

## Actividad 3.4

Representación plot de datasets,  
selección de características y  
entrenamiento de modelos



# ÍNDICE

<b>1. Descripción del Dataset.....</b>	<b>3</b>
Origen del Dataset.....	3
Propósito del Dataset.....	3
Análisis Inicial de las Características.....	3
Instalación del paquete.....	4
Configuración inicial.....	4
<b>2. Procesamiento de Datos.....</b>	<b>5</b>
Manejo de Valores Faltantes.....	5
Normalización de Datos.....	6
<b>3. Selección de Características.....</b>	<b>6</b>
3.1. Matriz de Correlación.....	6
Resultado.....	7
3.2. Matriz de Gráficos de Dispersión.....	7
Resultado.....	8
3.3. Selección con SelectKBest (Modelo entrenado).....	9
4. Reflexión sobre la elección de las características elegidas.....	9
<b>5. NaiveBayes - Sin utilizar Cross Validation y con Cross Validation.....</b>	<b>10</b>
5.1 Sin Cross Validation.....	10
5.2 Con Cross Validation.....	11
<b>6. Conclusión resultados obtenidos.....</b>	<b>11</b>
<b>7. Uso de Herramientas Adicionales.....</b>	<b>12</b>
Mutual Information para Selección de Características.....	12
Entrenar y evaluar el modelo con las características seleccionadas.....	13
Análisis de los resultados.....	13
<b>Referencias.....</b>	<b>13</b>
<b>Repositorios.....</b>	<b>13</b>

# 1. Descripción del Dataset

## Origen del Dataset

El conjunto de datos "Communities and Crime" proviene de un análisis realizado sobre comunidades en los Estados Unidos. Combina datos socioeconómicos del censo de 1990, datos policiales de la encuesta LEMAS de 1990, y datos de crímenes del informe UCR del FBI de 1995.

- Número de instancias: 1.994.
- Número de atributos: 128.
- Tareas asociadas: Regresión.
- Valores faltantes: Sí.
- Área de estudio: Social.

Descargar dataset: <https://archive.ics.uci.edu/dataset/183/communities+and+crime>

## Propósito del Dataset

El dataset se utiliza principalmente para modelar problemas relacionados con la tasa de criminalidad y explorar cómo diversos factores socioeconómicos y características comunitarias afectan los niveles de delincuencia.

## Análisis Inicial de las Características

El archivo de datos contiene valores numéricos y categóricos, separados por comas. Algunos campos están vacíos, lo que indica valores faltantes. Las características incluyen:

- Identificadores: Campos que identifican comunidades o regiones (e.g., nombres de ciudades, códigos numéricos).
- Factores demográficos: Proporciones relacionadas con la raza, la educación, y la pobreza.
- Indicadores económicos: Tasas de desempleo, ingresos, etc.
- Factores policiales: Información sobre presencia y desempeño de fuerzas policiales.
- Tasa de crímenes: Valores que pueden actuar como objetivos de predicción.



## Instalación del paquete

```
[1] 1 # IVAN FALCON MONZON
    2 # Instalación del paquete necesario para descargar y cargar el dataset
    3 !pip install ucimlrepo

Requirement already satisfied: ucimlrepo in /usr/local/lib/python3.11/dist-packages (0.0.7)
Requirement already satisfied: pandas>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from ucimlrepo) (2.2.2)
Requirement already satisfied: certifi>=2020.12.5 in /usr/local/lib/python3.11/dist-packages (from ucimlrepo) (2024.12.14)
Requirement already satisfied: numpy>=1.23.2 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.0.0->ucimlrepo) (1.26.4)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.0.0->ucimlrepo) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.0.0->ucimlrepo) (2024.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas>=1.0.0->ucimlrepo) (2025.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas>=1.0.0->ucimlrepo) (1.17.0)
```

