



Trabajo Práctico N° 4

Arboles de Decision

Presentado en la fecha: 21/10/2023

Hecho por: Huarca Brian

Nicolas Benitez

Facundo Rodriguez

Derlis Walter Hodge

Contents

Resumen	3
Sumario	4
0.1 Sumario	4
Objetivo	5
Desarrollo	6
1 Dataset	6
1.1 Informacion del dataset	6
1.2 Diccionario de datos	7
2 Dataset	9
2.1 Deteccion de outliers	9
3 RPart	10
3.1	10
3.2 Arbol de Decision	11
3.3 Analisis de Importancia de las Variables	12
3.4 Matriz de Confusion	13
4 RForest	14
4.1 Análisis con Random Forest	14
4.2 Resumen Matriz de Confusion	16
5 Tecnicas de balanceo de clases	17
5.1 Undersampling	17
5.2 Oversampling	17

Conclusión	19
Anexo	20

Resumen

Para la siguiente práctica de árboles de decisión utilizaremos un set de datos que contiene observaciones del programa Aprender 2016 - Primaria.

Sumario

0.1 Sumario

- Preparación del dataset (Valores Null, Outliers, Casteo de Datos).
- Clasificación de Variables.
- Discretización de Valores.
- Análisis del Árbol de Decisión.
- Análisis de Importancia de las Variables
- Análisis de Matriz de Confusión.
- Comparación y Análisis de Modelos RandomForest

Objetivo

El informe tiene por objetivo el analisis de datos para realizar una clasificacion en base a arboles de decision con los metodos RPart Y RForest. Haciendo un analisis de Variables para determinar una predictora, realizar si fuese necesario la discretizacion y la comparacion de distintos modelos con RForest en caso de visualizar altos niveles de Error OOB.

Dataset

1.1 Informacion del dataset

Es el dispositivo nacional de evaluación de los aprendizajes de los estudiantes y de sistematización de información acerca de algunas condiciones en las que ellos se desarrollan.

El documento se basa en una evaluación a jóvenes de una institución para medir su desempeño en sus clases de matemáticas y prácticas del lenguaje. Se nos presentan distintas preguntas realizadas a los alumnos como a qué tipo de institución asisten, en qué sector, si repitió de grado, entre otras.

1.2 Diccionario de datos

- **ap1** = ¿Cuántos años tenés?.
- **ap2** = ¿Sos varón o mujer?.
- **ap3** = ¿Fuiste al jardín de infantes?.
- **ap4** = ¿Alguna vez repetiste de grado?.
- **ap5** = ¿Cómo leés?.
- **ap6** = ¿Cómo escribís?.
- **ap7** = ¿Cómo resolvés problemas o cuentas de matemática?
- **ap8a** = ¿Los maestros y maestras pueden explicar sin que nadie los interrumpa?
- **ap8b** = ¿Los maestros y maestras se enojan con ustedes?
- **ap8c** = ¿Los maestros y maestras explican los temas hasta que los entendés?
- **ap9** = ¿Te llevás bien con tus compañeros?
- **ap10** = Si te dijeran que tenés que cambiarte de escuela, ¿cómo te sentirías?
- **ap11** = Cuando estás en la escuela, ¿estás contento?
- **ap12** = Cuando estás en la escuela, ¿estás aburrido?
- **ap13** = Cuando estás en la escuela, ¿te sentís incómodo?
- **ap14** = Resolver la prueba de Lengua, ¿te fue fácil o difícil?
- **ap15** = Resolver la prueba de Matemática, ¿te fue fácil o difícil?
- **ponder** = Factor de expansión
- **cod_{provincia}** = *Jurisdiccin*
- **lpondera** = Factor de expansión Lengua
- **lpuntaje** = Puntaje Lengua
- **ldesemp** = Nivel de Desempeño en Lengua

- **mpondera** = Factor de expansión Matemática
- **mpuntaje** = Puntaje Matemática
- **mdesemp** = Nivel de Desempeño en Matemática
- **sector** = Sector de Gestión
- **ambito** = Espacio donde se ubica Institucion Rural-Urbano
- **qvulneraa** = Cuartil de alumnos según el porcentaje de hogares en el radio de la escuela en estrato socioeconómico muy bajo Total País
- **qvulneraap** = Cuartil de alumnos según el porcentaje de hogares en el radio de la escuela en estrato socioeconómico muy bajo por jurisdicción
- **iclima** = Índice de Clima escolar
- **autoconl** = Autoconcepto en Lengua alumnos que participaron en evaluación de Lengua
- **autoconm** = Autoconcepto en Matemática alumnos que participaron en evaluación de Matemática

