

Proyecto MAT-281

Titanic Espacial

Martin Alonso Urrea Barros

Universidad Tecnica Federico Santa Maria

Lunes 4 de Diciembre 2023

Tabla de Contenidos

- 1 Definición del Problema
 - Contexto
 - Carga de Datos
- 2 Estadística Descriptiva
- 3 Visualización Descriptiva
- 4 Preprocesamiento
 - Ingeniería de Atributos
 - Pre-procesamiento
- 5 Modelos
 - Modelo para Datos Estandarizados
 - Modelo para Datos NO Estandarizados
- 6 Métricas y Análisis de Resultados
- 7 Conclusiones

Tabla de Contenidos

- 1 Definición del Problema
 - Contexto
 - Carga de Datos
- 2 Estadística Descriptiva
- 3 Visualización Descriptiva
- 4 Preprocesamiento
 - Ingeniería de Atributos
 - Pre-procesamiento
- 5 Modelos
 - Modelo para Datos Estandarizados
 - Modelo para Datos NO Estandarizados
- 6 Métricas y Análisis de Resultados
- 7 Conclusiones

Contexto

- La nave espacial Titanic con casi 13.000 pasajeros a bordo. En dirección a 3 exoplanetas.
- Tras chocar con una anomalía del espacio-tiempo, casi la mitad de los pasajeros fueron transportados a una dimensión alternativa.



Figura: Imagen del Titanic Espacial

Carga de Datos

Librerías:

- Numpy
- Pandas
- Matplotlib
- Seaborn
- Sklearn
- Missingno

De Sklearn:

- Modelos
- Métricas
- Preprocesamiento
- Impute

Datos a Manejar:

- Train
- Test

Tabla de Contenidos

- 1 Definición del Problema
 - Contexto
 - Carga de Datos
- 2 **Estadística Descriptiva**
- 3 Visualización Descriptiva
- 4 Preprocesamiento
 - Ingeniería de Atributos
 - Pre-procesamiento
- 5 Modelos
 - Modelo para Datos Estandarizados
 - Modelo para Datos NO Estandarizados
- 6 Métricas y Análisis de Resultados
- 7 Conclusiones

Estadística Descriptiva

Dimensiones:

- Entrenamiento → Train: (8693, 14)
- Prueba → Test: (4277, 13)

Columnas:

- PassengerId: object
- HomePlanet: object
- CryoSleep: object
- Cabin: object
- Destination: object
- Age: float64
- VIP: object

Columnas:

- RoomService: float64
- FoodCourt: float64
- ShoppingMall: float64
- Spa: float64
- VRDeck: float64
- Name: object
- Transported: bool

Tabla de Contenidos

- 1 Definición del Problema
 - Contexto
 - Carga de Datos
- 2 Estadística Descriptiva
- 3 Visualización Descriptiva**
- 4 Preprocesamiento
 - Ingeniería de Atributos
 - Pre-procesamiento
- 5 Modelos
 - Modelo para Datos Estandarizados
 - Modelo para Datos NO Estandarizados
- 6 Métricas y Análisis de Resultados
- 7 Conclusiones

Visualización Descriptiva

Según los datos del conjunto de Entrenamiento, aproximadamente la mitad de las personas fueron transportadas.