## Configuración inicial

```
1 # IVAN FALCON MONZON
2 # Cargar librerías necesarias
3 from ucimlrepo import fetch_ucirepo
4 import pandas as pd
5
6 # Obtener el dataset Communities and Crime
7 dataset = fetch_ucirepo(id=183) # ID del dataset en UCI Machine Learning Repository
8
9 # Convertir a DataFrame
10 X = pd.DataFrame(dataset.data.features, columns=dataset.data.headers[:-1])
11
12 # Convertir 'targets' a un array de numpy y aplanarlo
13 y = pd.Series(dataset.data.targets.to_numpy().flatten(), name=dataset.data.headers[-1])
14
15 # Mostrar primeras filas del dataset
16 print(X.head())
17 print(y.head())
```

	state	county	community	communityname	fold	population \
0	8	?	?	Lakewoodcity	1	0.19
1	53	?	?	Tukwilacity	1	0.00
2	24	?	?	Aberdeentown	1	0.00
3	34	5	81440	Willingborotownship	1	0.04
4	42	95	6096	Bethlehemtownship	1	0.01

	householdsize	racepctblack	racePctWhite	racePctAsian	...	\
0	0.33	0.02	0.90	0.12	...	
1	0.16	0.12	0.74	0.45	...	
2	0.42	0.49	0.56	0.17	...	
3	0.77	1.00	0.08	0.12	...	
4	0.55	0.02	0.95	0.09	...	

	PolicAveOTWorked	LandArea	PopDens	PctUsePubTrans	PolicCars	\
0	0.29	0.12	0.26	0.20	0.06	
1	?	0.02	0.12	0.45	?	
2	?	0.01	0.21	0.02	?	
3	?	0.02	0.39	0.28	?	
4	?	0.04	0.09	0.02	?	

## 2. Procesamiento de Datos

```
✓ 1 s [3] 1 # Librerías necesarias
2 import pandas as pd
3 import numpy as np
4 from sklearn.feature_selection import SelectKBest, f_classif
5 from sklearn.preprocessing import LabelEncoder, MinMaxScaler
6 from sklearn.impute import SimpleImputer
7 from sklearn.model_selection import train_test_split
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.metrics import accuracy_score
10 from ucimlrepo import fetch_ucirepo
```

### Manejo de Valores Faltantes

```
✓ 0 s Manejo de Valores Faltantes

1 # IVAN FALCON MONZON
2 # Manejo de valores faltantes
3 X.replace("?", np.nan, inplace=True) # Reemplaza '?' por NaN
4 X = X.apply(pd.to_numeric, errors='coerce') # Convierte todo a numérico
5
6 # Eliminar columnas con todos los valores NaN
7 X.dropna(axis=1, how='all', inplace=True)
8
9 # Imputar valores NaN con la mediana
10 imputer = SimpleImputer(strategy='median')
11 X = pd.DataFrame(imputer.fit_transform(X), columns=X.columns)
12
13 # Verificar si aún hay valores NaN después de la imputación
14 if X.isnull().sum().sum() > 0:
15     print("Advertencia: Aún quedan valores NaN en X después de la imputación.")
16     X.dropna(inplace=True) # Eliminar filas con NaN restantes como último recurso
17     y = y.loc[X.index] # Asegurar que y coincida con X
18
19 # Asegurar que X e y tienen las mismas filas
20 y = y.loc[X.index]
21
22 print("Valores faltantes tratados correctamente.")

⇒ Valores faltantes tratados correctamente.
```

## Normalización de Datos

### Normalización de Datos

```
[5] 1 # Normalización de datos
    2 scaler = MinMaxScaler()
    3 X_scaled = scaler.fit_transform(X)
    4 X = pd.DataFrame(X_scaled, columns=X.columns)
```

