

TIME TO GET UNDERWAY...

SOBREVIVIENDO AL TITANIC

A decorative graphic on the left side of the poster. It features a grid of light blue squares of varying sizes, some of which are partially cut off by the edges. Below the grid is a large, dark blue, stylized shape that resembles the hull of a ship, with a pointed top and a wavy bottom. The entire graphic is set against a dark blue background.

**PROYECTO FINAL ANALISIS DE DATOS
GUATEMALA DE LA ASUNCION 7 DE
MAYO 2023**



JULIO ANTHONY ENGELS RUIZ
COTO - 1284719

LIBRERIAS

PARA PODER RESOLVER ESTE PROBLEMA FUE NECESARIO INSTALAR LAS SIGUIENTES LIBRERÍAS PARA EXTRAER TEXTO, DIVIDIR EL CONJUNTO DE DATOS, IMPUTACIONES, Y ARBOLES DE REGRESIÓN.

```
7. # Importar librerías
8. library(dplyr)
9. library(ggplot2)
10. library(caret)
11. library(RColorBrewer)
12. library(stringr) # para extraer los títulos relevantes de names
13. library(caTools) # la utilidad para la función saveResult dividir el conjunto de datos en entrenamiento y prueba
14. library(caret) # la utilidad para árboles de regresión
15. library(rpart) # para árboles de regresión
16. library(MICE) # la librería MICE permite realizar la imputación de los datos mediante diferentes técnicas: Multivariate Imputations by Chained Equations (MICE).
17. library(kableExtra) # tablas de Markov para visualizar correlaciones de predictores
18.
19. #
```

DATOS

LUEGO SE NOS BRINDO EL DATASET LLAMADO DATA_TITANIC_PROYECTO.CSV ESTE DATASET ES PARA ENTRENAMIENTO Y VALIDACIÓN.

```
21. # Importar el conjunto de datos
22. titanic_data <- read.csv("data_titanic_proyecto.csv")
23. titanic_data
24.
25. #
```

Description: #784 x 12							
PassengerId	Name	Age	SibSp	Parch	Ticket	Fare	Cabin
8400	Braund, Mr. Owen Harris	22.00	1	0	A/5 21171	72500	
7009	Cummings, Mrs. John Bradley (Florence Briggs Thayer)	38.00	1	0	PC 17599	712833	C85
3622	Heikkinen, Miss. Laina	26.00	0	0	STON/O2 3101282	79250	
5596	Futrelle, Mrs. Jacques Heath (Lily May Peel)	35.00	1	0	113803	531000	C123
8403	Aller, Mr. William Henry	35.00	0	0	373450	80500	
8094	Moran, Mr. James	19.00	0	0	330977	84582	
7475	McCarthy, Mr. Timothy J	54.00	0	0	17463	518625	E46
7947	Palsson, Master. Gosta Leonard	2.00	3	1	349909	210750	
9549	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	27.00	0	2	347742	111333	
539	Nasser, Mrs. Nicholas (Adele Achem)	14.00	1	0	237736	300708	

1-10 of 784 rows | 1-8 of 12 columns

Previous 1 2 3 4 5 6 ... 79 Next

IMPUTACIÓN

CON EL MÉTODO COLSUMS SE VERIFICA QUE EN LA COLUMNA DE AGE HACEN FALTA DATOS.

```
30- [r]
31- colsums(is.na(titanic_data))
32-
33-
```

PassengerId	Name	Age	SibSp	Parch	Ticket	Fare
0	0	160	0	0	0	0
Cabin	Embarked	passenger_class	passenger_sex	passenger_survived		
0	0	0	0	0		

PARA IMPUTAR LOS VALORES FALTANTES SE HACE USO DE UNA REGRESION LINEAL LA CUAL VIENE DADA MICE. (MULTIVARIATE IMPUTATIONS BY CHAINED EQUATIONS)

```
34- [r]
35- #a imputa el valor a una variable temporal
36- imputed_data <- mice(titanic_data%select(Age,PassengerId), method = "norm.boot")
37-
38-
39-
40-
```

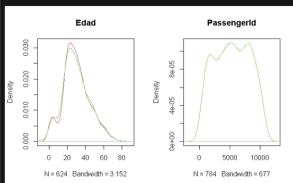
iter	imp	variable
1	1	Age
1	2	Age
1	3	Age
1	4	Age
1	5	Age
2	1	Age
2	2	Age
2	3	Age
2	4	Age
2	5	Age
3	1	Age
3	2	Age
3	3	Age
3	4	Age
3	5	Age
4	1	Age
4	2	Age
4	3	Age
4	4	Age
4	5	Age
5	1	Age
5	2	Age
5	3	Age
5	4	Age

ACÁ SE PUEDE OBSERVAR A LA IZQUIERDA UNA GRÁFICA DE COLOR ROJO QUE ES EL VALOR OBTENIDO DE UNA REGRESION LINEAL DE LOS VALORES IMPUTADOS Y LA GRAFICA DE COLOR VERDE QUE ESTA INTERCALADA ES LA FORMA EN LA QUE SE DISTRIBUYEN LOS DATOS CONOCIDOS.

```

42. # [r]
43. # Se utiliza la función de regresión lineal de la librería mice para imputar valores
44. data_impu_boot <- mice::complete(imputed_data)
45.
46. par(mfrow=c(1,2))
47. plot(density(titanic_data$Age, na.rm = T), col="red", main="Edad")
48. lines(density(data_impu_boot$Age), col="j")
49.
50. plot(density(titanic_data$PassengerId, na.rm = T), col="j", main="PassengerId")
51. lines(density(data_impu_boot$PassengerId), col="j")
52.
53.
54.

