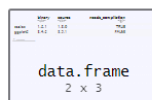


```

1 ---
2 title: "Actividad Integradora 2"
3 author: "Nallely Serna"
4 date: "r Sys.Date()"
5 output: word_document
6 ---
7
8 {r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10 options(repos = c(CRAN = "https://cran.rstudio.com/"))
11
12
13
14 {r}
15 # Librerías
16
17
18 #Algunas librerías no me funcionaron, así que no las utilicé
19 install.packages("dplyr")
20 install.packages("ggplot2")
21
22
23 library(dplyr)
24 library(ggplot2)
25
26
27

```



	binary <chr>	source <chr>	needs_compilation <lgl>
scales	1.2.1	1.3.0	TRUE
ggplot2	3.4.2	3.5.1	FALSE

```

29 # Preparación de la base de datos y análisis descriptivo
30
31 {r}
32 # Cargar datos
33 titanic <- read.csv("Titanic.csv")
34 titanic_test <- read.csv("Titanic_test.csv")
35
36 # Analizar datos faltantes
37 summary(titanic)
38
39 # Reemplazar valores faltantes
40 titanic$Age[is.na(titanic$Age)] <- median(titanic$Age, na.rm = TRUE)
41 titanic$Fare[is.na(titanic$Fare)] <- median(titanic$Fare, na.rm = TRUE)
42
43 # Eliminar la columna Cabin
44 titanic <- titanic[, !names(titanic) %in% "Cabin"]
45
46 # Análisis descriptivo
47 summary(titanic)
48 table(titanic$Survived, titanic$Pclass)
49 table(titanic$Survived, titanic$Sex)
50

```

```

PassengerId      Survived      Pclass      Name      Sex
Min.   :    1      Min.   :0.0000      Min.   :1.000      Length:1309      Length:1309
1st Qu.:   328      1st Qu.:0.0000      1st Qu.:2.000      Class :character      Class :character
Median :   655      Median :0.0000      Median :3.000      Mode  :character      Mode  :character
Mean   :   655      Mean   :0.3774      Mean   :2.295
3rd Qu.:   982      3rd Qu.:1.0000      3rd Qu.:3.000
Max.   :  1309      Max.   :1.0000      Max.   :3.000

      Age      SibSp      Parch      Ticket      Fare
Min.   : 0.17      Min.   :0.0000      Min.   :0.000      Length:1309      Min.   : 0.000
1st Qu.:21.00      1st Qu.:0.0000      1st Qu.:0.000      Class :character      1st Qu.: 7.896
Median :28.00      Median :0.0000      Median :0.000      Mode  :character      Median :14.454
Mean   :29.88      Mean   :0.4989      Mean   :0.385
3rd Qu.:39.00      3rd Qu.:1.0000      3rd Qu.:0.000
Max.   :80.00      Max.   :8.0000      Max.   :9.000

```

```

NA's :263
Cabin
Length:1309
Class :character
Mode :character

Embarked
Length:1309
Class :character
Mode :character

NA's :1

PassengerId      Survived      Pclass      Name      Sex
Min. : 1      Min. :0.0000      Min. :1.000      Length:1309      Length:1309
1st Qu.: 328      1st Qu.:0.0000      1st Qu.:2.000      Class :character      Class :character
Median : 655      Median :0.0000      Median :3.000      Mode :character      Mode :character
Mean : 655      Mean :0.3774      Mean :2.295
3rd Qu.: 982      3rd Qu.:1.0000      3rd Qu.:3.000
Max. :1309      Max. :1.0000      Max. :3.000

Age      SibSp      Parch      Ticket      Fare
Min. : 0.17      Min. :0.0000      Min. :0.000      Length:1309      Min. : 0.000
1st Qu.:22.00      1st Qu.:0.0000      1st Qu.:0.000      Class :character      1st Qu.: 7.896
Median :28.00      Median :0.0000      Median :0.000      Mode :character      Median : 14.454
Mean :29.50      Mean :0.4989      Mean :0.385
3rd Qu.:35.00      3rd Qu.:1.0000      3rd Qu.:0.000
Max. :80.00      Max. :8.0000      Max. :9.000
Embarked
Length:1309
Class :character
Mode :character

      1      2      3
0 137 160 518
1 186 117 191

female male
0      81 734
1     385 109

