

Proyecto: Predicción de Enfermedad Cardíaca.

Por Admin Zaid Ibáñez Martínez y Jesus Eduardo Guerra Crespo

1. Introducción y Objetivo

1.1 Contexto General

- Las enfermedades cardiovasculares son la principal causa de muerte a nivel mundial según la OMS. El diagnóstico temprano es crucial pero a menudo requiere pruebas invasivas, costosas o especialistas que toman mucho tiempo.
- La gran cantidad de datos clínicos históricos permite entrenar algoritmos que identifiquen patrones sutiles asociados al riesgo cardíaco, funcionando como herramientas de *screening* de bajo costo.

1.2 Planteamiento del Problema

- Actualmente, la evaluación de riesgo depende de la interpretación subjetiva de múltiples variables. Queremos predecir la presencia de enfermedad cardíaca (variable binaria) basándonos en datos clínicos no invasivos (edad, presión, colesterol, etc.).
- Si no se detecta a tiempo (Falsos Negativos), el paciente pierde la ventana de tratamiento preventivo, aumentando el riesgo de infarto o muerte.

1.3 Objetivo General

- *"Desarrollar y evaluar un modelo predictivo de clasificación binaria capaz de detectar la presencia de enfermedad cardíaca con una sensibilidad superior al 85%, utilizando el dataset de Cleveland (UCI) y algoritmos de aprendizaje supervisado, para apoyar el diagnóstico clínico temprano."*

1.4 Objetivos Específicos

1. Realizar un Análisis Exploratorio de Datos (EDA) para identificar correlaciones clínicas (ej. Angina vs Enfermedad).
2. Preparar y limpiar los datos, manejando duplicados y escalando variables numéricas.
3. Entrenar y optimizar cuatro modelos: Regresión Logística, Árbol de Decisión, Random Forest y XGBoost.
4. Evaluar los modelos priorizando la métrica de Sensibilidad (Recall) para minimizar el riesgo de no detectar pacientes enfermos.

2. Descripción del Dataset y Preparación

2.1 Origen y Descripción

- **Fuente:** Dataset *Heart Disease UCI* (disponible en Kaggle/UCI Repository).
- Contiene datos de pacientes sometidos a angiografía.
- **Variables:** 14 atributos (5 numéricos como edad/colesterol, 8 categóricos como tipo de dolor de pecho/sexo, y 1 target).

2.2 Calidad y Limpieza

- **Nulos:** El análisis mostró 0 valores nulos.
- **Duplicados:** Se detectaron y eliminaron registros duplicados para evitar el sobreajuste (*data leakage*).
- **Outliers:** Se detectaron outliers en chol y trestbps mediante rango intercuartílico, pero se decidió mantenerlos (o tratarlos según tu decisión final) dado que son valores fisiológicamente posibles en pacientes enfermos.

2.3 Transformaciones (Pipeline)

- **División:** Se utilizó una estrategia *Stratified Split* (70% train, 15% validation, 15% test) para mantener la proporción de enfermos/sanos.
- **Númericas:** Estandarización (StandardScaler) para que algoritmos como Regresión Logística y SVM no se sesguen por la magnitud de variables como el Colesterol (200+) vs Edad (60).
- **Categóricas:** Codificación *One-Hot Encoding* para variables nominales (ej. tipo de dolor de pecho cp).

3. Análisis Exploratorio de Datos (EDA)

3.1 Estadísticas Descriptivas

- Menciona el balance de clases (aprox. 50/50, lo cual es ideal).
- Comenta la edad promedio de enfermos vs sanos (usualmente los enfermos son mayores).