Dataset

2.1 Deteccion de outliers

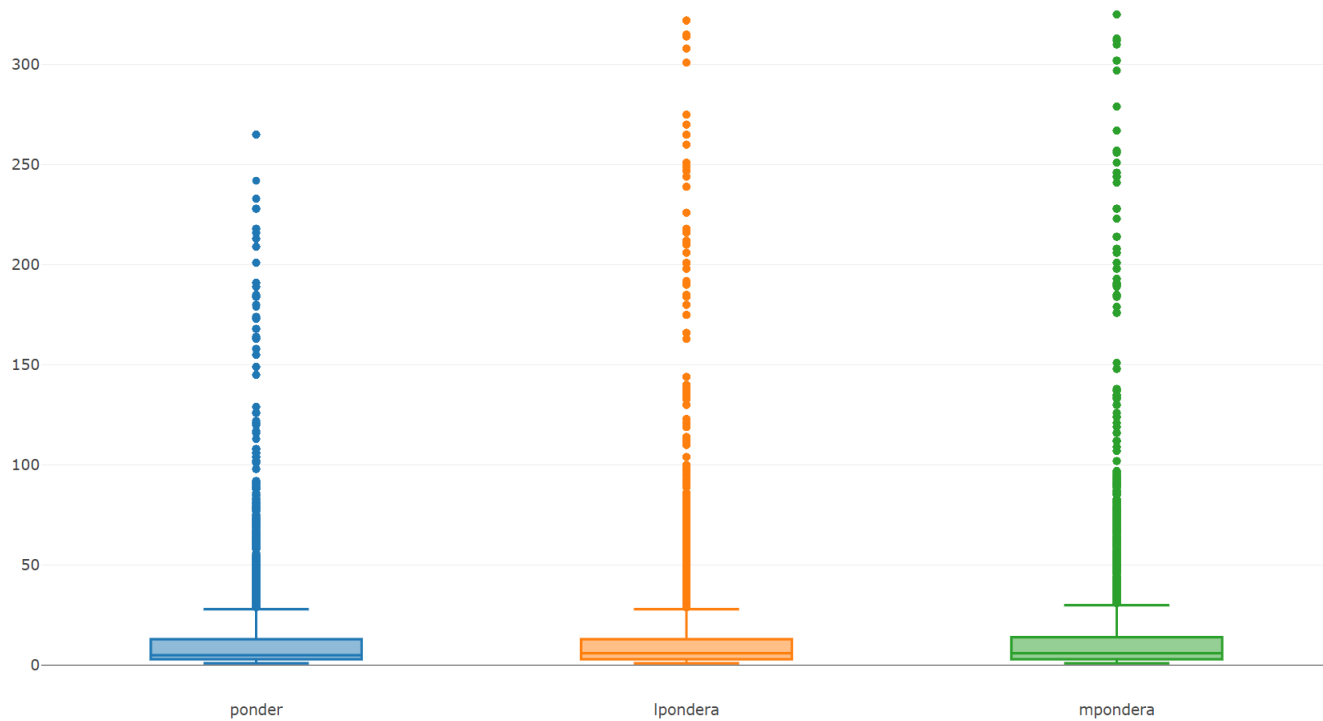


Figure 1: Boxplot

En este BoxPlot podemos ver una cantidad de Outliers Heterogeneos de las variables Ponder, Mpondera y Lpondera. Estas variables representan el factor de expansion general, el factor de expansion en practicas del lenguaje y en matematicas.

RPart

3.1

Separamos los datos en entrenamiento y testing (70/30%)

- entrenamiento: 34153 observaciones
- testing: 14636 observaciones

Con la función Rpart creamos un árbol.

El Rpart decide con qué variables se queda para la creación del árbol a partir de cuáles son las mas significativas.

Todo en concordancia con nuestro análisis previo. Así sucesivamente se van formando las ramas del árbol hasta quedar con las hojas/categorías finales.

3.2 Arbol de Decision

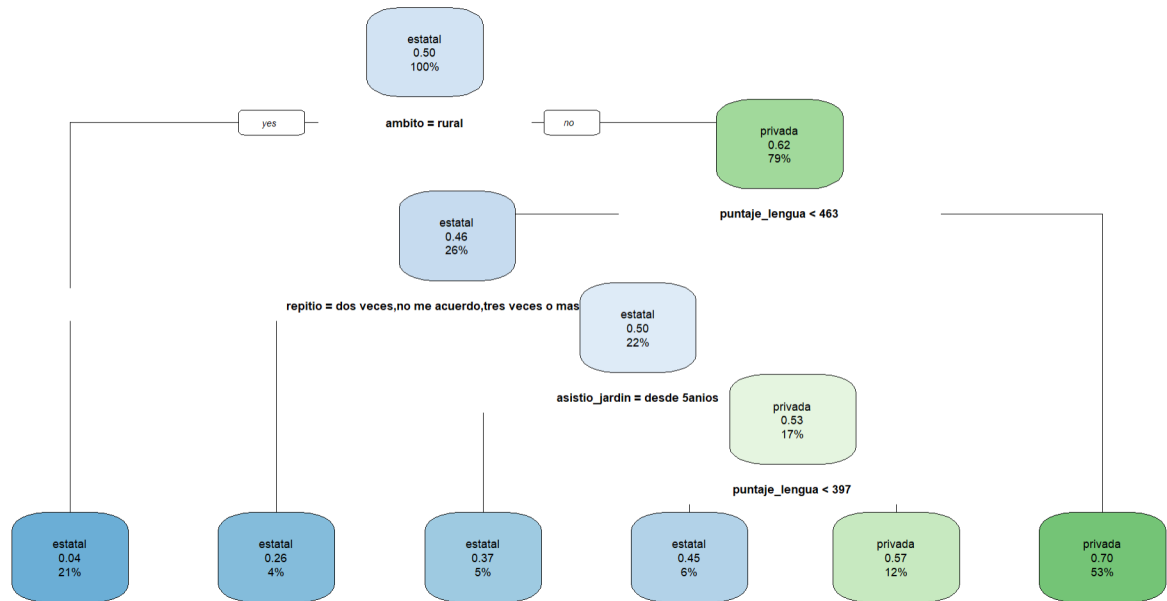


Figure 2

Aca podemos apreciar nuestro arbol de decision el cual nuestra variable "Predictora" a utilizar es la variable 'Sector', la cual indica si el alumno estudia en una institucion publica o privada, por lo cual nos da de resultado en nuestra raiz que un 50 del 100 de nuestra muestra total estudia en una institucion publica. Bajando un poco en detalles en el arbol notamos que la variable ambito es igual a Rural por ende si vamos a la Izquierda, osea por el lado del SI, vemos que la gran mayoria tiene chances de estudiar en un ambito rural y un 0.04 NO.

