



UNIVERSIDAD POLITÉCNICA METROPOLITANA DE HIDALGO

PROGRAMA EDUCATIVO DE MAESTRÍA EN INTELIGENCIA ARTIFICIAL

“ESTADÍSTICA APLICADA”

Reporte: Tarea 5

*Flores García Katherine
Itzel
(253220100)
Morales Hernández
Emmanuel
(253220003)*

26 de octubre de 2025

Índice

| | |
|---|-----------|
| 1. OBJETIVO | 4 |
| 2. INTRODUCCIÓN | 4 |
| 3. Desarrollo | 5 |
| 3.1. PARTE I | 5 |
| 3.1.1. Análisis del Dataset | 6 |
| 3.1.2. Separación de variables | 7 |
| 3.1.3. Limpieza de valores nulos | 7 |
| 3.1.4. Modelo de imputación | 8 |
| 3.1.5. Revisión del Dataset imputado | 9 |
| 3.1.6. Codificación | 9 |
| 3.1.7. Estandarización | 10 |
| 3.2. PARTE II | 11 |
| 3.2.1. Creación del modelo de regresión logística | 11 |
| 3.2.2. Predicción | 11 |
| 3.2.3. Coeficientes e Intercepto del Modelo | 11 |
| 3.2.4. Ecuación como String | 12 |
| 3.2.5. Interpretación de coeficientes | 13 |
| 3.3. PARTE III | 14 |
| 3.3.1. Extraer métricas | 15 |
| 3.3.2. Pseudo R-cuadrada (McFadden) | 15 |
| 3.3.3. Log-Likelihood | 16 |
| 3.3.4. LL-Null (Log-Likelihood Nulo) | 17 |
| 3.3.5. LLR p-value (Prueba de Razón de Verosimilitudes) | 18 |
| 3.4. PARTE IV | 20 |
| 3.4.1. Interpretación de coeficientes | 20 |
| 3.4.2. Interpretación de std err | 23 |
| 3.4.3. Variables significativas | 24 |
| 3.4.4. Variables con mayor poder explicativo (valor z) | 25 |
| 3.4.5. Intervalo de confianza del 95 % para la mejor variable | 26 |
| 3.4.6. Visualización de los Intervalos de Confianza | 27 |
| 3.5. PARTE V | 28 |
| 3.5.1. Métricas derivadas | 30 |
| 3.5.2. Mejor modelo con máximo 4 variables (criterio LLR p-value) | 30 |
| 3.5.3. Entrenar modelo reducido | 31 |
| 3.5.4. Creación de gráfica logística y diagrama de dispersión | 32 |
| 3.5.5. Precisión del Modelo | 33 |
| 3.6. PARTE VI | 35 |
| 3.6.1. Entrenamiento del modelo final con estas 4 variables | 37 |
| 4. CONCLUSIÓN | 39 |

Referencias**40****Índice de figuras**

| | | |
|-----|---|----|
| 1. | Análisis de la distribución del Dataset | 7 |
| 2. | Resultados del modelo de imputación | 9 |
| 3. | Resultados de codificación | 10 |
| 4. | Resultados de estandarización | 10 |
| 5. | Resultados de la regresión logística | 15 |
| 6. | Resultados del DataFrame | 20 |
| 7. | Coeficientes e intervalos de confianza al 95 % para las principales variables del modelo. | 28 |
| 8. | Matriz de confusión | 29 |
| 9. | Gráfica logística y diagrama de dispersión | 33 |
| 10. | Reporte Completo de Clasificación | 34 |
| 11. | Resultados de la regresión | 38 |

Índice de cuadros

| | | |
|----|--|----|
| 1. | Primeras filas del Heart Failure Prediction Dataset. | 5 |
| 2. | Resumen de las columnas, conteo de valores no nulos y tipos de datos del Heart Failure Prediction Dataset. | 6 |
| 3. | Resultados de la limpieza de datos | 8 |
| 4. | Resultados de la predicción | 11 |

1. OBJETIVO

Objetivo general

Construir y evaluar un modelo de **regresión logística** utilizando un dataset de Kaggle para identificar las variables que mejor explican una variable binaria dependiente y analizar la calidad del modelo.

Objetivos específicos

- Analizar y describir el dataset seleccionado, incluyendo todas las variables numéricas y la variable dependiente.
- Estimar la ecuación de regresión logística con todas las variables y analizar la interpretación de los coeficientes.
- Evaluar el modelo mediante métricas como pseudo R-cuadrada, Log-Likelihood, LL-Null y LLR p-value.
- Interpretar los resultados de los coeficientes, errores estándar, valores z y significancia estadística de las variables independientes.
- Construir la matriz de confusión, graficar la función logística junto con los datos y determinar la precisión del modelo.
- Seleccionar y construir un modelo optimizado con un máximo de cuatro variables independientes según criterios estadísticos.

2. INTRODUCCIÓN

En la actualidad, el análisis de datos se ha convertido en una herramienta esencial para la toma de decisiones en diversos ámbitos, desde la investigación científica hasta la industria y los negocios. La disponibilidad de plataformas como Kaggle permite a los equipos acceder a una amplia variedad de datasets que pueden ser utilizados para entrenar modelos predictivos y realizar análisis estadísticos avanzados.

En este trabajo, se propone utilizar un dataset seleccionado de Kaggle que contenga al menos diez variables numéricas y una variable binaria como dependiente. Se busca aplicar la regresión logística para comprender cómo las variables independientes influyen en la variable dependiente y evaluar la eficacia del modelo mediante métricas estadísticas y gráficas.

El código completo utilizado en este análisis puede visualizarse en el siguiente [link a Colab](#), donde se incluyen todos los pasos desde la carga del dataset hasta la construcción y evaluación del modelo de regresión logística.

Este análisis permitirá identificar las variables más relevantes y medir el desempeño predictivo del modelo, proporcionando información útil para la interpretación de los resultados.

La tabla a continuación muestra ejemplos representativos de los registros incluidos en el dataset.

| Age | Sex | ChestPainType | RestingBP | Cholesterol | FastingBS | RestingECG | ... |
|-----|-----|---------------|-----------|-------------|-----------|------------|-----|
| 40 | M | ATA | 140 | 289 | 0 | Normal | ... |
| 49 | F | NAP | 160 | 180 | 0 | Normal | ... |
| 37 | M | ATA | 130 | 283 | 0 | ST | ... |
| 48 | F | ASY | 138 | 214 | 0 | Normal | ... |
| 54 | M | NAP | 150 | 195 | 0 | Normal | ... |
| 39 | M | NAP | 120 | 339 | 0 | Normal | ... |
| 45 | F | ATA | 130 | 237 | 0 | Normal | ... |
| 54 | M | ATA | 110 | 208 | 0 | Normal | ... |
| ... | ... | ... | ... | ... | ... | ... | ... |

Cuadro 1: Primeras filas del Heart Failure Prediction Dataset.

3. Desarrollo

3.1. PARTE I

El dataset utilizado para esta actividad es el **Heart Failure Prediction Dataset**, obtenido el **18 de octubre de 2025** de la plataforma *Kaggle*, elaborado por **FEDESORIANO**[1]. El conjunto de datos está disponible en el siguiente enlace: [Heart Failure Prediction Dataset](#).

Este dataset está enfocado en enfermedades cardiovasculares, una de las **principales causas de mortalidad a nivel mundial**. Contiene **11 variables** que permiten predecir la presencia de insuficiencia cardíaca, incluyendo factores de riesgo como hipertensión, diabetes, edad, sexo y nivel de colesterol, entre otros.

El análisis de estas variables facilita la identificación de pacientes con **alto riesgo cardiovascular**, lo que permite aplicar estrategias de prevención y tratamiento oportunas. Las variables incluidas en el dataset son las siguientes:

- **Age:** Edad del paciente (numérica, en años).
- **Sex:** Sexo del paciente (categórica, M = masculino, F = femenino).
- **ChestPainType:** Tipo de dolor torácico (categórica):
 - ATA = angina típica,
 - NAP = angina no típica,
 - ASY = asintomático,
 - TA = dolor atípico.
- **RestingBP:** Presión arterial en reposo (numérica, mmHg).
- **Cholesterol:** Nivel de colesterol sérico (numérica, mg/dL).
- **FastingBS:** Glucosa en ayunas (binaria, 1 = >120 mg/dL, 0 = ≤120 mg/dL).
- **RestingECG:** Electrocardiograma en reposo (Normal / ST / LVH).

- **MaxHR:** Frecuencia cardíaca máxima alcanzada (latidos por minuto).
- **ExerciseAngina:** Angina inducida por ejercicio (Y = sí, N = no).
- **Oldpeak:** Depresión del segmento ST inducida por ejercicio.
- **ST_Slope:** Pendiente del segmento ST durante ejercicio (Up / Flat / Down).
- **HeartDisease:** Variable dependiente binaria (0 = no, 1 = sí).

3.1.1. Análisis del Dataset

Este análisis permite conocer rápidamente:

- El número de registros no nulos por columna y el tipo de dato ('df.info()').
- Estadísticas descriptivas básicas de las variables numéricas ('df.describe()').
- Las categorías únicas presentes en cada variable categórica, lo que facilita su comprensión y preparación para modelos predictivos.

```

1 # Resumen de la informacion general del dataset
2 display(df.info())
3
4 # Estadisticas descriptivas de las variables numericas
5 display(df.describe())
6
7 # Identificacion de columnas de tipo objeto (categoricas)
8 col_obj = df.select_dtypes(include='object').columns
9
10 # Mostrar valores unicos de cada columna categorica
11 for column in col_obj:
12     print(f"=== '{column}' ===")
13     display(df[column].unique())

```

| Columna | Non-Null Count | Dtype |
|----------------|----------------|---------|
| Age | 918 | int64 |
| Sex | 918 | object |
| ChestPainType | 918 | object |
| RestingBP | 918 | int64 |
| Cholesterol | 918 | int64 |
| FastingBS | 918 | int64 |
| RestingECG | 918 | object |
| MaxHR | 918 | int64 |
| ExerciseAngina | 918 | object |
| Oldpeak | 918 | float64 |
| ST_Slope | 918 | object |
| HeartDisease | 918 | int64 |

Cuadro 2: Resumen de las columnas, conteo de valores no nulos y tipos de datos del Heart Failure Prediction Dataset.

| | Age | RestingBP | Cholesterol | FastingBS | MaxHR | Oldpeak | HeartDisease |
|-------|------------|------------|-------------|------------|------------|------------|--------------|
| count | 918.000000 | 918.000000 | 918.000000 | 918.000000 | 918.000000 | 918.000000 | 918.000000 |
| mean | 53.510893 | 132.396514 | 198.799564 | 0.233115 | 136.809368 | 0.887364 | 0.553377 |
| std | 9.432617 | 18.514154 | 109.384145 | 0.423046 | 25.460334 | 1.066570 | 0.497414 |
| min | 28.000000 | 0.000000 | 0.000000 | 0.000000 | 60.000000 | -2.600000 | 0.000000 |
| 25% | 47.000000 | 120.000000 | 173.250000 | 0.000000 | 120.000000 | 0.000000 | 0.000000 |
| 50% | 54.000000 | 130.000000 | 223.000000 | 0.000000 | 138.000000 | 0.600000 | 1.000000 |
| 75% | 60.000000 | 140.000000 | 267.000000 | 0.000000 | 156.000000 | 1.500000 | 1.000000 |
| max | 77.000000 | 200.000000 | 603.000000 | 1.000000 | 202.000000 | 6.200000 | 1.000000 |

Figura 1: Análisis de la distribución del Dataset

3.1.2. Separación de variables

En esta etapa se realiza la separación entre las variables independientes (X) y la variable dependiente (y), preparándolas para el análisis de regresión logística:

```

1 # Separacion de variables independientes y dependiente
2 X = df.drop('HeartDisease', axis=1) # Variables independientes
3 y = df['HeartDisease']             # Variable dependiente
4
5 # Identificacion de columnas numericas
6 col_num = X.select_dtypes(include='number').columns

```

Este paso permite:

- Definir claramente cuál es la variable a predecir (HeartDisease).
- Seleccionar las variables que se utilizarán como predictores.
- Identificar las columnas numéricas para aplicar análisis estadísticos o transformaciones específicas.