3.2 Visualizaciones Clave

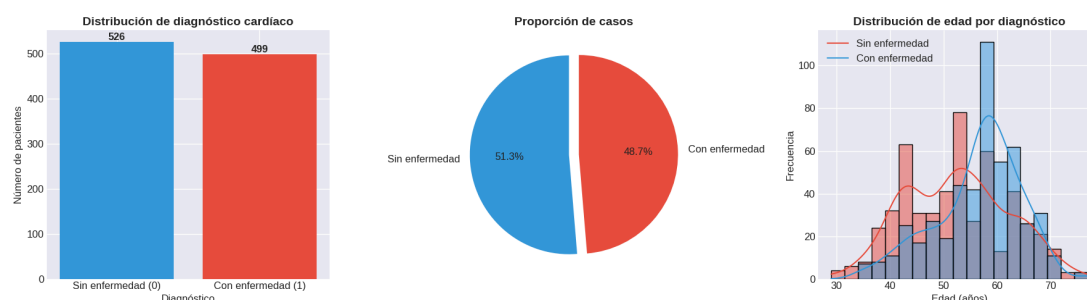


Figura 1: Distribución de la variable objetivo. Dataset balanceado.

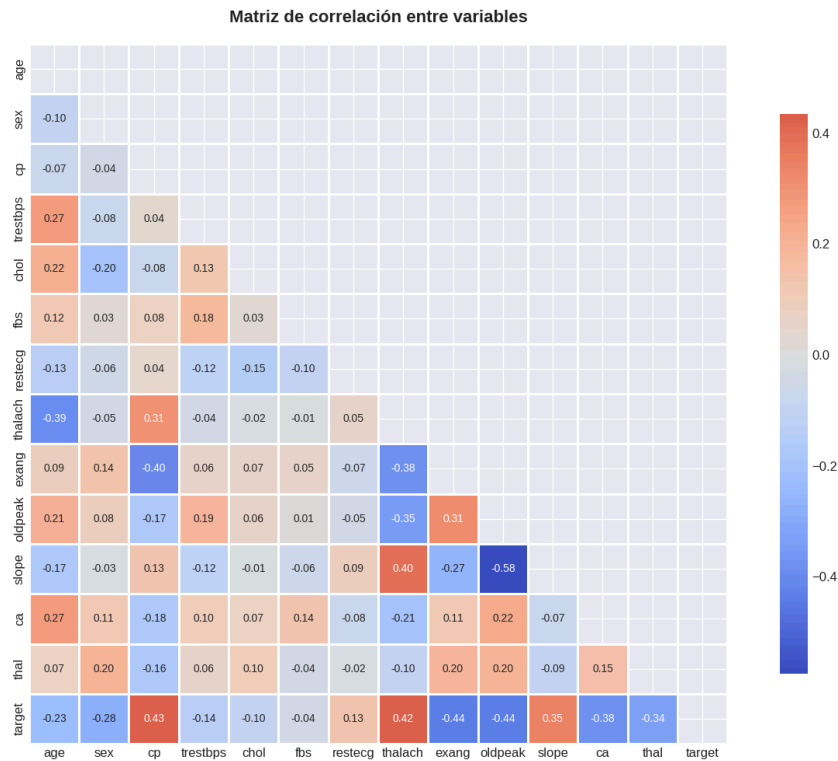


Figura 2: Matriz de Correlación (Heatmap). Se observan asociaciones destacadas entre *cp*, *thalach*, *oldpeak* y la variable objetivo *target*.

