



# Predicción de Enfermedades Cardiovasculares

Víctor Bórquez  
vitocode.ds

## Introducción

En este tutorial aprenderás a construir un modelo de **Machine Learning** que prediga el riesgo de enfermedades cardiovasculares usando un conjunto real de datos. Esto nos ayudará a entender cómo preparar datos, elegir un modelo adecuado, entrenarlo, evaluarlo y utilizarlo para hacer predicciones.

La predicción temprana de enfermedades cardiovasculares es fundamental para mejorar la prevención y el tratamiento. La **Inteligencia Artificial** y el análisis de datos ofrecen herramientas poderosas para apoyar esta tarea.

## Paso 1: Preparación

Para trabajar con los datos y modelos usaremos las librerías **pandas** para manipulación de datos, **scikit-learn** para Machine Learning y **matplotlib** y **seaborn** para visualización.

Instálalas si aún no las tienes con:

```
1 !pip install pandas scikit-learn matplotlib seaborn
```

## Paso 2: Carga y Exploración de Datos

Primero, sube el archivo `cardio_train.csv` y cárgalo en un `DataFrame` de **pandas**. Explorar los datos nos permite entender qué columnas contiene, qué tipo de datos hay y cómo están distribuidos.

```
1 import pandas as pd
2
3 data = pd.read_csv("cardio_train.csv", sep=';')
4 print(data.head())
5 print(data.info())
6 print(data['cardio'].value_counts(normalize=True))
```

La columna `cardio` es nuestra variable objetivo: indica si una persona tiene enfermedad cardiovascular (1) o no (0).

## Paso 3: Preprocesamiento de Datos

Antes de entrenar el modelo, limpiamos y preparamos los datos.

- Eliminamos la columna `id` porque no aporta información útil para predecir.
- Convertimos la edad de días a años para facilitar la interpretación.
- Eliminamos duplicados o filas con datos faltantes para evitar sesgos o errores.
- Codificamos las variables categóricas (como género, colesterol y glucosa) usando **One-Hot Encoding**, que crea columnas binarias para cada categoría. Esto es necesario porque los modelos matemáticos trabajan con números, no con etiquetas de texto.

```
1 # Limpieza
2 data = data.drop_duplicates()
3 data = data.dropna()
4 data.drop(['id'], axis=1, inplace=True)
5
6 # Transformacion
7 data['age'] = (data['age'] / 365).round().astype(int)
8
9 # One-Hot Encoding para variables categoricas
10 data = pd.get_dummies(data,
11                        columns=['gender', 'cholesterol', 'gluc'])
```

## Paso 4: Visualización de Datos

Visualizamos distribuciones y relaciones para entender mejor el conjunto.

- Histograma de la edad para ver qué rango domina.
- Mapa de calor de correlaciones para detectar variables que están relacionadas con la presencia de enfermedad.

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 sns.histplot(data['age'], kde=True)
5 plt.title('Distribuci n de la Edad')
6 plt.show()
7
8 corr = data.corr()
9 plt.figure(figsize=(6,6))
10 sns.heatmap(corr, cmap='coolwarm')
11 plt.title('Mapa de calor de correlaciones')
12 plt.show()
```

## Paso 5: Selección y Entrenamiento del Modelo

Elegimos la **Regresión Logística** porque es un modelo simple y fácil de entender que nos permite predecir la probabilidad de que un paciente tenga o no una enfermedad. Además, suele funcionar bien cuando queremos separar dos grupos (enfermos y sanos) y es rápido de entrenar.

- Separa tus datos en X (características) y y (objetivo).
- Normaliza las características para que tengan media 0 y desviación estándar 1, mejorando el entrenamiento.
- Divide en datos de entrenamiento y prueba.
- Entrena el modelo con los datos de entrenamiento.

```
1 from sklearn.preprocessing import StandardScaler
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LogisticRegression
4
5 X = data.drop('cardio', axis=1)
6 y = data['cardio']
7
8 scaler = StandardScaler()
9 X = scaler.fit_transform(X)
10
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
    =0.2, random_state=42)
12
13 clf = LogisticRegression(max_iter=1000)
14 clf.fit(X_train, y_train)
```

## Paso 6: Evaluación del Modelo

Evaluamos el desempeño con varias métricas:

- **Accuracy (Exactitud):** porcentaje de aciertos totales.
- **Precision (Precisión):** proporción de verdaderos positivos entre las predicciones positivas.
- **Recall (Sensibilidad):** proporción de verdaderos positivos detectados entre todos los positivos reales.
- **F1-Score:** medida que combina precisión y recall en un solo valor balanceado, destacando el equilibrio entre ambos.
- **Matriz de Confusión:** resumen de aciertos y errores clasificados por clase.

```
1 from sklearn.metrics import accuracy_score, classification_report,
  confusion_matrix
2
3 y_pred = clf.predict(X_test)
4
5 print("Accuracy:", accuracy_score(y_test, y_pred))
6 print("Reporte de Clasificación:\n", classification_report(y_test,
  y_pred))
7 print("Matriz de Confusión:\n", confusion_matrix(y_test, y_pred))
```

**Nota importante:** Un alto porcentaje de falsos negativos puede ser peligroso en medicina, porque significa no detectar a quienes realmente están enfermos. Aquí, es crucial balancear precisión y recall según el caso.

**¡Felicidades!** Has aprendido un flujo completo y fundamentado para construir un modelo de clasificación de machine learning desde datos crudos hasta predicción.

*Este tutorial es parte de vitocode.ds por Víctor Bórquez.*

