

# PRÁCTICA 2 - MACHINE LEARNING

## DOCUMENTACIÓN

Inteligencia Artificial

Yago Iglesias Díaz



# Índice

1. Contexto
  - 1.1. Descripción general del proyecto
2. Problema a resolver
  - 2.1. Cuestiones objetivo
  - 2.2. Variable fundamental
3. Estrategia seguida
  - 3.1. Análisis del conjunto de datos
4. Recursos usados
  - 4.1. Librerías
  - 4.2. Algoritmos
5. Objetivos esperados y logrados
  - 5.1. Objetivos esperados
  - 6.2. Objetivos logrados
6. Conclusión

# 1. Contexto

## 1.1 Descripción general del proyecto

El mundo desde hace ya varias décadas ha sido cada vez más monótono, debido a las comodidades que se van generando con los años ya sea por la televisión, el internet, los trabajos de oficina... Y todo esto a pesar de lo bueno que trae consigo, también tiene su parte negativa, así como las personas están cada vez más ocupadas en actividades sin moverse y no pueden dedicar tiempo a sí mismas, lo que genera un estrés físico y una sociedad sedentaria a la que no le merece la pena salir de su casa. Como resultado, la enfermedad cardíaca se ha convertido en uno de los factores más influyentes de muerte de hombres y mujeres.



Las enfermedades cardiovasculares (ECV) son la principal causa de muerte a nivel mundial, cobrando un estimado de 17,9 millones de vidas cada año, lo que representa el 31% de todas las muertes en todo el mundo. Cuatro de cada 5 muertes por ECV se deben a ataques cardíacos y accidentes cerebrovasculares, y un tercio de estas muertes ocurren prematuramente en personas menores de 70 años.

## 2. Problema a resolver

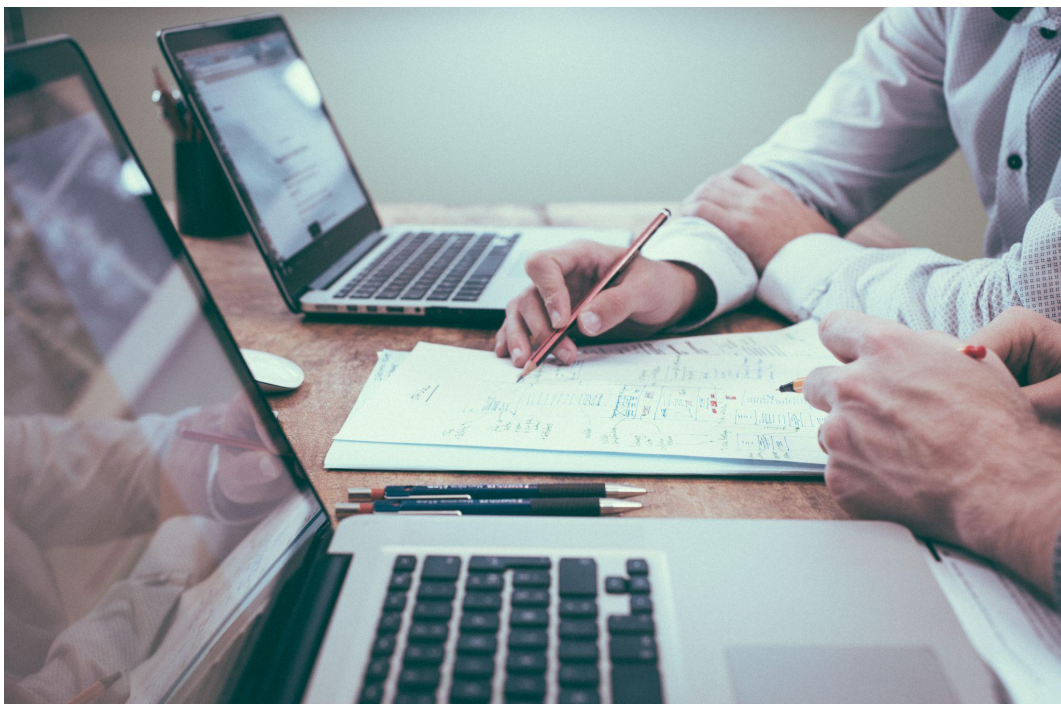
Las personas con enfermedad cardiovascular o que se encuentran en alto riesgo cardiovascular (debido a la presencia de uno o más factores de riesgo como hipertensión, diabetes, hiperlipidemia o enfermedad ya establecida) necesitan una detección y manejo precoces donde un modelo de aprendizaje automático puede ser de gran ayuda. De ahí que el principal objetivo sea predecir si alguien tiene un alto riesgo de ser diagnosticado como un paciente cardíaco.