Porcentaje de Transportados

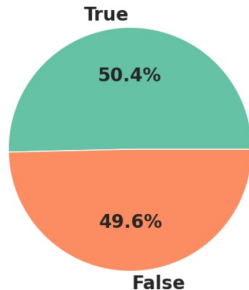


Figura: Porcentaje de Transportados

Visualización Descriptiva

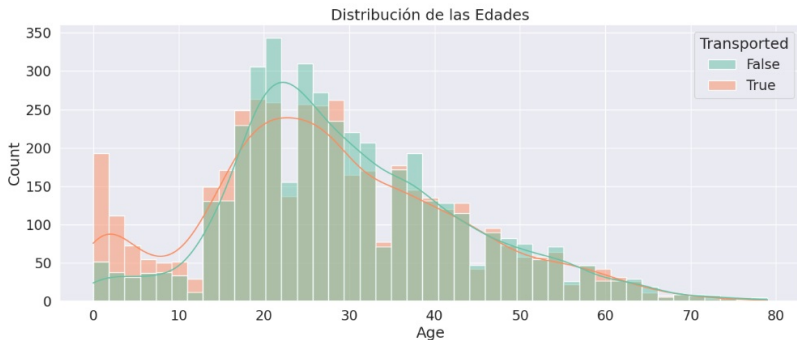


Figura: Distribucion de Edades

Visualización Descriptiva

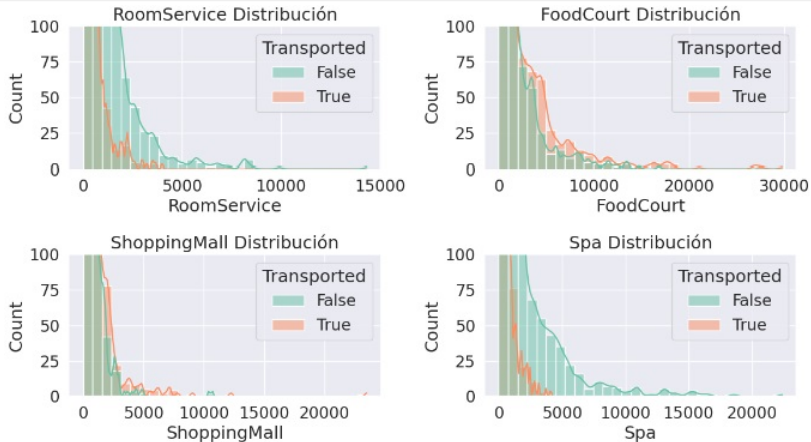


Figura: Distribucion de Gastos

Visualización Descriptiva

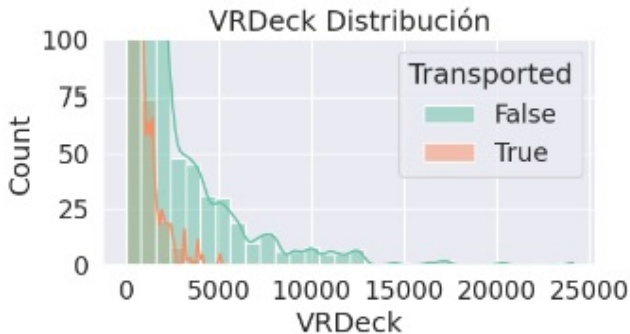


Figura: Distribucion de Gastos

Visualización Descriptiva

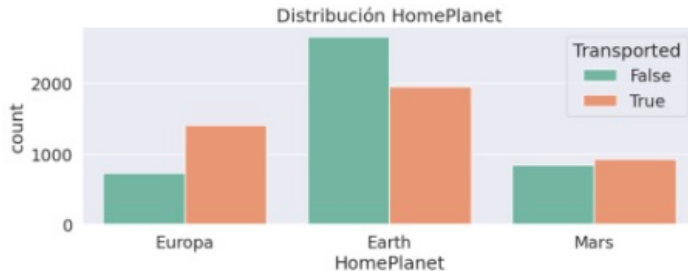


Figura: Distribucion de Lugar de Origen

Visualización Descriptiva



Figura: Distribucion de CryoSleep

Visualización Descriptiva



Figura: Distribucion de Destino

Visualización Descriptiva

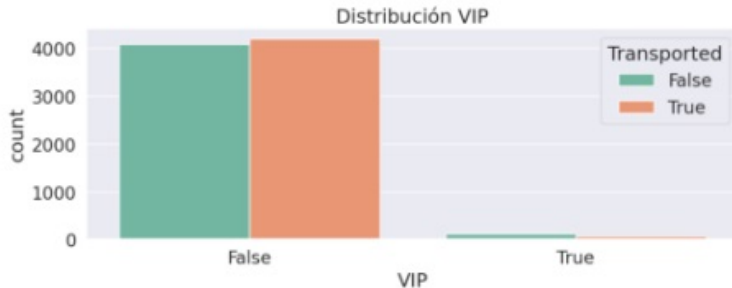


Figura: Distribucion de Pasajeros VIP

Tabla de Contenidos

- 1 Definición del Problema
 - Contexto
 - Carga de Datos
- 2 Estadística Descriptiva
- 3 Visualización Descriptiva