3.1.3. Limpieza de valores nulos

En esta etapa se realiza la separación entre las variables independientes (X) y la variable dependiente (y), preparándolas para el análisis de regresión logística:

```

1 # Separacion de variables independientes y dependiente
2 X = df.drop('HeartDisease', axis=1) # Variables independientes
3 y = df['HeartDisease']             # Variable dependiente
4
5 # Identificacion de columnas numericas
6 col_num = X.select_dtypes(include='number').columns

```

Este paso permite:

- Definir claramente cuál es la variable a predecir (HeartDisease).
- Seleccionar las variables que se utilizarán como predictores.
- Identificar las columnas numéricas para aplicar análisis estadísticos o transformaciones específicas.

| | |
|-------------|-----|
| | 0 |
| RestingBP | 1 |
| Cholesterol | 172 |

Cuadro 3: Resultados de la limpieza de datos

3.1.4. Modelo de imputación

Para manejar valores faltantes en el dataset, se utiliza un modelo de **imputación iterativa** (IterativeImputer) con un RandomForestRegressor como estimador. Este método permite predecir los valores faltantes basándose en las demás variables del dataset.

```

1      # Configurar IterativeImputer con RandomForestRegressor
2  imputer = IterativeImputer(
3      estimator=RandomForestRegressor(
4          n_estimators=100,          # Numero de arboles
5          max_depth=10,             # Profundidad para evitar overfitting
6          random_state=42,
7          n_jobs=-1                 # Usar todos los cores
8      ),
9      max_iter=10,                  # Iteraciones maximas
10     random_state=42,
11     initial_strategy='median',    # Inicializacion con mediana
12     imputation_order='ascending', # Menos faltantes primero
13     verbose=1
14 )

```

```

1      # Realizar la imputacion
2  X[columns_with_zeros] = imputer.fit_transform(X[columns_with_zeros])
3  # Redondear RestingBP y Cholesterol a enteros (variables discretas)
4  X[columns_with_zeros] = X[columns_with_zeros].round().astype(int)

```

Durante la ejecución del IterativeImputer con RandomForestRegressor, se obtuvieron los siguientes mensajes que indican la convergencia del algoritmo y los cambios en cada iteración:

```

1  [IterativeImputer] Completing matrix with shape (918, 2)
2  [IterativeImputer] Change: 79.09199999999998, scaled tolerance: 0.603
3  [IterativeImputer] Change: 98.03282084774568, scaled tolerance: 0.603
4  [IterativeImputer] Change: 32.56462456117043, scaled tolerance: 0.603
5  [IterativeImputer] Change: 8.521020423501142, scaled tolerance: 0.603
6  [IterativeImputer] Change: 0.6216106996030817, scaled tolerance: 0.603
7  [IterativeImputer] Change: 8.526512829121202e-14, scaled tolerance: 0.603
8  [IterativeImputer] Early stopping criterion reached.

```


Interpretación:

- Cada línea muestra el cambio promedio en los valores imputados durante cada iteración.
- El algoritmo detiene las iteraciones cuando el cambio se vuelve muy pequeño comparado con la tolerancia escalada ('scaled tolerance'), indicando que la imputación ha convergido.
- Esto asegura que los valores faltantes han sido estimados de manera estable y confiable.

3.1.5. Revisión del Dataset imputado

```

1 # Mostrar informaciOn
2 display(X.isnull().sum())
3 display(X.describe())
4

```

| | Age | RestingBP | Cholesterol | FastingBS | MaxHR | Oldpeak |
|-------|------------|------------|-------------|------------|------------|------------|
| count | 918.000000 | 918.000000 | 918.000000 | 918.000000 | 918.000000 | 918.000000 |
| mean | 53.510893 | 132.538126 | 245.295207 | 0.233115 | 136.809368 | 0.887364 |
| std | 9.432617 | 17.990127 | 54.526364 | 0.423046 | 25.460334 | 1.066570 |
| min | 28.000000 | 80.000000 | 85.000000 | 0.000000 | 60.000000 | -2.600000 |
| 25% | 47.000000 | 120.000000 | 213.000000 | 0.000000 | 120.000000 | 0.000000 |
| 50% | 54.000000 | 130.000000 | 241.000000 | 0.000000 | 138.000000 | 0.600000 |
| 75% | 60.000000 | 140.000000 | 271.750000 | 0.000000 | 156.000000 | 1.500000 |
| max | 77.000000 | 200.000000 | 603.000000 | 1.000000 | 202.000000 | 6.200000 |

Figura 2: Resultados del modelo de imputación

3.1.6. Codificación

Para preparar las variables categóricas para el modelo de regresión logística, se aplicaron técnicas de codificación según el tipo de variable:

- Variables **ordinales** (con más de dos categorías y un orden implícito): ChestPainType, RestingECG, ST_Slope, se codificaron usando `OrdinalEncoder`.
- Variables **nominales** (con dos categorías sin orden): Sex y ExerciseAngina, se codificaron con **One Hot Encoding** o mapeo a 0 y 1.

```

1      """TIPO DE VARIABLES CATEGORICAS
2  'Sex' ---> Nominal (2) ---> One Hot Encoding
3  'ChestPainType' ---> Ordinal (4)
4  'RestingECG' ---> Ordinal (3)
5  'ExerciseAngina' ---> Nominal (2)---> One Hot Encoding
6  'ST_Slope' ---> Ordinal (3)
7  """
8  # Codificacion Ordinal Encoder
9  encoder = OrdinalEncoder()
10 X['ChestPainType'] = encoder.fit_transform(X[['ChestPainType']])
11 X['RestingECG'] = encoder.fit_transform(X[['RestingECG']])
12 X['ST_Slope'] = encoder.fit_transform(X[['ST_Slope']])
13
14 # Codificacion One Hot Encoding
15 X['Sex'] = X['Sex'].map({'M': 1, 'F': 0})
16 X['ExerciseAngina'] = X['ExerciseAngina'].map({'Y': 1, 'N': 0})
17 X.head()

```

| | Age | Sex | ChestPainType | RestingBP | Cholesterol | FastingBS | RestingECG | MaxHR | ExerciseAngina | Oldpeak | ST_Slope |
|---|-----|-----|---------------|-----------|-------------|-----------|------------|-------|----------------|---------|----------|
| 0 | 40 | 1 | 1.0 | 140 | 289 | 0 | 1.0 | 172 | 0 | 0.0 | 2.0 |
| 1 | 49 | 0 | 2.0 | 160 | 180 | 0 | 1.0 | 156 | 0 | 1.0 | 1.0 |
| 2 | 37 | 1 | 1.0 | 130 | 283 | 0 | 2.0 | 98 | 0 | 0.0 | 2.0 |
| 3 | 48 | 0 | 0.0 | 138 | 214 | 0 | 1.0 | 108 | 1 | 1.5 | 1.0 |
| 4 | 54 | 1 | 2.0 | 150 | 195 | 0 | 1.0 | 122 | 0 | 0.0 | 2.0 |

Figura 3: Resultados de codificación

3.1.7. Estandarización

Para asegurar que todas las variables numéricas tengan la misma escala y evitar que algunas dominen el modelo de regresión logística, se aplica una **estandarización** utilizando `StandardScaler` de `scikit-learn`. Esto transforma los valores para que tengan media 0 y desviación estándar 1.

```

1      # Estandarizacion de valores StandardScaler
2  scaler = StandardScaler()
3  X[col_num] = scaler.fit_transform(X[col_num])
4  X.head()

```

| | Age | Sex | ChestPainType | RestingBP | Cholesterol | FastingBS | RestingECG | MaxHR | ExerciseAngina | Oldpeak | ST_Slope |
|---|-----------|-----|---------------|-----------|-------------|-----------|------------|-----------|----------------|-----------|----------|
| 0 | -1.433140 | 1 | 1.0 | 0.415002 | 0.801972 | -0.551341 | 1.0 | 1.382928 | 0 | -0.832432 | 2.0 |
| 1 | -0.478484 | 0 | 2.0 | 1.527329 | -1.198151 | -0.551341 | 1.0 | 0.754157 | 0 | 0.105664 | 1.0 |
| 2 | -1.751359 | 1 | 1.0 | -0.141161 | 0.691874 | -0.551341 | 2.0 | -1.525138 | 0 | -0.832432 | 2.0 |
| 3 | -0.584556 | 0 | 0.0 | 0.303769 | -0.574259 | -0.551341 | 1.0 | -1.132156 | 1 | 0.574711 | 1.0 |
| 4 | 0.051881 | 1 | 2.0 | 0.971166 | -0.922904 | -0.551341 | 1.0 | -0.581981 | 0 | -0.832432 | 2.0 |

Figura 4: Resultados de estandarización

3.2. PARTE II

Antes de construir el modelo, se realiza la **división de los datos** en conjuntos de entrenamiento y prueba, con el fin de evaluar el desempeño del modelo en datos no vistos durante el ajuste:

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42, shuffle = True)
```

3.2.1. Creación del modelo de regresión logística

Una vez preprocesadas y codificadas las variables, se ajusta un modelo de **regresión logística** utilizando el conjunto de entrenamiento[2, 3]:

```
1 from sklearn.linear_model import LogisticRegression
2
3 # Creación del modelo de regresión logística sin regularización
4 modelo = LogisticRegression(penalty=None)
5
6 # Ajuste del modelo con los datos de entrenamiento
7 modelo.fit(X_train, y_train)
```

3.2.2. Predicción

Una vez ajustado el modelo de regresión logística, se realizan predicciones sobre el conjunto de prueba para evaluar su desempeño:

```
1 # Prediccion de la variable dependiente para los datos de prueba
2 prediccion = modelo.predict(X_test)
3
4 # Convertir a DataFrame para visualizar los resultados
5 predicciones = pd.DataFrame(prediccion, columns=['Prediction'])
6
7 # Mostrar las primeras 3 predicciones
8 predicciones.head(3)
```

| | Predicción |
|---|------------|
| 0 | 0 |
| 1 | 1 |
| 2 | 1 |

Cuadro 4: Resultados de la predicción

3.2.3. Coeficientes e Intercepto del Modelo

A continuación se presentan los coeficientes estimados por el modelo de Regresión Logística junto con el intercepto. Estos valores se obtuvieron tras el ajuste del modelo con las variables estandarizadas:

```

1 # Obtener coeficientes e intercepto
2 coeficientes = modelo.coef_[0]
3 intercepto = modelo.intercept_[0]
4 feature_names = X_train.columns.tolist()
5
6 print("\n--- ECUACION DE REGRESION LOGISTICA ---")
7 print(f"\nIntercepto (\beta_0): {intercepto:.6f}")
8 print("\nCoeficientes (\beta_i):")
9 for i, (nombre, coef) in enumerate(zip(feature_names, coeficientes)):
10     print(f"    \beta_{i+1} ({nombre}): {coef:.6f}")

```

Los resultados obtenidos fueron los siguientes:

```

1 --- ECUACION DE REGRESION LOGISTICA ---
2
3 Intercepto (\beta_0): 1.662323
4
5 Coeficientes (\beta_i):
6     \beta_1 (Age): 0.148690
7     \beta_2 (Sex): 1.433789
8     \beta_3 (ChestPainType): -0.590140
9     \beta_4 (RestingBP): 0.045719
10    \beta_5 (Cholesterol): 0.101711
11    \beta_6 (FastingBS): 0.488293
12    \beta_7 (RestingECG): -0.130859
13    \beta_8 (MaxHR): -0.301389
14    \beta_9 (ExerciseAngina): 1.179116
15    \beta_10 (Oldpeak): 0.441164
16    \beta_11 (ST_Slope): -1.712709

```

Interpretación:

- El intercepto ($\beta_0 = 1.662323$) representa el logaritmo de las probabilidades de presentar enfermedad cardiaca cuando todas las variables independientes son cero (en su escala estandarizada).
- Los coeficientes (β_i) indican la dirección y magnitud del efecto de cada variable sobre la probabilidad de tener enfermedad cardiaca.
- Por ejemplo, el coeficiente de Sex ($\beta_2 = 1.433789$) sugiere que ser hombre incrementa las probabilidades de padecer enfermedad cardiaca.
- En cambio, variables como ChestPainType ($\beta_3 = -0.590140$) y ST_Slope ($\beta_{11} = -1.712709$) muestran una relación negativa con la presencia de la enfermedad, indicando que ciertos tipos de dolor de pecho y pendiente del segmento ST están asociados con menor riesgo.

3.2.4. Ecuación como String

Se construye la ecuación de regresión logística como log-odds, utilizando los coeficientes estimados:

```

1  # Construir la ecuacion como string
2  ecuacion_partes = [f"{intercepto:.4f}"]
3  i = 0
4  for nombre, coef in zip(feature_names, coeficientes):
5      i += 1
6      print(f"x{i} ---> {nombre}")
7      signo = "+" if coef >= 0 else "-"
8      ecuacion_partes.append(f"{signo}{coef:.4f}*x{i}")
9
10 ecuacion = " ".join(ecuacion_partes)
11 print(f"\n--- ECUACION COMPLETA (log-odds) ---")
12 print(f"log(p/(1-p)) = {ecuacion}")

```

Resultados:

```

1  x1 ---> Age
2  x2 ---> Sex
3  x3 ---> ChestPainType
4  x4 ---> RestingBP
5  x5 ---> Cholesterol
6  x6 ---> FastingBS
7  x7 ---> RestingECG
8  x8 ---> MaxHR
9  x9 ---> ExerciseAngina
10 x10 ---> Oldpeak
11 x11 ---> ST_Slope
12
13 --- ECUACION COMPLETA (log-odds) ---
14 log(p/(1-p)) = 1.6623 +0.1487*x1 +1.4338*x2 -0.5901*x3 +0.0457*x4 +0.1017*
    x5 +0.4883*x6 -0.1309*x7 -0.3014*x8 +1.1791*x9 +0.4412*x10 -1.7127*x11

```

3.2.5. Interpretación de coeficientes

En esta sección se realiza la interpretación de los coeficientes obtenidos del modelo de regresión logística. Los coeficientes indican el cambio en el **log-odds** de la variable dependiente (HeartDisease) por cada unidad de cambio en la variable independiente, manteniendo constantes las demás variables.

A continuación se presentan los resultados calculados para cada variable independiente.

```

1  print("\n--- INTERPRETACION DE COEFICIENTES ---")
2  print("""
3  Los coeficientes representan el cambio en el LOG-ODDS (logaritmo de las
4  probabilidades)
5  de la variable dependiente por cada unidad de cambio en la variable
6  independiente,
7  manteniendo todas las demás variables constantes.
8
9  - Coeficiente POSITIVO: aumenta la probabilidad del evento (HeartDisease
10  =1)
11 - Coeficiente NEGATIVO: disminuye la probabilidad del evento
12
13 Para interpretar en términos de ODDS RATIO (más intuitivo):

```

```

11     Odds Ratio = exp(coeficiente)
12
13     Si OR > 1: La variable aumenta las probabilidades del evento
14     Si OR < 1: La variable disminuye las probabilidades del evento
15     Si OR = 1: La variable no tiene efecto
16     """)
17
18 # Calcular Odds Ratios
19 print("\n--- ODDS RATIOS ---")
20 for nombre, coef in zip(feature_names, coeficientes):
21     odds_ratio = np.exp(coef)
22     print(f"    {nombre}: OR = {odds_ratio:.4f}")
23     if odds_ratio > 1:
24         print(f"        Incrementa las probabilidades en {(odds_ratio-1)
25             *100:.2f}%")
26     else:
27         print(f"        Reduce las probabilidades en {(1-odds_ratio)*100:.2f
28             }%")

```

3.3. PARTE III

Para este apartado, se ajustó un modelo de *Regresión Logística* utilizando la librería *statsmodels*, que permite obtener información estadística detallada del ajuste, incluyendo métricas como el *Pseudo R-cuadrado*, el *Log-Likelihood*, el *LL-Null* y el *LLR p-value* [5].

El modelo se entrenó agregando una constante (intercepto) a la matriz de predictores:

```

1     # Agregar constante (intercepto) a X_train para statsmodels
2     X_train_sm = sm.add_constant(X_train)
3     X_test_sm = sm.add_constant(X_test)
4
5     # Entrenar modelo con statsmodels para obtener estadísticas completas
6     modelo_statsmodels = sm.Logit(y_train, X_train_sm).fit(dis=0)
7     # Mostrar resumen estadístico del modelo
8     print(modelo_statsmodels.summary())

```

| Logit Regression Results | | | | | | |
|--------------------------|------------------|-------------------|-----------|-------|--------|--------|
| Dep. Variable: | HeartDisease | No. Observations: | 734 | | | |
| Model: | Logit | Df Residuals: | 722 | | | |
| Method: | MLE | Df Model: | 11 | | | |
| Date: | Tue, 21 Oct 2025 | Pseudo R-squ.: | 0.4732 | | | |
| Time: | 05:01:36 | Log-Likelihood: | -266.38 | | | |
| converged: | True | LL-Null: | -505.62 | | | |
| Covariance Type: | nonrobust | LLR p-value: | 1.250e-95 | | | |
| | coef | std err | z | P> z | [0.025 | 0.975] |
| const | 1.6636 | 0.433 | 3.846 | 0.000 | 0.816 | 2.511 |
| Age | 0.1493 | 0.130 | 1.146 | 0.252 | -0.106 | 0.405 |
| Sex | 1.4331 | 0.286 | 5.010 | 0.000 | 0.872 | 1.994 |
| ChestPainType | -0.5901 | 0.121 | -4.864 | 0.000 | -0.828 | -0.352 |
| RestingBP | 0.0460 | 0.121 | 0.381 | 0.703 | -0.191 | 0.283 |
| Cholesterol | 0.1017 | 0.129 | 0.791 | 0.429 | -0.151 | 0.354 |
| FastingBS | 0.4879 | 0.120 | 4.076 | 0.000 | 0.253 | 0.723 |
| RestingECG | -0.1306 | 0.184 | -0.710 | 0.478 | -0.491 | 0.230 |
| MaxHR | -0.3013 | 0.133 | -2.272 | 0.023 | -0.561 | -0.041 |
| ExerciseAngina | 1.1790 | 0.255 | 4.615 | 0.000 | 0.678 | 1.680 |
| Oldpeak | 0.4408 | 0.135 | 3.272 | 0.001 | 0.177 | 0.705 |
| ST_Slope | -1.7134 | 0.225 | -7.610 | 0.000 | -2.155 | -1.272 |

Figura 5: Resultados de la regresión logística

3.3.1. Extraer métricas

En esta sección se analizan los principales indicadores estadísticos derivados del modelo de regresión logística, utilizando la librería statsmodels. Estos valores permiten evaluar la calidad del ajuste, la capacidad explicativa del modelo y la significancia global de sus parámetros.

```

1 pseudo_r2 = modelo_statsmodels.prsquared # McFadden's Pseudo R-dos
2 log_likelihood = modelo_statsmodels.llf # Log-Likelihood del modelo
  completo
3 ll_null = modelo_statsmodels.llnull # Log-Likelihood del modelo nulo
4 llr = modelo_statsmodels.llr # Likelihood Ratio (LR)
5 llr_pvalue = modelo_statsmodels.llr_pvalue # p-value del LR test

```

3.3.2. Pseudo R-cuadrada (McFadden)

```

1 # III.1 - Pseudo R-dos
2 print(f"\n===== III.1: Pseudo R-cuadrada (McFadden) =====")
3 print(f"Valor: {pseudo_r2:.4f}")
4 print("""
5 INTERPRETACION:
6 El Pseudo R-dos de McFadden mide la mejora proporcional del modelo
  completo

```

```

7 respecto al modelo nulo (solo con intercepto).
8
9 Fórmula: Pseudo R_dos = 1 - (Log-Likelihood / LL-Null)
10
11 Valores de referencia (McFadden, 1974):
12 - 0.2 a 0.4: ajuste EXCELENTE
13 - 0.1 a 0.2: ajuste BUENO
14 - < 0.1: ajuste POBRE
15
16 NO es equivalente al R_dos de regresión lineal. Valores más bajos son
17 normales
18 y aceptables en regresión logística.
19 """
20 if pseudo_r2 >= 0.2:
21     print(f" Con {pseudo_r2:.4f}, tu modelo tiene un AJUSTE EXCELENTE")
22 elif pseudo_r2 >= 0.1:
23     print(f" Con {pseudo_r2:.4f}, tu modelo tiene un AJUSTE BUENO")
24 else:
25     print(f" Con {pseudo_r2:.4f}, tu modelo tiene un AJUSTE POBRE")

```

Obteniendo:

```

1 ===== III.1: Pseudo R-cuadrada (McFadden) =====
2 Valor: 0.4732
3
4 INTERPRETACION:
5 El Pseudo R_dos de McFadden mide la mejora proporcional del modelo
6 completo
7 respecto al modelo nulo (solo con intercepto).
8
9 Fórmula: Pseudo R_dos = 1 - (Log-Likelihood / LL-Null)
10
11 Valores de referencia (McFadden, 1974):
12 - 0.2 a 0.4: ajuste EXCELENTE
13 - 0.1 a 0.2: ajuste BUENO
14 - < 0.1: ajuste POBRE
15
16 NO es equivalente al R_dos de regresión lineal. Valores más bajos son
17 normales
18 y aceptables en regresión logística.
19
20 Con 0.4732, tu modelo tiene un AJUSTE EXCELENTE

```

El valor de 0.4732 indica que el modelo presenta un **ajuste excelente**, lo que sugiere que las variables predictoras explican una proporción importante de la variabilidad en la probabilidad de presentar enfermedad cardíaca.

3.3.3. Log-Likelihood

```

1 # III.2 - Log-Likelihood
2 print(f"\n===== III.2: Log-Likelihood =====")
3 print(f"Valor: {log_likelihood:.4f}")
4 print("""

```



```

5  INTERPRETACION:
6  El Log-Likelihood (LL) es el logaritmo natural de la función de
   verosimilitud.
7  Mide qué tan bien el modelo se ajusta a los datos.
8
9  - Valores MAS CERCANOS A 0 indican MEJOR ajuste
10 - Valores MAS NEGATIVOS indican PEOR ajuste
11 - El LL por si solo NO es interpretable; se usa para COMPARAR modelos
12
13 Un modelo perfecto tendria LL = 0, pero esto es imposible en la practica.
14 """
15 print(f" Tu modelo tiene un Log-Likelihood de {log_likelihood:.4f}")

```

Obteniendo:

```

1  ===== III.2: Log-Likelihood =====
2  Valor: -266.3836
3
4  INTERPRETACION:
5  El Log-Likelihood (LL) es el logaritmo natural de la función de
   verosimilitud.
6  Mide qué tan bien el modelo se ajusta a los datos.
7
8  - Valores MAS CERCANOS A 0 indican MEJOR ajuste
9  - Valores MAS NEGATIVOS indican PEOR ajuste
10 - El LL por sí solo NO es interpretable; se usa para COMPARAR modelos
11
12 Un modelo perfecto tendría LL = 0, pero esto es imposible en la práctica.
13
14 Tu modelo tiene un Log-Likelihood de -266.3836

```

En este caso, el valor de -266.3836 indica un buen ajuste general, ya que es considerablemente mayor (menos negativo) que el del modelo nulo, reflejando una mejora significativa en la capacidad explicativa del modelo.

3.3.4. LL-Null (Log-Likelihood Nulo)

```

1  # III.3 - LL-Null
2  print(f"\n===== III.3: LL-Null (Log-Likelihood Nulo) =====")
3  print(f"Valor: {ll_null:.4f}")
4  print("""
5  INTERPRETACION:
6  LL-Null es el Log-Likelihood del modelo "nulo" o baseline.
7  El modelo nulo SOLO incluye el intercepto (sin variables independientes).
8
9  Este valor sirve como punto de referencia para evaluar si tu modelo
   completo
10 aporta información útil para predecir la variable dependiente.
11
12 Si tu modelo completo tiene un LL mucho mejor (más cercano a 0) que LL-
   Null,
13 significa que tus variables independientes S ayudan a explicar la
   variable dependiente.

```

```

14 """
15 mejora = ll_null - log_likelihood
16 print(f"      Mejora del modelo completo vs nulo: {mejora:.4f}")

```

Obteniendo:

```

1      ===== III.3: LL-Null (Log-Likelihood Nulo) =====
2 Valor: -505.6156
3
4 INTERPRETACION:
5 LL-Null es el Log-Likelihood del modelo "nulo" o baseline.
6 El modelo nulo SOLO incluye el intercepto (sin variables independientes).
7
8 Este valor sirve como punto de referencia para evaluar si tu modelo
9   completo
10  aporta información útil para predecir la variable dependiente.
11
12 Si tu modelo completo tiene un LL mucho mejor (más cercano a 0) que LL-
13   Null,
14  significa que tus variables independientes SI ayudan a explicar la
15   variable dependiente.
16
17 Mejora del modelo completo vs nulo: -239.2320

```

Valor obtenido: LL-Null = -505.6156

El *LL-Null* corresponde al logaritmo de la verosimilitud de un modelo sin variables independientes (solo con intercepto). Sirve como punto de comparación para evaluar cuánto mejora el modelo completo.

$$\text{Mejora} = \text{LL-Null} - \text{Log-Likelihood} = -239.2320$$

Esto indica una mejora sustancial, lo que confirma que las variables independientes aportan información significativa para predecir la enfermedad cardíaca.

3.3.5. LLR p-value (Prueba de Razón de Verosimilitudes)

```

1 # III.4 - LLR p-value
2 print(f"\n===== III.4: LLR p-value (Likelihood Ratio Test) =====
3   ")
4 print(f"Valor: {llr_pvalue:.10f}")
5 print(f"Estadístico LR: {llr:.4f}")
6 print("""
7 INTERPRETACION:
8 El LLR p-value prueba la hipótesis nula de que el modelo nulo (solo
9   intercepto)
10  es tan bueno como el modelo completo (con todas las variables).
11
12 Hipótesis:
13   H0: El modelo nulo es suficiente (las variables NO aportan información)
14   H1: El modelo completo es mejor (las variables SI aportan información)
15
16 Criterio de decision (alpha = 0.05):

```

```

15 - Si p-value < 0.05: RECHAZAMOS  $H_0$  - El modelo completo es
    significativamente mejor
16 - Si p-value  $\geq$  0.05: NO rechazamos  $H_0$  - El modelo no es
    significativamente mejor
17 """
18 if llr_pvalue < 0.001:
19     print(f"- Con p-value = {llr_pvalue:.2e} (< 0.001), el modelo es
    ALTAMENTE SIGNIFICATIVO")
20     print("    Las variables independientes explican de manera
    significativa la variable dependiente")
21 elif llr_pvalue < 0.05:
22     print(f"- Con p-value = {llr_pvalue:.4f} (< 0.05), el modelo es
    SIGNIFICATIVO")
23 else:
24     print(f"- Con p-value = {llr_pvalue:.4f} ( $\geq$  0.05), el modelo NO es
    significativo")

```

Obteniendo:

```

1  ===== III.4: LLR p-value (Likelihood Ratio Test) =====
2  Valor: 0.0000000000
3  Estadístico LR: 478.4641
4
5  INTERPRETACION:
6  El LLR p-value prueba la hipótesis nula de que el modelo nulo (solo
    intercepto)
7  es tan bueno como el modelo completo (con todas las variables).
8  Hipótesis:
9       $H_0$ : El modelo nulo es suficiente (las variables NO aportan información)
10      $H_1$ : El modelo completo es mejor (las variables SI aportan información)
11
12  Criterio de decisión (alpha = 0.05):
13     - Si p-value < 0.05: RECHAZAMOS  $H_0$  - El modelo completo es
        significativamente mejor
14     - Si p-value  $\geq$  0.05: NO rechazamos  $H_0$  - El modelo no es
        significativamente mejor
15
16 - Con p-value = 1.25e-95 (< 0.001), el modelo es ALTAMENTE SIGNIFICATIVO
17     Las variables independientes explican de manera significativa la
        variable dependiente

```

Resultados:

Estadístico LR = 478.4641, $p\text{-value} = 1.25 \times 10^{-95}$

El *LLR p-value* contrasta dos hipótesis:

H_0 : El modelo nulo (sin variables) es suficiente

H_1 : El modelo completo es significativamente mejor

Con un valor de $p < 0.001$, se rechaza contundentemente H_0 , concluyendo que el modelo completo es altamente significativo. Esto demuestra que las variables independientes incluidas explican de manera estadísticamente significativa la variable dependiente.

3.4. PARTE IV

Se analizan los resultados obtenidos del modelo de regresión logística completo, con base en los valores de la tabla generada mediante `modelo_statsmodels.summary()` y los coeficientes extraídos en un DataFrame personalizado.

```

1 resultados_df = pd.DataFrame({
2     'Variable': ['Intercepto'] + feature_names,
3     'coef': modelo_statsmodels.params.values,
4     'std_err': modelo_statsmodels.bse.values,
5     'z': modelo_statsmodels.tvalues.values,
6     'p_value': modelo_statsmodels.pvalues.values,
7     'CI_lower': modelo_statsmodels.conf_int()[0].values,
8     'CI_upper': modelo_statsmodels.conf_int()[1].values
9 })
10
11 print("\n--- TABLA COMPLETA DE RESULTADOS ---")
12 print(resultados_df.to_string(index=False))

```

Obteniendo:

```

--- TABLA COMPLETA DE RESULTADOS ---
  Variable    coef  std_err         z    p_value  CI_lower  CI_upper
Intercepto  1.663575  0.432525  3.846189  1.199690e-04  0.815841  2.511309
Age         0.149275  0.130274  1.145851  2.518568e-01 -0.106058  0.404608
Sex         1.433095  0.286038  5.010147  5.438835e-07  0.872470  1.993720
ChestPainType -0.590115  0.121322 -4.864054  1.150056e-06 -0.827901 -0.352329
RestingBP    0.046043  0.120790  0.381180  7.030700e-01 -0.190701  0.282787
Cholesterol  0.101749  0.128710  0.790525  4.292211e-01 -0.150519  0.354017
FastingBS    0.487935  0.119700  4.076314  4.575527e-05  0.253327  0.722543
RestingECG   -0.130590  0.183870 -0.710227  4.775633e-01 -0.490969  0.229790
MaxHR        -0.301256  0.132578 -2.272285  2.306933e-02 -0.561105 -0.041407
ExerciseAngina 1.179007  0.255460  4.615228  3.926649e-06  0.678314  1.679699
Oldpeak      0.440785  0.134703  3.272266  1.066891e-03  0.176772  0.704799
ST_Slope    -1.713367  0.225141 -7.610191  2.736916e-14 -2.154635 -1.272098

```

Figura 6: Resultados del DataFrame

3.4.1. Interpretación de coeficientes

```

1     # IV.1 - Interpretación de coeficientes
2 print(f"\n--- IV.1: Interpretación de cada coeficiente ---")
3 print("""
4 El coeficiente (coef) representa el cambio en el LOG-ODDS por cada unidad
5 de cambio en la variable independiente (manteniendo otras variables
6     constantes).
7
8 En datos ESTANDARIZADOS (como los tuyos):
9     - El coeficiente representa el cambio en log-odds por cada DESVIACION
10        ESTANDAR
11        de cambio en la variable
12     - Esto permite COMPARAR directamente la importancia de las variables

```

```

11     - Variables con |coeficiente| mayor tienen mayor impacto en la predicción
12     """
13
14     for idx, row in resultados_df.iterrows():
15         if row['Variable'] == 'Intercepto':
16             continue
17         var = row['Variable']
18         coef = row['coef']
19         odds_ratio = np.exp(coef)
20
21         print(f"\n{var}:")
22         print(f"    Coeficiente: {coef:.6f}")
23         print(f"    Odds Ratio: {odds_ratio:.6f}")
24
25         if coef > 0:
26             print(f"    Efecto POSITIVO: aumentar 1 desviación estándar en {var}
27                 ")
28             print(f"    multiplica las probabilidades de HeartDisease por {
29                 odds_ratio:.4f}")
30         else:
31             print(f"    Efecto NEGATIVO: aumentar 1 desviación estándar en {var}
32                 ")
33             print(f"    multiplica las probabilidades de HeartDisease por {
34                 odds_ratio:.4f}")

```

Teniendo como resultados del análisis:

```

1     --- IV.1: Interpretación de cada coeficiente ---
2
3     El coeficiente (coef) representa el cambio en el LOG-ODDS por cada unidad
4     de cambio en la variable independiente (manteniendo otras variables
5     constantes).
6
7     En datos ESTANDARIZADOS (como los tuyos):
8     - El coeficiente representa el cambio en log-odds por cada DESVIACION
9       ESTANDAR
10      de cambio en la variable
11      - Esto permite COMPARAR directamente la importancia de las variables
12      - Variables con |coeficiente| mayor tienen mayor impacto en la predicción
13
14      Age:
15          Coeficiente: 0.149275
16          Odds Ratio: 1.160992
17          - Efecto POSITIVO: aumentar 1 desviación estándar en Age
18            multiplica las probabilidades de HeartDisease por 1.1610
19
20      Sex:
21          Coeficiente: 1.433095
22          Odds Ratio: 4.191652
23          - Efecto POSITIVO: aumentar 1 desviación estándar en Sex
24            multiplica las probabilidades de HeartDisease por 4.1917

```

