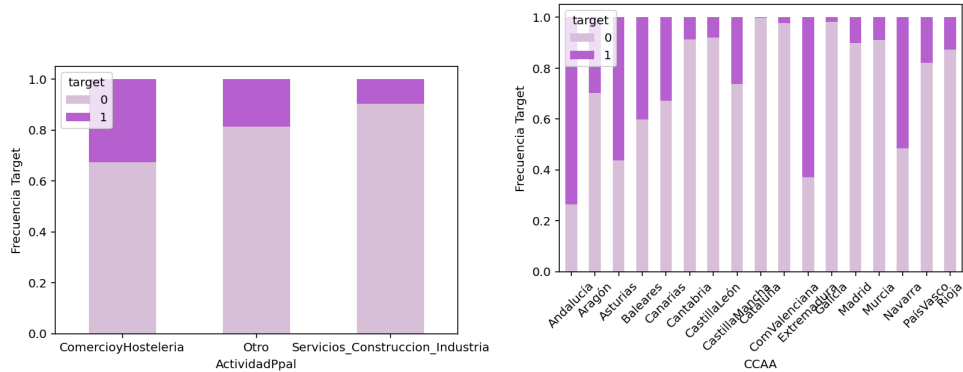


Figura 6: Gráficos de mosaico para relación entre la variable objetivo binaria **Izquierda** con **ActividadPpal** (izqa) y **CCAA** (dcha).

Los gráficos de caja pueden servir para comprobar cómo de bien una variable numérica describe cualquier variable objetivo. Lo que se busca, al igual que con los gráficos de mosaico, son dos variables que muestren diferencias claras entre los dos gráficos de caja para los dos valores 1 y 0 de **Izquierda**. En el caso de la variable objetivo binaria se han elegido las columnas **ServicesUnemploymentPtge** y **UnemployLess25_Ptge**, es decir, los porcentajes de desempleo en el sector de servicios (con coeficiente V de Cramer pequeño según la figura 5), y de desempleo en habitantes menores de 25 años (con V de Cramer alta).

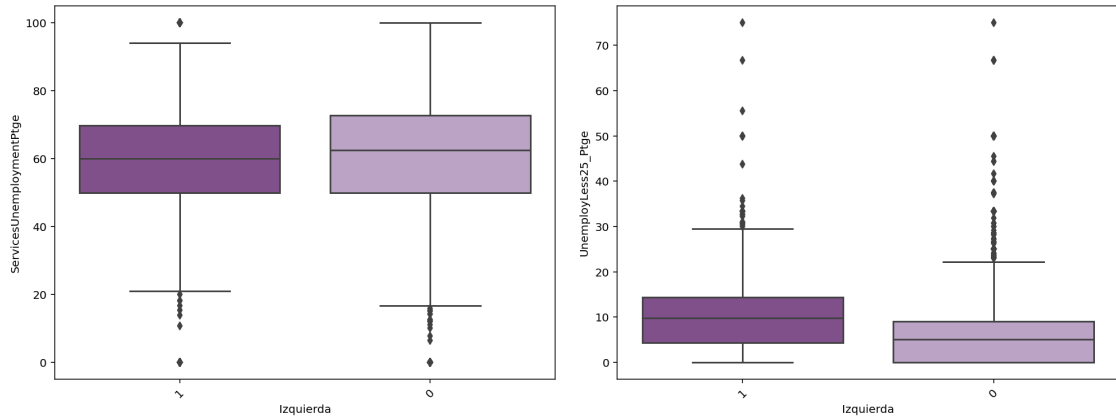


Figura 7: Diagrama de caja para **Izquierda** en función de **ServicesUnemploymentPtge** y **UnemployLess25_Ptge**

Al igual que con los gráficos anteriores, lo que se busca con los histogramas son variables que difieran mucho para poder entrenar el modelo de regresión. En este caso las elegidas son:

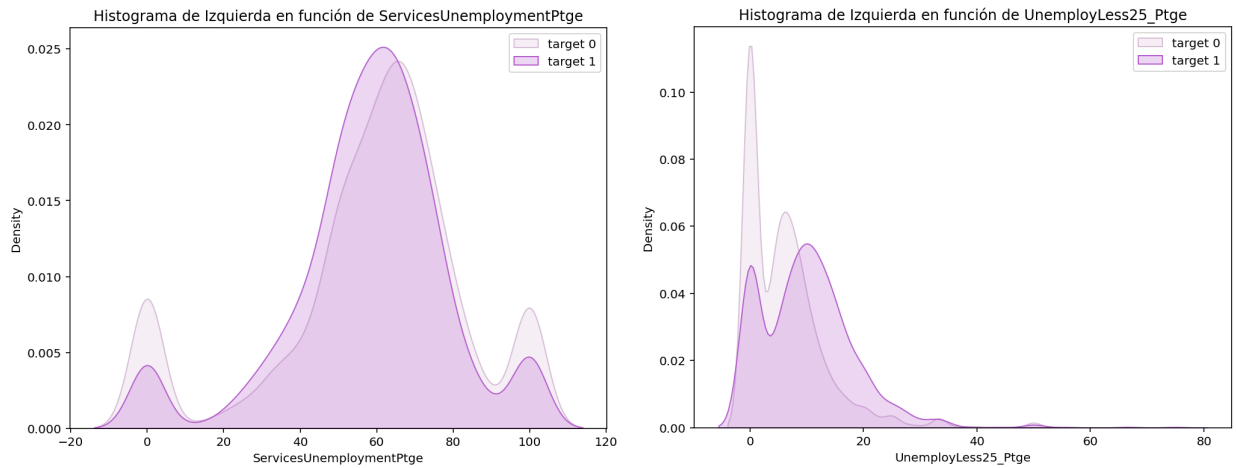


Figura 8: Histograma de **Izquierda** en función de **ServicesUnemploymentPtge** y **UnemployLess25_Ptge**.