3.3 Analisis de Importancia de las Variables

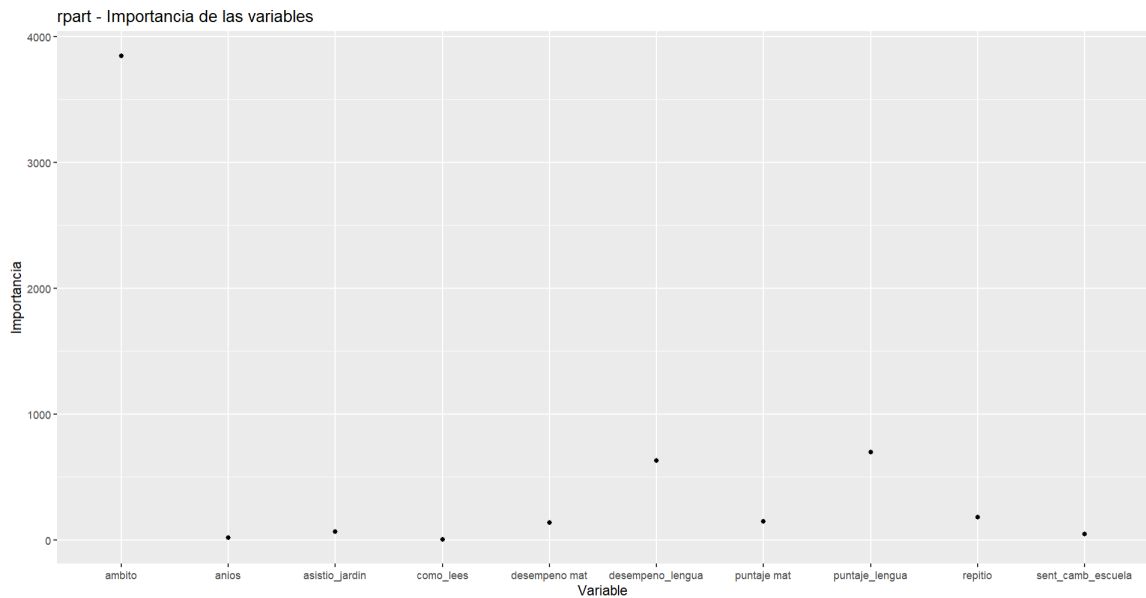


Figure 3

Aca podemos observar que entre las variables con mayor importancia estan las variables que preguntan si los alumnos alguna vez repitieron, cuando comenzaron el jardin, su puntaje y desempeño en matematicas y practicas del lenguaje. Pero la que mas resalta entre ellas es el ambito escolar en el cual ellos estudian.

3.4 Matriz de Confusion

Aquí podemos visualizar lo que es nuestra Matriz de Confusion sobre el modelo RPart realizado, en la cual se puede apreciar nuestra Taza de Exactitud (Accuracy) es del 0.72, un P-Value = $2.2e-16$ que se encuentra por debajo del 0.05 por lo cual se rechaza la hipótesis nula, una sensibilidad del 0.82 y una especificidad del 0.66.

```
> tabla_mc_aprender
Confusion Matrix and Statistics

      Predicho
Actual  estatal privada
estatal  4224   3161
privada   888   6363

      Accuracy : 0.7234
      95% CI   : (0.716, 0.7306)
 No Information Rate : 0.6507
 P-Value [Acc > NIR] : < 2.2e-16

      Kappa : 0.4482

 Mcnemar's Test P-Value : < 2.2e-16

      Sensitivity : 0.8263
      Specificity : 0.6681
   Pos Pred Value : 0.5720
   Neg Pred Value : 0.8775
      Prevalence : 0.3493
   Detection Rate : 0.2886
 Detection Prevalence : 0.5046
   Balanced Accuracy : 0.7472

      'Positive' Class : estatal

> |
```

Figure 4

RForest

4.1 Análisis con Random Forest

Para realizar un modelo mas profundo y buscar mejorar la exactitud de la categorización procedimos a utilizar Random Forest, obteniendo los siguientes resultados.

	edad	sexo	grado			
15 años	:7202	Femenino :17229	1er año/10mo grado nivel Polimodal o 3er año nivel Secundario:7			
14 años	:7020	Masculino:14636	2do año/11vo grado nivel Polimodal o 4to año nivel Secundario:7			
16 años	:6738		3er año/12vo grado nivel Polimodal o 5to año nivel Secundario:5			
13 años	:5854		8vo grado nivel Primario/Polimodal o 1er año nivel Secundario:4			
17 años	:4773		9no grado nivel Primario/Polimodal o 2do año nivel Secundario :6			
18 años o mas:	201					
(other)	: 77					
	dias_con_hambre	ult_12meses_atac_fisi	ult_12meses_part_pelea	lesion_grave		
Algunas veces:	2743	Ninguna :27076	Ninguna :25422	Ninguna	:21559	
Casi siempre :	341	1 vez : 2648	1 vez : 3472	1 vez	: 6138	
Nunca	:22093	2 o 3 veces : 1192	2 o 3 veces : 1758	2 o 3 veces	: 2893	
Rara vez	: 6568	4 o 5 veces : 334	4 o 5 veces : 462	4 o 5 veces	: 670	
Siempre	: 120	12 o mas veces: 329	12 o mas veces: 366	6 o 7 veces	: 224	
		6 o 7 veces : 145	6 o 7 veces : 196	12 o mas veces:	208	
		(other) : 141	(Other) : 189	(Other)	: 173	

Figure 5

consideraste_suicidio	plan_de_suicidio	veces_intentaste_suici	amigos_cercanos	edad_f
No:25433	No:26838	0 veces :27833	0 : 1604	Nunca probÃ© cigarrillos
Si: 6432	Si: 5027	1 vez : 2470	1 : 2565	12 o 13 aÃ±os
		2 o 3 veces : 1022	2 : 5281	14 o 15 aÃ±os
		4 o 5 veces : 270	3 o mÃ¡s:22415	10 u 11 aÃ±os
		6 o mas veces: 270		16 o 17 aÃ±os
				8 o 9 aÃ±os
				(other)

Figure 6

edad_fumador	ult_30dias_fumaste	dejas
Nunca prob� cigarrillos:19819	0 d�as :26085	No
12 o 13 a��os : 4124	1 o 2 d�as : 2450	No fum� cigarrillos durante los �ltimos 12 mese
14 o 15 a��os : 3984	10 a 19 d�as: 617	Nunca fum� cigarrillos
10 u 11 a��os : 1581	20 a 29 d�as: 396	S�
16 o 17 a��os : 1109	3 a 5 d�as : 961	
8 o 9 a��os : 754	6 a 9 d�as : 664	
(Other) : 494	Los 30 d�as : 692	
edad_alcohol	dias_alcohol	veces_tomaste_alcohol
12 o 13 a��os :8558	0 d�as :14732	0 veces :20174
Nunca tom� alcohol m��s que unos pocos sorbos:8392	1 o 2 d�as : 8390	1 o 2 veces : 7149
14 o 15 a��os :7923	10 a 19 d�as: 1748	10 o m��s veces: 1569
10 o 11 a��os :3028	20 a 29 d�as: 514	3 a 9 veces : 2973
8 o 9 a��os :1393	3 a 5 d�as : 3756	
7 a��os o menos :1325	6 a 9 d�as : 2403	
(other) :1246	Los 30 d�as : 322	