```

```

51
52 # Haz una partición de los datos (70-30) para el entrenamiento y la validación. Revisa la proporción
53 de sobrevivientes para la partición y la base original.
54 ```{r}
55
56 # Configurar semilla para reproducibilidad
57 set.seed(42)
58
59 # Crear índices para dividir los datos (70% para entrenamiento)
60 train_index <- sample(1:nrow(titanic), size = 0.7 * nrow(titanic))
61
62 # Dividir en datos de entrenamiento y validación
63 train_data <- titanic[train_index, ]
64 validation_data <- titanic[-train_index, ]
65
66 # Revisar la proporción de sobrevivientes
67 original_survival_rate <- mean(titanic$Survived)
68 train_survival_rate <- mean(train_data$Survived)
69 validation_survival_rate <- mean(validation_data$Survived)
70
71 # Mostrar las proporciones
72 original_survival_rate
73 train_survival_rate
74 validation_survival_rate
75
76
77

```

```

[1] 0.3773873
[1] 0.3766376
[1] 0.3791349

```

```

79 #Modelo logístico y selección de predictores
80
81 {r}
82 train_data <- na.omit(train_data)
83
84 # Crear un modelo logístico inicial
85 modelo_inicial <- glm(Survived ~ Pclass + Sex + Age + SibSp + Parch + Fare + Embarked,
86                       data = train_data, family = binomial)
87
88 # Ver resumen del modelo
89 summary(modelo_inicial)
90
91 # Selección de modelo por AIC
92 modelo_final <- step(modelo_inicial, direction = "both")
93
94 # Resumen del modelo final
95 summary(modelo_final)
96
97 # Ajustar un segundo modelo
98 modelo_alternativo <- glm(Survived ~ Pclass + Sex + Age, data = train_data, family = binomial)
99
100 # Comparar con el modelo final
101 summary(modelo_alternativo)
102 AIC(modelo_final, modelo_alternativo)
103

```



R Console

data.frame  
2 x 2

Description: df [2 x 2]

	df <dbl>	AIC <dbl>
modelo_final	6	710.3632
modelo_alternativo	4	714.9462

Call:  
glm(formula = Survived ~ Pclass + Sex + Age, family = binomial,  
data = train\_data)

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.5896	-0.5447	-0.3707	0.5195	2.4787

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	4.270745	0.499368	8.552	< 2e-16 ***
Pclass	-0.905588	0.129473	-6.994	2.66e-12 ***
Sexmale	-3.529521	0.204364	-17.271	< 2e-16 ***
Age	-0.023842	0.008412	-2.834	0.00459 **

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 1211.55 on 914 degrees of freedom  
Residual deviance: 706.95 on 911 degrees of freedom  
AIC: 714.95

Number of Fisher Scoring iterations: 5

```

104
105 El modelo tiene un AIC de 710.36, que es relativamente bajo, lo que sugiere que el modelo es
    adecuado, aunque no es el modelo más ajustado posible. El AIC se utiliza para comparar la calidad del
    ajuste entre diferentes modelos, por lo que siempre es importante considerar si podemos mejorar este
    valor mediante la inclusión de más variables o transformaciones.

```

```

106
107 # Análisis del modelo
108
109 ```{r}
110 # Desviación residual y nula
111 desviacion_residual <- modelo_final$deviance
112 desviacion_nula <- modelo_final$null.deviance
113 desviacion_explicada <- 1 - (desviacion_residual / desviacion_nula)
114
115 # Prueba de razón de verosimilitud
116 anova(modelo_inicial, modelo_final, test = "Chisq")
117
118
119 # Interpretación de los coeficientes
120 exp(coef(modelo_final))
121
122
123 ```