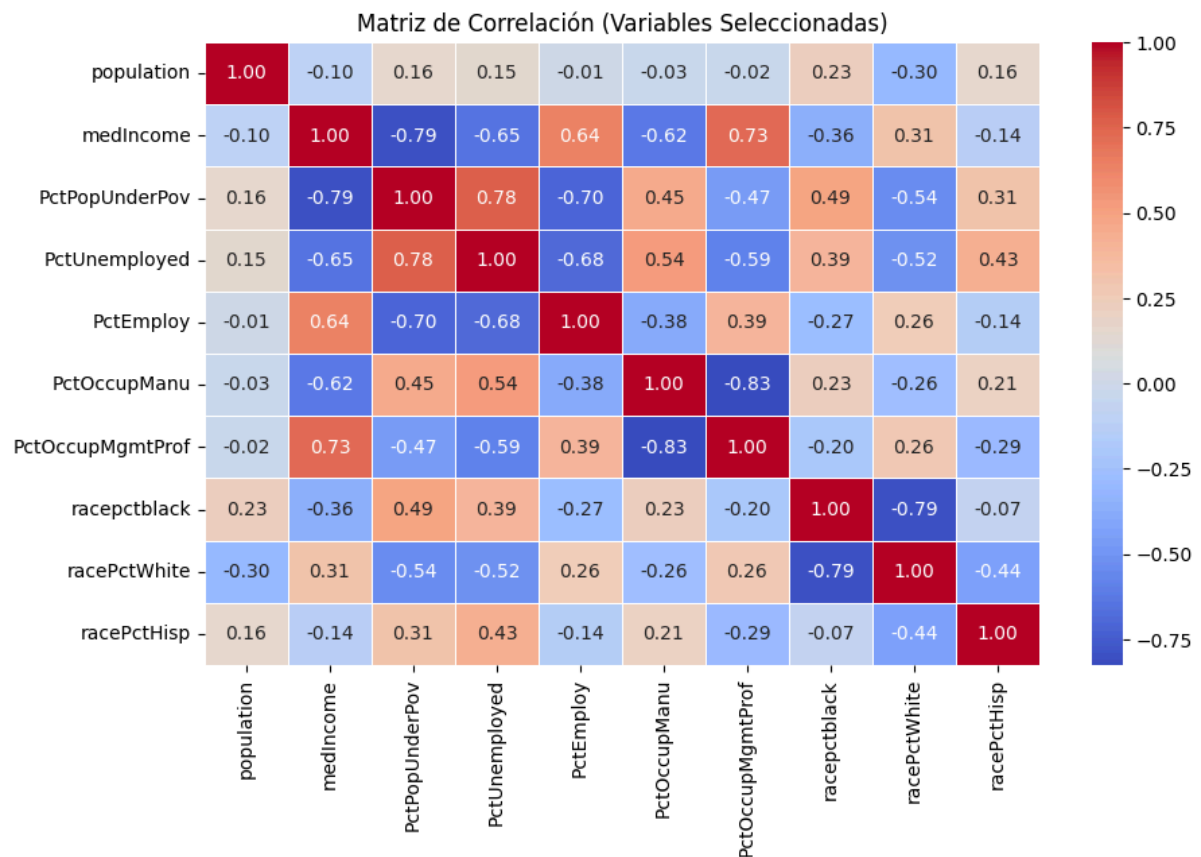
## 3. Selección de Características

### 3.1. Matriz de Correlación

#### 3.1. Matriz de Correlación

```
1 # IVAN FALCON MONZON
2 # Importamos las librerías necesarias
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5
6 # Seleccionamos un subconjunto de características relevantes para reducir la cantidad de datos en la matriz
7 variables_seleccionadas = [
8     "population", "medIncome", "PctPopUnderPov", "PctUnemployed",
9     "PctEmploy", "PctOccupManu", "PctOccupMgmtProf", "racepctblack",
10    "racePctWhite", "racePctHisp"
11 ]
12
13 # Verificamos que las variables seleccionadas existen en X
14 X_sub = X[variables_seleccionadas].copy()
15
16 # Configuramos la figura para la matriz de correlación
17 plt.figure(figsize=(10, 6)) # Ajustamos el tamaño para mejor visualización
18
19 # Generamos el mapa de calor de correlaciones con las variables seleccionadas
20 sns.heatmap(
21     X_sub.corr(), # Calculamos la matriz de correlación de las variables seleccionadas
22     annot=True, # Mostramos los valores de correlación en el gráfico
23     cmap="coolwarm", # Esquema de colores para mejor interpretación
24     fmt=".2f", # Mostramos valores con 2 decimales
25     linewidths=0.5 # Ajustamos el espacio entre celdas
26 )
27
28 # Agregamos el título al gráfico
29 plt.title("Matriz de Correlación (Variables Seleccionadas)")
30
31 # Mostramos el gráfico
32 plt.show()
```