3.3. Variable continua.

Por otro lado, para la variable objetivo continua (**Izda_Pct** se genera la siguiente imagen:

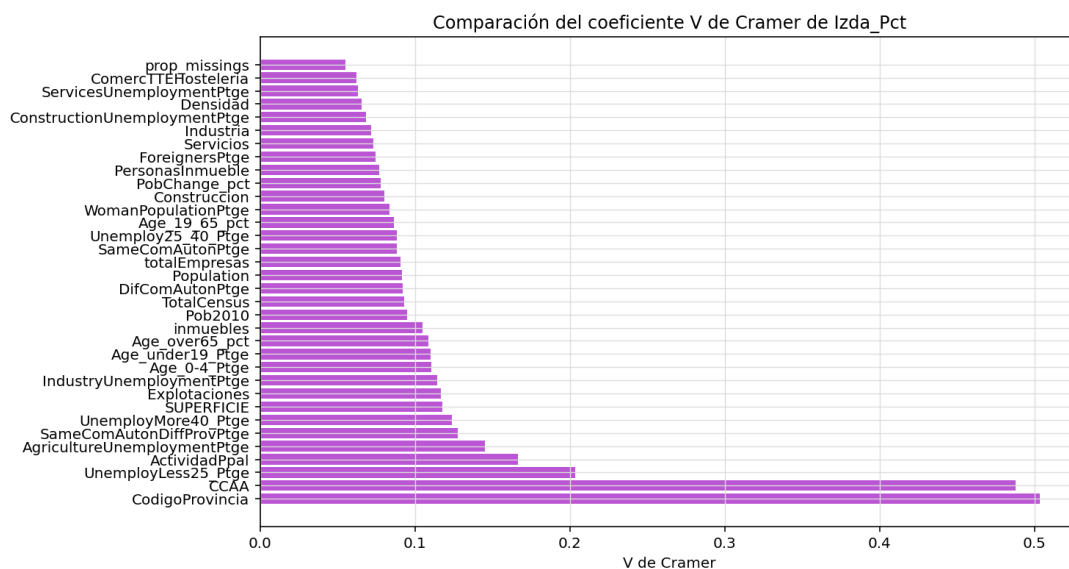


Figura 9: Coeficientes V de Cramer con respecto a la variable continua **Izda_Pct**.

De la misma forma, se buscan variables en la parte alta y baja del gráfico de la figura 9 para poder entrenar el modelo de regresión lineal. En este caso, al tratarse de una variable numérica, se obtienen gráficos con numerosas distribuciones debido a los distintos valores posibles de las otras variables tal y como se ve a continuación.

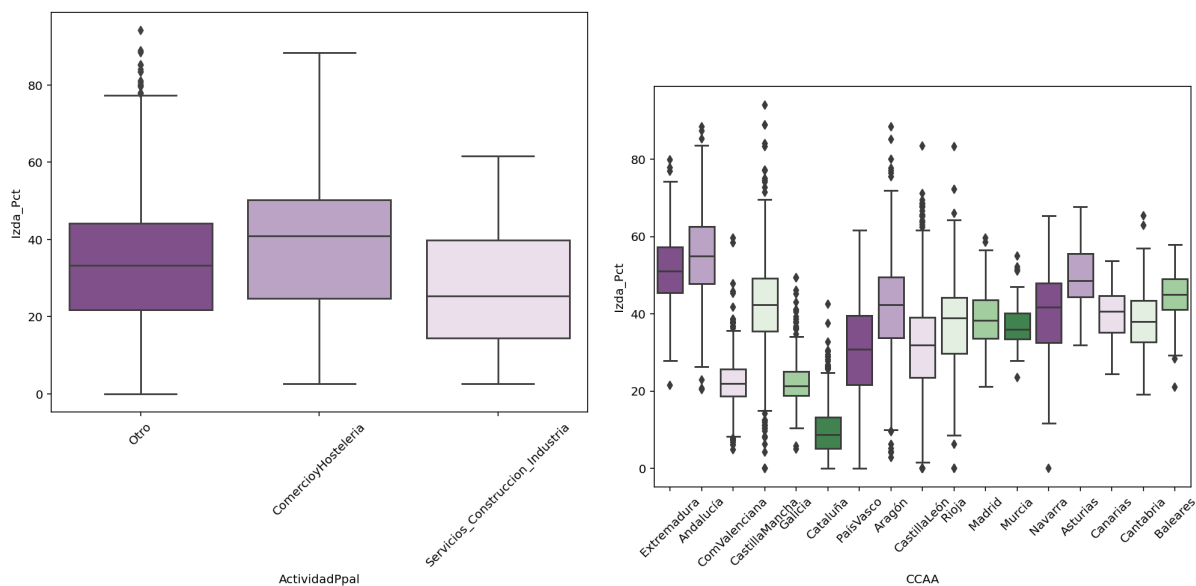


Figura 10: Diagrama de caja para **Izquierda** en función de **ActividadPpal** y **CCAA**.

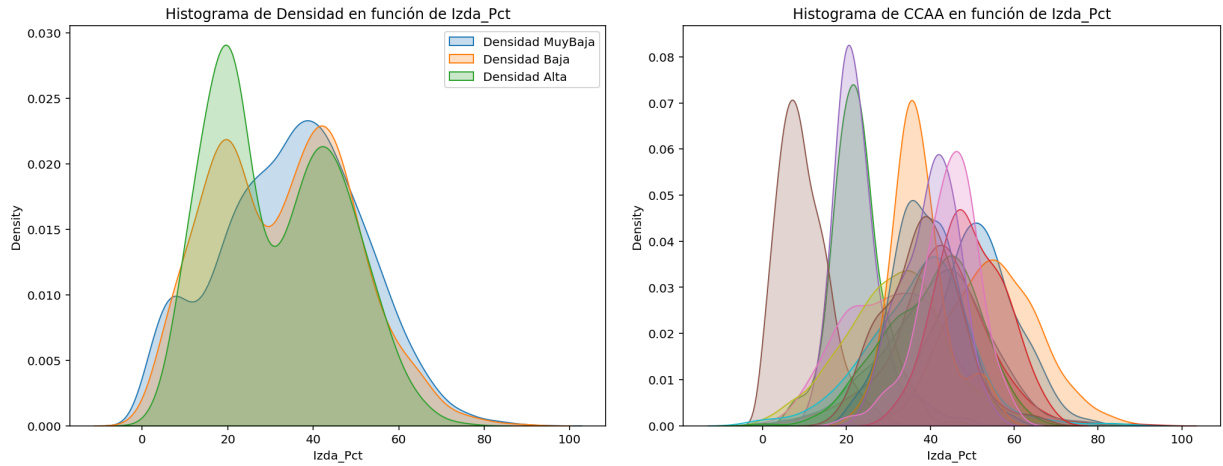


Figura 11: Diagrama de caja para **Izquierda** en función de

Esta matriz nos indica el grado de correlación entre dos variables, lo cual es de gran ayuda para la selección manual.