```



SE PUEDE DECIR QUE SE OBTIENE UN VALOR DE IMPUTACION DE LA EDAD BASTANTE ACEPTABLE.

EN LA GRÁFICA DE LA DERECHA SE PUEDE OBSERVAR LA VARIABLE PASSENGERID Y COMO ESTA VARIABLE NO TIENE NINGUN DATO FALTANTE NO IMPUTA NINGÚN VALOR.

```

56. # [r]
57. # Se utiliza para imputar valores el método por valor medio
58. imputed_data <- mice(titanic_data[,c(Age, PassengerId)], method = "mean")
59.
60.
61.

```

```

iter imp variable
1 1 Age
1 2 Age
1 3 Age
1 4 Age
1 5 Age
2 1 Age
2 2 Age
2 3 Age
2 4 Age
2 5 Age
3 1 Age
3 2 Age
3 3 Age
3 4 Age
3 5 Age
4 1 Age
4 2 Age
4 3 Age
4 4 Age
4 5 Age
5 1 Age
5 2 Age
5 3 Age
5 4 Age
5 5 Age

```

AHORA SE COMPARA CON EL MÉTODO PARA IMPUTAR DE VALOR PROMEDIO DE LOS VALORES EXISTENTES



COMO SE PUEDE OBSERVAR CON EL VALOR PROMEDIO NO SE OBTIENE UNA BUENA DISTRIBUCIÓN DE DATOS IMPUTADOS.

```

76- # [r]
77 # De acuerdo al mejor desempeño en imputación se eligió el metodo de regresión lineal
78 titanic_data <- titanic_data %>% mutate(Age = Data_impu_boot$Age, PassengerId = Data_impu_boot$PassengerId)
79
80 titanic_data
81
82

```

PassengerId	Name	Age	SibSp	Parch	Ticket	Fare	Cabin
8429	Brand, Mr. Owen Harris	22.00000000	1	0	A/5 21171	7.2500	
7009	Cummings, Mrs. John Bradley (Florence Briggs Thayer)	38.00000000	1	0	PC 17599	71.2833	C85
3622	Heikkinen, Miss. Laina	26.00000000	0	0	STON/O2 3101282	7.9250	
5595	Futrelle, Mrs. Jacques Heath (Lily May Peel)	35.00000000	1	0	113803	53.1000	C123
8403	Allen, Mr. William Henry	35.00000000	0	0	378450	8.0500	
8094	Moran, Mr. James	52.98559762	0	0	310877	8.4583	
7475	McCarthy, Mr. Timothy J	54.00000000	0	0	17463	51.8625	E46
7947	Palsson, Master. Gosta Leonard	2.00000000	3	1	349309	21.0750	
9549	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	27.00000000	0	2	347742	11.1333	
539	Nasser, Mrs. Nicholas (Adele Achem)	14.00000000	1	0	237736	30.0708	

1-10 of 784 rows | 1-8 of 12 columns

DE ACUERDO AL MEJOR DESEMPEÑO EN IMPUTACIÓN SE ELIGIÓ EL MÉTODO DE REGRESIÓN LINEAL.

```

84- # [r]
85 colSums(is.na(titanic_data))
86-

```

PassengerId	Name	Age	SibSp	Parch	Ticket	Fare
0	0	0	0	0	0	0
Cabin	Embarked	passenger_class	passenger_sex	passenger_survived	0	0

CON ESTO SE COMPRUEBA QUE YA TODAS LAS COLUMNAS NO TIENEN VALORES NA

FEATURE ENGINEERING

SE TRANSFORMA LA COLUMNA EMBARKED, PASSENGER_CLASS A FACTOR Y LA COLUMNA PASSENGER_SURVIVED A INTEGER

```
88. [r]
89. #limpiar y transformar los datos
90. titanic_data$embarked <- as.factor(titanic_data$embarked)#transformar columna embarked en un factor(variable categorica)
91. titanic_data$passenger_class <- as.factor(titanic_data$passenger_class)
92. titanic_data$passenger_survived <- as.integer(titanic_data$passenger_survived == "Y") #comparo cada valor de la columna passenger_survived con "Y" TRUE si valor es igual a "Y" y FALSE "N" con as.integer TRUE pasa a 1 y FALSE pasa a 0
93. titanic_data$passenger_sex <- as.integer(titanic_data$passenger_sex == "M") #comparo cada valor de la columna passenger_sex con "Y" TRUE si valor es igual a "Y" y FALSE "N" con as.integer TRUE pasa a 1 y FALSE pasa a 0
94.
```

ACÁ EXTRAIGO EL TITULO DE LA COLUMNA NAMES Y CREO UNA COLUMNA TITULOS

```
96. [r]
97. #extraer los titulos de la columna Names y crear una nueva columna Titulos
98. titanic_data$titulos <- str_extract(titanic_data$Name, "(Mr|Mrs|Miss|Master|Don|Dr|Rev|Wm|Ms|Major|Lady|Sir|Title|Col|Capt|the Countess)\.")
99.
100. #separar las columnas en 6 subconjuntos con dplyr
101. titanic_data$titulos[is.na(titanic_data$titulos)] <- "otro"
102.
103. #mostrar los cambios realizados
104. head(titanic_data)
105.
```

SibSp	Parch	Ticket	Fare	Cabin	Embarked	passenger_class	passenger_sex	passenger_survived	Titulos
1	0	A/5 21171	7.2500	S	Lower	1	0	Mr.	
1	0	PC 17599	71.2833	C85	C	Upper	0	1	Mrs.
0	0	STON/O2 3101282	7.9250	S	Lower	0	1	Miss.	
1	0	113803	53.1000	C123	S	Upper	0	1	Mrs.
0	0	373450	8.0500	S	Lower	1	0	Mr.	
0	0	310877	8.4583	Q	Lower	1	0	Mr.	