## Resultado



## 3.2. Matriz de Gráficos de Dispersión

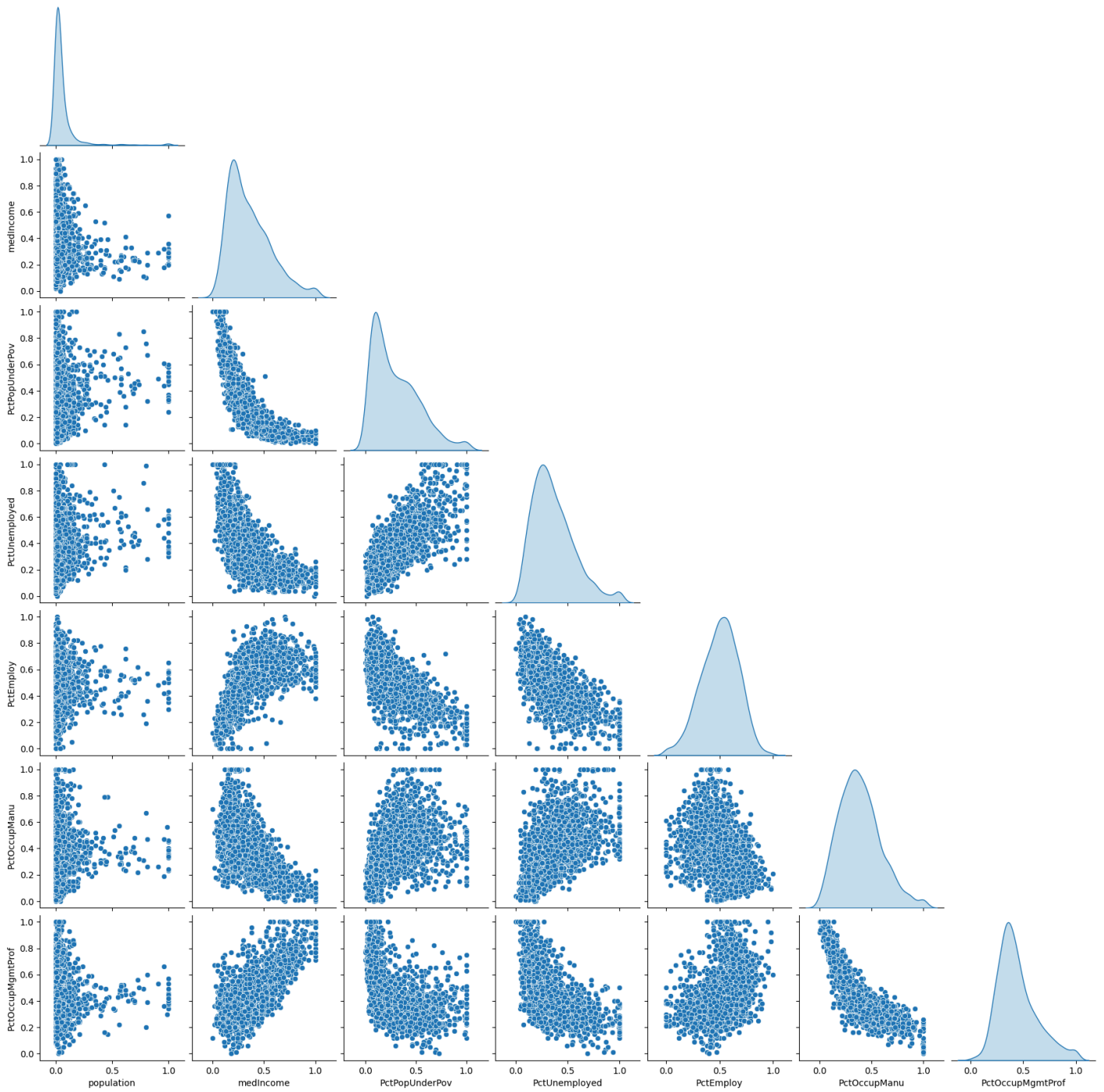
```

3.2. Matriz de Gráficos de Dispersión

1 # IVAN FALCON MONZON
2 # Importamos las librerías necesarias
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5
6 # Seleccionamos un subconjunto de variables para reducir el tiempo de ejecución
7 variables_seleccionadas = [
8     "population", "medIncome", "PctPopUnderPov", "PctUnemployed",
9     "PctEmploy", "PctOccupManu", "PctOccupMgmtProf"
10 ]
11
12 # Creamos un nuevo DataFrame solo con las variables seleccionadas
13 X_sub = X[variables_seleccionadas].copy()
14
15 # Generamos el gráfico de dispersión con las variables seleccionadas
16 sns.pairplot(X_sub, comar=True, diag_kind="kde") # Usamos densidad en la diagonal para mejor visualización
17
18 # Mostramos el gráfico
19 plt.show()