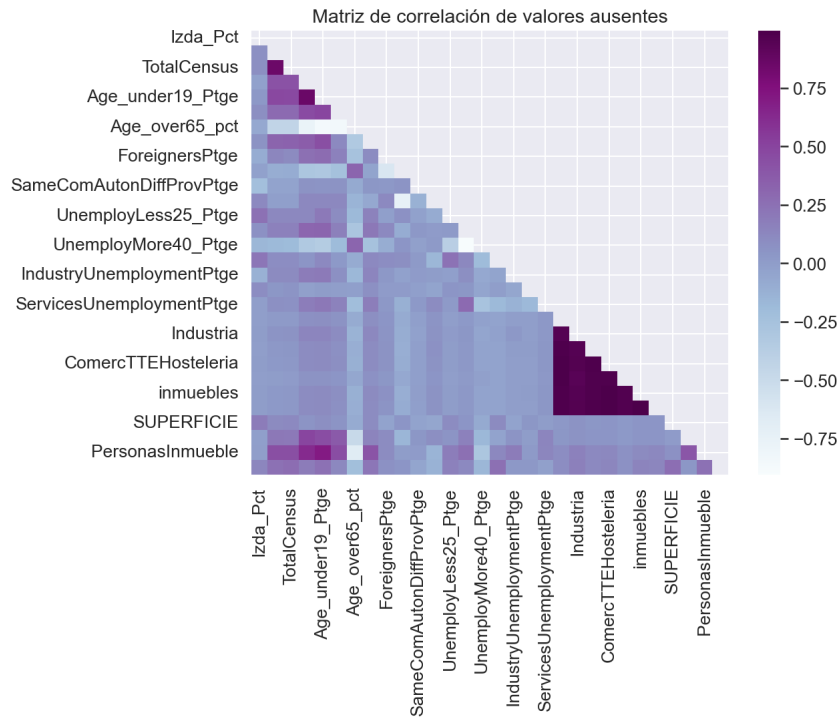


Figura 12: Matriz de correlación.

4. Modelos de regresión lineal.

En este apartado se muestra como crear los modelos de regresión lineal para poder predecir la variable objetivo continua (**Izda_Pct**). Para ello se hace uso de la librería **FuncionesMineria** y de los métodos vistos en clase, como la validación cruzada o los modelos con métricas AIC o BIC.

4.1. Modelo de regresión lineal. Selección de variables clásica.

En este apartado se emplean los métodos *Stepwise selection*, *Backward selection* y *Forward selection*. Estos métodos de selección de variables crean modelos y van añadiendo variables (*Forward*) o quitándolas (*Backward*) del modelo inicial mientras se comprueba la bondad del mismo mediante los criterios AIC (Criterio de Información de Akaike) o BIC (Criterio de Información Bayesiano), que permiten ajustar la complejidad del mismo para evitar seleccionar variables poco útiles o que den información no necesaria. El método *Stepwise* combina los dos anteriores para un desarrollo más eficiente. Lectura de los datos y asignación de variables:

```
1 with open('datosTarea_cont.pickle', 'rb') as f:
2     todo = pickle.load(f)
3     # Identificacion de variables:
4     varObjCont = todo['Izda_Pct']
5     todo = todo.drop('Izda_Pct', axis = 1)
6     var_cont = todo.select_dtypes(include = ['int', 'float']).columns.tolist()
7     var_cont_sin_transf = ['CodigoProvincia', ..., 'Explotaciones']
8     var_categ = [var for var in todo.columns.tolist() if var not in var_cont]
```

Listing 12: Modelos de regresión lineal. Selección de variables.

En este caso se han construido 6 modelos de regresión lineal mediante tres métodos diferentes, cada uno con ambas métricas estudiadas. En este primer caso se han considerado interacciones entre las variables. El código para la creación de los modelos se puede comprobar en el siguiente listado.

```
1 x_train, x_test, y_train, y_test = train_test_split(todo, varObjCont,
2     test_size = 0.2, random_state = 1234567)
3 interacciones = var_cont_sin_transf
4 interacciones_unicas = list(combinations(interacciones, 2))
5
6 # MODELO STEPWISE AIC.
7 modeloStepAIC_int = lm_stepwise(y_train, x_train, var_cont_sin_transf,
8     var_categ, interacciones_unicas, 'AIC')
9 modeloStepAIC_int['Modelo'].summary()
10 r_train_stepAIC_int = Rsq(modeloStepAIC_int['Modelo'], y_train,
11     modeloStepAIC_int['X'])
12 x_test_modeloStepAIC_int = crear_data_modelo(x_test, modeloStepAIC_int['
13     Variables']['cont'], modeloStepAIC_int['Variables']['categ'], modeloStepAIC_int
14     ['Variables']['inter'])
15
16 r_test_stepAIC_int = Rsq(modeloStepAIC_int['Modelo'], y_test,
17     x_test_modeloStepAIC_int)
18 param_stepAIC_int = len(modeloStepAIC_int['Modelo'].params)
19
20 # MODELO STEPWISE BIC.
21 modeloStepBIC_int = lm_stepwise(y_train, x_train, var_cont_sin_transf,
22     var_categ, interacciones_unicas, 'BIC')
23 modeloStepBIC_int['Modelo'].summary()
24 r_train_stepBIC_int = Rsq(modeloStepBIC_int['Modelo'], y_train,
25     modeloStepBIC_int['X'])
26 x_test_modeloStepBIC_int = crear_data_modelo(x_test, modeloStepBIC_int['
27     Variables']['cont'], modeloStepBIC_int['Variables']['categ'], modeloStepBIC_int
28     ['Variables']['inter'])
```

```

19     r_test_stepBIC_int = Rsq(modeloStepBIC_int['Modelo'], y_test,
20     x_test_modeloStepBIC_int)
21     param_stepBIC_int = len(modeloStepBIC_int['Modelo'].params)
22
23     # MODELO BACKWARD AIC CON EL PRIMER CONJUNTO DE INTERACCIONES.
24     modeloBackAIC_int = lm_backward(y_train, x_train, var_cont_sin_transf,
25     var_categ, interacciones_unicas, 'AIC')
26     modeloBackAIC_int['Modelo'].summary()
27     r_train_backAIC_int = Rsq(modeloBackAIC_int['Modelo'], y_train,
28     modeloBackAIC_int['X'])
29     x_test_modeloBackAIC_int = crear_data_modelo(x_test, modeloBackAIC_int['
30     Variables']['cont'], modeloBackAIC_int['Variables']['categ'], modeloBackAIC_int
31     ['Variables']['inter'])
32     r_test_backAIC_int = Rsq(modeloBackAIC_int['Modelo'], y_test,
33     x_test_modeloBackAIC_int)
34     param_backAIC_int = len(modeloBackAIC_int['Modelo'].params)
35
36     # MODELO BACKWARD BIC CON EL PRIMER CONJUNTO DE INTERACCIONES.
37     modeloBackBIC_int = lm_backward(y_train, x_train, var_cont_sin_transf,
38     var_categ, interacciones_unicas, 'BIC')
39     modeloBackBIC_int['Modelo'].summary()
40     r_train_backBIC_int = Rsq(modeloBackBIC_int['Modelo'], y_train,
41     modeloBackBIC_int['X'])
42     x_test_modeloBackBIC_int = crear_data_modelo(x_test, modeloBackBIC_int['
43     Variables']['cont'], modeloBackBIC_int['Variables']['categ'], modeloBackBIC_int
44     ['Variables']['inter'])
45     r_test_backBIC_int = Rsq(modeloBackBIC_int['Modelo'], y_test,
46     x_test_modeloBackBIC_int)
47     param_backBIC_int = len(modeloBackBIC_int['Modelo'].params)
48
49     # MODELO FORWARD AIC CON EL PRIMER CONJUNTO DE INTERACCIONES.
50     modeloForwAIC_int = lm_stepwise(y_train, x_train, var_cont_sin_transf,
51     var_categ, interacciones_unicas, 'AIC')
52     modeloForwAIC_int['Modelo'].summary()
53     r_train_forwAIC_int = Rsq(modeloForwAIC_int['Modelo'], y_train,
54     modeloForwAIC_int['X'])
55
56     x_test_modeloForwAIC_int = crear_data_modelo(x_test, modeloForwAIC_int['
57     Variables']['cont'], modeloForwAIC_int['Variables']['categ'], modeloForwAIC_int
58     ['Variables']['inter'])
59     r_test_forwAIC_int = Rsq(modeloForwAIC_int['Modelo'], y_test,
60     x_test_modeloForwAIC_int)
61     param_forwAIC_int = len(modeloForwAIC_int['Modelo'].params)
62
63     # MODELO FORWARD BIC CON EL PRIMER CONJUNTO DE INTERACCIONES.
64     modeloForwBIC_int = lm_stepwise(y_train, x_train, var_cont_sin_transf,
65     var_categ, interacciones_unicas, 'BIC')
66     modeloForwBIC_int['Modelo'].summary()
67     r_train_forwBIC_int = Rsq(modeloForwBIC_int['Modelo'], y_train,
68     modeloForwBIC_int['X'])
69     x_test_modeloForwBIC_int = crear_data_modelo(x_test, modeloForwBIC_int['
70     Variables']['cont'], modeloForwBIC_int['Variables']['categ'], modeloForwBIC_int
71     ['Variables']['inter'])
72     r_test_forwBIC_int = Rsq(modeloForwBIC_int['Modelo'], y_test,
73     x_test_modeloForwBIC_int)
74     param_forwBIC_int = len(modeloForwBIC_int['Modelo'].params)