La insuficiencia cardíaca es un evento común causado por las enfermedades cardiovasculares y este conjunto de datos contiene once características que pueden usarse para predecir una posible enfermedad cardíaca. Tanto los modelos de aprendizaje lineal como automático se utilizan para predecir la insuficiencia cardíaca en función de diversos datos como entradas.

### 2.1 Cuestiones objetivo

En este estudio, se intentará resolver las siguientes cuestiones:

- ¿Cuál es la probabilidad de que una persona en específico pueda tener este problema?
- ¿Cuáles son los principales indicadores de que una persona pueda padecer una de estas enfermedades?
- ¿Qué políticas o estrategias pueden adoptarse a partir de los resultados para prevenir las enfermedades cardiovasculares?



## 2.2 Variable fundamental

Dado que se dispone de gran variedad de datos recolectados de personas referentes a la salud cardiovascular, se pretende predecir lo expuesto con anterioridad. Para ello, hay una variable objetivo, la cual es "HeartDisease" que funciona de la siguiente manera, si una persona tiene insuficiencia cardíaca se representa con un 1 y si no la tiene se hace con un 0. Esta variable es el objetivo del estudio, averiguar la probabilidad de padecer problemas del corazón.

## 3. Estrategia seguida

La estrategia seguida se ha compuesto de varios puntos elegidos estratégicamente para ir poco a poco, asentando los conocimientos y la cantidad de datos que se han proporcionado. Para ello, primeramente se ha desarrollado un análisis individual para las características del conjunto de datos, y más tarde se ha expuesto la información relevante en un pequeño resumen. Por otro lado, después ya se ha comenzado con el desarrollo de los algoritmos de machine learning como mas adelante se explicará. Y por último lugar, se llevó a cabo el desempeño de las conclusiones que se han obtenido de hacer dicho estudio.

### 3.2 Análisis del Conjunto de Datos

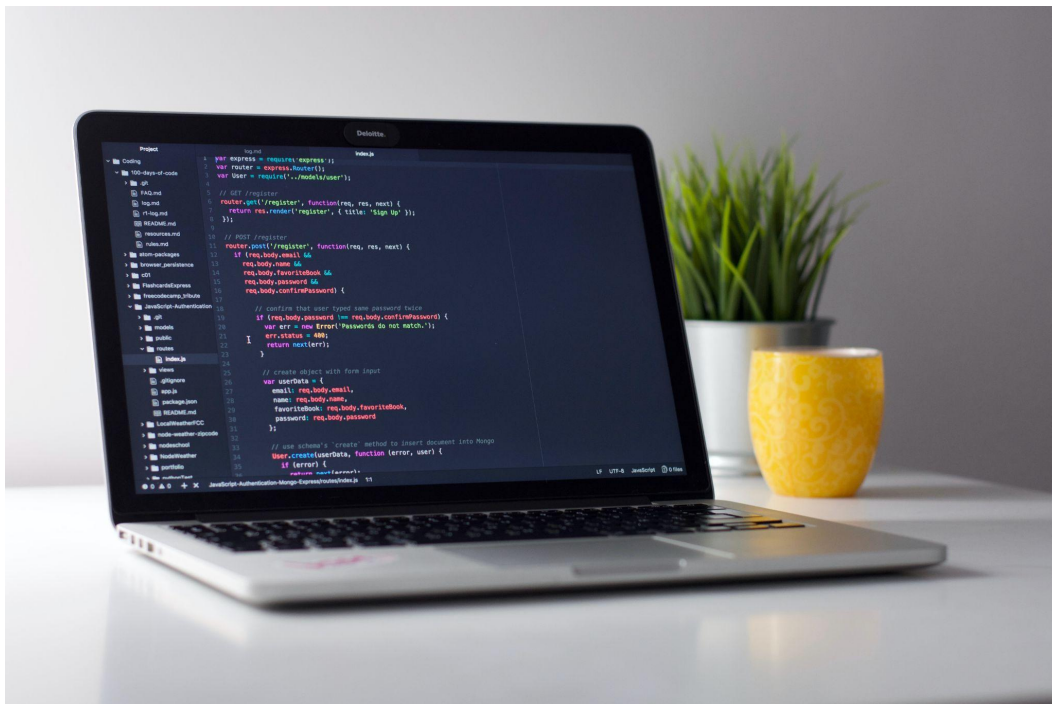
En este caso de estudio, se cuenta con datos extraídos de sitio web Kaggle, donde a través de esta url "<https://www.kaggle.com/fedesoriano/heart-failure-prediction>" se han podido adquirir la información necesaria para desarrollar este análisis, contando con 918 personas con diversos datos sobre ellos. Y todos ellos serán usados para la predicción de que una persona pueda tener problemas de insuficiencia cardíaca.