Figure 7

problem_entorno_alcohol	edad_drogas	veces_maria	ult_30dias_consu_maria
0 veces :28298	Nunca us� Drogas:27451	0 veces :27727	0 veces :29281
1 o 2 veces : 2459	14 o 15 a��os : 2142	1 o 2 veces : 1378	1 o 2 veces : 1214
10 o m��s veces: 422	12 o 13 a��os : 1055	10 a 19 veces : 532	10 a 19 veces : 291
3 a 9 veces : 686	16 o 17 a��os : 872	20 veces o m��s: 1128	20 veces o m��s: 353
	10 o 11 a��os : 175	3 a 9 veces : 1100	3 a 9 veces : 726
	7 a��os o menos : 86		
	(Other) : 84		
veces_consu_metanfet	tuvistesexo	edadsexo_lra_vez	cuantas
0 veces :31306	No:18952	Nunca tuve relaciones sexuales:19994	1 persona
1 o 2 veces : 360	S�:12913	14 a��os : 3297	2 personas
10 a 19 veces : 37		15 a��os : 3231	3 personas
20 veces o m��s: 48		13 a��os : 1942	4 personas
3 a 9 veces : 114		16 o 17 a��os : 1798	5 personas
		12 a��os : 919	6 o m��s personas
		(Other) : 684	Nunca tuve relaciones sexuales

Figure 8

4.2 Resumen Matriz de Confusion

```
> emse_red_rf_pred_mc
Confusion Matrix and Statistics

      Predicho
Actual   No   Si
No  7767  284
Si   629  879

      Accuracy : 0.9045
      95% CI : (0.8984, 0.9103)
No Information Rate : 0.8783
P-Value [Acc > NIR] : 3.559e-16

      Kappa : 0.6037

McNemar's Test P-Value : < 2.2e-16

      Sensitivity : 0.9251
      Specificity : 0.7558
      Pos Pred Value : 0.9647
      Neg Pred Value : 0.5829
      Prevalence : 0.8783
      Detection Rate : 0.8125
      Detection Prevalence : 0.8422
      Balanced Accuracy : 0.8404

      'Positive' Class : No
```

Figure 9

En esta matriz nuestros casos positivos se representan con los casos Negativos como Verdaderos Positivos y los valores Positivos como Falsos Negativos. Observamos que Nuestra Exactitud del modelo es del 0.90 y nuestro P-Value es igual a 3.559e-16, Nuestra Sensibilidad es del 0.92 y la especificidad es del 0.75.

Técnicas de balanceo de clases

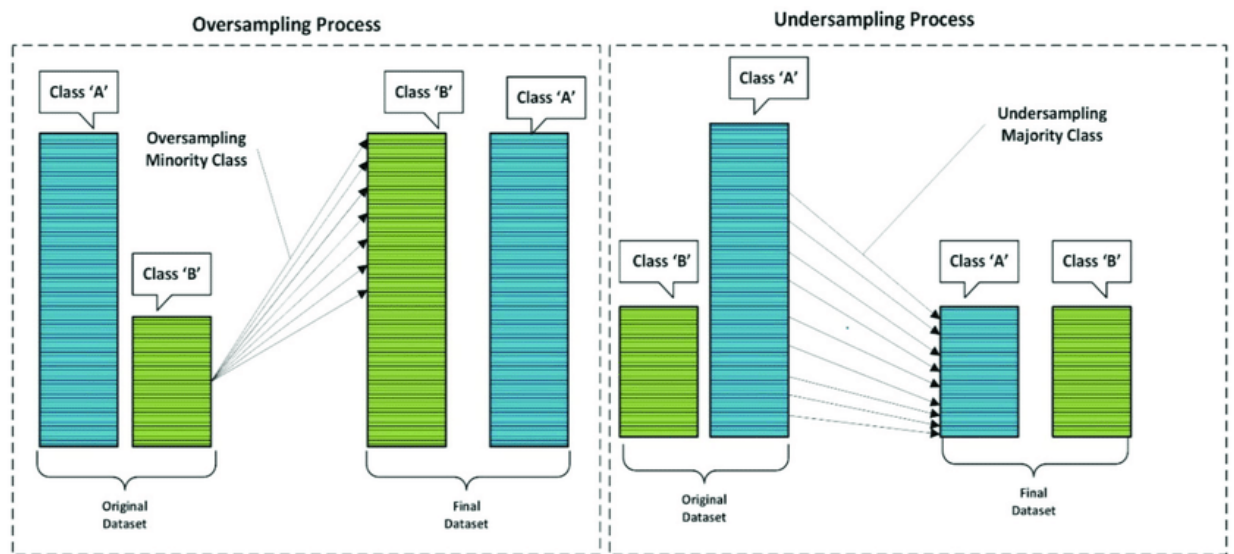
5.1 Undersampling

Es una técnica que consiste en mantener todos los datos de la clase de menor frecuencia y reducir la cantidad de los de la clase de mayor frecuencia, haciendo que las observaciones del conjunto tengan datos con la variable objetivo balanceada.

5.2 Oversampling

Es una técnica que consiste en aumentar el número de registros de la clase con menor frecuencia hasta que la base de datos tenga un número equilibrado entre las clases de la variable objetivo. Para aumentar la cantidad de registros, podemos duplicar aleatoriamente los registros de la clase con menos frecuencia. Sin embargo, esto hará que mucha información sea idéntica, lo que puede afectar el modelo.

Una ventaja de esta técnica es que no se pierde ninguna información de los registros que tenían la clase con mayor frecuencia. Esto hace que el conjunto de datos tenga muchos registros para alimentar los algoritmos de aprendizaje automático. A su vez, el tiempo de almacenamiento y procesamiento crece significativamente y existe la posibilidad de sobreajustar los datos que se han duplicado. Este sobreajuste ocurre cuando el modelo se vuelve muy bueno para predecir los resultados de los datos de entrenamiento, pero no generaliza bien los datos nuevos.