```

## Resultado





### 3.3. Selección con SelectKBest (Modelo entrenado)

```
3.3. Selección con SelectKBest (Modelo entrenado)

1 # Aplicar SelectKBest
2 k = min(5, X.shape[1]) # Asegurar que k no sea mayor que el número de características
3 selector = SelectKBest(score_func=f_classif, k=k)
4 X_new = selector.fit_transform(X, y)
5 selected_features = X.columns[selector.get_support()]
6 X_selected = pd.DataFrame(X_new, columns=selected_features)
7
8 # Clasificación binaria de y
9 y_clas = (y > y.median()).astype(int)
10
11 # División en conjuntos de entrenamiento y prueba
12 X_train, X_test, y_train, y_test = train_test_split(X_selected, y_clas, test_size=0.2, random_state=42)
13
14 # Entrenar el modelo
15 model = LogisticRegression()
16 model.fit(X_train, y_train)
17 y_pred = model.predict(X_test)
18
19 # Evaluación
20 print("Accuracy:", accuracy_score(y_test, y_pred))

Accuracy: 0.8320802005012531
```

## 4. Reflexión sobre la elección de las características elegidas

### Matriz de Correlación:

- Proporciona una visión clara de las relaciones lineales entre las variables seleccionadas. Las variables socioeconómicas, como el ingreso medio y el porcentaje de desempleo, muestran relaciones fuertes entre sí, lo que sugiere que estas características están vinculadas a fenómenos socioeconómicos comunes, como el nivel de pobreza y las tasas de criminalidad.

### Matriz de Gráficos de Dispersión (Pairplot):

- El pairplot revela visualmente las interacciones entre las variables seleccionadas. Permite identificar relaciones no lineales y patrones complejos que no son tan evidentes en la matriz de correlación. Este gráfico es útil para detectar posibles agrupamientos o tendencias que podrían no estar capturados por simples estadísticas de correlación.

### Selección con SelectKBest:

- Identifica las características más relevantes basadas en su relación estadística con la variable objetivo. Al seleccionar las 5 características más significativas, mejora la eficiencia del modelo al reducir el número de variables y enfocar el modelo en las características que realmente importan para la predicción, optimizando el rendimiento y evitando el sobreajuste.

## 5. NaiveBayes – Sin utilizar Cross Validation y con Cross Validation.

En este apartado, he puesto un nuevo hiperparámetro para comprobar su funcionamiento:

`var_smoothing` es un hiperparámetro de Gaussian Naive Bayes que controla la cantidad de suavizado aplicado a la estimación de la varianza de cada característica. Su propósito es evitar problemas numéricos cuando la varianza estimada es muy pequeña o cercana a cero, lo que podría hacer que las probabilidades calculadas sean inestables.

