



## Proyecto Titanic

Este proyecto trabajé en el análisis del famoso *Titanic* dataset, el cual incluye información sobre los pasajeros del transatlántico que naufragó en 1912. El dataset contiene variables como la clase del pasajero, el género, la edad, la cantidad de familiares a bordo. Busqué explorar cómo estas variables influyen en la supervivencia de los pasajeros. Además, el análisis explora patrones como la regla de "mujeres y niños primero", investigando si realmente se cumplió y cómo varió la tasa de supervivencia según el género y la edad.

El enfoque del proyecto combina un análisis exploratorio de datos (EDA) con la construcción de un modelo predictivo utilizando un algoritmo de Random Forest. Se realiza una limpieza profunda de los datos, se identifican valores nulos y duplicados, y se transforman las variables categóricas en valores numéricos para facilitar el entrenamiento del modelo. También se generan gráficos descriptivos para visualizar la distribución de edades y la supervivencia en función del género, destacando patrones que surgieron durante el hundimiento del Titanic.

Este proyecto es el de un estudiante en formación en *Data Science* y *Machine Learning*, con sólidos conocimientos en análisis de datos y programación en Python. Tengo experiencia trabajando con librerías como *pandas*, *scikit-learn* y *seaborn*, enfocado en desarrollar mis habilidades en la implementación de modelos predictivos y la interpretación de datos históricos.

El proyecto refleja mi capacidad para resolver problemas reales mediante el uso de datos y modelos estadísticos.

## Contaré como trabaje el proyecto paso a paso

### **Paso 1: Mostrar primeras líneas**

Comencé mi proyecto cargando el dataset del Titanic en un DataFrame de pandas. Para asegurarme de que los datos se habían cargado correctamente, utilicé el método `head()` para visualizar las primeras filas del dataset. También verifiqué las dimensiones del dataset con `shape` para conocer la cantidad de filas y columnas que contenía.

### **Paso 2: Valores duplicados y nulos**

A continuación, analicé la calidad de los datos. Utilicé `duplicated().sum()` para contar cuántas filas duplicadas había en el dataset y `isnull().sum()` para identificar cuántos valores faltantes había en cada columna. Esta etapa fue crucial para la limpieza del dataset antes de realizar un análisis más profundo.

### **Paso 3: Valores únicos**

Luego, realicé un análisis de los valores únicos en las columnas. Usé `nunique()` para contar cuántos valores únicos había en cada columna y `unique()` para mostrarlos. Esta etapa me permitió entender mejor la estructura de los datos y las características de cada variable.

### **Paso 4: Análisis exploratorio de datos (EDA)**

Después, realicé un análisis exploratorio de datos para comprender mejor las relaciones entre las características de los pasajeros y su probabilidad de supervivencia. Creé gráficos de barras y histogramas utilizando `matplotlib` y `seaborn` para visualizar la distribución de la edad y la supervivencia de los pasajeros.

### **Paso 5: Preparación de las columnas predictoras y objetivo**

Procedí a preparar el dataset para el modelado eliminando las columnas que no serían útiles para la predicción, como "Cabin", "Fare", "Ticket" y "Name". Utilicé `drop()` para hacerlo. Luego, separé las variables predictoras (X) de la variable objetivo (y).

### **Paso 6: Conversión de datos y creación del modelo**

Para manejar las columnas categóricas, utilicé `OrdinalEncoder` para convertirlas en numéricas. Además, rellené los valores nulos con `SimpleImputer`. Después, dividí el dataset en conjuntos de entrenamiento y prueba usando `train_test_split`, lo que me permitió validar el modelo en datos no vistos.

### **Paso 7: Creación del modelo y predicción**

Creé un modelo de clasificación utilizando `RandomForestClassifier`. Entrené el modelo con los datos de

entrenamiento y realicé predicciones sobre el conjunto de prueba. Para medir el rendimiento del modelo, utilicé `accuracy_score` para calcular la precisión del modelo.

### Paso 8: Predicción con nuevos datos

Finalmente, implementé una función que permite realizar predicciones introduciendo los datos de un nuevo pasajero. Esta función acepta parámetros correspondientes a las variables predictoras y utiliza el modelo previamente entrenado para predecir si el pasajero habría sobrevivido o no.

The screenshot shows a Jupyter Notebook titled "titanic\_notebook.ipynb" in Visual Studio Code. The notebook is divided into three sections:

- Librerías para el manejo de datos:** This section contains the following code:

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

from sklearn.metrics import accuracy_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import OrdinalEncoder
from sklearn.model_selection import train_test_split
```
- Cargo el Dataset:** This section contains the following code:

```
# Cargo el Dataset
dataset = pd.read_csv('titanic.csv')
```
- Muestro las primeras 5 líneas:** This section contains the following code:

```
dataset.sample(5)
```

The output of the third section shows a sample of 5 rows from the Titanic dataset:

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	
314	315	0	2	Hart, Mr. Benjamin	male	43.0	1	1	F.C.C. 13529	26.2500	NaN	S
468	469	0	3	Scanlan, Mr. James	male	NaN	0	0	36209	7.7250	NaN	Q
18	19	0	3	Vander Planke, Mrs. Julius (Emelia Maria Vande...)	female	31.0	1	0	345763	18.0000	NaN	S
566	567	0	3	Stoytcheff, Mr. Ili	male	19.0	0	0	349205	7.8958	NaN	S
391	392	1	3	Jansson, Mr. Carl Olof	male	21.0	0	0	350034	7.7958	NaN	S

The bottom of the image shows the Windows taskbar with the date 7/10/2024 and time 23:03.

```
titanic_prediccion.py - Proyecto Titanic - Visual Studio Code

titanic_prediccion.py > ...
1 import pandas as pd
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.impute import SimpleImputer
4 from sklearn.preprocessing import OrdinalEncoder
5 from sklearn.model_selection import train_test_split
6
7
8
9
10 # ----- Cargar y explorar el dataset -----
11 # Cargar el dataset del Titanic
12 df_titanic = pd.read_csv('titanic.csv')
13
14
15 # Verificar si hay valores duplicados
16 duplicados = df_titanic.duplicated().sum()
17 print(f"Cantidad de filas duplicadas: {duplicados}")
18
19
20 # Eliminar filas duplicadas si existen
21 if duplicados > 0:
22     df_titanic = df_titanic.drop_duplicates()
23     print("Duplicados eliminados.")
24
25
26 # Verificar si hay valores nulos en cada columna
27 print("Valores nulos por columna antes de preprocesar:\n", df_titanic.isnull().sum())
28
29
30
31
32 # ----- Preprocesamiento de datos -----
33 # Eliminar columnas irrelevantes para el análisis o predicción
34 df_titanic = df_titanic.drop(columns=["Cabin", "Fare", "Ticket", "Name"])
35
36
37
38 # Separar las variables predictoras (X) de la variable objetivo (y)
39 X = df_titanic.drop("Survived", axis=1)
40 y = df_titanic["Survived"]
41
42
43
44
45 # ----- Conversión de columnas categóricas a numéricas -----
46 # Identificar las columnas categóricas (variables no numéricas)
47 columnas_categoricas = X.select_dtypes(include=["object"]).columns
48
49
```