Conclusión

Para empezar realizamos varios análisis, un summary inicial para observar la distribución de las variables

Contando con un amplio número de observaciones y de variables se pudo obtener una predicción de clasificación aceptable.

A la vez afianzamos nuestro dominio sobre la herramienta R.

Anexo

```
1
2
3
4 #####
5 # Arboles de desicion del dataset Aprender 2016.
6 #
7 # Creado: 19/10/2023
8 # Version: 1.0
9 # Autor: Grupo 4
10
11 #####
12
13 ##### Importamos las biliotecas #####
14 library(plyr)
15 library(rpart)
16 library(rpart.plot)
17 library(caret)
18 library(ggplot2)
19 library(gridExtra)
20 library(tidyverse)
21 library(rsample)
22 library(e1071)
23 library(GGally)
24 library(data.table)
25 library(DT)
26 library(readr)
27
28 library(dplyr)
29 library(tidyr)
30 library(corrplot)
31 #library(rms)
32 library(MASS)
33 library(e1071)
34 library(ROCR)
35 library(gplots)
36 library(pROC)
37
38 library(randomForest)
39 library(ggpubr)
40 library(plotly)
41
42 ##### Lectura y verificacion de la informacion leida #####
```

```

43 # Seteamos el entorno de desarrollo.
44 setwd(paste0(getwd(), "/" ))
45
46 aprender <- read_delim("./UNO_EDD/TP4_Arboles_Decision/aprender2016-primaria-3.csv",
47                       delim = ";", escape_double = FALSE, locale = locale(decimal_mark =
48                               ", ", grouping_mark = "."), trim_ws = TRUE)
49 View(aprender)
50 dim(aprender)
51 nrow(aprender)
52
53 head(aprender)
54
55
56 # Verificamos si hay valores nulos.
57 sapply(aprender, function(x) sum(is.na(x)))
58
59 # complete.cases() devuelve un vector con los valores completos y que sean no nulos.
60 aprender_clean <- aprender[complete.cases(aprender), ]
61
62 # Verificamos si hay valores nulos.
63 sapply(aprender_clean, function(x) sum(is.na(x)))
64
65 summary(aprender_clean)
66 aprender_clean$mpondera
67 ##### Normalizacion y analisis #####
68 # Grafico dinamico de boxplot para identificacion de outliers
69 grafico.boxplot <- plot_ly(y=aprender_clean$mpondera, type="box", name="ponder") %>%
70   add_trace(y=aprender_clean$lpondera, type="box", name="lpondera") %>%
71   add_trace(y=aprender_clean$mpondera, type="box", name="mpondera")
72
73 grafico.boxplot
74
75 names(aprender_clean)
76
77 aprender_clean = subset(aprender_clean, select = c(Ap1, Ap2, Ap3, Ap4, Ap5, Ap6, Ap7,
78           Ap8a, Ap8b, Ap8c, Ap9, Ap10, Ap11,
79           Ap12, Ap13, Ap14, Ap15, lpuntaje,
80           ldesemp, mpuntaje, mdesemp, sector,
81           ambito, iclima, autoconl, autoconm
82 ))
83
84
85 # Normalizamos campos.
86 aprender_clean$Ap1 <- as.factor(mapvalues(aprender_clean$Ap1,
87                                           from=c(1,2,3,4,5, -1,-9),
88                                           to=c("7 a os o menos", "8 a os", "9
89           a os", "10 a os", "11 a os o mas", "En blanco", "No respondio")))
90
91 aprender_clean$Ap2 <- as.factor(mapvalues(aprender_clean$Ap2,
92                                           from=c(1,2, -1,-9),
93                                           to=c("varon", "mujer", "En blanco", "No respondio
94           ")))

```

[illegible]

```

139         to=c("siempre o casi siempre", "a veces", "nunca
    o casi nunca", "En blanco", "No respondio")))
140
141
142 aprender_clean$Ap12 <- as.factor(mapvalues(aprender_clean$Ap12,
143         from=c(1,2,3, -1,-9),
144         to=c("siempre o casi siempre", "a veces", "nunca
    o casi nunca", "En blanco", "No respondio")))
145
146
147 aprender_clean$Ap13 <- as.factor(mapvalues(aprender_clean$Ap13,
148         from=c(1,2,3, -1,-9),
149         to=c("siempre o casi siempre", "a veces", "nunca
    o casi nunca", "En blanco", "No respondio")))
150
151 aprender_clean$Ap14 <- as.factor(mapvalues(aprender_clean$Ap14,
152         from=c(1,2,3,4, -1,-9),
153         to=c("muy facil", "un poco facil", "un poco
    dificil", "dificil", "En blanco", "No respondio")))
154
155
156 aprender_clean$Ap15 <- as.factor(mapvalues(aprender_clean$Ap15,
157         from=c(1,2,3,4, -1,-9),
158         to=c("muy facil", "un poco facil", "un poco
    dificil", "dificil", "En blanco", "No respondio")))
159
160 aprender_clean$lde Kemp <- as.factor(mapvalues(aprender_clean$lde Kemp,
161         from=c(1,2,3,4),
162         to=c("por debajo del nivel basico", "basico", "
    satisfactorio", "avanzado")))
163
164 aprender_clean$mdes Kemp <- as.factor(mapvalues(aprender_clean$mdes Kemp,
165         from=c(1,2,3,4),
166         to=c("por debajo del nivel basico", "basico"
    , "satisfactorio", "avanzado")))
167
168 aprender_clean$sector <- as.factor(mapvalues(aprender_clean$sector,
169         from=c(1,2,3),
170         to=c("estatal", "privada", "sin datos")))
171
172 aprender_clean$ambito <- as.factor(mapvalues(aprender_clean$ambito,
173         from=c(1,2,3),
174         to=c("urbano", "rural", "sin datos")))
175
176 aprender_clean$iclima <- as.factor(mapvalues(aprender_clean$iclima,
177         from=c(1,2,3,-1),
178         to=c("bajo", "medio", "alto", "en blanco"))
    )
179
180 aprender_clean$autoconl <- as.factor(mapvalues(aprender_clean$autoconl,
181         from=c(1,2,3,-1),
182         to=c("bajo", "medio", "alto", "en blanco")))
183
184 aprender_clean$autoconm <- as.factor(mapvalues(aprender_clean$autoconm,
185         from=c(1,2,3,-1),