```

#### Analysis of Deviance Table

```

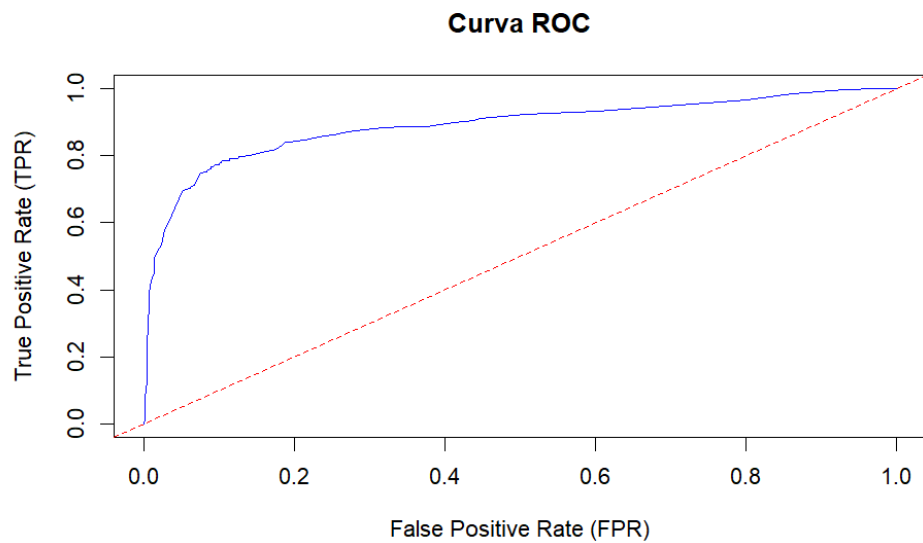
Model 1: Survived ~ Pclass + Sex + Age + SibSp + Parch + Fare + Embarked
Model 2: Survived ~ Pclass + Sex + Age + SibSp + Fare
  Resid. Df Resid. Dev Df Deviance Pr(>Chi)
1      906      694.46
2      909      698.36 -3   -3.9001  0.2725
(Intercept)      Pclass      Sexmale      Age      SibSp      Fare
62.63548725  0.45577156  0.02755284  0.97414712  0.75243074  1.00322799

```

```

124
125 # Predicción y evaluación en datos de entrenamiento
126
127 ```{r}
128 # Generar predicciones (probabilidades)
129 predicciones_train <- predict(modelo_final, train_data, type = "response")
130
131 # Crear un rango de umbrales
132 umbrales <- seq(0, 1, by = 0.01)
133
134 # Inicializar vectores para TPR y FPR
135 tpr <- numeric(length(umbrales)) # True Positive Rate
136 fpr <- numeric(length(umbrales)) # False Positive Rate
137
138 # Calcular TPR y FPR para cada umbral
139 for (i in seq_along(umbrales)) {
140   umbral <- umbrales[i]
141   pred_clase <- ifelse(predicciones_train >= umbral, 1, 0)
142
143   # Matriz de confusión
144   tp <- sum(pred_clase == 1 & train_data$Survived == 1) # Verdaderos Positivos
145   fp <- sum(pred_clase == 1 & train_data$Survived == 0) # Falsos Positivos
146   fn <- sum(pred_clase == 0 & train_data$Survived == 1) # Falsos Negativos
147   tn <- sum(pred_clase == 0 & train_data$Survived == 0) # Verdaderos Negativos
148
149   # TPR y FPR
150   tpr[i] <- tp / (tp + fn) # Sensibilidad
151   fpr[i] <- fp / (fp + tn) # 1 - Especificidad
152 }
153
154 # Graficar la Curva ROC
155 plot(fpr, tpr, type = "l", col = "blue", xlab = "False Positive Rate (FPR)",
156      ylab = "True Positive Rate (TPR)", main = "Curva ROC")
157 abline(0, 1, col = "red", lty = 2) # Línea diagonal (clasificación aleatoria)
158
159 # Calcular el Área Bajo la Curva (AUC) de forma manual
160 auc <- sum(diff(fpr)) * (head(tpr, -1) + tail(tpr, -1)) / 2
161 auc