6 rows | 5.14 of 13 columns

```

111. # Se crea un nuevo atributo con Cabin a pesar que tiene bastantes datos faltantes vamos a probar donde no hay numero de cabina se tomara como No y donde
112. # hay se tomara como Yes
113. titanic_data$Hascabinnum <- ifelse(titanic_data$Cabin != "", "Yes", "No")
114.
115.
116.
117. # Se crea una variable dummy que es Ticket_Class
118.
119.
120. titanic_data$Ticket <- as.character(titanic_data$Ticket)
121. titanic_data$Ticket_class <- ifelse(titanic_data$Ticket != "", substr(titanic_data$Ticket, 1, 3), "")
122. titanic_data$Ticket_class <- as.factor(titanic_data$Ticket_class)
123. # Eliminar las columnas en numeros enteros
124. titanic_data$Ticket_class <- as.integer(titanic_data$Ticket_class)
125.
126.
127.
128. # Se crea una variable dummy que es Deck a partir de la variable cabin
129.
130. titanic_data$Deck <- as.character(titanic_data$Cabin)
131. titanic_data$Deck <- ifelse(titanic_data$Cabin == "", "No", substr(titanic_data$Cabin, 3, 3))
132. titanic_data$Deck <- as.factor(titanic_data$Deck)
133. titanic_data$Deck <- as.integer(titanic_data$Deck)
134.
135.

```

SE CREAN VARIABLES DUMMIES HASCABINNUM (RELACIONA CON LA VARIABLE CABIN), TICKET_CLASS(RELACIONA CON TICKET), DECK(SE RELACIONA CON CABIN).

```

241. # Se crea un conjunto de datos sin columnas innecesarias
242. titanic_data_clean <- titanic_data %>% select(-Ticket)
243. titanic_data_clean
244.

```

PassengerId	Name	Age	SibSp	Parch	Fare	Cabin	Embarked
8400	Braund, Mr. Owen Harris	22.00000000	1	0	7.2500		S
7009	Cummings, Mrs. John Bradley (Florence Briggs Thayer)	38.00000000	1	0	71.2833	C85	C
3662	Heikkinen, Miss. Laina	26.00000000	0	0	7.9250		S
5196	Futrelle, Mrs. Jacques Heath (Lily May Peel)	35.00000000	1	0	53.1000	C123	S
6403	Allen, Mr. William Henry	35.00000000	0	0	8.0500		S
8094	Moran, Mr. James	52.9698782	0	0	8.4583		Q
7475	McCarthy, Mr. Timothy J	54.00000000	0	0	51.8625	E46	S
7047	Palsson, Master. Gosta Leonard	2.00000000	3	1	21.0750		S
9540	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	27.00000000	0	2	11.1333		S
539	Nasser, Mrs. Nicholas (Adele Achem)	14.00000000	1	0	30.0708		C

1-10 of 784 rows | 1-8 of 15 columns

ELIMINO LAS COLUMNAS INNECESARIAS

```

249. # Se crean variables dummy comparando variables categoricas en variables numericas binarias para las variables categoricas
250. titanic_data_clean <- titanic_data_clean %>%
251. mutate(
252.   Embarked_S = as.integer(Embarked == "S"),
253.   Embarked_C = as.integer(Embarked == "C"),
254.   Embarked_Q = as.integer(Embarked == "Q"),
255.   Age = as.integer(Age),
256.   Fare = as.integer(Fare)
257. )
258. titanic_data_clean
259.
260.
261.

```

passenger_class	passenger_sex	passenger_survived	Titulos	Hascabinnum	Ticket_class	Deck	Embarked_S	Embarked_C	Embarked_Q
Lower	1	0	Mr	No	10	9	1	0	0
Upper	0	1	Mrs	Yes	14	9	0	1	0
Lower	0	1	Mrs	No	15	9	1	0	0
Upper	0	1	Mrs	Yes	1	3	1	0	0
Lower	1	0	Mr	No	3	9	1	0	0
Lower	1	0	Mr	No	3	9	0	0	1
Upper	1	0	Mr	Yes	1	5	1	0	0
Lower	1	0	Master	No	3	9	1	0	0
Lower	0	1	Mrs	No	3	9	1	0	0
Middle	0	1	Mrs	No	2	9	0	1	0

1-10 of 784 rows | 9-18 of 18 columns

CREO VARIABLES DUMMY DE LA COLUMNA EMBARKED(EMBARKED_S, EMBARKED_C, EMBARKED_Q)

```

154. [r]
155 # eliminar las variables originales categoricas
156 titanic_data_clean <- titanic_data_clean %>%
157   select(-Embarked, -cabin, -name)
158 titanic_data_clean
159.