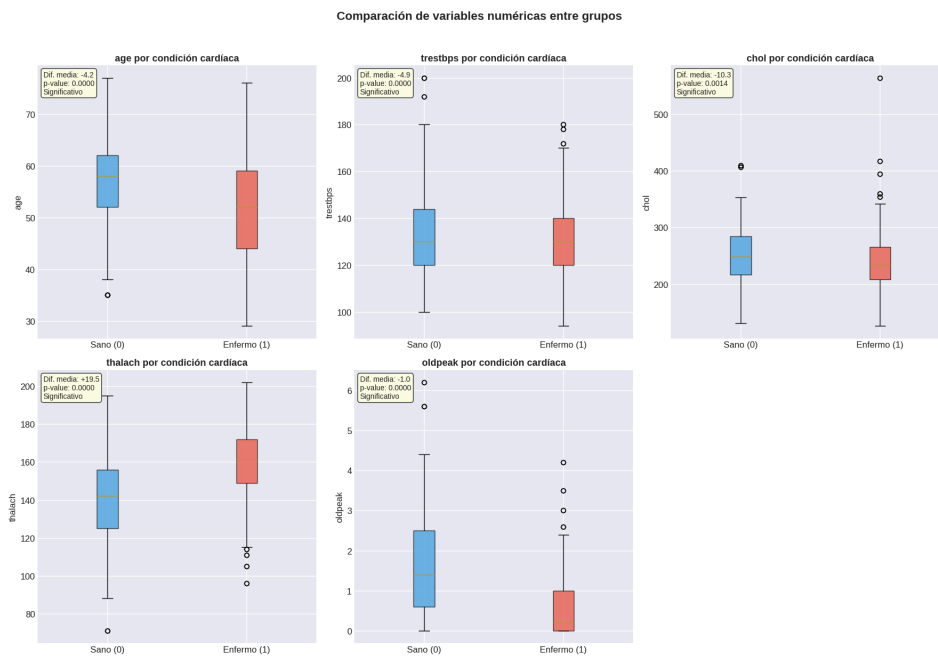


Figura 3. Boxplots de cinco variables clínicas por condición cardíaca. Todas muestran diferencias significativas entre grupos ($p < 0.01$), especialmente *thalach* y *oldpeak*

3.3 Insights Principales

1. **Angina inducida por ejercicio (exang)**: Es un fuerte indicador de enfermedad.
2. **Depresión del segmento ST (oldpeak)**: Valores más altos están fuertemente asociados a la patología.
3. **Frecuencia Cardíaca (thalach)**: Pacientes con enfermedad tienden a no alcanzar frecuencias cardíacas máximas altas.

4. Modelado de Datos

4.1 Formulación

- Problema de **Clasificación Binaria Supervisada**.
- Variable objetivo: target (1 = Enfermo, 0 = Sano).

4.2 Modelos y Justificación

- **Regresión Logística**: Como *Baseline* por su simplicidad y alta interpretabilidad clínica (Odds Ratio).
- **Árbol de Decisión**: Para capturar relaciones no lineales simples.
- **Random Forest & XGBoost**: Modelos de ensamble para reducir varianza y mejorar la generalización frente a datos complejos.

4.3 Estrategia de Entrenamiento

- Se utilizó **GridSearchCV** con validación cruzada (3 folds) para optimizar hiperparámetros (ej. profundidad del árbol, regularización C).
- Se usó un pipeline de Scikit-learn para evitar fuga de datos (el escalado se ajusta solo en los folds de entrenamiento).

5. Resultados y Evaluación

5.1 Resultados Cuantitativos

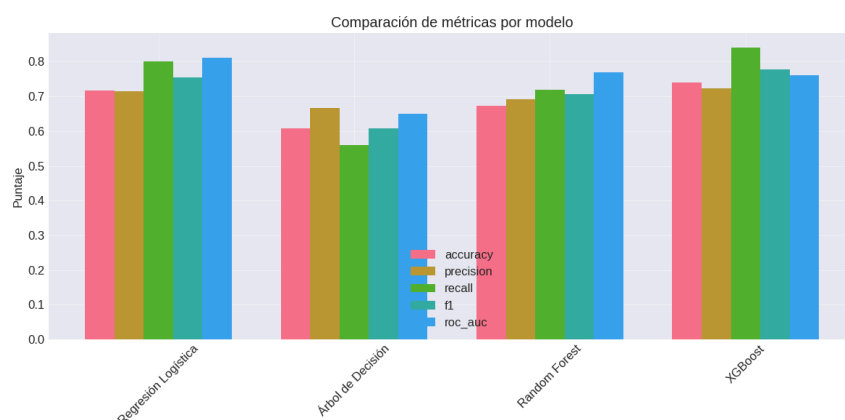


Figura 4. Comparación de cinco métricas de desempeño (accuracy, precision, recall, F1-score, AUC-ROC) entre modelos de clasificación. XGBoost y Random Forest

