



# Competencia en Plataforma de Ciencia de Datos: Kaggle (Caso Integrador 2024)

# Trabajo práctico en forma Individual

1\_ Primero estudiar y realizar la ejercitación de los videos de las Clases 1, 2 y 3



\* Videos de Clases y Material Importante para estudiar ✍

2\_ Resolver el mini caso y mandar la solución a la plataforma kaggle

## Clasificador de Riesgo de Enfermedad Coronaria

Riesgo de Enfermedad Coronaria

**Link de la competencia:**

<https://www.kaggle.com/t/4746c883bac144fdb1ed235a7832f87>

**(\*ENTREGAR EL DÍA: 30/05/2024 hasta las 11:00(am) hay tiempo)**

**\*\*\*Recordar hacer las prácticas en forma INDIVIDUAL\*\*\***

**\*\*\*La evaluación es individual\*\*\***

# Para Leer los archivos: train, test

**\*1\_ Pasar los archivos de la competencia a Drive**



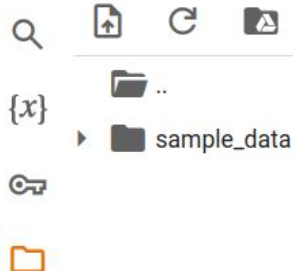
datos\_train.csv  
nuevas\_instancias\_clasificar.csv



Prueba\_Competencia\_Enviar\_Datos.ipynb

Archivo Editar Ver Insertar Entorno de ejecución Hi

Archivos



**\*3\_ Conectar con Drive**

**\*4\_ Ejecutar**

```
from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

**\*5\_ Copiar ruta de drive y colocarla en pd.read\_csv()**

```
pandas as pd
datos1 =
pd.read_csv('/content/drive/MyDrive/Gestion_2024/ruta_en_drive/da
tos_train.csv')
datos1.head(2)
```

**\*2\_ En colab seleccionar**



# Reglas y Puntaje

## GAMIFICATION

Es el proceso de utilizar mecánicas de juego para incentivar los deseos naturales de las personas:



**1\_ Se debe trabajar en forma individual**

**2\_ La evaluación es individual**

**3\_ Hay tiempo hasta el 30/05/2024 11:00 a.m. para mandar la solución del caso a kaggle.**

**4\_ Se permiten 4 intentos de envío por día en la plataforma de ciencia de datos.**

**5\_ Cumplir con el flujo de trabajo que se muestra en las filminas siguientes.**

# Flujo de Trabajo para Caso Integrador Kaggle

1\_ Aplicar lo visto en las clases: 1,2 y 3.

2\_ Seguir el método



\* Videos de Clases y Material Importante para estudiar



1\_Clases\_2024.zip



2\_material\_de\_lectura\_práctica\_U4.zip



3\_tener\_en\_cuenta.zip

## Método



Pipelines

Utilizar Distintos Algoritmos de Minería de Datos. Definir línea base.

Optimización de Hiper Parámetros

Calibrar: Poner a competir a todos los métodos y posibles parámetros para determinar cual da el menor error

Aplicar métodos de consenso y de potenciación

1\_ Importar las librerías

2\_ Obtener los datos y Analizarlos

3\_ Realizar el preprocesamiento de los Datos

4\_ Dividir el Conjunto de Datos

5\_ Crear el Modelo

6\_ Entrenar, Calibrar y Validar el Modelo

7\_ Probar el Modelo

8\_ Realizar Predicciones

# Flujo de Trabajo para Caso Integrador Kaggle

**La Notebook utilizada para enviar la solución a la competencia de Kaggle debe tener en cuenta:**

**Análisis y exploración de datos.**

**El preprocesamiento:**

- Variables numéricas y categóricas**

- Atributos importantes**

- Valores nulos, faltantes o atípicos**

- Desbalanceo de clases**

- (Estudiar la clase 3)**

**División de los datos. Si es posible utilizar validación cruzada.**

**Métodos y algoritmos utilizados para entrenar y obtener el modelo. Comparar varios modelos utilizando métricas.**

**Optimización de hiperparámetros. Calibración del modelo.**

**Validación y prueba del modelo**

**Utilizar las siguientes métricas: Matriz de confusión, classification report(recall, F1), kappa, curva roc.**

**Integrar la teoría con la práctica**

**Evitar el Overfitting y el Underfitting**