```

passenger.class	passenger.sex	passenger.survived	Titulos	HasCabinNum	Ticket.class	Deck	Embarked.S	Embarked.C	Embarked.Q
Lower	1	0	Mr	No	10	9	1	0	0
Upper	0	1	Mrs	Yes	14	3	0	1	0
Lower	0	1	Miss	No	15	9	1	0	0
Upper	0	1	Mrs	Yes	1	3	1	0	0
Lower	1	0	Mr	No	3	9	1	0	0
Lower	1	0	Mr	No	3	9	0	0	1
Upper	1	0	Mr	Yes	1	5	1	0	0
Lower	1	0	Master	No	3	9	1	0	0
Lower	0	1	Mrs	No	3	9	1	0	0
Middle	0	1	Mrs	No	2	9	0	1	0

1-10 of 784 rows | 6-15 of 15 columns

ELIMINO LAS VARIABLES ORIGINALES CATEGÓRICAS

```

174. [r]
175 #convertir los titulos en factores
176 #obtener los niveles unicos de los titulos
177 titulos_unicos <- unique(titanic_data_clean$titulos)
178
179 #convertir los titulos en factores utilizando los niveles unicos
180 titanic_data_clean$titulos <- factor(titanic_data_clean$titulos, levels = titulos_unicos)
181
182 #convertir los factores en numeros enteros
183 titanic_data_clean$titulos <- as.integer(titanic_data_clean$titulos)
184
185 titanic_data_clean
186.

```

passenger.class	passenger.sex	passenger.survived	Titulos	HasCabinNum	Ticket.class	Deck	Embarked.S	Embarked.C	Embarked.Q
Lower	1	0	1 No	No	10	9	1	0	0
Upper	0	1	2 Yes	Yes	14	3	0	1	0
Lower	0	1	3 No	No	15	9	1	0	0
Upper	0	1	2 Yes	Yes	1	3	1	0	0
Lower	1	0	1 No	No	3	9	1	0	0
Lower	1	0	1 No	No	3	9	0	0	1
Upper	1	0	1 Yes	Yes	1	5	1	0	0
Lower	1	0	4 No	No	3	9	1	0	0
Lower	0	1	2 No	No	3	9	1	0	0
Middle	0	1	2 No	No	2	9	0	1	0

1-10 of 784 rows | 6-15 of 15 columns

CONVIERTO LA COLUMNA DE TITULOS EN CATEGÓRICAS

```

187. [r]
188 #convertir class pasajeros en factores
189 #obtener los niveles unicos de los titulos
190 pasajeros_unicos <- unique(titanic_data_clean$passenger.class)
191
192 #convertir los titulos en factores utilizando los niveles unicos
193 titanic_data_clean$passenger.class <- factor(titanic_data_clean$passenger.class, levels = pasajeros_unicos)
194
195 #convertir los factores en numeros enteros
196 titanic_data_clean$passenger.class <- as.integer(titanic_data_clean$passenger.class)
197
198 #mostrar un resumen de la columna Titulos
199 summary(titanic_data_clean$passenger.class)
200
201
202 titanic_data_clean
203.

```

passenger.class	passenger.sex	passenger.survived	Titulos	HasCabinNum	Ticket.class	Deck	Embarked.S	Embarked.C	Embarked.Q
1	1	0	1 No	No	10	9	1	0	0
2	0	1	2 Yes	Yes	14	3	0	1	0
1	0	1	3 No	No	15	9	1	0	0
2	0	1	2 Yes	Yes	1	3	1	0	0
1	1	0	1 No	No	3	9	1	0	0
1	1	0	1 No	No	3	9	0	0	1
2	0	1	Yes	Yes	1	5	1	0	0
1	1	0	4 No	No	3	9	1	0	0
1	0	1	2 No	No	3	9	1	0	0

CONVIERTO LA COLUMNA DE PASSENGER_CLASS EN CATEGÓRICAS


```

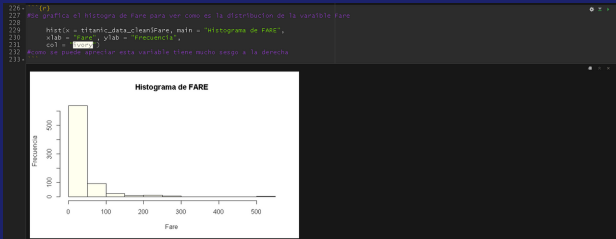
209. # [r]
210. # Se crea una nueva variable dummy llamada hasCabinNum
211. # Convertir hasCabinNum en factores
212. # obtener las cabinas únicas de las cabinas
213. cabinas_unicas <- unique(titanic_data_clean$hasCabinNum)
214.
215. # Convertir las cabinas en factores utilizando las cabinas únicas
216. titanic_data_clean$hasCabinNum <- factor(titanic_data_clean$hasCabinNum, levels = cabinas_unicas)
217.
218. # Convertir los factores en números enteros
219. titanic_data_clean$hasCabinNum <- as.integer(titanic_data_clean$hasCabinNum)
220.
221. titanic_data_clean
222.