```
25 ChestPainType:
26   Coeficiente: -0.590115
27   Odds Ratio: 0.554263
28   - Efecto NEGATIVO: aumentar 1 desviación estándar en ChestPainType
29     multiplica las probabilidades de HeartDisease por 0.5543
30
31 RestingBP:
32   Coeficiente: 0.046043
33   Odds Ratio: 1.047119
34   - Efecto POSITIVO: aumentar 1 desviación estándar en RestingBP
35     multiplica las probabilidades de HeartDisease por 1.0471
36
37 Cholesterol:
38   Coeficiente: 0.101749
39   Odds Ratio: 1.107105
40   - Efecto POSITIVO: aumentar 1 desviación estándar en Cholesterol
41     multiplica las probabilidades de HeartDisease por 1.1071
42
43 FastingBS:
44   Coeficiente: 0.487935
45   Odds Ratio: 1.628950
46   - Efecto POSITIVO: aumentar 1 desviación estándar en FastingBS
47     multiplica las probabilidades de HeartDisease por 1.6289
48
49 RestingECG:
50   Coeficiente: -0.130590
51   Odds Ratio: 0.877578
52   - Efecto NEGATIVO: aumentar 1 desviación estándar en RestingECG
53     multiplica las probabilidades de HeartDisease por 0.8776
54
55 MaxHR:
56   Coeficiente: -0.301256
57   Odds Ratio: 0.739888
58   - Efecto NEGATIVO: aumentar 1 desviación estándar en MaxHR
59     multiplica las probabilidades de HeartDisease por 0.7399
60
61 ExerciseAngina:
62   Coeficiente: 1.179007
63   Odds Ratio: 3.251143
64   - Efecto POSITIVO: aumentar 1 desviación estándar en ExerciseAngina
65     multiplica las probabilidades de HeartDisease por 3.2511
66
67 Oldpeak:
68   Coeficiente: 0.440785
69   Odds Ratio: 1.553927
70   - Efecto POSITIVO: aumentar 1 desviación estándar en Oldpeak
71     multiplica las probabilidades de HeartDisease por 1.5539
72
73 ST_Slope:
74   Coeficiente: -1.713367
75   Odds Ratio: 0.180258
76   - Efecto NEGATIVO: aumentar 1 desviación estándar en ST_Slope
77     multiplica las probabilidades de HeartDisease por 0.1803
```

3.4.2. Interpretación de std err

```

1      # IV.2 - Interpretación de std err
2  print(f"\n--- IV.2: Interpretación del Error Estándar (std err) ---")
3  print("""
4  El error estándar mide la VARIABILIDAD o INCERTIDUMBRE del coeficiente
      estimado.
5
6  - Error estándar DIMINUTO: el coeficiente es más PRECISO y CONFIABLE
7  - Error estándar GRANDE: el coeficiente tiene más INCERTIDUMBRE
8
9  El error estándar se usa para:
10     1. Calcular el estadístico z:  $z = \text{coef} / \text{std\_err}$ 
11     2. Construir intervalos de confianza
12     3. Evaluar la precisión de las estimaciones
13 """)
14
15 print("\nVariables ordenadas por precisión (menor error estándar):")
16 resultados_ordenados = resultados_df[resultados_df['Variable'] != '
      Intercepto'].sort_values('std_err')
17 print(resultados_ordenados[['Variable', 'coef', 'std_err']].to_string(
      index=False))

```

Obteniendo:

```

1      --- IV.2: Interpretación del Error Estándar (std err) ---
2
3  El error estándar mide la VARIABILIDAD o INCERTIDUMBRE del coeficiente
      estimado.
4
5  - Error estándar DIMINUTO: el coeficiente es más PRECISO y CONFIABLE
6  - Error estándar GRANDE: el coeficiente tiene más INCERTIDUMBRE
7
8  El error estándar se usa para:
9     1. Calcular el estadístico z:  $z = \text{coef} / \text{std\_err}$ 
10     2. Construir intervalos de confianza
11     3. Evaluar la precisión de las estimaciones
12
13
14 Variables ordenadas por precisión (menor error estándar):
15     Variable      coef  std_err
16     FastingBS    0.487935  0.119700
17     RestingBP    0.046043  0.120790
18     ChestPainType -0.590115  0.121322
19     Cholesterol  0.101749  0.128710
20     Age          0.149275  0.130274
21     MaxHR        -0.301256  0.132578
22     Oldpeak      0.440785  0.134703
23     RestingECG   -0.130590  0.183870
24     ST_Slope     -1.713367  0.225141
25     ExerciseAngina 1.179007  0.255460
26     Sex          1.433095  0.286038

```

3.4.3. Variables significativas

```

1      # IV.3 - Variables significativas
2  print(f"\n--- IV.3: Variables significativas (alpha = 0.05) ---")
3  print("""
4  Un variable es SIGNIFICATIVA si su p-value < 0.05 (nivel de significancia
    común).
5
6  Esto significa que hay menos del 5% de probabilidad de que el coeficiente
    observado
7  sea debido al azar (rechazamos  $H_0: \beta = 0$ ).
8
9  Variables significativas: REALMENTE explican parte de la variable
    dependiente
10 Variables NO significativas: NO aportan información estadísticamente
    relevante
11 """)
12
13 variables_significativas = resultados_df[
14     (resultados_df['Variable'] != 'Intercepto') &
15     (resultados_df['p_value'] < 0.05)
16 ].sort_values('p_value')
17
18 print(f"\nVariables SIGNIFICATIVAS (p < 0.05):")
19 if len(variables_significativas) > 0:
20     print(variables_significativas[['Variable', 'coef', 'p_value']].
21         to_string(index=False))
22     print(f"\n {len(variables_significativas)} de {len(feature_names)}
23     variables SI explican la variable dependiente")
24 else:
25     print("NINGUNA variable es significativa")
26
27 variables_no_significativas = resultados_df[
28     (resultados_df['Variable'] != 'Intercepto') &
29     (resultados_df['p_value'] >= 0.05)
30 ]
31
32 print(f"\nVariables NO SIGNIFICATIVAS (p ≥ 0.05):")
33 if len(variables_no_significativas) > 0:
34     print(variables_no_significativas[['Variable', 'coef', 'p_value']].
35         to_string(index=False))
36     print("\n Estas variables podrían ser candidatas para eliminación")
37 else:
38     print("Todas las variables son significativas")

```

Obteniendo:

```

1      --- IV.3: Variables significativas (alpha = 0.05) ---
2
3  Un variable es SIGNIFICATIVA si su p-value < 0.05 (nivel de significancia
    común).
4
5  Esto significa que hay menos del 5% de probabilidad de que el coeficiente
    observado
6  sea debido al azar (rechazamos  $H_0: \beta = 0$ ).

```



```

7
8 Variables significativas: REALMENTE explican parte de la variable
  dependiente
9 Variables NO significativas: NO aportan información estadísticamente
  relevante
10
11
12 Variables SIGNIFICATIVAS (p < 0.05):
13     Variable      coef      p_value
14     ST_Slope -1.713367 2.736916e-14
15     Sex      1.433095 5.438835e-07
16     ChestPainType -0.590115 1.150056e-06
17     ExerciseAngina 1.179007 3.926649e-06
18     FastingBS      0.487935 4.575527e-05
19     Oldpeak        0.440785 1.066891e-03
20     MaxHR          -0.301256 2.306933e-02
21
22 - 7 de 11 variables SI explican la variable dependiente
23
24 Variables NO SIGNIFICATIVAS (p ≥ 0.05):
25     Variable      coef      p_value
26     Age      0.149275 0.251857
27     RestingBP 0.046043 0.703070
28     Cholesterol 0.101749 0.429221
29     RestingECG -0.130590 0.477563
30
31 - Estas variables podrían ser candidatas para eliminación

```

3.4.4. Variables con mayor poder explicativo (valor z)

```

1     # IV.4 - Variables con mayor poder explicativo (valor z)
2 print(f"\n--- IV.4: Variables con mayor poder explicativo (según |z|) ---"
3       )
4 print("""
5 El estadístico z (valor z) mide cuántas desviaciones estándar el
6 coeficiente
7 está alejado de cero.
8
9 - |z| GRANDE: La variable tiene un efecto FUERTE y SIGNIFICATIVO
10 - |z| DIMINUTO: La variable tiene un efecto DEBIL o NO SIGNIFICATIVO
11
12 Valores de referencia:
13     |z| > 2.576 - p < 0.01 (muy significativo)
14     |z| > 1.960 - p < 0.05 (significativo)
15     |z| < 1.960 - p ≥ 0.05 (no significativo)
16 """)
17
18 top_variables = resultados_df[resultados_df['Variable'] != 'Intercepto'].
19     copy()
20 top_variables['abs_z'] = np.abs(top_variables['z'])
21 top_variables = top_variables.sort_values('abs_z', ascending=False)
22
23 print("\nVariables ordenadas por |z| (poder explicativo):")

```

```

21 print(top_variaciones[['Variable', 'coef', 'z', 'abs_z', 'p_value']].head
    (10).to_string(index=False))
22
23 top_2 = top_variaciones.head(2)
24 print(f"\n- Las 2 variables que MEJOR explican la variable dependiente son
    :")
25 for idx, row in top_2.iterrows():
26     print(f"    {idx+1}. {row['Variable']} (|z| = {row['abs_z']:.4f}, p = {
        row['p_value']:.6f})")

```

Obtenemos:

```

1     --- IV.4: Variables con mayor poder explicativo (según |z|) ---
2
3 El estadístico z (valor z) mide cuántas desviaciones estándar el
4 coeficiente
5 está alejado de cero.
6
7 - |z| GRANDE: La variable tiene un efecto FUERTE y SIGNIFICATIVO
8 - |z| DIMINUTO: La variable tiene un efecto DEBIL o NO SIGNIFICATIVO
9
10 Valores de referencia:
11 |z| > 2.576 - p < 0.01 (muy significativo)
12 |z| > 1.960 - p < 0.05 (significativo)
13 |z| < 1.960 - p ≥ 0.05 (no significativo)
14
15 Variables ordenadas por |z| (poder explicativo):
16     Variable      coef      z      abs_z      p_value
17     ST_Slope -1.713367 -7.610191 7.610191 2.736916e-14
18     Sex      1.433095  5.010147 5.010147 5.438835e-07
19     ChestPainType -0.590115 -4.864054 4.864054 1.150056e-06
20     ExerciseAngina 1.179007  4.615228 4.615228 3.926649e-06
21     FastingBS      0.487935  4.076314 4.076314 4.575527e-05
22     Oldpeak      0.440785  3.272266 3.272266 1.066891e-03
23     MaxHR      -0.301256 -2.272285 2.272285 2.306933e-02
24     Age      0.149275  1.145851 1.145851 2.518568e-01
25     Cholesterol  0.101749  0.790525 0.790525 4.292211e-01
26     RestingECG -0.130590 -0.710227 0.710227 4.775633e-01
27
28 - Las 2 variables que MEJOR explican la variable dependiente son:
29 12. ST_Slope (|z| = 7.6102, p = 0.000000)
30 3. Sex (|z| = 5.0101, p = 0.000001)

```

3.4.5. Intervalo de confianza del 95 % para la mejor variable

```

1     # IV.5 - Intervalo de confianza del 95% para la mejor variable
2 print(f"\n--- IV.5: Intervalo de Confianza 95% para la mejor variable ---"
3     )
4 mejor_variable = top_variaciones.iloc[0]
5 var_name = mejor_variable['Variable']
6 coef_val = mejor_variable['coef']

```