- 4 Preprocesamiento**
 - Ingeniería de Atributos
 - Pre-procesamiento
- 5 Modelos
 - Modelo para Datos Estandarizados
 - Modelo para Datos NO Estandarizados
- 6 Métricas y Análisis de Resultados
- 7 Conclusiones

Ingeniería de Atributos

Creacion de Nuevas Variables:

Reinterpretamos 'PassengerId' que viene en el siguiente formato: gggg-pp

- 'Group_Size'
- "Travelling_Solo"

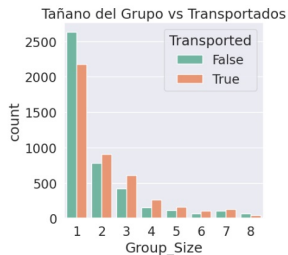


Figura: Tamaño del Grupo vs Transportados

Ingeniería de Atributos

Creación de Nuevas Variables:

Reinterpretamos 'cabin' que viene en el siguiente formato: deck/num//side

- "Cabin_Deck"
- "Cabin_Number"
- "Cabin_Side"

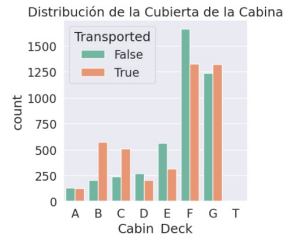


Figura: Distribución de la Cubierta de la Cabina

Ingeniería de Atributos

Creación de Nuevas Variables:

- Distribuimos las Cabinas por Regiones del 1 al 6
- Distribuimos las Edades por Rangos Etarios
- Fusionamos los Gastos