```

PassengerId	Age	Sex	SibSp	Parch	Fare	passenger.class	passenger.sex	passenger.survived	Titulos	HasCabinNum
8400	22	1	0	7	1	1	1	0	1	1
7009	38	1	0	71	2	0	1	2	2	2
9622	26	0	0	7	1	0	1	3	1	1
1506	35	1	0	53	2	0	1	2	2	2
8403	35	0	0	8	1	1	0	1	1	1
8034	52	0	0	8	1	1	0	1	1	1
7475	54	0	0	51	2	1	0	1	2	2
7947	2	3	1	21	1	1	0	4	1	1
9549	27	0	2	11	1	0	1	2	1	1
539	14	1	0	30	3	0	1	2	1	1

1-10 of 784 rows | 1-10 of 15 columns

CONVIERTO LA COLUMNA DE HASCABINNUM EN CATEGÓRICAS

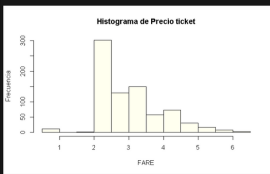


SE GRAFICA EL HISTOGRAMA DE FARE PARA VER COMO ES LA DISTRIBUCION DE LA VARIABLE FARE, COMO SE PUEDE APRECIAR ESTA VARIABLE TIENE MUCHO SESGO A LA DERECHA

```

218: ##(r)
219:
220: #se grafica el histograma de fare para ver como es la distribucion de la variable fare
221:
222: hist(x = log(titanic_data_clean$fare), main = "Histograma de Precio ticket",
223:      xlab = "fare", ylab = "Frecuencia",
224:      col = "yellow")
225: # aplicando logaritmos naturales a los datos de la variable se ve que mejora la distribucion
226:
227:

```



APLICANDO LOGARITMOS NATURALES A LOS DATOS DE LA VARIABLE SE VE QUE MEJORA LA DISTRIBUCIÓN

```

251: ##(r)
252:
253: #se transforman los datos aplicando logaritmos
254: log1 <- log(titanic_data_clean$fare)
255: titanic_data_clean$log1 <- log1
256:
257: titanic_data_clean
258:
259:

```

#	passenger_sex	passenger_survived	Titulos	HasCabinNum	Ticket.class	Deck	Embarked.S	Embarked.C	Embarked.Q	Fare_log
1	1	0	1	1	10	9	1	0	0	2.1972246
0	1	2	2	14	3	0	1	0	0	4.2804594
0	1	3	1	15	9	1	0	0	0	2.1972246
0	1	2	2	1	3	1	0	0	0	4.0073332
1	0	1	1	1	3	9	1	0	0	2.3025851
1	0	1	1	3	9	0	0	1	2	2.3025851
1	0	1	2	1	5	1	0	0	0	3.9702919
1	0	4	1	3	9	1	0	0	0	3.1354342
0	1	2	1	3	9	1	0	0	0	2.5643434
0	1	2	1	2	9	0	1	0	0	3.4657359

ACÁ SE SUMA A LA COLUMNA FARE UN 2 PARA EVITAR VALORES INDEFINIDOS DE LOGARITMO NATURAL CUANDO FARE TENGA UN VALOR 0 Y SE CREA LA COLUMNA FARE_LOG CON ESTE CALCULO

```

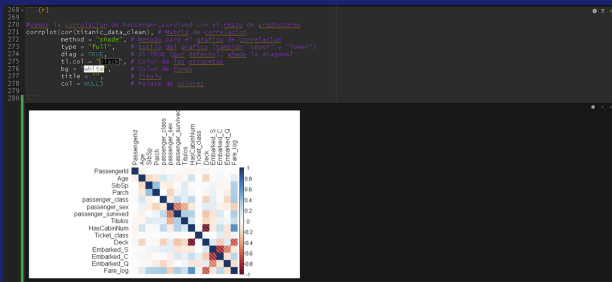
261. # Correlacion
262. titanic_data_clean <- titanic_data_clean %>%
263.   select(-fare)
264.   titanic_data_clean
265.
266.

```

PassengerId	Age	Sex	Survived	Embarked	Class	Survived	Titles	HasCabinNum	Ticket_class
8400	22	1	0	1	1	0	0	1	10
7009	38	1	0	2	1	0	1	2	2
3622	26	0	0	1	0	1	3	1	14
5586	35	1	0	2	0	1	2	2	1
8403	35	0	0	1	1	0	1	1	3
8094	52	0	0	1	1	0	1	1	3
7475	54	0	0	2	1	0	1	2	1
7947	2	3	1	1	1	0	4	1	3
9549	27	0	2	1	0	1	2	1	3
539	14	1	0	3	0	1	2	1	2

1-10 of 784 rows | 1-10 of 15 columns

ELIMINO LA COLUMNA FARE



SE GRÁFICA LA MATRIZ DE CORRELACIONES PARA LA VARIABLE PASSENGER_SURVIVED DONDE SE PUEDE APRECIAR CON QUE OTRAS VARIABLES TIENE CORRELACIÓN SEA POSITIVA O NEGATIVA ESTO SERVIRÁ PARA DEFINIR LOS PREDICTORES QUE SERVIRÁN PARA ENTRENAR EL MODELO

```

Dec: 4 3 the right, agreed:119, adju:159, (0.92115)
PARA ENTRENAR EL MODELO DE ARBOLES
DE REGRESIÓN LOGÍSTICA SE HACE USO
DEL DATASET TRAIN_DATA Y ADEMÁS POR
CORRELACIÓN SE DEFINEN LOS
PREDICTORES QUE SERVIRÁN PARA
PREDECIR A LA VARIABLE OBJETIVO
PASSENGER SURVIVED