```

7  ci_lower = mejor_variable['CI_lower']
8  ci_upper = mejor_variable['CI_upper']
9
10 print(f"\nMejor variable: {var_name}")
11 print(f"Coeficiente: {coef_val:.6f}")
12 print(f"Intervalo de Confianza 95%: [{ci_lower:.6f}, {ci_upper:.6f}]")
13
14 print("""
15  INTERPRETACION:
16  El intervalo de confianza del 95% indica que estamos 95% seguros de que el
17  VERDADERO valor del coeficiente poblacional está dentro de este rango.
18
19  Si el intervalo NO contiene el 0, la variable es significativa (p < 0.05).
20  Cuanto más ESTRECHO el intervalo, más PRECISA es la estimación.
21  """)
22
23 if ci_lower > 0 and ci_upper > 0:
24     print(f"- El IC no contiene 0 y es completamente POSITIVO")
25     print(f"    {var_name} tiene un efecto POSITIVO significativo y
26           confiable")
27 elif ci_lower < 0 and ci_upper < 0:
28     print(f"- El IC no contiene 0 y es completamente NEGATIVO")
29     print(f"    {var_name} tiene un efecto NEGATIVO significativo y
30           confiable")
31 else:
32     print(f"- El IC contiene 0: la variable NO es significativa")

```

Obtenemos:

```

1  --- IV.5: Intervalo de Confianza 95% para la mejor variable ---
2
3  Mejor variable: ST_Slope
4  Coeficiente: -1.713367
5  Intervalo de Confianza 95%: [-2.154635, -1.272098]
6
7  INTERPRETACION:
8  El intervalo de confianza del 95% indica que estamos 95% seguros de que el
9  VERDADERO valor del coeficiente poblacional está dentro de este rango.
10
11 Si el intervalo NO contiene el 0, la variable es significativa (p < 0.05).
12 Cuanto más ESTRECHO el intervalo, más PRECISA es la estimación.
13
14 - El IC no contiene 0 y es completamente NEGATIVO
15   ST_Slope tiene un efecto NEGATIVO significativo y confiable

```

3.4.6. Visualización de los Intervalos de Confianza

La Figura 7 presenta los coeficientes estimados para las cinco variables más relevantes, junto con sus intervalos de confianza al 95 %. En donde Los puntos representan los coeficientes estimados ($\hat{\beta}$). Las barras horizontales indican el intervalo de confianza. La línea roja vertical ($x = 0$) representa el punto de “sin efecto”

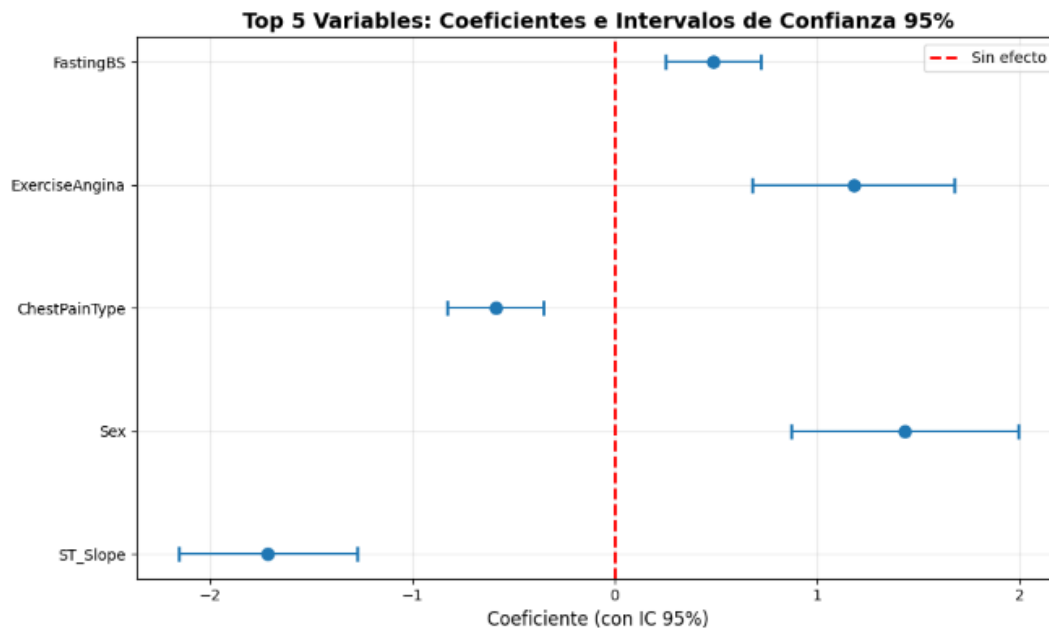


Figura 7: Coeficientes e intervalos de confianza al 95 % para las principales variables del modelo.

Interpretación: Ninguna de las variables seleccionadas cruza la línea de referencia, lo que confirma su significancia estadística. Las variables **ST_Slope** y **Sex** presentan los intervalos más alejados de cero, siendo las de mayor peso explicativo en la probabilidad de enfermedad cardíaca.

3.5. PARTE V

En esta sección se evalúa el desempeño del modelo de regresión logística mediante la construcción de la matriz de confusión. Esta matriz permite comparar las predicciones del modelo (*Predicho*) con los valores reales (*Actual*) de la variable dependiente (HeartDisease[5][1]).

```

1  # Predicciones en conjunto de prueba
2  y_pred = modelo.predict(X_test)
3
4  # Crear matriz de confusión
5  cm = confusion_matrix(y_test, y_pred)
6
7  print("\nMatriz de Confusión:")
8  print(cm)
9  print("""
10 Estructura:
11
12           Predicho
13      0      1
14 Actual 0  [TN] [FP]
15          1  [FN] [TP]
16
17 TN (True Negative): Correctamente predicho como NO enfermedad
18 TP (True Positive): Correctamente predicho como SI enfermedad
19 FN (False Negative): Incorrectamente predicho como NO (era SI)
20 FP (False Positive): Incorrectamente predicho como SI (era NO)

```

```

20 """
21
22 tn, fp, fn, tp = cm.ravel()
23 print(f"True Negatives (TN): {tn}")
24 print(f"False Positives (FP): {fp}")
25 print(f"False Negatives (FN): {fn}")
26 print(f"True Positives (TP): {tp}")
27
28 # Visualizar matriz de confusión
29 fig, ax = plt.subplots(figsize=(8, 6))
30 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=['No
    Enfermedad', 'Enfermedad'])
31 disp.plot(cmap='Blues', ax=ax, values_format='d')
32 plt.title('Matriz de Confusión', fontsize=14, fontweight='bold')

```

La matriz de confusión tiene la siguiente estructura:

```

1   Matriz de Confusión:
2   [[68  9]
3    [17 90]]
4
5   True Negatives (TN): 68
6   False Positives (FP): 9
7   False Negatives (FN): 17
8   True Positives (TP): 90
9   Text(0.5, 1.0, 'Matriz de Confusión')

```

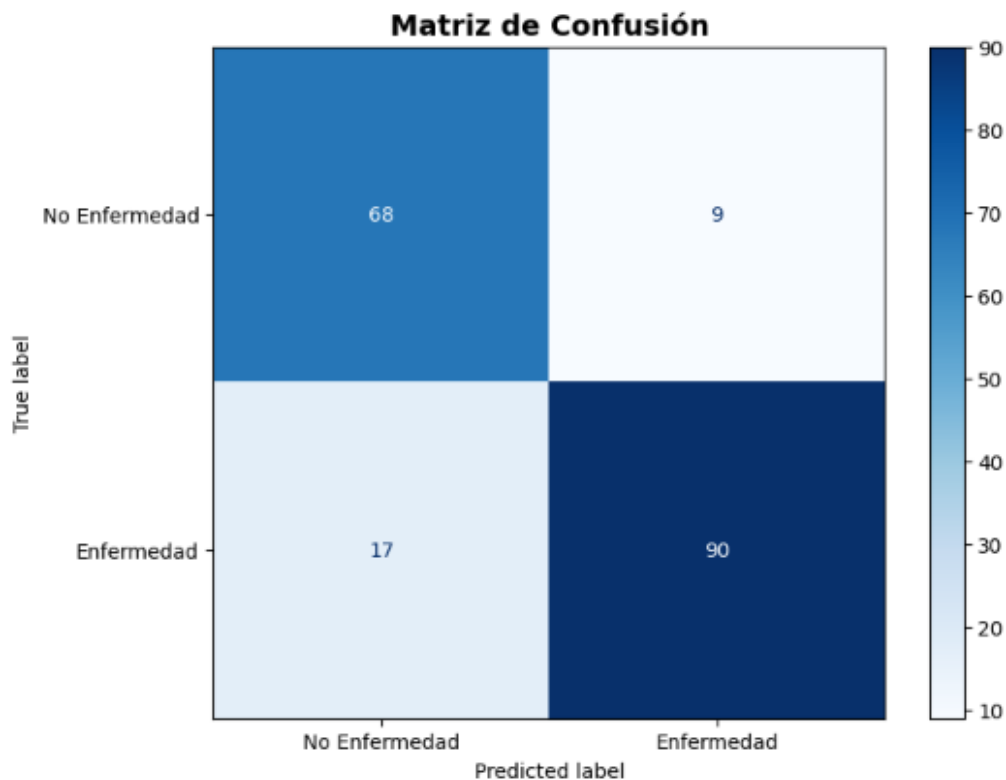


Figura 8: Matriz de confusión

3.5.1. Métricas derivadas

```

1 precision = tp / (tp + fp) if (tp + fp) > 0 else 0
2 recall = tp / (tp + fn) if (tp + fn) > 0 else 0
3 f1_score = 2 * (precision * recall) / (precision + recall) if (precision +
4   recall) > 0 else 0
5 specificity = tn / (tn + fp) if (tn + fp) > 0 else 0
6
7 print(f"\nMétricas de Evaluación:")
8 print(f" Precision (Precisión): {precision:.4f}")
9 print(f" Recall (Sensibilidad): {recall:.4f}")
10 print(f" F1-Score: {f1_score:.4f}")
11 print(f" Specificity (Especificidad): {specificity:.4f}")

```

En donde sus resultados de evaluación son:

```

1 Métricas de Evaluación:
2 Precision (Precisión): 0.9091
3 Recall (Sensibilidad): 0.8411
4 F1-Score: 0.8738
5 Specificity (Especificidad): 0.8831

```

3.5.2. Mejor modelo con máximo 4 variables (criterio LLR p-value)

En esta etapa se busca identificar las variables más significativas estadísticamente dentro del modelo de Regresión Logística. Para ello, se utiliza el criterio del p -value obtenido del test de razón de verosimilitudes (**Likelihood Ratio Test, LLR**). Las variables con menor p -value son aquellas que tienen una mayor evidencia de influencia sobre la variable dependiente (HeartDisease), es decir, son las más relevantes para el modelo.

A continuación se presenta el fragmento de código utilizado para seleccionar las variables más significativas:

```

1 # Seleccionar top 4 variables por p-value (excluyendo intercepto)
2 top_4_pvalue = resultados_df[resultados_df['Variable'] != 'Intercepto'].
3   nsmallest(4, 'p_value')
4 top_4_vars = top_4_pvalue['Variable'].tolist()
5
6 print(f"\nTop 4 variables seleccionadas (menor p-value):")
7 for i, row in top_4_pvalue.iterrows():
8     print(f" {row['Variable']}: p-value = {row['p_value']:.6f}")

```

Los resultados obtenidos fueron los siguientes:

