

```

1 #####
2 # TITANIC #
3 # Pepa Montero #
4 #####
5
6 # PAQUETES UTILIZADOS
7 install.packages("xtable")
8 library(xtable)
9 install.packages("e1071")
10 library(e1071)
11
12 #           Análisis de los datos (ap. 3 del pdf)
13
14 setwd("") # directorio donde se encuentran los datos
15
16 # Primero guardamos los datos en una variable
17 data <- load(file = 'titanic_train.rda')
18 attach(titanic_train)
19
20 # Visualizamos el dataset
21 View(titanic_train)
22
23 # LIMPIEZA DEL DATASET
24 # Eliminar columnas innecesarias
25 titanic_train$Ticket <- NULL
26 titanic_train$PassengerId <- NULL
27 titanic_train$Name <- NULL
28
29 # Ver huecos vacíos
30 # 1. Cabin
31 table(Cabin)
32 # vemos que 687 pasajeros toman el valor " "
33 # eliminamos la variable Cabin
34 titanic_train$Cabin <- NULL
35 # 2. Age
36 colSums(is.na(titanic_train))
37 # vemos que 177 pasajeros no tienen valor de edad
38 # lo solucionaremos más adelante
39
40 # Variable Embarked
41 table(Embarked)
42 # Vemos que 644 pasajeros son de Southampton
43 titanic_train$Embarked <- NULL
44
45 # Volvemos a visualizar los datos, ahora solo con las variables a estudiar
46 View(titanic_train)
47
48
49 #           Análisis unidimensional
50
51
52 # SUPERVIVIENTES
53 # Tablas de frecuencias
54 table_survived <- table(Survived)
55 frec_survived <- prop.table(table_survived)
56 percent_survived <- frec_survived*100
57 # Resultado: 61.62% de muertos y 38.38% de supervivientes
58 # Exportar tablas para latex
59 table_survived_latex <- xtable(table_survived)
60 frec_survived_latex <- xtable(frec_survived)
61 # Pie chart
62 percent_0 <- toString(round(percent_survived[1], 2))
63 percent_1 <- toString(round(percent_survived[2], 2))
64 pie(percent_survived, labels=c(paste("Fallecidos:", percent_0, "%"),
65                                paste("Supervivientes:", percent_1, "%")),
66      col=c("red", "lightblue"), main="Survived")
67
68 # SEXO
69 # Tablas de frecuencias
70 table_sex <- table(Sex)
71 frec_sex <- prop.table(table_sex)
72 percent_sex <- frec_sex*100
73 # Resultado: 35.24% de mujeres y 64.76% de hombres

```

```

74 # Exportar tablas para latex
75 table_sex_latex <- xtable(table_sex)
76 frec_sex_latex <- xtable(frec_sex)
77 # Pie chart
78 percent_female <- toString(round(percent_sex[1], 2))
79 percent_male <- toString(round(percent_sex[2], 2))
80 pie(percent_sex, labels=c(paste("Mujeres:", percent_female, "%"),
81                             paste("Hombres:", percent_male, "%")),
82      col=c("yellow", "magenta"), main="Sex")
83
84 # CLASE
85 # Tablas de frecuencias
86 table_class <- table(Pclass)
87 frec_class <- prop.table(table_class)
88 percent_class <- frec_class*100
89 # Resultado: 24.24% Clase 1, 20.65% Clase 2, 55.10% Clase 3
90 # Exportar tablas para latex
91 table_class_latex <- xtable(table_class)
92 frec_class_latex <- xtable(frec_class)
93 # Diagrama de barras
94 percent_class1 <- toString(round(percent_class[1], 2))
95 percent_class2 <- toString(round(percent_class[2], 2))
96 percent_class3 <- toString(round(percent_class[3], 2))
97 barplot(percent_class, legend.text=c("Clase 1", "Clase 2", "Clase 3"),
98        col=c("yellow", "green", "blue"), main="Clase")
99
100 # EDAD
101 # .....
102 # Resolver datos NA
103 # Sobreescribimos los NA con el valor de la media
104 summary(Age)
105 # Nos devuelve que la media es 29.70
106 titanic_train$Age[is.na(titanic_train$Age)] <- 29.70 # sobreescribimos
107 Age = titanic_train$Age
108 # De esta forma la media se mantiene
109 # .....
110 # Convertir en franjas de edad
111 Age_bands <- cut(Age, 10)
112 # .....
113 # Tablas de frecuencias
114 table_age <- table(Age_bands)
115 frec_age <- prop.table(table_age)
116 # Exportar tablas para latex
117 table_age_latex <- xtable(table_age)
118 frec_age_latex <- xtable(frec_age)
119 # Histograma
120 breaks <- c(0.34, 8.38, 16.3, 24.3, 32.3, 40.2, 48.2, 56.1, 64.1, 72, 80)
121 hist(Age, breaks=breaks, freq = T, col="lightgreen", include.lowest = T,
122      right = T, main = "Age", xlab=NULL, ylab = "Frecuencia")
123 # Datos
124 summary(Age)
125 mean_age <- mean(Age)
126 var_age <- var(Age)
127 dt_age <- sd(Age)
128 # La media de edad de los pasajeros es 29.70
129 # El rango intercuartílico nos dice que el 50% de los pasajeros estaban
130 # entre los 22 y los 35 años
131 # La varianza es 169.05 y la desviación típica es 13
132 # .....
133 # Coeficiente de variación
134 cv_age <- dt_age / mean_age * 100 # = 43.78
135 # Coeficiente de asimetría de Fisher
136 b3_age <- moment(Age, order = 3, center = TRUE)
137 af_age <- b3_age / dt_age**3 # = 0.43
138
139 # FAMILIARES A BORDO
140 Familiares = SibSp + Parch
141 # Tablas de frecuencias
142 table_fam = table(Familiares)
143 frec_fam = prop.table(table_fam)
144 # Observamos que el 60% de los pasajeros viajaban solos
145 # Exportar tablas para latex
146 table_fam_latex <- xtable(table_fam)