```

```

122. #####r)
123. #Entonces la supervivencia utilizando el modelo entrenado en el conjunto de validacion
124. predictions <- predict(model, newdata = validation_data, type = "matrix")
125. predicted_survived <- ifelse(predictions > 0.5, "1", "0")
126.
127.
128. #Calculo la precision del modelo en el conjunto de validacion
129. accuracy <- sum(predicted_survived == validation_data$passenger_survived) / length(validation_data$passenger_survived)
130. accuracy
131.
132.
[1] 0

```

UNA VEZ EL MODELO ENTRENADO SE PROCEDE A PREDECIR LOS VALORES DE LA VARIABLE PASSENGER_SURVIVED, USANDO EL METODO PREDICT Y PARA ELLO SE USA EL DATASET VALIDATION_DATA.

```

134. #####r)
135. #Entonces para calcular precision, recall y f1 score
136. precision <- function(matrix) {
137.   # True positive
138.   tp <- matrix[2, 2]
139.   # False positive
140.   fp <- matrix[1, 2]
141.   return (tp / (tp + fp))
142. }
143.
144. #recall <- function(matrix) {
145.   # True positive
146.   tp <- matrix[2, 2] # False positive
147.   fn <- matrix[2, 1]
148.   return (tp / (tp + fn))
149. }
150. #Cuando una matriz de confusion
151. conf_matrix <- table(predicted_survived, validation_data$passenger_survived)
152.
153.
154. prec <- precision(conf_matrix)
155. prec
156. rec <- recall(conf_matrix)
157. rec
158.
[1] 0.6393443
[1] 0.7990909
159.
160. #####r)
161. #F1 score
162. f1 <- 2 * ((prec * rec) / (prec + rec))
163. f1
164.
[1] 0.6724138

```

CON ESTAS MÉTRICAS SIRVEN PARA MEDIR QUE TAN BIEN ESTA PREDICIENDO EL MODELO CON DATOS DE ENTRENAMIENTO.

ESTAS METRICAS SE OBTUVIERON DE LOS DIFERENTES INTENTOS Y MODELOS DE MACHINE LEARNING DE APRENDIZAJE SUPERVISADO.

	F1 SCORE	PRECISION	RECALL
Entrenamiento fallido	-	-	-
Entrenamiento fallido	-	-	-
Entrenamiento fallido	-	-	-
Entrenamiento fallido	-	-	-
Entrenamiento y validación 60% - 40 %	[1] 0.7420814	[1] 0.7130435	[1] 0.7735849
Entrenamiento y validación 60% - 40 %	[1] 0.7420814	[1] 0.7130435	[1] 0.7735849
Entrenamiento y validación 80% - 20 %	[1] 0.733945	[1] 0.7017544	[1] 0.7692308
Entrenamiento y validación 90% - 10 %	[1] 0.7666667	[1] 0.71875	[1] 0.8214286
Quitando Embarked_S, Embarked_C, Embarked_Q y agregando columna Embarked entrenamiento y validación 90% - 10%	[1] 0.7868852	[1] 0.75	[1] 0.8275862
Normalizando la columna Fare log 90% - 10%	[1] 0.7868852	[1] 0.75	[1] 0.8275862
Aplicando Árboles de regresión logística 90% - 10%	[1] 0.7741935	[1] 0.75	[1] 0.8
Aplicando arboles de decisión con variables Dummies	[1] 0.8	[1] 0.7333333	[1] 0.88
Aplicando regresión logística	[1] 0.7857143	[1] 0.7333333	[1] 0.8461538
Aplicando regresión logística	[1] 0.7857143	[1] 0.7333333	[1] 0.8461538
Aplicando la normalización 80% - 20%	[1] 0.789916	[1] 0.7704918	[1] 0.8103448
Aplicando arboles de decisión 90% - 10%	[1] 0.7586207	[1] 0.6875	[1] 0.8461538
Aplicando RandomForest	[1] 0.75	[1] 0.7	[1] 0.8076923
Aplicando regresión logística	[1] 0.7333333	[1] 0.6875	[1] 0.7857143
Aplicando regresión logística	[1] 0.6896552	[1] 0.6666667	[1] 0.7142857
Aplicando regresión logística	[1] 0.7142857	[1] 0.6666667	[1] 0.7692308
Aplicando regresión logística	[1] 0.7213115	[1] 0.7333333	[1] 0.7096774
Aplicando regresión logística	[1] 0.6785714	[1] 0.6333333	[1] 0.7307692
Aplicando regresión logística	[1] 0.6785714	[1] 0.6333333	[1] 0.7307692
Aplicando arboles de regresión logística	[1] 0.6785714	[1] 0.6333333	[1] 0.7307692
Aplicando arboles de regresión logística	[1] 0.6724138	[1] 0.6393443	[1] 0.7090909
Aplicando arboles de regresión logística	[1] 0.6724138	[1] 0.6393443	[1] 0.7090909
Aplicando arboles de regresión logística	[1] 0.6724138	[1] 0.6393443	[1] 0.7090909

EL INTENTO FINAL CON MAYOR SCORE EN KAGGLE

PRUEBAS DEL MODELO DE APRENDIZAJE CON DATOS DE EVALUACIÓN

LUEGO SE NOS BRINDO EL DATASET LLAMADO EVALUATION_DATASET.CSV ESTE DATASET NO CONTO CON LA COLUMNA DE PASSENGER_SURVIVED

```
373: #AQUI TERMINA EL ENTRENAMIENTO DEL MODELO DE REGRESION LOGISTICA BINOMIAL
374:
375: print(r)
376: #Importo el conjunto de datos
377: evaluation_data = read.csv("evaluation_dataset.csv")
378: evaluation_data
379:
```