```

Listing 13: Creación de modelos de regresión lineal por selección de variables clásica.

Los resultados más importantes de estos modelos se recogen en la tabla 1. De esta tabla se puede concluir que los mejores modelos son aquellos que emplean los métodos Forward y Stepwise

con métrica BIC, ya que tienen un coeficiente R^2 parecido y el menor número de parámetros. Es decir, estos dos modelos serán los mejores para predecir la variable objetivo continua.

Método	Métrica	R^2_{train}	R^2_{test}	Parámetros
Backward	AIC	0.6625	0.6364	83
Backward	BIC	0.6582	0.6335	65
Forward	AIC	0.6620	0.6377	79
Forward	BIC	0.6566	0.6343	64
Stepwise	AIC	0.6620	0.6377	79
Stepwise	BIC	0.6566	0.6343	64

Tabla 1: Resultados de R^2 en entrenamiento y test para distintos métodos y métricas

El código siguiente ejecuta un proceso de validación cruzada que permite comparar los modelos:

```

1  results = pd.DataFrame({
2      'Rsquared': []
3      , 'Resample': []
4      , 'Modelo': []})
5  for rep in range(20):
6      # Realiza validación cruzada en cuatro modelos diferentes y almacena sus
6      # R-squared en listas separadas
7
8      modelo_stepAIC_int = validacion_cruzada_lm(
9          5
10         , x_train
11         , y_train
12         , modeloStepAIC_int['Variables']['cont']
13         , modeloStepAIC_int['Variables']['categ']
14         , modeloStepAIC_int['Variables']['inter']
15     )
16     modelo_stepBIC_int = validacion_cruzada_lm(
17         5
18         , x_train
19         , y_train
20         , modeloStepBIC_int['Variables']['cont']
21         , modeloStepBIC_int['Variables']['categ']
22         , modeloStepBIC_int['Variables']['inter']
23     )
24     modelo_backAIC_int = validacion_cruzada_lm(
25         5
26         , x_train
27         , y_train
28         , modeloBackAIC_int['Variables']['cont']
29         , modeloBackAIC_int['Variables']['categ']
30         , modeloBackAIC_int['Variables']['inter']
31     )
32     modelo_backBIC_int = validacion_cruzada_lm(
33         5
34         , x_train
35         , y_train
36         , modeloBackBIC_int['Variables']['cont']
37         , modeloBackBIC_int['Variables']['categ']
38         , modeloBackBIC_int['Variables']['inter']
39     )
40     modelo_forwaAIC_int = validacion_cruzada_lm(
41         5
42         , x_train
43         , y_train

```

```

44     , modeloForwAIC_int['Variables']['cont']
45     , modeloForwAIC_int['Variables']['categ']
46     , modeloForwAIC_int['Variables']['inter']
47 )
48 modelo_forwBIC_intt = validacion_cruzada_lm(
49     5
50     , x_train
51     , y_train
52     , modeloForwBIC_int['Variables']['cont']
53     , modeloForwBIC_int['Variables']['categ']
54     , modeloForwBIC_int['Variables']['inter']
55 )
56
57 results_rep = pd.DataFrame({
58     'Rsquared': modelo_stepAIC_int + modelo_stepBIC_int +
59     modelo_backAIC_int + modelo_backBIC_int + modelo_forwAIC_int +
60     modelo_forwBIC_intt
61     , 'Resample': ['Rep' + str((rep + 1))] * 5 * 6
62     , 'Modelo': [1] * 5 + [2] * 5 + [3] * 5 + [4] * 5 + [5] * 5 + [6] * 5
63 })
64 results = pd.concat([results, results_rep], axis = 0)

```

Listing 14: Modelos de regresión lineal. Validación cruzada.