```



```

164
165 ~~~{r}
166
167 # Crear la matriz de confusión manualmente
168 table(Predicción = pred_clase_train, Real = train_data$Survived)
169
170 # Calcular precisión, sensibilidad, especificidad y F1
171 tp <- sum(pred_clase_train == 1 & train_data$Survived == 1) # Verdaderos positivos
172 tn <- sum(pred_clase_train == 0 & train_data$Survived == 0) # Verdaderos negativos
173 fp <- sum(pred_clase_train == 1 & train_data$Survived == 0) # Falsos positivos
174 fn <- sum(pred_clase_train == 0 & train_data$Survived == 1) # Falsos negativos
175
176 precision <- tp / (tp + fp)
177 sensibilidad <- tp / (tp + fn)
178 especificidad <- tn / (tn + fp)
179 f1 <- 2 * (precision * sensibilidad) / (precision + sensibilidad)
180
181 # Mostrar los resultados
182 cat("Precisión: ", precision, "\n")
183 cat("Sensibilidad: ", sensibilidad, "\n")
184 cat("Especificidad: ", especificidad, "\n")
185 cat("F1: ", f1, "\n")
186
187 ~~~

```

	Real	
Predicción	0	1
0	510	74
1	61	270

```

Precisión: 0.81571
Sensibilidad: 0.7848837
Especificidad: 0.8931699
F1: 0.8

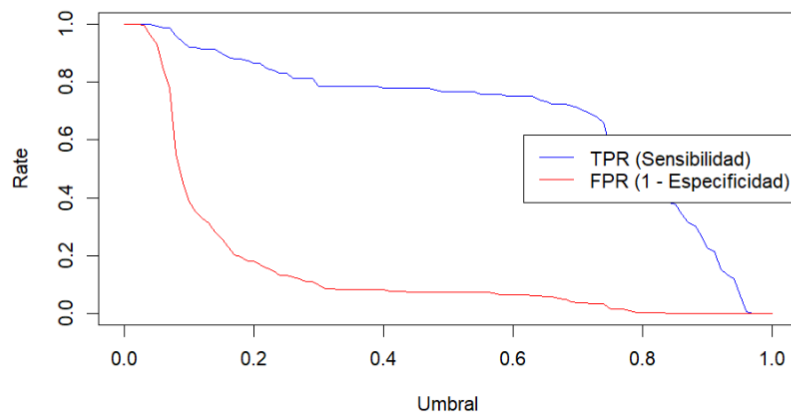
```

```

188 #Validación y conclusión
189 
190 
191 {r}
192 # Generar predicciones (probabilidades) para los datos de validación
193 predicciones_validation <- predict(modelo_final, validation_data, type = "response")
194 
195 # Crear un rango de umbrales
196 umbrales <- seq(0, 1, by = 0.01)
197 
198 # Inicializar vectores para TPR, FPR y el criterio de Youden
199 tpr <- numeric(length(umbrales)) # True Positive Rate
200 fpr <- numeric(length(umbrales)) # False Positive Rate
201 youden_index <- numeric(length(umbrales)) # Índice de Youden
202 
203 # Calcular TPR, FPR y Youden para cada umbral
204 for (i in seq_along(umbrales)) {
205   umbral <- umbrales[i]
206   pred_clase <- ifelse(predicciones_validation >= umbral, 1, 0)
207 
208   # Matriz de confusión
209   tp <- sum(pred_clase == 1 & validation_data$Survived == 1) # Verdaderos Positivos
210   fp <- sum(pred_clase == 1 & validation_data$Survived == 0) # Falsos Positivos
211   fn <- sum(pred_clase == 0 & validation_data$Survived == 1) # Falsos Negativos
212   tn <- sum(pred_clase == 0 & validation_data$Survived == 0) # Verdaderos Negativos
213 
214   # TPR y FPR
215   tpr[i] <- tp / (tp + fn) # Sensibilidad
216   fpr[i] <- fp / (fp + tn) # 1 - Especificidad
217 
218   # Índice de Youden
219   youden_index[i] <- tpr[i] - fpr[i]
220 }
221 
222 # Encontrar el umbral óptimo (el máximo índice de Youden)
223 umbral_optimo <- umbrales[which.max(youden_index)]
224 cat("El umbral óptimo es:", umbral_optimo, "\n")
225 
226 # Graficar TPR y FPR vs. Umbrales
227 plot(umbrales, tpr, type = "l", col = "blue", ylim = c(0, 1),
228      ylab = "Rate", xlab = "Umbral", main = "TPR y FPR vs. Umbral")
229 lines(umbrales, fpr, col = "red")
230 legend("right", legend = c("TPR (Sensibilidad)", "FPR (1 - Especificidad)"),
231       col = c("blue", "red"), lty = 1)
232 
233 

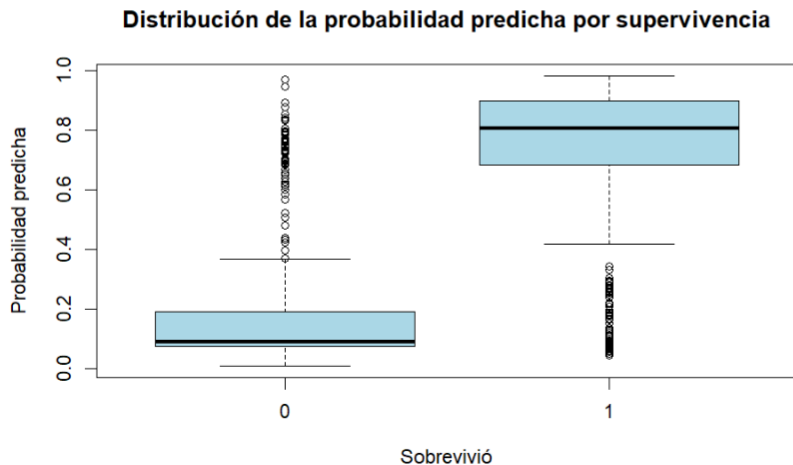
```