PassengerId	Name	Age	Sex	Survived	Parch	Ticket	Fare	Cabin	Embarked
6182	Ali, Mr. William	25.00	0	0	50TON/O.Q.	3101312	7.0500	S	
1777	Kamei, Mr. Abraham (David Lohan)	25.00	0	0	374887		7.2100	S	
1224	Sjöblom, Miss Anna Sofia	18.00	0	0	3101265		7.4900	S	
2151	Rice, Master. George Hugh	8.00	4	1	382652		29.1250	Q	
3813	Dean, Master. Bertram Vere	1.00	1	2	CA. 2115		20.5700	S	
4072	Cuppenham, Mr. Benjamin	46.00	0	0	PC.17593		79.0000	B82 B84	C
9943	Keane, Mr. Andrew "Andy"	NA	0	0	12460		7.7100	Q	
5058	Casell, Mr. Alfred	16.00	0	0	219855		26.0000	S	
947	Sage, Miss Stella Anna	NA	0	2	CA. 2343		69.5100	S	
5928	Hoyt, Mr. William Fisher	NA	0	0	PC.17600		30.6900	C	

110 of 107 rows (1-9 of 11 columns)

Resized 2 3 4 5 6 11 Next

LUEGO DE IMPORTAR ESTE DATASET DE EVALUATION SE APLICO EL MISMO PROCEDIMIENTO DE ANALISIS DE DATOS, LIMPIEZA DE DATOS, IMPUTACION DE DATOS PARA DEJARLO DE LA MISMA FORMA QUE EL DATASET DE ENTRENAMIENTO Y CON ELLO SE BUSCA HOMOGENISAR LOS DATOS EN AMBOS ARCHIVOS. CON ESTE DATASET DE EVALUATION SE TIENE COMO OBJETIVO EVALUAR EL MODELO CON DATOS TOTALMENTE AGENOS A LOS QUE SIRVIERON EN EL ENTRENAMIENTO DEL MODELO DE APRENDIZAJE SUPERVISADO.

```

614- """(r)
615- #predicciones a base de modelo utilizando el modelo entrenado en el conjunto de evaluacion
616- predictions = predict(model_test, newdata = evaluation_data_clean, type = "matrix")
617- predicted_survived_test = (ifelse(predictions > 0.5, "Y", "N"))
618- predicted_survived_test
619-
620-
621-
622- """(r)
623- #se convierten los valores predichos a valores de 0/1
624- passenger_survived <- as.integer(predicted_survived_test=="Y") #poner 1 cada valor de la columna passenger_survived con "Y" TRUE si valor es igual a "Y"
625- #false 0 si no lo es
626- passenger_survived
627-
628-
629-
630-
631-
632-
633-
634-
635-
636-
637-
638-
639-
640-
641-
642-
643-
644-
645-
646-
647-
648-
649-
650-
651-
652-
653-
654-
655-
656-
657-
658-
659-
660-
661-
662-
663-
664-
665-
666-
667-
668-
669-
670-
671-
672-
673-
674-
675-
676-
677-
678-
679-
680-
681-
682-
683-
684-
685-
686-
687-
688-
689-
690-
691-
692-
693-
694-
695-
696-
697-
698-
699-
700-
701-
702-
703-
704-
705-
706-
707-
708-
709-
710-
711-
712-
713-
714-
715-
716-
717-
718-
719-
720-
721-
722-
723-
724-
725-
726-
727-
728-
729-
730-
731-
732-
733-
734-
735-
736-
737-
738-
739-
740-
741-
742-
743-
744-
745-
746-
747-
748-
749-
750-
751-
752-
753-
754-
755-
756-
757-
758-
759-
760-
761-
762-
763-
764-
765-
766-
767-
768-
769-
770-
771-
772-
773-
774-
775-
776-
777-
778-
779-
780-
781-
782-
783-
784-
785-
786-
787-
788-
789-
790-
791-
792-
793-
794-
795-
796-
797-
798-
799-
800-
801-
802-
803-
804-
805-
806-
807-
808-
809-
810-
811-
812-
813-
814-
815-
816-
817-
818-
819-
820-
821-
822-
823-
824-
825-
826-
827-
828-
829-
830-
831-
832-
833-
834-
835-
836-
837-
838-
839-
840-
841-
842-
843-
844-
845-
846-
847-
848-
849-
850-
851-
852-
853-
854-
855-
856-
857-
858-
859-
860-
861-
862-
863-
864-
865-
866-
867-
868-
869-
870-
871-
872-
873-
874-
875-
876-
877-
878-
879-
880-
881-
882-
883-
884-
885-
886-
887-
888-
889-
890-
891-
892-
893-
894-
895-
896-
897-
898-
899-
900-
901-
902-
903-
904-
905-
906-
907-
908-
909-
910-
911-
912-
913-
914-
915-
916-
917-
918-
919-
920-
921-
922-
923-
924-
925-
926-
927-
928-
929-
930-
931-
932-
933-
934-
935-
936-
937-
938-
939-
940-
941-
942-
943-
944-
945-
946-
947-
948-
949-
950-
951-
952-
953-
954-
955-
956-
957-
958-
959-
960-
961-
962-
963-
964-
965-
966-
967-
968-
969-
970-
971-
972-
973-
974-
975-
976-
977-
978-
979-
980-
981-
982-
983-
984-
985-
986-
987-
988-
989-
990-
991-
992-
993-
994-
995-
996-
997-
998-
999-
1000-