### ¿Cómo funciona?

-Pequeños valores de `var_smoothing` → La varianza se estima casi sin modificación, pero puede generar divisiones por valores muy cercanos a cero, causando problemas numéricos.

-Valores mayores de `var_smoothing` → Se añade una pequeña constante a la varianza de cada característica, suavizando la probabilidad y evitando sobreajustes.

### 5.1 Sin Cross Validation

```
5.1 Sin Cross Validation

1 # IVAN FALCON MONZON
2 # Librerías necesarias
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import accuracy_score
6 import numpy as np
7
8 # Convertir la variable objetivo en clases binarias
9 threshold = np.median(y)
10 y_categorico = (y > threshold).astype(int)
11
12 # Dividir los datos en entrenamiento y prueba
13 X_train, X_test, y_train, y_test = train_test_split(X, y_categorico, test_size=0.2, random_state=42)
14
15 # Ajuste manual de var_smoothing
16 best_accuracy = 0
17 best_smoothing = None
18
19 for smoothing in np.logspace(-10, -3, 20): # Se prueban valores en un rango logarítmico
20     model = GaussianNB(var_smoothing=smoothing)
21     model.fit(X_train, y_train)
22     y_pred = model.predict(X_test)
23     accuracy = accuracy_score(y_test, y_pred)
24
25     if accuracy > best_accuracy:
26         best_accuracy = accuracy
27         best_smoothing = smoothing
28
29 print(f"Mejor precisión con ajuste manual: {best_accuracy} usando var_smoothing={best_smoothing}")

Mejor precisión con ajuste manual: 0.7694235588972431 usando var_smoothing=0.001
```

## 5.2 Con Cross Validation

```
5.2 Con Cross Validation

1 # IVAN FALCON MONZON
2 # Librerías necesarias
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.model_selection import cross_val_score
5 import numpy as np
6
7 # Convertir la variable objetivo en clases binarias
8 threshold = np.median(y)
9 y_categorico = (y > threshold).astype(int)
10
11 # Lista de valores de var_smoothing a probar
12 smoothing_values = np.logspace(-10, -3, 20)
13
14 best_accuracy = 0
15 best_smoothing = None
16
17 for smoothing in smoothing_values:
18     model = GaussianNB(var_smoothing=smoothing)
19     scores = cross_val_score(model, X, y_categorico, cv=5)
20     mean_score = scores.mean()
21
22     if mean_score > best_accuracy:
23         best_accuracy = mean_score
24         best_smoothing = smoothing
25
26 print("Mejor precisión media con Cross Validation: {best_accuracy} usando var_smoothing={best_smoothing}")

Mejor precisión media con Cross Validation: 0.7507399151144192 usando var_smoothing=0.001
```

## 6. Conclusión resultados obtenidos

Los resultados obtenidos muestran una leve diferencia entre entrenar el modelo con una única partición de datos (sin validación cruzada) y utilizar validación cruzada para evaluar su desempeño en múltiples particiones.

- Sin validación cruzada, el modelo alcanzó una precisión de 0.7694 tras ajustar `var_smoothing`, lo que representa una mejora respecto al valor inicial de 0.7468.
- Con validación cruzada, la precisión media mejoró de 0.7422 a 0.7507. Aunque esta precisión es ligeramente menor que la obtenida sin validación cruzada, proporciona una estimación más estable y confiable, ya que el modelo fue evaluado en distintas particiones del conjunto de datos.

### Conclusiones clave:

- Si queremos la mejor precisión absoluta en un solo entrenamiento es mejor entrenar sin validación cruzada, pero con el riesgo de sobreajuste (el modelo puede adaptarse demasiado a una división específica de los datos).
- Si queremos un mejor rendimiento real del modelo y evitar sobreajustes, la validación cruzada es la mejor opción, ya que permite medir la variabilidad del modelo en diferentes subconjuntos de datos.