TPR y FPR vs. Umbral



El umbral óptimo es: 0.44

```
236 {r}
237 # Crear una variable para la supervivencia como factor
238 train_data$Survived <- as.factor(train_data$Survived)
239
240 # Crear el gráfico de cajas
241 boxplot(predicciones_train ~ Survived,
242         data = train_data,
243         col = "lightblue",
244         xlab = "Sobrevivió",
245         ylab = "Probabilidad predicha",
246         main = "Distribución de la probabilidad predicha por supervivencia")
247
248
```



```
249 {r}
250 # Predicción para datos de prueba
251 predicciones_test <- predict(modelo_final, titanic_test, type = "response")
252 pred_clase_test <- ifelse(predicciones_test >= umbral_optimo, 1, 0)
253
254 # Agregar las predicciones a los datos de prueba
255 titanic_test$Predicted_Survived <- pred_clase_test
256 head(titanic_test)
257
258
259
```

Description: df [6 x 12]

	PassengerId <int>	Pclass <int>	Name <chr>	Sex <chr>	Age <dbl>	SibSp <int>	Parch <int>
1	892	3	Kelly, Mr. James	male	34.5	0	0
2	893	3	Wilkes, Mrs. James (Ellen Needs)	female	47.0	1	0
3	894	2	Myles, Mr. Thomas Francis	male	62.0	0	0
4	895	3	Wirz, Mr. Albert	male	27.0	0	0
5	896	3	Hirvonen, Mrs. Alexander (Helga E Lindqvist)	female	22.0	1	1
6	897	3	Svensson, Mr. Johan Cervin	male	14.0	0	0

6 rows | 1-8 of 12 columns

El análisis de regresión logística aplicado al conjunto de datos del Titanic muestra que las variables más influyentes en la supervivencia de los pasajeros son la clase social, el sexo y la edad. Estos hallazgos son consistentes con los eventos históricos conocidos sobre el hundimiento del Titanic. El modelo tiene un buen ajuste, pero aún existen áreas de mejora, especialmente en lo que respecta a las variables menos significativas.