Y los resultados son:

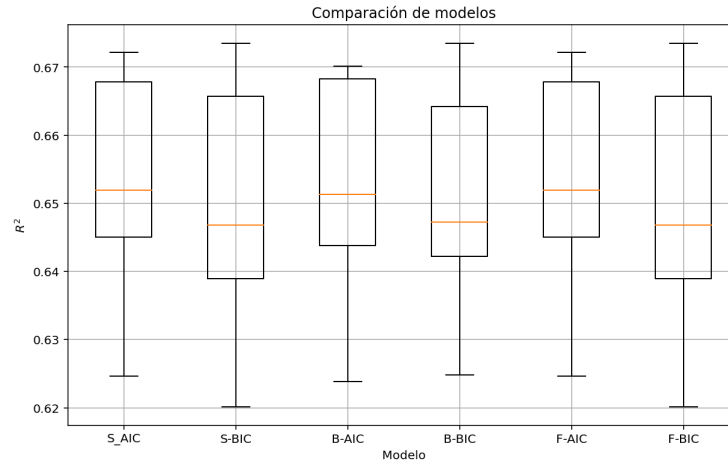


Figura 13: Diagrama de cajas de los coeficientes R^2 de los modelos seleccionados.

Modelo	$\overline{R^2}$	$\sigma(R^2)$	Parámetros
S-AIC	0.6523	0.0171	79
S-BIC	0.6490	0.0192	64
B-AIC	0.6515	0.0171	83
B-BIC	0.6504	0.0172	65
F-AIC	0.6523	0.0171	79
F-BIC	0.6490	0.0192	64

Tabla 2: Valores de promedio y desviación típica de R^2 y número de parámetros para distintos modelos

Por lo que el modelo ganador es el *Stepwise* BIC, ya que tiene el menor número de parámetros y los promedios de R^2 son similares.

4.2. Modelo de regresión lineal. Selección de variables aleatoria.

En este caso se va a ejecutar un bucle de 30 iteraciones. En cada una se hace una partición aleatoria de los datos de entrenamiento $\mathbf{x}_{\text{train}}$ (de los cuales se destina el 70 % para el nuevo entrenamamiento y 30 % para las pruebas). Una vez divididos se desarrolla un modelo con método de selección paso a paso (*Stepwise*) y según el criterio bayesiano (BIC). Es decir, en cada iteración se genera un modelo nuevo, y del conjunto de modelos se obtienen todas las combinaciones de variables que potencialmente puedan mejorar las predicciones. En esta selección se tienen en cuenta las interacciones anteriores.

```
1 variables_seleccionadas = {
2     'Formula': [],
3     'Variables': []}
4 for x in range(30):
5     print('----- iter: ' + str(x))
6     x_train2, x_test2, y_train2, y_test2 = train_test_split(x_train, y_train,
7 test_size = 0.3, random_state = 1234567 + x)
8     modelo = lm_stepwise(y_train2.astype(int), x_train2, var_cont, var_categ,
9 interacciones_unicas, 'BIC')
10    variables_seleccionadas['Variables'].append(modelo['Variables'])
11    variables_seleccionadas['Formula'].append(sorted(modelo['Modelo'].model.
12 exog_names))
13 variables_seleccionadas['Formula'] = list(map(lambda x: '+'.join(x),
14 variables_seleccionadas['Formula']))
15 frecuencias = Counter(variables_seleccionadas['Formula'])
16 frec_ordenada = pd.DataFrame(list(frecuencias.items()), columns = ['Formula',
17 'Frecuencia'])
18 frec_ordenada = frec_ordenada.sort_values('Frecuencia', ascending = False).
19 reset_index()
20 var_1 = variables_seleccionadas['Variables'][variables_seleccionadas['Formula']
21 ].index(frec_ordenada['Formula'][0])]
```

Listing 15: Selección de variables aleatoria.

Gracias a este proceso, se pueden emplear las tres fórmulas que más veces han aparecido de forma aleatoria para el desarrollo de otros 3 modelos, que se comparan a continuación con el modelo ganador de selección clásica (S-BIC):

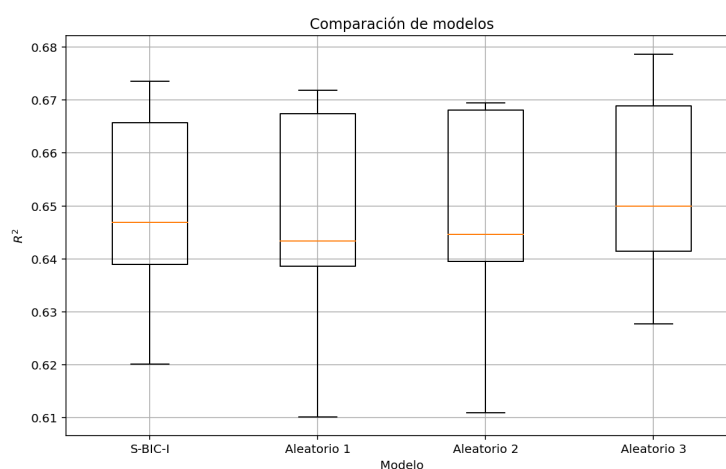


Figura 14: Comparación entre el modelo ganador de selección clásica y tres modelos aleatorios.

Y los resultados son:

Modelo	$\overline{R^2}$	$\sigma(R^2)$	Parámetros
S-BIC	0.6490	0.0192	64
Aleatorio 1	0.6462	0.0224	33
Aleatorio 2	0.6465	0.0216	35
Aleatorio 3	0.6533	0.0185	62

Tabla 3: Resultados de R^2 promedio y desviación estándar para distintos modelos.

Donde es evidente que el **primer modelo aleatorio** es el ganador por el principio de Parsimonia, ya que necesita muchos menos parámetros (aproximadamente la mitad) que los modelos S-BIC, aleatorio 2 y aleatorio 3, aunque el modelo aleatorio 2 es un buen candidato.

4.3. Modelo de regresión lineal. Conclusiones.

El resumen del modelo ganador es el siguiente. Se puede observar que el modelo ganador tiene un coeficiente F de 374,7. Se puede observar que, mientras que un aumento en la mayoría de variables supone una predicción menor para la variable **Izquierda**, hay algunas que indican aumento, como **AgricultureUnemploymentPtge** o **UnemployLess25_Ptge**. Con respecto al valor del estadístico p , la mayoría de variables son relevantes al modelo pues este sólo supera el 5 % en dos ocasiones (ambas para factores relativos a interacciones entre variables).