Este conjunto de datos, creado por el usuario "fedesoriano", contiene un total de 35.92 kB de almacenamiento de información en formato CSV, 145 códigos públicos creados por la comunidad en base a este, y fue subido a la plataforma Kaggle en septiembre del año 2021.



## 4. Recursos usados

En este estudio se han utilizado alguna variedad de recursos para la correcta realización de este, ya que sin estos recursos, como pueden ser librerías y algoritmos de machine learning no es posible hacer el análisis. En primer lugar se presentarán las librerías y mas tarde los algoritmos utilizados.

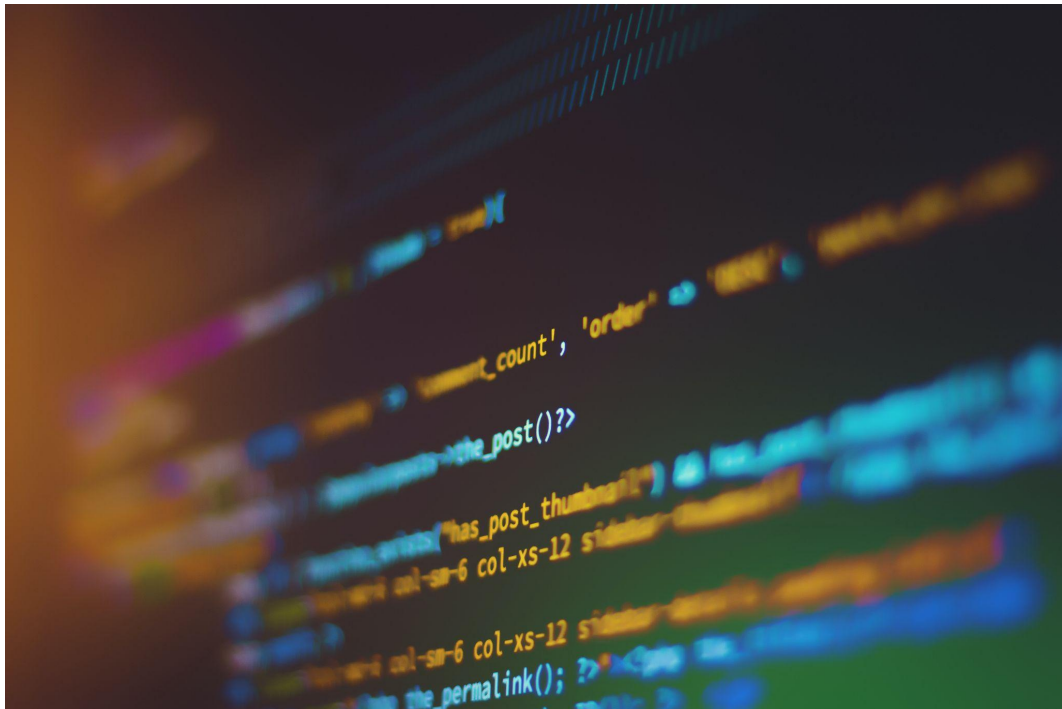


### 4.1 Librerías

En cuanto a las librerías utilizadas para el funcionamiento del código han destacado las siguientes:

- **Pandas:** es una librería de Python especializada en el manejo y análisis de estructuras de datos. Con él se han realizado visualizaciones gráficas que nos permiten determinar entre otros la dependencia y el grado de correlación entre variables del modelo.
- **Numpy:** da soporte para crear vectores y matrices grandes multidimensionales, junto con una gran colección de funciones matemáticas de alto nivel para operar con ellas.
- **Sklearn:** es un conjunto de rutinas escritas en Python para hacer análisis predictivo, que incluyen clasificadores, algoritmos de clusterización, etc. Está basada en NumPy, SciPy y matplotlib, de forma que es fácil reaprovechar el código que use estas librerías.