## 7. Uso de Herramientas Adicionales

### Mutual Information para Selección de Características

- Este método mide cuánta información aporta cada característica para predecir la variable objetivo.

```
7. Uso de Herramientas Adicionales

Mutual Information para Selección de Características

• Este método mide cuánta información aporta cada característica para predecir la variable objetivo.

[11] 1 # IVAN FALCON MONZON
      2 # Librerías necesarias
      3 from sklearn.feature_selection import mutual_info_classif
      4 import pandas as pd
      5
      6 # Calcular la información mutua entre las características y la variable objetivo
      7 mi_scores = mutual_info_classif(X, y_categorico) # Usar la variable ya binarizada
      8
      9 # Convertir los resultados en un DataFrame para análisis
     10 mi_df = pd.DataFrame({'Característica': X.columns, 'MI Score': mi_scores})
     11
     12 # Ordenar de mayor a menor relevancia
     13 mi_df = mi_df.sort_values(by='MI Score', ascending=False)
     14
     15 # Seleccionar las 5 mejores características
     16 top_features = mi_df.head(5)['Característica'].values
     17 print("Mejores características según Información Mutua:", top_features)
     18
     19 # Filtrar el dataset con las mejores características
     20 X_selected = X[top_features]
```

Mejores características según Información Mutua: ['PctKids2Par' 'NumIlleg' 'PctIlleg' 'PctFam2Par' 'PctYoungKids2Par']

### ¿Por qué es útil este método?

- No depende de un modelo específico, solo analiza la relación entre las características y la variable objetivo.
- Funciona bien con Naive Bayes, ya que el algoritmo asume independencia entre variables y Mutual Information mide cuánta información aporta cada una.
- Evita errores de compatibilidad como los que surgen con RFE (método que utilice al principio) y Naive Bayes.

## Entrenar y evaluar el modelo con las características seleccionadas

```
Entrenar y evaluar el modelo con las características seleccionadas

1 # IVAN FALCON MONZON
2 # Librerías necesarias
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.model_selection import cross_val_score
5
6 # Crear el modelo Naive Bayes
7 model = GaussianNB(var_smoothing=0.001) # Usamos el mejor hiperparámetro encontrado antes
8
9 # Aplicar validación cruzada con las características seleccionadas
10 scores = cross_val_score(model, X_selected, y_categorico, cv=5)
11
12 # Mostrar resultados
13 print("Precisión en cada iteración de Cross Validation:", scores)
14 print("Precisión media con Cross Validation (nuevas características):", scores.mean())
```

Precisión en cada iteración de Cross Validation: [0.77192982 0.77694236 0.79197995 0.78947368 0.77386935]  
Precisión media con Cross Validation (nuevas características): 0.7808390322540018

Antes, con todas las características, la mejor precisión media con Cross Validation era 0.7507. Ahora, después de la selección de características con Información Mutua, la precisión media subió a 0.7808, lo que indica que el modelo está generalizando mejor.

### Análisis de los resultados

- Eliminación de ruido: Seleccionar solo las características más relevantes ayuda a reducir el impacto de variables irrelevantes o redundantes.
- Mejor generalización: La precisión se volvió más estable en diferentes particiones, evitando que el modelo se sobreajuste a características poco útiles.
- Eficiencia mejorada: Ahora el modelo entrena con menos variables, lo que mejora la velocidad y simplicidad sin perder precisión

### Referencias

[https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.GaussianNB.html](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html)

### Repositorios

-Google Colab: <https://colab.research.google.com/drive/1IRyGEFcQ2fCRYhPw3zjgYg4eMmcNQkHx?usp=sharing>

-Github: [https://github.com/IvanFalconMonzon/SNS\\_ACT3\\_4\\_IvanFalconMonzon.git](https://github.com/IvanFalconMonzon/SNS_ACT3_4_IvanFalconMonzon.git)