```

```

147 freq_fam_latex <- xtable(freq_fam)
148 # Datos
149 median(Familiares) # = 0
150 # Diagrama de puntos
151 plot(freq_fam, type = "o", main = "Familiares",
152       xlab = "n° de familiares a bordo", ylab = "f. relativa",
153       col = "darkred")
154
155 # PRECIO DEL TICKET
156 summary(Fare)
157 # Diagrama de cajas y bigotes
158 boxplot(Fare, col="lightpink", main="Fare")
159 # Medidas de dispersión
160 mean_fare = mean(Fare)
161 dt_fare = sd(Fare) # = 49.69
162 b3_fare = moment(Fare, order = 3, center = TRUE)
163 cv_fare = dt_fare / mean_fare * 100 # = 154.3
164 af_fare = b3_fare / dt_fare**3 # = 4.77
165
166
167 # Análisis bidimensional
168
169
170 # SEXO VS. SUPERVIVENCIA
171 # Tablas de frecuencias
172 table2_sex_surv <- table(Sex, Survived)
173 freq2_sex_surv <- prop.table(table2_sex_surv)
174 percent2_sex_surv <- prop.table(table2_sex_surv, 1)*100
175 # Exportar tablas para latex
176 # xtable(table2_sex_surv)
177 # xtable(percent2_sex_surv)
178 # Barras apiladas
179 barplot(table2_sex_surv, names.arg = c("Fallecidos", "Supervivientes"),
180         col = c("lightpink", "lightblue"), main = "Sex vs. Survived",
181         legend.text = c("Mujeres", "Hombres"))
182 # Test chi-cuadrado
183 chi_SS <- chisq.test(table2_sex_surv)
184
185 # EDAD VS. SUPERVIVENCIA
186 # Tablas de frecuencias
187 table2_age_surv <- table(Survived, Age_bands)
188 freq2_age_surv <- prop.table(table2_age_surv)
189 percent2_age_surv <- prop.table(table2_age_surv, 1)*100
190 # Exportar tablas para latex
191 # xtable(table(Age_bands, Survived))
192 # xtable(prop.table(table(Age_bands, Survived), 1)*100)
193 # Barras apiladas
194 barplot(table2_age_surv, names.arg = breaks[1:10],
195         xlab = "Edad mínima de cada intervalo", main = "Age vs. Survived",
196         legend.text = c("Fallecidos", "Supervivientes"))
197 # Test chi-cuadrado
198 chi_AS <- chisq.test(table2_age_surv)
199 # Al ejecutarlo recibimos un error,
200 # por lo que podría no ser concluyente
201
202 # FAMILIARES VS. SUPERVIVENCIA
203 # Tablas de frecuencias
204 table2_fam_surv <- table(Survived, Familiares)
205 freq2_fam_surv <- prop.table(table2_fam_surv)
206 percent2_fam_surv <- prop.table(table2_fam_surv, 1)*100
207 # Exportar tablas para latex
208 # xtable(table(Familiares, Survived))
209 # xtable(prop.table(table(Familiares, Survived), 1)*100)
210 # Barras agrupadas
211 barplot(table2_fam_surv, xlab = "n° de familiares a bordo",
212         main = "Familiares vs. Survived", beside = T,
213         legend.text = c("Fallecidos", "Supervivientes"))
214 # Convertir Familiares en una variable binaria
215 Familiares_bin <- Familiares
216 Familiares_bin[Familiares_bin > 0] <- 1
217 table2_famb_surv <- table(Survived, Familiares_bin)
218 barplot(table2_famb_surv, names.arg = c("No", "Sí"), xlab = "Tiene familiares a bordo",
219         main = "Familiares vs. Survived", beside = T,

```

```

220         legend.text = c("Fallecidos", "Supervivientes"))
221 # Test chi-cuadrado
222 chi_FS <- chisq.test(table2_famb_surv)
223
224 # CLASS VS. SUPERVIVENCIA
225 # Class vs. Fare
226 boxplot(Fare ~ Pclass)
227 # Vemos que claramente están relacionadas
228 # Luego podemos continuar sólo fijándonos en la clase
229 # Tablas de frecuencias
230 table2_class_surv <- table(Pclass, Survived)
231 freq2_class_surv <- prop.table(table2_class_surv)
232 percent2_class_surv <- prop.table(table2_class_surv, 1)*100
233 # Exportar tablas para latex
234 # xtable(table2_class_surv)
235 # xtable(percent2_class_surv)
236 # Barras apiladas
237 barplot(table2_class_surv, names.arg = c("Fallecidos", "Supervivientes"),
238         col = c("lightpink", "lightblue", "lightgreen"), main = "Pclass vs. Survived",
239         legend.text = c("Clase 1", "Clase 2", "Clase 3"))
240 # Test chi-cuadrado
241 chi_CS <- chisq.test(table2_class_surv)
242
243 # FARE VS. AGE
244 # Diagrama de dispersión
245 plot(Age, Fare)
246 # Covarianza
247 cov(Age, Fare) # = 59.16
248 # Correlación
249 cor(Age, Fare) # = 0.09

```