```



```

186                                     to=c("bajo", "medio", "alto", "en blanco"))
187 aprender_clean %>% distinct(Ambito)
188
189 names(aprender_clean)
190
191 aprender %>% distinct(ambito)
192 names(aprender_clean) <- c("años",
193                             "genero",
194                             "asistio_jardin",
195                             "repetio",
196                             "como_lees",
197                             "como_escribis",
198                             "como_resolv_mat",
199                             "interrum_maestros",
200                             "maestros_se_enojan",
201                             "maestros_explican",
202                             "relacion_companeros",
203                             "sent_camb_escuela",
204                             "contento_escuela",
205                             "aburrido",
206                             "incomodo",
207                             "examen_lengua",
208                             "examen_mat",
209                             "puntaje_lengua",
210                             "desempeno_lengua",
211                             "puntaje_mat",
212                             "desempeno_mat",
213                             "sector",
214                             "ambito",
215                             "clima_escolar",
216                             "autoconcepto_lengua",
217                             "autoconcepto_matematica"
218 )
219
220 attach(aprender_clean)
221 aprender_clean %>% distinct (sector)
222 #aprender_clean_cordoba <- subset(aprender_clean, aprender_clean$Provincia == "Cordoba")
223 #head(aprender_clean_cordoba)
224
225 ##### Modelado - Arbol de desicion #####
226 set.seed(1250) # Numero inicial del cual comenzara a generar una secuencia aleatoria.
227 split_train_test <- createDataPartition(aprender_clean$sector, p=0.7, list=FALSE) #p =
    porcentaje de los datos a usar, list = Si el resultado devolviera una lista o una
    matriz.
228 dtrain<- aprender_clean[split_train_test,]
229 dtest<- aprender_clean[-split_train_test,]
230
231 sqldf('
232 SELECT asistio_jardin, count(*)
233 FROM dtrain
234 WHERE 'puntaje lengua' < 463
235 GROUP BY asistio_jardin
236 ORDER BY 2 DESC
237 ')
238

```

```

239 dtrain$'puntaje lengua'
240
241 sqldf('
242 SELECT DISTINCT repitio
243 FROM dtrain
244
245 ')
246 #11.853 < 463 --> 33%
247 #22.299 > 463 --> 66%
248 #34153
249
250 tr_fit <- rpart(sector ~., data = dtrain, method="class") # Indicamos que deseamos un
    arbol de clasificacion, tambien podemos armar un arbol de regresion.
251 tr_fit # Nuestro arbol obtenido.
252 rpart.plot(tr_fit, tweak = 1.6) # Graficamos el arbol.
253
254 prp(tr_fit,
255     type = 2, #Especifica que el cada nodo quede etiquetado, y que el split quede debajo de
        cada nodo
256     extra = 104, #Muestra la probabilidad de cada clase en el nodo
257     nn = TRUE, #Etiqueta el nro de nodo
258     fallen.leaves = TRUE, #Muestra los nodos hojas abajo de todo el grafico
259     faclen = 10, #Se utiliza para abreviar el nombre de las clases en 4caracteres
260     varlen = 10, #Se utiliza para abreviar el nombre de las variables en 8caracteres
261     shadow.col = "gray")
262
263 #Importancia de las variables.
264 qplot(x = names(tr_fit$variable.importance), y=tr_fit$variable.importance,
265       xlab="Variable", ylab="Importancia", main="rpart - Importancia de las variables")
266
267 #
268 #####
269 ##### 05 - TEST Y VALIDACION DEL MODELO
270 #####
271 #
272 #####
273
274 #EVALUACION DE CALIDAD DEL MODELO: MATRIZ DE CONFUSION:
275
276 #5.1 - CREAMOS EL OBJETO DE PREDICCION SOBRE EL CUAL SE HAR LA MATRIZ DE CONFUSION
277 aprender_rpart_pred <- predict(tr_fit, dtest, type="class")
278 aprender_rpart_pred
279 #5.2 - MATRIZ DE CONFUSION
280 tabla_mc_aprender <- confusionMatrix(table(dtest$sector, aprender_rpart_pred,
281     dnn = c("Actual", "Predicho")))
282
283 tabla_mc_aprender
284
285 #5.3 - DIAGRAMA ROC: MIDE EL RENDIMIENTO: NOS DICE QUE TAN BIEN PUEDE CLASIFICAR EL
    MODELO ENTRES 2 CLASES
286 #5.3 - SE GENERA EL INDICE DE PROBABILIDAD EN LUGAR DE LAS CANTIDADES DE CLASIFICACIONES
287 #aprender_rpart_pred2 <- predict(tr_fit, dtest, type = "prob")
288 #aprender_rpart_pred2
289 #5.3.1 - CON QUE PROB CLASIFICA?