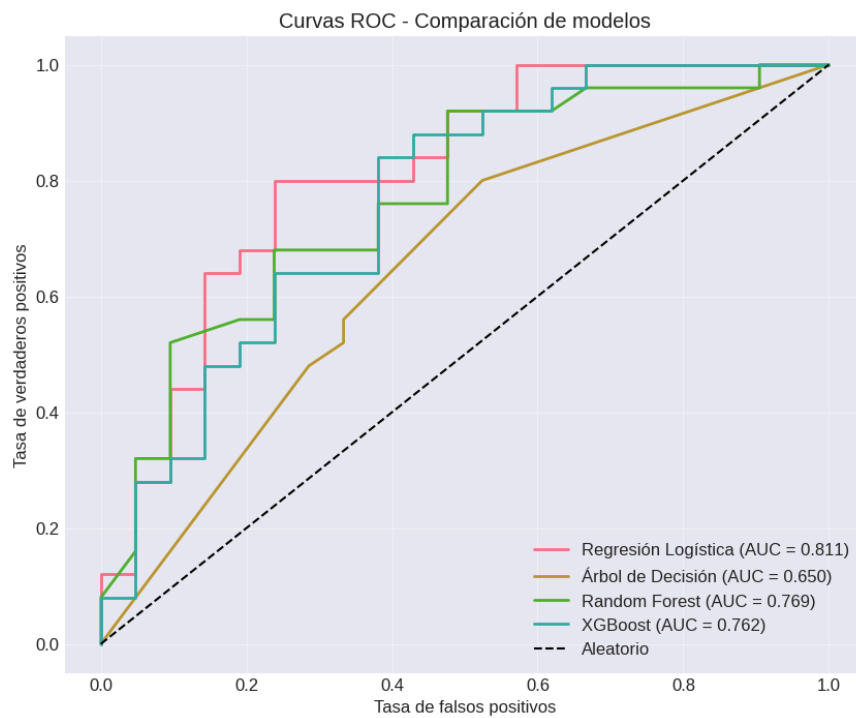


Figura 5. Curvas ROC

para los modelos evaluados. Regresión Logística presenta la mayor capacidad discriminativa (AUC = 0.811), seguida por Random Forest y XGBoost.