Ingeniería de Atributos

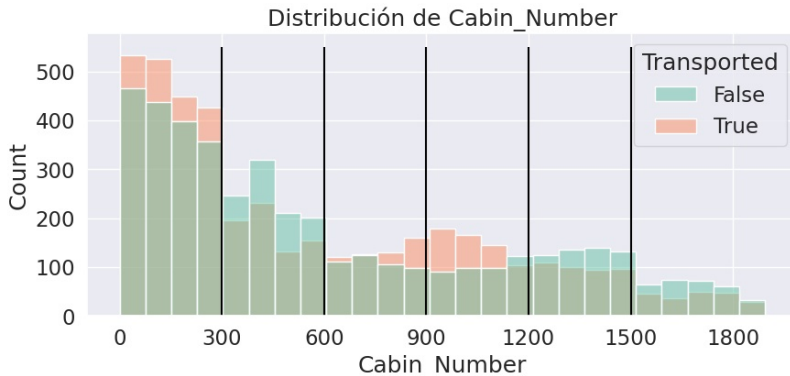


Figura: Distribucion de Regiones de la Cabina

Ingeniería de Atributos

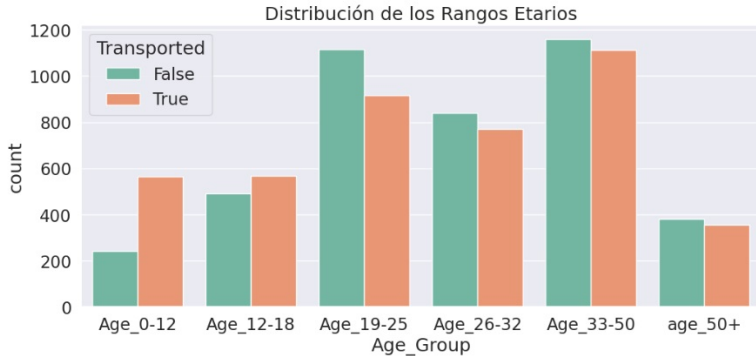


Figura: Distribucion de Rangos Etarios

Ingeniería de Atributos

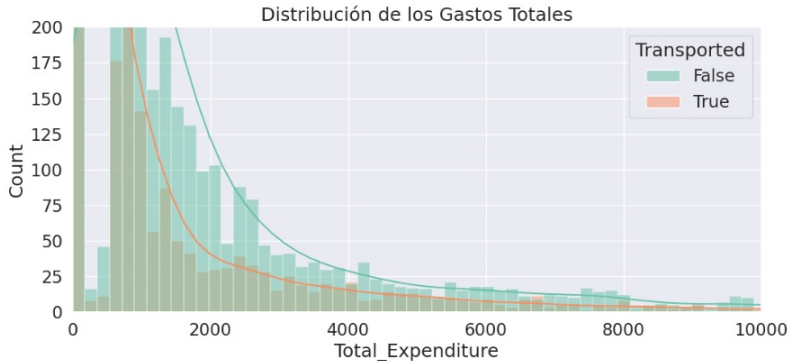


Figura: Gastos Totales

Pre-procesamiento

Información:

- Cambiamos funciones de Tipo Objeto a Tipo Booleano
- Cambiamos funciones de Tipo Int a Tipo Float

Ademas: (usando missingno)

- Se completaron los datos faltantes con lo mas frecuente (datos categoricos)
- Se completaron los datos faltantes con la mediana (datos numericos)

Pre-procesamiento

Información:

- Aplicamos una Transformación Logarítmica a las variables de Gasto
- Hacemos One Hot Encoding para las variables categóricas nominales.
- Hacemos LabelEncoding para las variables categóricas ordinales.

Además:

También pasando a Booleano los datos que falten

Tabla de Contenidos

- 1 Definición del Problema
 - Contexto
 - Carga de Datos
- 2 Estadística Descriptiva
- 3 Visualización Descriptiva
- 4 Preprocesamiento
 - Ingeniería de Atributos
 - Pre-procesamiento
- 5 **Modelos**
 - Modelo para Datos Estandarizados
 - Modelo para Datos NO Estandarizados
- 6 Métricas y Análisis de Resultados
- 7 Conclusiones

Modelo para Datos Estandarizados

Pasos:

- Estandarizamos con StandardScaler
- Train Test Split
- Creamos una función que nos entregara las métricas

Modelos:

- Regresión Logística
- Support Vector Machine

Métricas y Análisis de Resultados

Modelo de Regresión Logística:

- Accuracy_Score del conjunto de Entrenamiento es: 77.86
- Accuracy_Score del conjunto de Testeo es: 77.17
- Precision Score es: 0.75
- Recall Score es: 0.80
- F1 Score es: 0.78

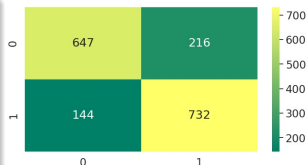


Figura: Matriz de Confusion

Métricas y Análisis de Resultados

Modelo de SVM:

- Accuracy_Score del conjunto de Entrenamiento es: 81.82
- Accuracy_Score del conjunto de Testeo es: 79.70
- Precision Score es: 0.79
- Recall Score es: 0.79
- F1 Score es: 0.79



Figura: Matriz de Confusion

Modelo para Datos NO Estandarizados

Pasos:

- Train Test Split
- Creamos una función que nos entregara las métricas

Modelos:

- Random Forest
- Gradient Boosting

Métricas y Análisis de Resultados

Modelo de Random Forest:

- Accuracy_Score del conjunto de Entrenamiento es: 98.51
- Accuracy_Score del conjunto de Testeo es: 80.56
- Precision Score es: 0.82
- Recall Score es: 0.77
- F1 Score es: 0.80



Figura: Matriz de Confusion

Métricas y Análisis de Resultados

Modelo de Gradient Boosting:

- Accuracy_Score del conjunto de Entrenamiento es: 82.05
- Accuracy_Score del conjunto de Testeo es: 79.29
- Precision Score es: 0.77
- Recall Score es: 0.83
- F1 Score es: 0.80

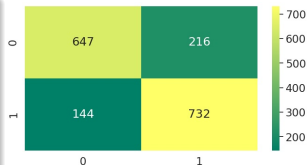
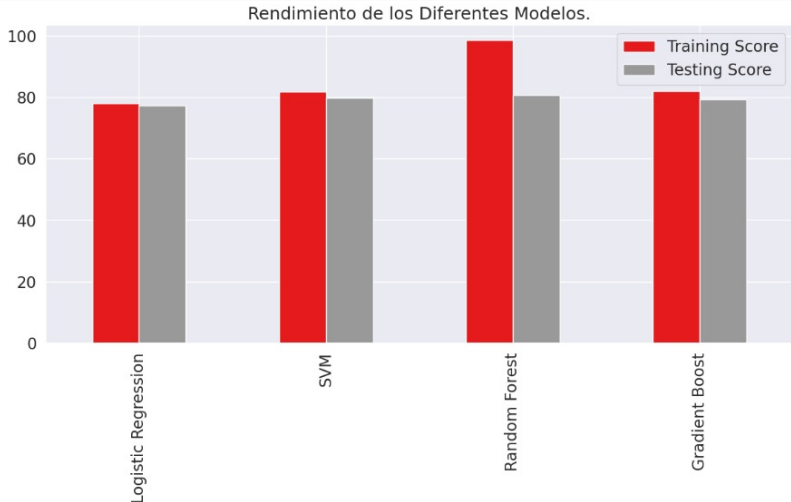


Figura: Matriz de Confusion

Tabla de Contenidos

- 1 Definición del Problema
 - Contexto
 - Carga de Datos
- 2 Estadística Descriptiva
- 3 Visualización Descriptiva
- 4 Preprocesamiento
 - Ingeniería de Atributos
 - Pre-procesamiento
- 5 Modelos
 - Modelo para Datos Estandarizados
 - Modelo para Datos NO Estandarizados
- 6 Métricas y Análisis de Resultados
- 7 Conclusiones

Métricas y Análisis de Resultados



Métricas y Análisis de Resultados

Modelo	Training Score	Testing Score
Regresión Logística	77.868852	77.170788
SVM	81.823411	79.700978
Random Forest	98.518838	80.563542
Gradient Boost	82.053494	79.298447

Cuadro: Tabla de Resultados

Tabla de Contenidos

- 1 Definición del Problema
 - Contexto
 - Carga de Datos
- 2 Estadística Descriptiva
- 3 Visualización Descriptiva
- 4 Preprocesamiento
 - Ingeniería de Atributos
 - Pre-procesamiento
- 5 Modelos
 - Modelo para Datos Estandarizados
 - Modelo para Datos NO Estandarizados
- 6 Métricas y Análisis de Resultados
- 7 Conclusiones

Conclusiones

Conclusiones Finales:

- Pocas funciones directamente trabajables
- Construimos nuevas funciones relevantes (Evitando perdida de informacion)
- Utilizamos tecnicas de Machine Learning para realizar predicciones
- Comparamos las predicciones y seleccionamos los mejores rendimientos:
 - 1 Random Forest
 - 2 Support-Vector Machine
 - 3 Gradient Boost
- Finalmente, se usa 'Stacking Model' para predecir los datos de prueba.