```

```
286 #head(aprender_rpart_pred2)
287
288 #aprender_rpart_pred2_roc <- prediction(aprender_rpart_pred2[,2], dtest[, "sector"])
289 #aprender_rpart_pred2_roc_perf <- performance(aprender_rpart_pred2_roc, "tpr", "fpr")
290 #plot(aprender_rpart_pred2_roc_perf)
291
292 aprender_rpart_pred2_roc
293 aprender_rpart_pred2_roc_perf
294
295 #
296 #####
297 ##### PAQUETES NECESARIOS
298 #####
299
300 install.packages("ROCR")
301
302 library(TTR)
303 library(tsibble)
304 library(readr)
305 library(sqldf)
306 library(tidyverse)
307 library(stringi)
308 library(stringr)
309 library(ROCR)
310 library(plyr)
311 library(caret)
312 library(gridExtra)
313 library(dplyr)
314 library(tidyr)
315 library(corrplot)
316 library(ggplot2)
317
318 library(rpart)
319 library(pROC)
320 library(MASS)
321 library(e1071)
322 library(ggpubr)
323 library(rsample)
324 library(e1071)
325 library(GGally)
326 library(data.table)
327 library(DT)
328 library(ROCR)
329 library(gplots)
330 library(randomForest)
331 library(rpart.plot)
332 #
333 #####
334 ##### 01 - IMPORTAR LOS DATOS
335 #####
```

```
332 # #####
333 emse <- read.csv("./UNO_EDD/TP4_Arboles_Decision/emse_datosabiertos/EMSE_DatosAbiertos.
      csv")
334 str(emse)
335 View(emse)
336 names(emse)
337
338 emse_red <- emse[, (4:144)]
339 names(emse_red)
340 rm(emse_red)
341 emse_red = subset(emse_red, select = c(texto_q1,
342                                       texto_q2,
343                                       texto_q3,
344                                       texto_q6,
345                                       texto_q15,
346                                       texto_q16,
347                                       texto_q17,
348                                       texto_q18,
349                                       texto_q19,
350                                       texto_q22,
351                                       texto_q23,
352                                       texto_q24,
353                                       texto_q25,
354                                       texto_q26,
355                                       texto_q27,
356                                       texto_q28,
357                                       texto_q29,
358                                       texto_q31,
359                                       texto_q34,
360                                       texto_q35,
361                                       texto_q38,
362                                       texto_q39,
363                                       texto_q40,
364                                       texto_q41,
365                                       texto_q42,
366                                       texto_q43,
367                                       texto_q44,
368                                       texto_q45,
369                                       texto_q46,
370                                       texto_q49,
371                                       texto_q51,
372                                       texto_q53,
373                                       texto_q55,
374                                       texto_q56,
375                                       texto_q57,
376                                       texto_q59,
377                                       texto_q60,
378                                       texto_q61,
379                                       texto_q62,
380                                       texto_q65,
381                                       texto_q66,
382                                       texto_q68,
383                                       texto_q73,
```

```

384         texto_q74,
385         texto_q76,
386         texto_q79
387     ))
388 #
389 #####
390 ##### 02 - ETL DEL MODELO
391 #####
392 #
393 #####
394
395 #borrar:
396 #2.1.1 - RENOMBRAR COLUMNAS
397 names(emse_red) <- c("edad",
398                     "sexo",
399                     "grado",
400                     "dias_con_hambre",
401                     "ult_12meses_atac_fisi",
402                     "ult_12meses_part_pelea",
403                     "lesion_grave",
404                     "lesion_mas_seria",
405                     "causa_lesion_seria",
406                     "sentirse_solo",
407                     "preocupacion",
408                     "consideraste_suicidio", #pred
409                     "plan_de_suicidio", #pred
410                     "veces_intentaste_suici",
411                     "amigos_cercanos",
412                     "edad_fumador",
413                     "ult_30dias_fumaste",
414                     "dejaste_fumar",
415                     "edad_alcohol",
416                     "dias_alcohol",
417                     "veces_tomaste_acohol",
418                     "problem_entorno_alcohol",
419                     "edad_drogas",
420                     "veces_maria",
421                     "ult_30dias_consu_maria",
422                     "veces_consu_metanfet",
423                     "tuviste_sexo",
424                     "edad_sexo_1ra_vez",
425                     "cuantas_perso_tuv_sex",
426                     "ult_7dias_activ_fisica",
427                     "vas_a_edu_fisica",
428                     "ult_30dias_faltaste",
429                     "padres_verfic_tarea",
430                     "ult_30dias_padres_atend_salud",
431                     "ult_30dias_padres_sabian_activi",
432                     "educacion_padre",
433                     "educacion_madre",
434                     "ult_7dias_frutas",
435                     "ult_7dias_verd",
436                     "comida_grasa",
437                     "intimidacion_en_esc",

```

```

434         "intimidacion-en-int",
435         "que_bebida_tomas",
436         "con_quien_tomas",
437         "si_te_ofrecen_tomas",
438         "quedaste_embarazada"
439     )
440
441 #2.1.2 CORRECCION DE CARACTERES ESPECIALES
442 #CAMBIO DE VALORES POR CONFLICTO EN APLICACION DE FUNCIONES CON LOS MISMOS
443 names(emse_red)
444 str(emse_red)
445 library(stringi)
446 attach(emse_red)
447 emse_red$veces_intentaste_suici
448 emse_red %>% distinct(veces_intentaste_suici)
449
450 emse_red$veces_intentaste_suici <- stri_replace_all_regex(emse_red$veces_intentaste_suici
451     ,
452     pattern=c(' ',' ',' ',' ','>','<'),
453     replacement=c('E','e','a','mayor a ',
454     'ni',' ','o','u','er','i','menor a '),
455     vectorize=FALSE)
456
457 emse_red$causa_lesion_seria <- stri_replace_all_regex(emse_red$causa_lesion_seria,
458     pattern=c(' ',' ',' ',' ','>','<'),
459     replacement=c('E','e','a','mayor a ',
460     'ni',' ','o','u','er','i','menor a '),
461     vectorize=FALSE)
462
463 emse_red$causa_lesion_seria <- stri_replace_all_regex(emse_red$causa_lesion_seria,
464     pattern=c(' ',' ',' ',' ','>','<'),
465     replacement=c('E','e','a','mayor a ',
466     'ni',' ','o','u','er','i','menor a '),
467     vectorize=FALSE)
468
469 emse_red$lesion_mas_seria <- stri_replace_all_regex(emse_red$lesion_mas_seria,
470     pattern=c(' ',' ',' ',' ','>','<'),
471     replacement=c('E','e','a','mayor a ',
472     'ni',' ','o','u','er','i','menor a '),
473     vectorize=FALSE)
474
475 emse_red$lesion_grave <- stri_replace_all_regex(emse_red$lesion_grave,
476     pattern=c(' ',' ',' ',' ','>','<'),
477     replacement=c('E','e','a','mayor a ',
478     'ni',' ','o','u','er','i','menor a '),
479     vectorize=FALSE)
480
481 ##
482 emse_red$edad <- stri_replace_all_regex(emse_red$edad,
483     pattern=c(' ',' ',' ',' ','>','<'),
484     replacement=c('E','e','a','mayor a ',
485     'ni',' ','o','u','er','i','menor a '),
486     vectorize=FALSE)