```

1 Top 4 variables seleccionadas (menor p-value):
2 ST_Slope: p-value = 0.000000
3 Sex: p-value = 0.000001
4 ChestPainType: p-value = 0.000001
5 ExerciseAngina: p-value = 0.000004

```

Las variables con menor p -value son las que más contribuyen a explicar la probabilidad de padecer una enfermedad cardíaca. En este caso, ST_Slope, Sex, ChestPainType y ExerciseAngina presentan una alta significancia estadística ($p < 0.001$), lo que indica que tienen un fuerte efecto

sobre el modelo. Estos resultados son consistentes con los obtenidos previamente mediante el criterio de Pseudo R^2 , lo que refuerza la confiabilidad de la selección.

3.5.3. Entrenar modelo reducido

Una vez identificadas las cuatro variables más significativas según el criterio del p -value (ST_Slope, Sex, ChestPainType y ExerciseAngina), se procede a entrenar un modelo reducido de Regresión Logística. El objetivo es comparar su desempeño con el modelo completo (que incluye todas las variables) mediante métricas como el Pseudo R^2 , el p -value del LLR (Likelihood Ratio Test) y el Log-Likelihood.

El siguiente bloque de código muestra el proceso de entrenamiento y comparación de ambos modelos:

```

1     X_train_reduced = X_train[top_4_vars]
2     X_test_reduced = X_test[top_4_vars]
3
4     X_train_reduced_sm = sm.add_constant(X_train_reduced)
5     modelo_reducido = sm.Logit(y_train, X_train_reduced_sm).fit(dis=0)
6
7     print(f"\n--- Comparación de Modelos ---")
8     print(f"Modelo Completo ({len(feature_names)} variables):")
9     print(f"    Pseudo R_dos: {pseudo_r2:.4f}")
10    print(f"    LLR p-value: {llr_pvalue:.2e}")
11    print(f"    Log-Likelihood: {log_likelihood:.4f}")
12
13    print(f"\nModelo Reducido (4 variables):")
14    print(f"    Pseudo R_dos: {modelo_reducido.prsquared:.4f}")
15    print(f"    LLR p-value: {modelo_reducido.llr_pvalue:.2e}")
16    print(f"    Log-Likelihood: {modelo_reducido.llf:.4f}")
17
18    if modelo_reducido.llr_pvalue < llr_pvalue:
19        print("\n El modelo reducido tiene un LLR p-value MEJOR (más pequeño)"
20            )
21    else:
22        print("\n El modelo completo mantiene un LLR p-value mejor")

```

Los resultados obtenidos fueron los siguientes:

```

1
2     --- Comparación de Modelos ---
3     Modelo Completo (11 variables):
4         Pseudo R_dos: 0.4732
5         LLR p-value: 1.25e-95
6         Log-Likelihood: -266.3836
7
8     Modelo Reducido (4 variables):
9         Pseudo R_dos: 0.4272
10        LLR p-value: 3.38e-92
11        Log-Likelihood: -289.6169
12
13    El modelo completo mantiene un LLR p-value mejor

```

El modelo reducido, con solo cuatro variables, mantiene un desempeño razonablemente alto ($R^2 = 0.4272$) y un p -value del LLR aún muy significativo (3.38×10^{-92}). No obstante, el modelo completo presenta un ajuste ligeramente mejor ($R^2 = 0.4732$) y un menor p -value, indicando que explica una mayor proporción de la variabilidad en los datos. Aun así, el modelo reducido es preferible en contextos donde se busca un balance entre simplicidad e interpretabilidad, sacrificando solo una pequeña parte del poder explicativo.

3.5.4. Creación de gráfica logística y diagrama de dispersión

Para visualizar la relación entre la variable más influyente del modelo y la probabilidad predicha de padecer enfermedad cardíaca, se genera una gráfica logística acompañada de un diagrama de dispersión. Este gráfico permite observar cómo el modelo ajusta la función sigmoide (*curva logística*) sobre los datos observados, y cómo la probabilidad estimada cambia a medida que la variable independiente varía.

```

1      # Seleccionar la variable más importante
2      mejor_var_idx = feature_names.index(var_name)
3
4      # Obtener valores de esa variable
5      X_plot = X_test.iloc[:, mejor_var_idx].values if isinstance(X_test, pd.
        DataFrame) else X_test[:, mejor_var_idx]
6      y_true = y_test.values if isinstance(y_test, pd.Series) else y_test
7
8      # Calcular probabilidades predichas
9      y_proba = modelo.predict_proba(X_test)[: , 1]
10
11     # Ordenar para la curva logística
12     idx_sort = np.argsort(X_plot)
13     X_sorted = X_plot[idx_sort]
14     y_proba_sorted = y_proba[idx_sort]
15
16     # Crear gráfico
17     fig, ax = plt.subplots(figsize=(12, 7))
18
19     # Diagrama de dispersión (con jitter para visibilidad)
20     jitter = 0.05
21     y_jitter = y_true + np.random.normal(0, jitter, size=len(y_true))
22     ax.scatter(X_plot, y_jitter, alpha=0.3, s=30, label='Datos observados',
        color='steelblue')
23
24     # Curva logística
25     ax.plot(X_sorted, y_proba_sorted, color='red', linewidth=3, label='Curva
        Logística')
26
27     # Línea de decisión
28     ax.axhline(y=0.5, color='green', linestyle='--', linewidth=2, label='
        Umbral de decisión (0.5)')
29
30     ax.set_xlabel(f'{var_name} (estandarizado)', fontsize=12)
31     ax.set_ylabel('Probabilidad de HeartDisease=1', fontsize=12)
32     ax.set_title(f'Regresión Logística: {var_name} vs Probabilidad de
        Enfermedad Cardíaca',

```



```

33         fontsize=14, fontweight='bold')
34 ax.legend(fontsize=10)
35 ax.grid(True, alpha=0.3)
36 ax.set_ylim(-0.1, 1.1)
37 plt.tight_layout()

```

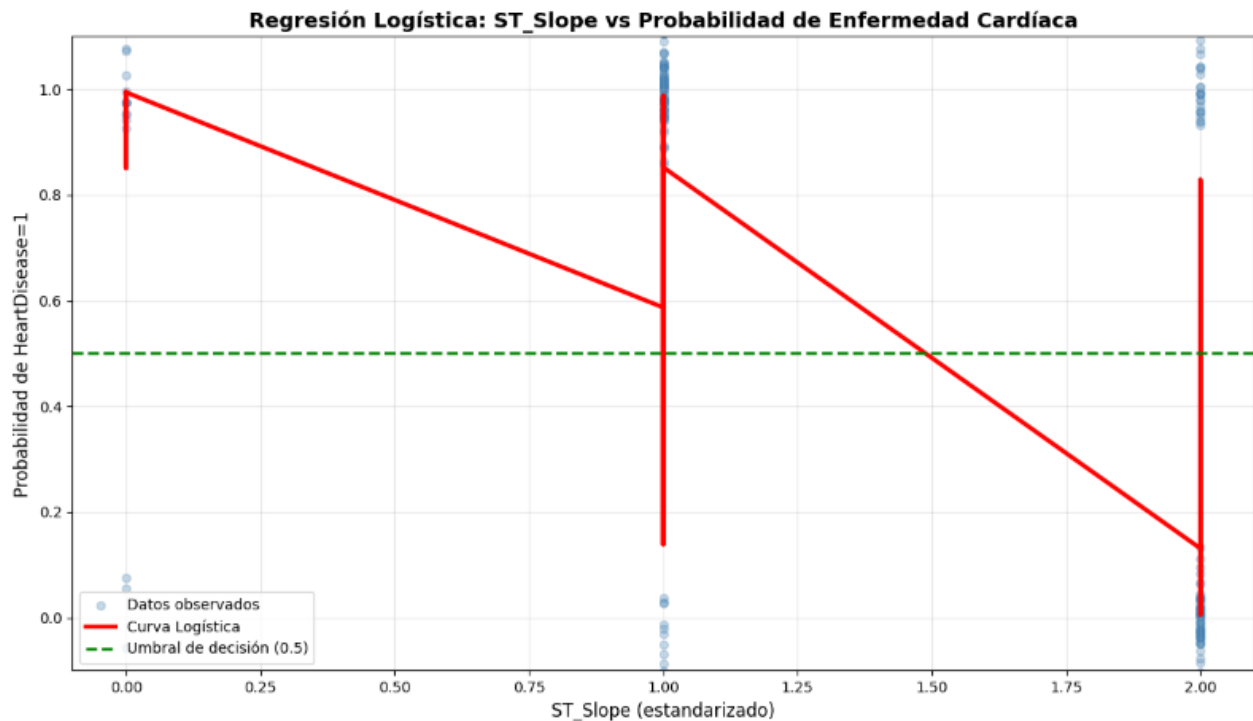


Figura 9: Gráfica logística y diagrama de dispersión

3.5.5. Precisión del Modelo

La precisión (*accuracy*) es una métrica fundamental para evaluar el rendimiento de un modelo de clasificación. Indica el porcentaje de predicciones correctas realizadas por el modelo sobre el total de observaciones, considerando tanto los aciertos en la clase positiva como en la negativa. Matemáticamente se expresa como:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

donde:

- *TP*: verdaderos positivos
- *TN*: verdaderos negativos
- *FP*: falsos positivos
- *FN*: falsos negativos

A continuación se muestra el código empleado para calcular la precisión y generar un reporte completo de clasificación:

```

1     accuracy = accuracy_score(y_test, y_pred)
2     print(f"\nPrecisión (Accuracy): {accuracy:.4f} ({accuracy*100:.2f}%)")
3
4     print("""
5     INTERPRETACIÓN:
6     La precisión (accuracy) es el porcentaje de predicciones correctas sobre
       el total.
7
8     Accuracy = (TP + TN) / (TP + TN + FP + FN)
9
10    Valores de referencia:
11    > 0.90: EXCELENTE
12    0.80 - 0.90: MUY BUENO
13    0.70 - 0.80: BUENO
14    0.60 - 0.70: ACEPTABLE
15    < 0.60: POBRE
16
17    IMPORTANTE: En datasets desbalanceados, la accuracy puede ser engañosa.
18    Siempre revisar precision, recall y F1-score.
19    """)
20
21    if accuracy >= 0.90:
22        print(f"    Con {accuracy:.2%}, tu modelo tiene precisión EXCELENTE")
23    elif accuracy >= 0.80:
24        print(f"    Con {accuracy:.2%}, tu modelo tiene precisión MUY BUENA")
25    elif accuracy >= 0.70:
26        print(f"    Con {accuracy:.2%}, tu modelo tiene precisión BUENA")
27    elif accuracy >= 0.60:
28        print(f"    Con {accuracy:.2%}, tu modelo tiene precisión ACEPTABLE")
29    else:
30        print(f"    Con {accuracy:.2%}, tu modelo tiene precisión POBRE")
31
32    # Reporte completo de clasificación
33    print("\n--- Reporte Completo de Clasificación ---")
34    print(classification_report(y_test, y_pred, target_names=['No Enfermedad',
       'Enfermedad']))

```

| --- Reporte Completo de Clasificación --- | | | | |
|---|-----------|--------|----------|---------|
| | precision | recall | f1-score | support |
| No Enfermedad | 0.80 | 0.88 | 0.84 | 77 |
| Enfermedad | 0.91 | 0.84 | 0.87 | 107 |
| accuracy | | | 0.86 | 184 |
| macro avg | 0.85 | 0.86 | 0.86 | 184 |
| weighted avg | 0.86 | 0.86 | 0.86 | 184 |

Figura 10: Reporte Completo de Clasificación

Interpretación: El modelo logra clasificar correctamente aproximadamente el 85.87 % de las observaciones, lo que se considera un desempeño **muy bueno** según las referencias comunes ($0.80 \leq \text{Accuracy} < 0.90$).

Dado que la precisión puede ser engañosa en conjuntos de datos desbalanceados, se complementa este análisis con otras métricas como la *precisión por clase (precision)*, la *sensibilidad (recall)* y el *F1-score*, las cuales permiten evaluar de forma más completa el comportamiento del modelo en cada categoría.

3.6. PARTE VI

En esta sección se realiza la selección del modelo de regresión logística reducido, considerando un máximo de 4 variables independientes. El criterio de selección utilizado es el **Pseudo R²**, que mide la proporción de variabilidad explicada por el modelo.

Dado que el número de variables es relativamente grande, se implementa un enfoque de **forward selection** para agregar iterativamente la variable que más incrementa el Pseudo R² en cada paso. Este método es computacionalmente más eficiente que la búsqueda exhaustiva de todas las combinaciones posibles.

El procedimiento consiste en:

1. Comenzar con ninguna variable seleccionada.
2. En cada paso, agregar la variable que más incrementa el Pseudo R².
3. Repetir hasta completar 4 variables.

A continuación se presenta el bloque de código utilizado para este análisis.