```

SE APLICA LAS PREDICCIONES CON EL DATASET DE EVALUATION_DATA_CLEAN PARA DETERMINAR LAS PREDICCIONES OBTENIDAS SON PROBABILIDADES DE OCURENCIA DEL EVENTO POR LO QUE SE USA LA FUNCION PARA CLASIFICAR TODA PROBABILIDAD > 0.5 SE CONSIDERA COMO VERDADERO Y TODA PROBABILIDAD < 0.5 SE CONSIDERA COMO FALSO.


```

628 # []
629 PassengerId <- evaluation_data_clean$PassengerId
630 pred$clones_data <- cbind(PassengerId, passenger_survived)
631 pred$clones_data
632
633

```

```

  PassengerId passenger_survived
[1,]         1192              0
[2,]         1777              0
[3,]         5224              1
[4,]         2151              0
[5,]         1811              1
[6,]         4072              0
[7,]         9943              0
[8,]         5058              0
[9,]          947              0
[10,]        5928              0
[11,]        5637              0
[12,]        1930              1
[13,]        1842              0
[14,]        9216              1
[15,]        7517              1
[16,]       10302              1
[17,]        3915              0
[18,]        2733              0
[19,]        1915              1
[20,]        7398              1
[21,]        7475              0
[22,]        5766              0
[23,]        8660              0
[24,]        1628              1
[25,]        5037              0
[26,]        7419              1
[27,]        6217              0
[28,]        4992              0
[29,]       10674              0
[30,]        8950              0

```

SE CONSTRUYE UN DATASET CON EL FORMATO PARA QUE KAGGLE PUEDA EVALUAR LA PRUEBA.

```

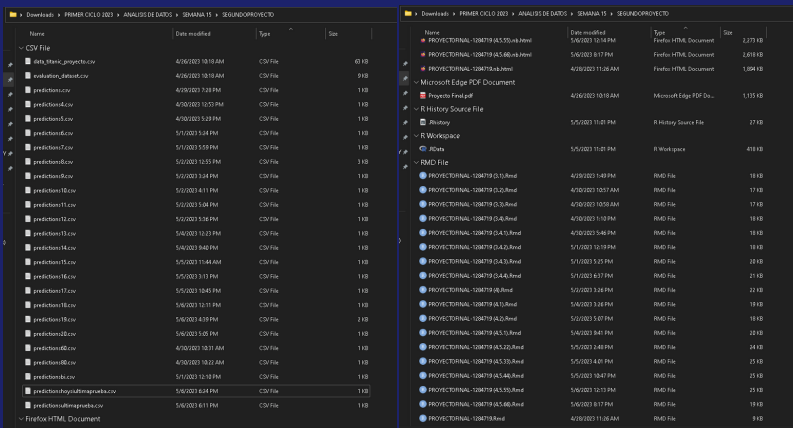
[82,]        5377              0
[83,]        9818              1
[84,]        3889              0
[85,]        4054              0
[86,]       10650              1
[87,]        6133              0
[88,]       1275              1
[89,]       1221              0
[90,]       9817              0
[91,]       2702              0
[92,]       2212              1
[93,]       4919              0
[94,]       6189              0
[95,]       6139              0
[96,]       9769              1
[97,]       1318              0
[98,]       2228              0
[99,]       8488              1
[100,]      10206              0
[101,]       9603              0
[102,]       6608              0
[103,]       981              0
[104,]      10484              1
[105,]       2769              1
[106,]       1934              0
[107,]       1016              0

```

SIEMPRE VERIFICANDO LA CANTIDAD DE FILAS CON LOS DATASET ORIGINALES PARA NO TENER MAS O MENOS DATOS.

```
635 # [r]
636
637 #m la variable predicciones_data estas las predicciones utilizo sep = ";" para especificar un separador personalizado
638 write.csv(predicciones_data, "prediccioneshoysultimaprueba.csv", sep = ";", row.names = FALSE)
639
640
641
Warning: attempt to set 'sep' ignored
```

PARA FINALIZAR SE GUARDA ESTE DATASET DE LAS PREDICCIONES OBTENIDAS POR EL MODELO EN UN FORMATO DE CSV SEPARADO POR COMA



DEJO EVIDENCIA DE TODAS LAS PRUEBAS DE LOS MODELOS QUE UTILICE PARA ESTA SOLUCION AL PROBLEMA.

COMENTARIO FINAL

CON MI MODELO DE APRENDIZAJE AUTOMATICO SUPERVISADO DE ARBOLES DE REGRESION LOGISTICA PUEDE OBSERVAR QUE LAS METRICAS OBTENIDAS EN EL ENTRENAMIENTO ERAN MAYORES A LAS OBTENIDAS CUANDO SE ENTRENA EL MODELO CON DATOS NUEVOS ESTO ME HACE PENSAR EN QUE EXISTE UN OVERFITTING CONSIDERO QUE ESO HACE QUE EL RESULTADO OBTENIDO EN KAGGLE NO SEA DEL TODO ACEPTABLE YA QUE EL MODELO NO APRENDE LO SUFICIENTEMENTE BIEN PARA ELLO HACE FALTA MAYOR ANALISIS DE DATOS POR QUE ESTO PUEDE CAUSAR DICHOS EFECTOS.