- SciPy: contiene módulos para optimización, álgebra lineal, integración, interpolación, funciones especiales, FFT, procesamiento de señales y de imagen, resolución de ODEs y otras tareas para la ciencia e ingeniería.
- Matplotlib: es una librería de Python especializada en la creación de gráficos en dos dimensiones.
- IPython: es un shell interactivo que añade funcionalidades extra al modo interactivo incluido con Python, como resaltado de líneas y errores mediante colores.
- Plotly: es otro paquete de gráficos en R para crear gráficos interactivos con calidad de publicación.
- Os: Este módulo provee una manera versátil de usar funcionalidades dependientes del sistema operativo.



## 4.2 Algoritmos

En cuanto a los algoritmos implementados en este estudio se han utilizado primero una serie de algoritmos de base (utilizando hiperparámetros de fábrica) antes de pasar a soluciones más sofisticadas.

Más tarde se realiza el desarrollo de los algoritmos considerados en este análisis: Regresión logística, Bosque aleatorio, SVM y KNN.

- La regresión logística: es un algoritmo de clasificación que se utiliza para predecir la probabilidad de una variable dependiente categórica. Permite decir que la presencia de un factor de riesgo aumenta la probabilidad de un resultado dado un porcentaje específico.
- Los Bosques Aleatorios: son un algoritmo de aprendizaje supervisado que, como ya se puede ver en su nombre, crea un bosque y lo hace de alguna manera aleatorio. Para decirlo en palabras simples: el Bosque Aleatorio crea múltiples árboles de decisión y los combina para obtener una predicción más precisa y estable.
- Support vector machine (SVM): es un algoritmo de aprendizaje supervisado que se utiliza en muchos problemas de clasificación y regresión.
- KNN: trabaja buscando las distancias entre una consulta y todos los ejemplos en los datos, seleccionando el número especificado ejemplos (K) más cercanos a la consulta, luego vota por la etiqueta más frecuente (en el caso de la clasificación) o promedia las etiquetas (en el caso de la regresión).

## 5. Objetivos esperados y logrados

### 5.1 Objetivos esperados

Los objetivos esperados en este proyecto eran varios, como ya se mencionó anteriormente en el punto uno.

Estas metas del estudio destacaban por no ser muy ambiciosas, queriendo predecir la probabilidad de que una persona en específico pueda tener enfermedades cardiovasculares, así como averiguar cuáles son los principales indicadores de que una persona pueda padecer una de estas enfermedades. Por otro lado, también se pretendía averiguar qué políticas o estrategias pueden adoptarse a partir de los resultados para prevenir estos problemas cardiovasculares.

### 5.2 Objetivos logrados

En cuanto a los objetivos logrados, se puede afirmar que se han alcanzado casi la totalidad de los esperados, teniendo algunos huecos en la información minúsculos, como puede ser el caso de las políticas o estrategias a adoptar para prevenir las enfermedades, donde en el estudio se han expuesto algunas, pero no respaldadas lo suficientes por estudios científicos.



Pero por otro lado se han conseguido una gran cantidad de objetivos, así como averiguar la probabilidad de que una persona en específico pueda tener enfermedades cardiovasculares, o saber cuáles son los principales indicadores de que una persona pueda padecer una de estas enfermedades, como por ejemplo la prueba del esfuerzo, la cual es el mayor indicador según los algoritmos desarrollados de que una persona pueda llegar a tener una enfermedad cardíaca.



## 6. Conclusión

El estudio ha resultado exitoso y útil para personas que le interese el tema o lo necesite para un fin de salud. Pero también hay que manifestar que el análisis está muy lejos de ser perfecto y a medida que se vayan generando más datos sobre personas diagnosticadas (sobre las nuevas personas enfermas y sanas), el algoritmo puede volver a entrenarse utilizando los datos adicionales y, en teoría, generar predicciones más precisas para identificar a los enfermos del corazón.

Debe elaborarse un "Plan de prevención ante las enfermedades del corazón" estratégico para cada grupo de categorías de riesgo. Además de los pasos sugeridos en el documento de Python para cada característica del conjunto de datos.

Como parte de la solución, pueden iniciarse eventos o anuncios publicitarios por parte del gobierno o empresas para concienciar a las personas a cuidarse más o realizarse con más frecuencia las pruebas/características comentadas anteriormente. Asimismo, una mayor implicación por todas las personas en cuanto a la salud se refiere vendría perfecto para combatir a una sociedad que poco a poco se está volviendo más sedentaria.