```

```

477 replacement=c('E','e','a','mayor a ','ni','',' 'o
    ', 'u', 'er', 'i', 'menor a '),
478 vectorize=FALSE)
479
480 emse_red$grado <- stri_replace_all_regex(emse_red$grado,
481 pattern=c(' ',' ',' ','='>', ' ',' ','
    ', ' ',' ',' ','='<'),
482 replacement=c('E','e','a','mayor a ','ni','',' '
    'o', 'u', 'er', 'i', 'menor a '),
483 vectorize=FALSE)
484
485 emse_red$sult_12meses_part_pelea <- stri_replace_all_regex(emse_red$sult_12meses_part_pelea
486 ',
    pattern=c(' ',' ',' ','='>', ' ',' ','
    ', ' ',' ',' ','='<'),
487 replacement=c('E','e','a','mayor
    a ','ni','',' 'o', 'u', 'er', 'i', 'menor a '),
488 vectorize=FALSE)
489
490 emse_red$sult_12meses_atac_fisi <- stri_replace_all_regex(emse_red$sult_12meses_atac_fisi,
491 pattern=c(' ',' ',' ','='>', ' ',' ','
    ', ' ',' ',' ','='<'),
492 replacement=c('E','e','a','mayor a ','ni','','
    'o', 'u', 'er', 'i', 'menor a '),
493 vectorize=FALSE)
494
495 #2.1.3 - TRATA DE VALORES NAs
496 View(emse_red)
497 str(emse_red)
498
499 emse_red[emse_red == ""] <- NA
500 emse_red[emse_red == "Dato perdido"] <- NA
501
502 #CONTAR NAs
503 View(summarise_all(emse_red, funs(sum(is.na(.)))))
504
505 #ELIMINAR TODOS LOS NAs
506 emse_red <- na.omit(emse_red)
507
508 View(summarise_all(emse_red, funs(sum(is.na(.)))))
509
510 emse_red <- data.frame(lapply(emse_red, as.factor))
511 str(emse_red)
512 #
    #####
513 ##### 03 - DIVIDIMOS EN TRAINEE Y TEST
    #####
514 #
    #####
515 #ACLARACION: EN RF, NO NECESARIAMENTE DIVIDIR LOS DATOS PARA MODELO, YA QUE EL PROPIO
    MODELO
516 #YA TIENE INCORPORADO EN EL ALGORITMO LA PARTICION.
517 set.seed(2018)

```

```

518 #plan_de_suicidio
519 #consideraste_suicidio
520
521 split_train_test <- createDataPartition(emse_red$plan_de_suicidio, p=0.7, list=FALSE) #p
    = porcentaje de los datos a usar, list = Si el resultado devolviera una lista o una
    matriz.
522 dtrain <- emse_red[split_train_test,]
523 dtest <- emse_red[-split_train_test,]
524
525 sqldf('
526 SELECT plan_de_suicidio, count(*)
527 FROM dtrain
528 GROUP BY plan_de_suicidio
529 ORDER BY 2 DESC
530 ')
531 #      plan_de_suicidio cantidad
532 #1                No      18787
533 #2                Si       3519
534 #
    #####
535 ##### 04 - CONSTRUCCION DEL MODELO
    #####
536 #
    #####
537 names(emse_red)
538 emse_red_rf <- randomForest(x = dtrain[, -13], #VARIABLES INDEPENDIENTES
539                             y = dtrain$plan_de_suicidio, #VARIABLE DEPENDIENTE
540                             ntree = 500, #CANTIDAD DE ARBOLES. POR DEFECTO SON 500
541                             keep.forest = TRUE) #TRUE = ME DEJA LOS ARBOLES INTERMEDIOS.
    OJO QUE ES COSTOSO, Y EN LA REALIDAD NO ES RECOMENDABLE
542
543 #KEEP.FOREST: POR DEFECTO, EL MODELO NO SE QUEDA CON LOS ARBOLES QUE FUE UTILIZANDO PARA
    CONSTRUIR EL MODELO.
544 #POR ENDE DICHOS ARBOLES NO SE USAN PARA PREDECIR SOBRE LOS CUALES SE VA CLASIFICAR. SIN
    EMBARGO, SE PUEDE FORZAR.
545
546 #
    #####
547 ##### MATRIZ DE CONFUSION #####
548 #
    #####
549 #NOTA: PREDICT CATALOGA/CLASIFICA; PREDICTION HACE UNA PREDICION PORCENTUAL
550 #4.1 - SE PREDICE LA BONDAD DEL MODELO: SE PREDICE EN BASE AL MODELO ACTUAL, LA
    CLASIFICACION DE LOS DATOS QUE NO INGRESARON AL MODELO
551 emse_red_rf_pred <- predict(emse_red_rf, dtest, type = "class")
552
553 #4.2 - MATRIZ DE CONFUSION:
554 emse_red_rf_pred_mc <- confusionMatrix(table(dtest$plan_de_suicidio, emse_red_rf_pred,
555                                             dnn = c("Actual", "Predicho")))
556
557 emse_red_rf_pred_mc

```



```
558
559 #CONCLUSION 1: Balanced Accuracy : 0.8404 NOS INDICA QUE EL MODELO PRESENTA UN BALANCEO
    DE CLASES
560 #EN DONDE LA CLASE A PRESENTA MAYOR CANTIDAD DE CASOS QUE LA OTRA CLASE.
561
562 #4.3 - PREDICR DE PROBABILIDADES: SE PREDICE CON QUE PROBABILIDAD EL MODELO PUEDE
    CLASIFICAR SOBRE LOS
563 #DATOS DE TEST
564 #banknote_probs <- predict(banknote_rf, banknote[-training.ids,], type = "prob")
565 #head(banknote_probs)
566
567 #4.4 - ROC: SOBRE LOS EXITOS,
568 #banknote_probs_pred <- prediction(banknote_probs[,2], banknote[-training.ids,"class"])
569 #banknote_probs_pred_perf <- performance(banknote_probs_pred, "tpr", "fpr") #TPR: TRUE
    POSITIVE RATE; FALSE POSITIVE RATE
570 #plot(banknote_probs_pred_perf)
571
572 #CONCLUSION: RF ES COSTOSO COMPUTANCIONALMENTE, PERO CLASIFICA MUY BIEN.
```