Dep. Variable:	Izda_Pct	R-squared:	0.650			
Model:	OLS	Adj. R-squared:	0.648			
Method:	Least Squares	F-statistic:	374.7			
Date:	Thu, 20 Feb 2025	Prob (F-statistic):	0.00			
Time:	22:16:44	Log-Likelihood:	-23995.			
No. Observations:	6493	AIC:	4.806e+04			
Df Residuals:	6460	BIC:	4.828e+04			
Df Model:	32					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	34.0268	2.446	13.909	0.000	29.231	38.822
raiz4DifComAutonPtge	10.3711	0.721	14.391	0.000	8.958	11.784
sqrtxAgricultureUnemploymentPtge	1.9167	0.254	7.537	0.000	1.418	2.415
Age_19_65_pct	0.1849	0.022	8.430	0.000	0.142	0.228
sqrtxUnemployLess25_Ptge	6.0633	0.720	8.420	0.000	4.652	7.475
logxPopulation	0.6471	0.129	5.013	0.000	0.394	0.900
raiz4Explotaciones	-5.5464	0.746	-7.432	0.000	-7.009	-4.083
ForeignersPtge	-0.1147	0.019	-6.066	0.000	-0.152	-0.078
logxinmuebles	0.8258	0.154	5.356	0.000	0.524	1.128
logxConstructionUnemploymentPtge	0.1921	0.038	5.035	0.000	0.117	0.267
sqrtxWomanPopulationPtge	2.2544	0.574	3.925	0.000	1.129	3.380
CCAA_Aragón	-9.3148	0.659	-14.133	0.000	-10.607	-8.023
CCAA_Asturias	-6.0731	1.325	-4.582	0.000	-8.671	-3.475
CCAA_Baleares	-8.7071	1.439	-6.051	0.000	-11.528	-5.886
CCAA_Canarias	-12.7730	1.299	-9.836	0.000	-15.319	-10.227
CCAA_Cantabria	-16.7808	1.225	-13.700	0.000	-19.182	-14.380
CCAA_CastillaLeón	-18.7859	0.562	-33.422	0.000	-19.888	-17.684
CCAA_CastillaMancha	-11.5004	0.629	-18.273	0.000	-12.734	-10.267
CCAA_Cataluña	-42.8001	0.622	-68.781	0.000	-44.020	-41.580
CCAA_ComValenciana	-30.2972	0.668	-45.328	0.000	-31.608	-28.987
CCAA_Extremadura	-4.7275	0.716	-6.603	0.000	-6.131	-3.324
CCAA_Galicia	-30.4978	0.827	-36.889	0.000	-32.118	-28.877
CCAA_Madrid	-15.6571	0.987	-15.857	0.000	-17.593	-13.721
CCAA_Murcia	-15.5734	1.613	-9.654	0.000	-18.736	-12.411
CCAA_Navarra	-12.9485	0.839	-15.433	0.000	-14.593	-11.304
CCAA_PaísVasco	-23.1846	0.845	-27.422	0.000	-24.842	-21.527
CCAA_Rioja	-14.2929	0.999	-14.312	0.000	-16.251	-12.335
sqrtxUnemployLess25_Ptge_ActividadPpal_Otro	4.3986	0.804	5.471	0.000	2.823	5.975
sqrtxUnemployLess25_Ptge_ActividadPpal_Servicios_Construccion_Industria	-0.6720	2.652	-0.253	0.800	-5.870	4.527
UnemployLess25_Ptge_ActividadPpal_Otro	-0.6543	0.070	-9.292	0.000	-0.792	-0.516
UnemployLess25_Ptge_ActividadPpal_Servicios_Construccion_Industria	-0.2955	0.332	-0.891	0.373	-0.945	0.354
sqrtxUnemployLess25_Ptge_Age_under19_Ptge	-0.6679	0.066	-10.087	0.000	-0.798	-0.538
UnemployLess25_Ptge_Age_under19_Ptge	0.0449	0.005	8.836	0.000	0.035	0.055
Omnibus:	394.304	Durbin-Watson:	2.013			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	825.039			
Skew:	0.413	Prob(JB):	7.00e-180			
Kurtosis:	4.539	Cond. No.	3.69e+03			

Figura 15: Modelo de regresión lineal ganador.

5. Modelos de regresión logística.

En esta sección se repite todo el proceso del apartado 4. Para ello se emplean funciones de **FuncionesMineria** parecidas pero adaptadas a la variable objetivo binaria. El proceso es el mismo: selección de variables clásica, selección aleatoria y comparación. Mediante validación cruzada se elegirán los mejores modelos para predecir la variable **Izquierda**.

```
1 with open('datosTarea_bin.pickle', 'rb') as f:
2     todo = pickle.load(f)
3     varObjBin = todo['Izquierda']
4     todo = todo.drop('Izquierda', axis = 1)
5     pd.DataFrame({
6         'n': varObjBin.value_counts(),
7         '%': varObjBin.value_counts(normalize = True)})
8     var_cont = todo.select_dtypes(include = ['int', 'float']).columns.tolist()
9     var_cont_sin_transf = ['Population', 'TotalCensus', ..., 'Explotaciones']
10    var_cont_transf = [var for var in var_cont if var not in var_cont_sin_transf]
11    var_categ = [var for var in todo.columns.tolist() if var not in var_cont]
12    todo = todo.drop(var_cont_transf, axis = 1)
13    # Partici n de datos. La variable respuesta se pone como num rica:
14    x_train, x_test, y_train, y_test = train_test_split(todo, varObjBin, test_size
15    = 0.2, random_state = 1234567)
16    y_train, y_test = y_train.astype(int), y_test.astype(int)
```

Listing 16: Modelo de regresión logística. Asignación de varibales y partición de datos.