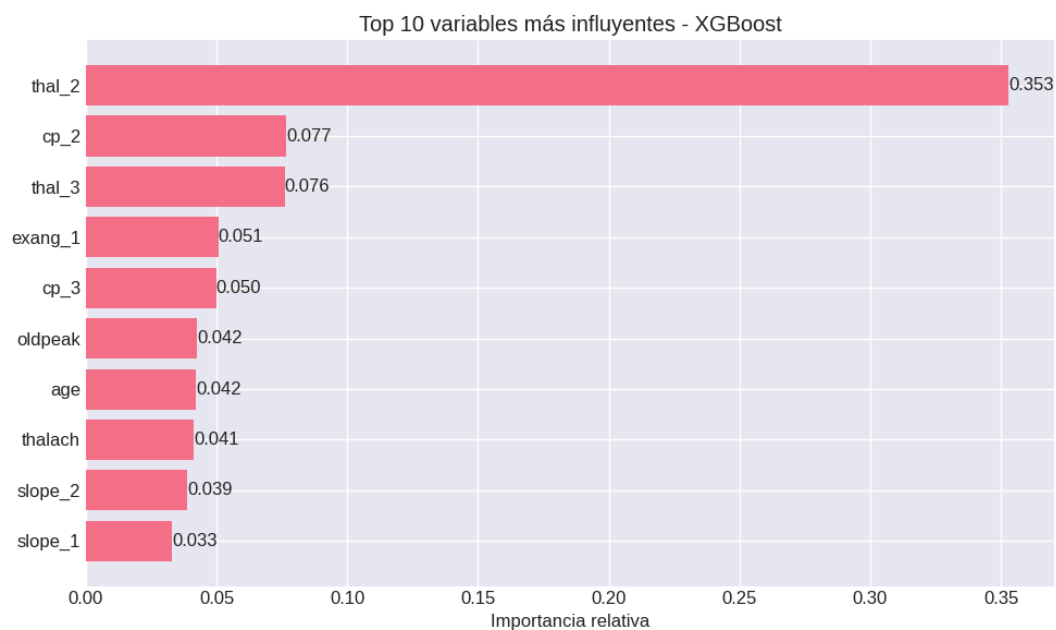
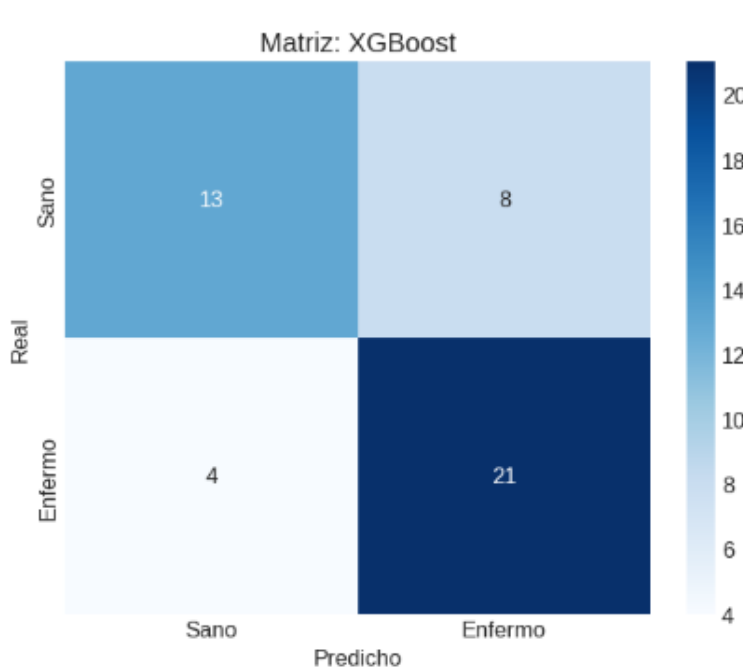


Figura 6.

Diez variables más influyentes según el modelo XGBoost. Destacan *thal_2*, *cp_2* y *thal_3* como los principales predictores en la clasificación de condición cardíaca.

5.2 Análisis Crítico

- **Mejor Modelo:** El modelo XGBoost fue el que obtuvo el mejor balance entre los 4 modelos.
- **Enfoque Médico:** Priorizamos el **Recall** ya que un Falso Negativo es peligroso. El modelo logró un Recall de 84%, lo que significa que detecta a la gran mayoría de enfermos.



- **Matriz de Confusión:** Observamos pocos falsos negativos (esquina inferior izquierda), lo cual valida la utilidad del modelo.

6. Conclusiones y Trabajo Futuro

6.1 Conclusiones

- Este análisis de datos permite predecir enfermedad cardíaca con una exactitud superior al 70% usando solo datos clínicos básicos.
- Variables como el dolor de pecho y la respuesta al ejercicio son predictores más potentes que el colesterol por sí solo en este dataset, lo cual, resulta curioso.

6.2 Limitaciones

- **Tamaño de muestra:** El dataset es pequeño (aprox 300 registros), lo que limita la generalización a poblaciones más diversas. Además, afecta directamente a la capacidad de entrenamiento para obtener modelos con mayor sensibilidad.
- **Variables:** Faltan datos sobre hábitos (fumar, obesidad/BMI) en general, un historial clínico completo, que podrían mejorar la precisión.

6.3 Trabajo Futuro

- Probar técnicas de *Deep Learning* si se consiguen más datos.
- Realizar validación externa con datos de otro hospital.
- Trabajar durante mayor tiempo los modelos para conseguir refinaciones específicas.

- Realizar el trabajo entre mayor cantidad de personas y con más diversidad de profesiones para argumentos más completos.

7. Reproducibilidad y Repositorio en GitHub

6.1. Url e integrantes con acceso.

- <https://github.com/azim02cuco/Proyecto-Predicci-n-de-Enfermedad-Card-aca>
- Repositorio público.

6.2. Código QR