```
1 print("""
2 ESTRATEGIA:
3 Probaremos todas las combinaciones posibles de 4 variables y
4 seleccionaremos
5 la combinacion que maximice el Pseudo R_dos.
6
7 NOTA: Con muchas variables, esto puede ser computacionalmente costoso.
8 Usaremos un enfoque iterativo (forward selection) como alternativa
9 eficiente.
10 """)
11
12 # Si tienes muchas variables, esto puede tardar. Limitar a forward
13 selection
14 if len(feature_names) > 10:
15     print("\n Numero grande de variables detectado.")
16     print(" Usando FORWARD SELECTION en lugar de búsqueda exhaustiva.")
17
18     # Forward Selection
19     selected_vars = []
20     remaining_vars = feature_names.copy()
21     best_r2_history = []
22
23     for step in range(4):
```

```

22     best_r2 = -np.inf
23     best_var = None
24
25     for var in remaining_vars:
26         test_vars = selected_vars + [var]
27         X_temp = sm.add_constant(X_train[test_vars])
28         try:
29             modelo_temp = sm.Logit(y_train, X_temp).fit(dis=0)
30             if modelo_temp.prsquared > best_r2:
31                 best_r2 = modelo_temp.prsquared
32                 best_var = var
33         except:
34             continue
35
36     if best_var:
37         selected_vars.append(best_var)
38         remaining_vars.remove(best_var)
39         best_r2_history.append(best_r2)
40         print(f"\nPaso {step+1}: Agregada '{best_var}' Pseudo R_dos =
           {best_r2:.4f}")
41
42     top_4_r2_vars = selected_vars
43
44 else:
45     # Búsqueda exhaustiva (solo si pocas variables)
46     print("\nBuscando la mejor combinación de 4 variables...")
47     best_r2 = -np.inf
48     best_combination = None
49
50     for combo in combinations(feature_names, 4):
51         X_temp = sm.add_constant(X_train[list(combo)])
52         try:
53             modelo_temp = sm.Logit(y_train, X_temp).fit(dis=0)
54             if modelo_temp.prsquared > best_r2:
55                 best_r2 = modelo_temp.prsquared
56                 best_combination = combo
57         except:
58             continue
59
60     top_4_r2_vars = list(best_combination)
61
62 print(f"\n--- Mejor Combinacion de 4 Variables (Pseudo R_dos) ---")
63 for i, var in enumerate(top_4_r2_vars, 1):
64     print(f"    {i}. {var}")

```

Los resultados de la selección muestran que las cuatro variables que mejor explican la presencia de enfermedad cardíaca según el criterio de Pseudo R^2 son:

```

1     --- Mejor Combinacion de 4 Variables (Pseudo R_dos) ---
2     1. ST_Slope
3     2. ExerciseAngina
4     3. Sex
5     4. ChestPainType

```

3.6.1. Entrenamiento del modelo final con estas 4 variables

```

1      X_train_r2 = sm.add_constant(X_train[top_4_r2_vars])
2  X_test_r2 = sm.add_constant(X_test[top_4_r2_vars])
3  modelo_final_r2 = sm.Logit(y_train, X_train_r2).fit(dis=0)
4
5  print(modelo_final_r2.summary())
6
7  print(f"\n--- Comparación Final ---")
8  print(f"\nModelo Completo ({len(feature_names)} variables):")
9  print(f"   Pseudo R_dos: {pseudo_r2:.4f}")
10 print(f"   AIC: {modelo_statsmodels.aic:.4f}")
11 print(f"   BIC: {modelo_statsmodels.bic:.4f}")
12
13 print(f"\nModelo con 4 vars (criterio p-value):")
14 print(f"   Pseudo R_dosR_dosR_dos: {modelo_reducido.prsquared:.4f}")
15 print(f"   AIC: {modelo_reducido.aic:.4f}")
16 print(f"   BIC: {modelo_reducido.bic:.4f}")
17 print(f"   Variables: {top_4_vars}")
18
19 print(f"\nModelo con 4 vars (criterio Pseudo R_dosR_dos):")
20 print(f"   Pseudo R_dos: {modelo_final_r2.prsquared:.4f}")
21 print(f"   AIC: {modelo_final_r2.aic:.4f}")
22 print(f"   BIC: {modelo_final_r2.bic:.4f}")
23 print(f"   Variables: {top_4_r2_vars}")
24
25 print("""
26 CRITERIOS DE SELECCION:
27
28 1. Pseudo R_dos: Mayor es mejor (más varianza explicada)
29 2. AIC (Akaike Information Criterion): Menor es mejor
30 3. BIC (Bayesian Information Criterion): Menor es mejor
31
32 AIC y BIC penalizan modelos con más parámetros. Son útiles para evitar
33   overfitting.
34 BIC penaliza más fuertemente que AIC.
35
36 RECOMENDACION:
37 - Si buscas INTERPRETABILIDAD: modelo con 4 variables
38 - Si buscas MAXIMA PRECISION: modelo completo (si no hay overfitting)
39 - Usa BIC si tienes muchos datos y quieres un modelo parsimonioso
40 """)
41 # Determinar mejor modelo
42 if modelo_final_r2.bic < modelo_reducido.bic and modelo_final_r2.bic <
43     modelo_statsmodels.bic:
44     print("\n MEJOR MODELO: 4 variables seleccionadas por Pseudo R_dos (
45         menor BIC)")
46 elif modelo_reducido.bic < modelo_statsmodels.bic:
47     print("\n MEJOR MODELO: 4 variables seleccionadas por p-value (menor
48         BIC)")
49 else:
50     print("\n MEJOR MODELO: Modelo completo (menor BIC)")

```

Los resultados del modelo de la regresión logística es:

| Logit Regression Results | | | | | | |
|--------------------------|------------------|-------------------|-----------|-------|--------|--------|
| Dep. Variable: | HeartDisease | No. Observations: | 734 | | | |
| Model: | Logit | Df Residuals: | 729 | | | |
| Method: | MLE | Df Model: | 4 | | | |
| Date: | Tue, 21 Oct 2025 | Pseudo R-squ.: | 0.4272 | | | |
| Time: | 05:01:39 | Log-Likelihood: | -289.62 | | | |
| converged: | True | LL-Null: | -505.62 | | | |
| Covariance Type: | nonrobust | LLR p-value: | 3.384e-92 | | | |
| | coef | std err | z | P> z | [0.025 | 0.975] |
| const | 1.9466 | 0.385 | 5.057 | 0.000 | 1.192 | 2.701 |
| ST_Slope | -2.1776 | 0.206 | -10.585 | 0.000 | -2.581 | -1.774 |
| ExerciseAngina | 1.4802 | 0.232 | 6.392 | 0.000 | 1.026 | 1.934 |
| Sex | 1.5483 | 0.271 | 5.721 | 0.000 | 1.018 | 2.079 |
| ChestPainType | -0.6133 | 0.114 | -5.395 | 0.000 | -0.836 | -0.390 |

Figura 11: Resultados de la regresión

```

1      --- Comparación Final ---
2
3      Modelo Completo (11 variables):
4          Pseudo R_doS: 0.4732
5          AIC: 556.7672
6          BIC: 611.9493
7
8      Modelo con 4 vars (criterio p-value):
9          Pseudo R_doS: 0.4272
10         AIC: 589.2337
11         BIC: 612.2263
12         Variables: ['ST_Slope', 'Sex', 'ChestPainType', 'ExerciseAngina']
13
14     Modelo con 4 vars (criterio Pseudo R_do):
15         Pseudo R_doS: 0.4272
16         AIC: 589.2337
17         BIC: 612.2263
18         Variables: ['ST_Slope', 'ExerciseAngina', 'Sex', 'ChestPainType']
19
20     CRITERIOS DE SELECCION:
21
22     1. Pseudo R_dos: Mayor es mejor (más varianza explicada)
23     2. AIC (Akaike Information Criterion): Menor es mejor
24     3. BIC (Bayesian Information Criterion): Menor es mejor
25
26     AIC y BIC penalizan modelos con más parámetros. Son útiles para evitar
27     overfitting.
28     BIC penaliza más fuertemente que AIC.
29
30     RECOMENDACION:

```

```
30 - Si buscas INTERPRETABILIDAD: modelo con 4 variables
31 - Si buscas MAXIMA PRECISION: modelo completo (si no hay overfitting)
32 - Usa BIC si tienes muchos datos y quieres un modelo parsimonioso
33
34
35 MEJOR MODELO: Modelo completo (menor BIC)
```

4. CONCLUSIÓN

El presente trabajo permitió aplicar la **regresión logística** para modelar una variable binaria a partir de múltiples variables numéricas y categóricas. Mediante la interpretación de los coeficientes, errores estándar, valores z y p -values, se identificaron las variables independientes que ejercen mayor influencia sobre la probabilidad de presentar enfermedad cardíaca, destacando especialmente **ST_Slope**, **Sex** y **ExerciseAngina** como las más significativas.

La construcción de la matriz de confusión y la evaluación del modelo mediante métricas como **pseudo R-cuadrada**, **Log-Likelihood** y **LLR p-value** permitió comprobar que el modelo presenta un buen ajuste y una precisión (*accuracy*) del 85.87 %, considerada **muy buena**. Asimismo, la selección de modelos optimizados con un máximo de cuatro variables independientes evidenció la importancia de elegir las variables más explicativas, logrando modelos más eficientes sin sacrificar la capacidad predictiva.

Este ejercicio resalta cómo los datasets disponibles en plataformas como Kaggle[1, 3] pueden ser utilizados de manera práctica para la formación de modelos predictivos confiables, proporcionando información valiosa para la toma de decisiones basada en datos y reforzando la comprensión de técnicas estadísticas avanzadas aplicadas en problemas reales de salud.

Referencias

- [1] Kaggle. (n.d.). *Heart Disease UCI dataset*. <https://www.kaggle.com/ronitf/heart-disease-uci>
- [2] Hosmer, D. W., Lemeshow, S., & Sturdivant, R. X. (2013). *Applied logistic regression* (3rd ed.). Wiley.
- [3] James, G., Witten, D., Hastie, T., & Tibshirani, R. (2021). *An Introduction to Statistical Learning: with Applications in R* (2nd ed.). Springer.
- [4] Kuhn, M., & Johnson, K. (2019). *Feature engineering and selection: A practical approach for predictive models*. CRC Press.
- [5] McFadden, D. (1974). Conditional logit analysis of qualitative choice behavior. In P. Zarembka (Ed.), *Frontiers in Econometrics* (pp. 105–142). Academic Press.