5.1. Modelo de regresión logística. Selección de variables clásica.

Para la obtención de modelos por selección clásica se ejecuta el código siguiente. En este caso no se contemplan modelos de selección hacia adelante (*Forward*) porque, tal y como se ha visto en el apartado 4.1 proporcionan los mismos resultados que los modelos con selección paso a paso. El código para crear estos modelos es el siguiente, donde no se han tenido en cuenta interacciones debido a problemas de capacidad de computación:

```
1 interacciones_unicas = []
2
3 # MODELO STEPWISE AIC.
4 modeloStepAIC = glm_stepwise(y_train, x_train, var_cont_sin_transf, var_categ,
5 [], 'AIC')
6 summary_glm(modeloStepAIC['Modelo'], y_train, modeloStepAIC['X'])
7 r_train_stepAIC = pseudoR2(modeloStepAIC['Modelo'], modeloStepAIC['X'],
8 y_train)
9 x_test_modeloStepAIC = crear_data_modelo(x_test, modeloStepAIC['Variables']['
10 cont'], modeloStepAIC['Variables']['categ'], modeloStepAIC['Variables']['inter'
11 ])
12 r_test_stepAIC = pseudoR2(modeloStepAIC['Modelo'], x_test_modeloStepAIC,
13 y_test)
14 param_stepAIC = len(modeloStepAIC['Modelo'].coef_[0])
15
16 # MODELO STEPWISE BIC.
17 modeloStepBIC = glm_stepwise(y_train, x_train, var_cont_sin_transf, var_categ,
18 [], 'BIC')
19 summary_glm(modeloStepBIC['Modelo'], y_train, modeloStepBIC['X'])
20 r_train_stepBIC = pseudoR2(modeloStepBIC['Modelo'], modeloStepBIC['X'],
21 y_train)
22 x_test_modeloStepBIC = crear_data_modelo(x_test, modeloStepBIC['Variables']['
23 cont'], modeloStepBIC['Variables']['categ'], modeloStepBIC['Variables']['inter'
24 ])
```

```

16     r_test_stepBIC = pseudoR2(modeloStepBIC['Modelo'], x_test_modeloStepBIC,
17     y_test)
18     param_stepBIC = len(modeloStepBIC['Modelo'].coef_[0])
19
20     # MODELO BACKWARD AIC.
21     modeloBackAIC = glm_backward(y_train, x_train, var_cont_sin_transf, var_categ,
22     [], 'AIC')
23     summary_glm(modeloBackAIC['Modelo'], y_train, modeloBackAIC['X'])
24     r_train_backAIC = pseudoR2(modeloBackAIC['Modelo'], modeloBackAIC['X'],
25     y_train)
26     x_test_modeloBackAIC = crear_data_modelo(x_test, modeloBackAIC['Variables']['
27     cont'], modeloBackAIC['Variables']['categ'], modeloBackAIC['Variables']['inter'
28     ])
29     r_test_backAIC = pseudoR2(modeloBackAIC['Modelo'], x_test_modeloBackAIC,
30     y_test)
31     param_backAIC = len(modeloBackAIC['Modelo'].coef_[0])
32
33     # MODELO BACKWARD BIC.
34     modeloBackBIC = glm_backward(y_train, x_train, var_cont_sin_transf, var_categ,
35     [], 'BIC')
36     summary_glm(modeloBackBIC['Modelo'], y_train, modeloBackBIC['X'])
37     r_train_backBIC = pseudoR2(modeloBackBIC['Modelo'], modeloBackBIC['X'],
38     y_train)
39     x_test_modeloBackBIC = crear_data_modelo(x_test, modeloBackBIC['Variables']['
40     cont'], modeloBackBIC['Variables']['categ'], modeloBackBIC['Variables']['inter'
41     ])
42     r_test_backBIC = pseudoR2(modeloBackBIC['Modelo'], x_test_modeloBackBIC,
43     y_test)
44     param_backBIC = len(modeloBackBIC['Modelo'].coef_[0])

```

Listing 17: Modelo de regresión logística. Creación de modelos por selección de variables clásica.

De nuevo, se genera una tabla con los resultados de cada modelo en función del proceso y la métrica empleados. :

Método	Métrica	R^2_{train}	R^2_{test}	Parámetros
Backward	AIC	0.6564	0.6126	73
Backward	BIC	0.6547	0.6203	65
Stepwise	AIC	0.6549	0.6340	68
Stepwise	BIC	0.6538	0.6336	63

Tabla 4: R^2 en modelos básicos de regresión logística.

A continuación se realiza un proceso de validación cruzada para comprobar qué modelo es mejor, obteniendo así una tabla con los promedios y las desviaciones típicas de los coeficientes R^2 .

```

1     results = pd.DataFrame({
2         'AUC': [],
3         'Resample': [],
4         'Modelo': []})
5
6     for rep in range(20):
7         modelo_StepAIC = validacion_cruzada_glm(5, x_train, y_train, modeloStepAIC
8         ['Variables']['cont'], modeloStepAIC['Variables']['categ'], modeloStepAIC['
9         Variables']['inter'])
10        modelo_StepBIC = validacion_cruzada_glm(5, x_train, y_train, modeloStepBIC
11        ['Variables']['cont'], modeloStepBIC['Variables']['categ'], modeloStepBIC['
12        Variables']['inter'])
13        modelo_BackAIC = validacion_cruzada_glm(5, x_train, y_train, modeloBackAIC
14        ['Variables']['cont'], modeloBackAIC['Variables']['categ'], modeloBackAIC['
15        Variables']['inter'])

```

```

10     modelo_BackBIC= validacion_cruzada_glm(5, x_train, y_train, modeloBackBIC[
    'Variables']['cont'], modeloBackBIC['Variables']['categ'], modeloBackBIC['
    Variables']['inter'])
11
12     results_rep = pd.DataFrame({
13         'AUC': modelo_StepAIC + modelo_StepBIC + modelo_BackAIC +
    modelo_BackBIC
14         , 'Resample': ['Rep' + str((rep + 1))*5*4
15         , 'Modelo': [1]*5 + [2]*5 + [3]*5 + [4]*5})
16     results = pd.concat([results, results_rep], axis = 0)

```

Listing 18: Modelo de regresión logística. Validación cruzada por selección de variables clásica.

Y los resultados obtenidos son los siguientes:

Modelo	$\overline{R^2}$	$\sigma(R^2)$	Parámetros
S-AIC	0.8625	0.0116	76
S-BIC	0.8625	0.0116	76
B-AIC	0.8642	0.0115	75
B-BIC	0.8642	0.0115	75

Tabla 5: Valores de promedio y desviación típica de R^2 y número de parámetros para distintos modelos

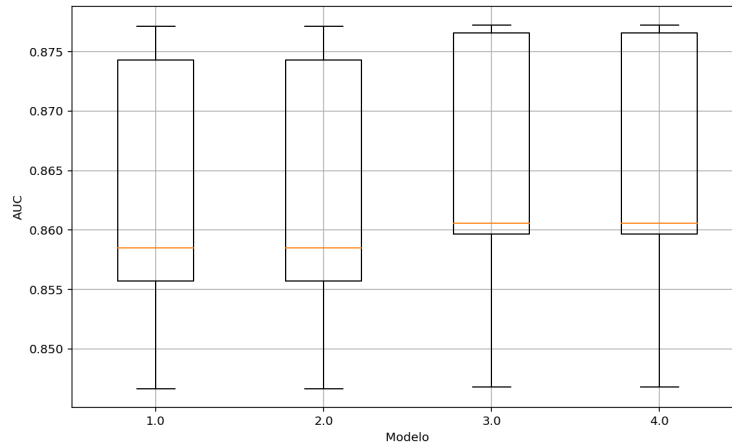


Figura 16: Gráfico de cajas para los modelos de regresión logística creados por selección de variables clásica.

Por lo que se puede afirmar que el modelo ganador es aquel que se ha desarrollado con el método *Backward* y la métrica BIC. Este modelo se comparará con tres obtenidos mediante selección aleatoria de variables en el siguiente apartado.

5.2. Modelo de regresión logística. Selección de variables aleatoria.

El código para la selección aleatoria es el siguiente:

```

1     variables_seleccionadas = {
2         'Formula': [],
3         'Variables': [],
4         'Interaccion': []}
5
6     for x in range(20):
7         print('----- iter: ' + str(x))
8         x_train2, x_test2, y_train2, y_test2 = train_test_split(x_train, y_train,
    test_size = 0.3, random_state = 1234567 + x)

```

```

9     modelo = glm_stepwise(y_train2.astype(int), x_train2, var_cont, var_categ,
10     interacciones_unicas, 'BIC')
11     variables_seleccionadas['Variables'].append(modelo['Variables'])
12     variables_seleccionadas['Formula'].append(sorted(modelo['X'].columns))
13     variables_seleccionadas['Formula'] = list(map(lambda x: '+'.join(x),
14     variables_seleccionadas['Formula']))
15
16     frecuencias = Counter(variables_seleccionadas['Formula'])
17     frec_ordenada = pd.DataFrame(list(frecuencias.items()), columns = ['Formula',
18     'Frecuencia'])
19     frec_ordenada = frec_ordenada.sort_values('Frecuencia', ascending = False).
20     reset_index()
21
22     var_1 = variables_seleccionadas['Variables'][variables_seleccionadas['Formula']
23     ].index(frec_ordenada['Formula'][0])]
24     var_2 = variables_seleccionadas['Variables'][variables_seleccionadas['Formula']
25     ].index(frec_ordenada['Formula'][1])]
26     var_3 = variables_seleccionadas['Variables'][variables_seleccionadas['Formula']
27     ].index(frec_ordenada['Formula'][2])]

```

Listing 19: Modelo de regresión logística. Selección de variables aleatoria.

De nuevo se ejecuta una validación cruzada del modelo ganador por selección clásica y los modelos aleatorios más fiables, obteniendo el siguiente gráfico de cajas y la tabla de coeficientes estadísticos:

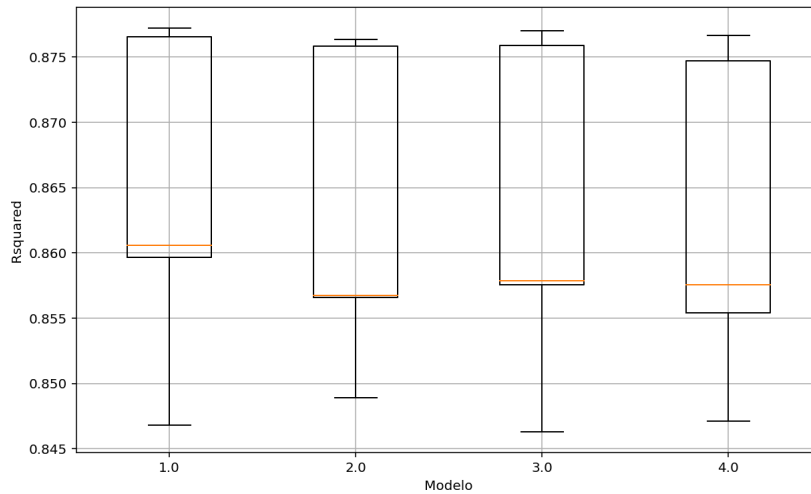


Figura 17: Validación cruzada para modelos de regresión logística aleatorios.

Modelo	$\overline{R^2}$	$\sigma(R^2)$	Parámetros
B-BIC	0.8625	0.0116	75
Aleatorio 1	0.8625	0.0116	74
Aleatorio 1	0.8642	0.0115	75
Aleatorio 3	0.8642	0.0115	74

Tabla 6: Valores de promedio y desviación típica de R^2 y número de parámetros para distintos modelos

Por lo que se elige como modelo ganador al **aleatorio 3**.


```

['Contrastes':
  Variable Estimate ... p value signif
0 (Intercept) 0.974633 ... nan
1 Explotaciones -0.000955 ... 0.000008 ***
2 xAgricultureUnemploymentPtge 0.190752 ... 0.000002 ***
3 PobChange_pct -0.014227 ... 0.000666 ***
4 xForeignersPtge -0.277276 ... 0.000000 ***
.. ...
69 CCAA_Madrid -0.668397 ... nan
70 CCAA_Murcia -0.688903 ... nan
71 CCAA_Navarra 0.171877 ... nan
72 CCAA_PaisVasco -0.648264 ... nan
73 CCAA_Rioja -0.794362 ... nan

[74 rows x 6 columns],
'BondadAjuste': LLK AIC BIC
0 -2297.09114 4598.18228 4611.73924}

```

Figura 18: Resumen del modelo de regresión logística ganador.

5.3. Modelo de regresión logística. Conclusiones.

El resumen del modelo ganador es el siguiente, donde se puede ver el efecto de todas las variables independientes sobre la variable objetivo binaria **Izquierda**.

6. Anexo.

En esta sección se incluyen tablas, gráficos o celdas de código adicionales que pueden servir para un mejor entendimiento de la práctica.

Variable	Izda_Pct	Izquierda
CodigoProvincia	0.1253	0.1609
CCAA	0.4876	0.5614
Population	0.0897	0.1322
TotalCensus	0.0946	0.1426
Age_0-4_Ptge	0.1107	0.1308
Age_under19_Ptge	0.1101	0.1401
Age_19_65_pct	0.0855	0.1227
Age_over65_pct	0.1086	0.1435
WomanPopulationPtge	0.0840	0.1227
ForeignersPtge	0.0737	0.1017
SameComAutonPtge	0.0886	0.0845
SameComAutonDiffProvPtge	0.1280	0.1060
DifComAutonPtge	0.0918	0.0308
UnemployLess25_Ptge	0.2037	0.2688
Unemploy25_40_Ptge	0.0883	0.1282
UnemployMore40_Ptge	0.1238	0.1855
AgricultureUnemploymentPtge	0.1448	0.1998
IndustryUnemploymentPtge	0.1118	0.1141
ConstructionUnemploymentPtge	0.0676	0.0649
ServicesUnemploymentPtge	0.0633	0.0842
totalEmpresas	0.0910	0.1158
Industria	0.0679	0.0134
Construccion	0.0829	0.0250
ComercTTEHosteleria	0.0625	0.0175
Servicios	0.0716	0.0085
ActividadPpal	0.1200	0.1759
inmuebles	0.1047	0.1471
Pob2010	0.0948	0.1346
SUPERFICIE	0.1178	0.1604
Densidad	0.0666	0.0320
PobChange_pct	0.0777	0.1149
PersonasInmueble	0.0773	0.0927
Explotaciones	0.1165	0.1761
prop_missings	0.0552	0.0608

Tabla 7: Coeficientes V de Cramer para las variables objetivo.