



## Laboratorio 5: Árboles de Decisión

Integrantes: Nicole Reyes  
Diego Águila  
Curso: Análisis de Datos  
Profesor: Max Chacón  
Ayudante: Francisco Muñoz

# Tabla de contenidos

<b>1. Introducción</b>	<b>1</b>
1.1. Objetivos . . . . .	1
<b>2. Obtención del árbol</b>	<b>2</b>
2.1. Pre Procesamiento . . . . .	2
2.1.1. Eliminación de valores vacíos . . . . .	2
2.1.2. Eliminación de variables . . . . .	2
2.2. Categorización de los datos . . . . .	3
2.3. Árbol de decisión . . . . .	3
<b>3. Análisis de los resultados</b>	<b>4</b>
3.1. Comparación con Reglas de Asociación . . . . .	5
<b>4. Anexos</b>	<b>7</b>
<b>Bibliografía</b>	<b>13</b>

## Índice de figuras

1.	Árbol obtenido. . . . .	3
----	-------------------------	---

## Índice de cuadros

1.	Reglas obtenidas del árbol de decisión. . . . .	4
2.	Entropía de las variables. . . . .	4
3.	Reglas obtenidas en la experiencia 3. . . . .	5

# 1. Introducción

Obtener conocimiento de un conjunto de datos es fundamental al momento de realizar un estudio o análisis con respecto a estos, existen diversas técnicas y estrategias en la minería de datos para lograrlo, la obtención de un árbol de decisión es una de las técnicas conocidas para lograr esto y permite establecer un modelo predictivo el cual, muchas veces se utiliza para representar visualmente decisiones y toma de decisiones o en el caso de minería de datos, se crea un modelo que predice el valor de una variable de destino en función de diversas variables de entrada.

Esta quinta experiencia de laboratorio consiste en, mediante el proceso de minería de datos y a través de la utilización de un árbol de decisión obtener conocimiento de la base de datos asignada en la primera experiencia ”*Allhypo*” correspondiente al hipotiroidismo.

## 1.1. Objetivos

1. Profundizar el entendimiento de los árboles de decisión como medida de obtención de conocimiento.
2. Corroborar la teoría vista en clases mediante la aplicación del árbol de decisión en el lenguaje estadístico R, haciendo uso del paquete ”*C50*”.
3. Obtener conocimiento de la base de datos ”*Allhypo*” a través de esta técnica de minería de datos.

A continuación se define cómo se trabaja la base de datos para obtener el árbol de decisión y finalmente se realiza una comparación entre los resultados obtenidos del árbol de decisión con respecto a los resultados de las experiencias anteriores, concretamente, la experiencia número 3 correspondiente a reglas de asociación.

## 2. Obtención del árbol

### 2.1. Pre Procesamiento

Al igual que en las experiencias anteriores es fundamental elaborar una limpieza a la base de datos. En esta oportunidad la limpieza está compuesta por las siguientes etapas:

#### 2.1.1. Eliminación de valores vacíos

Dentro de la base de datos, existen celdas con valores nulos (representados con el símbolo "?"), por lo que se opta por eliminar todos los sujetos de la muestra que contengan este símbolo para alguna variable ya que no aportan nada al análisis.

#### 2.1.2. Eliminación de variables

Dentro del conjunto de datos usados, hay columnas que no formarán parte del estudio, debido a que no aportan información importante para este pudiendo estropear el análisis. Las variables que son eliminadas son:

- Variable TBG: Como toda su columna representa valores vacíos, no entrega información relevante debido a que se puede inferir a que no se midió la presencia de esta proteína en ningún sujeto.
- Variables de verificación de medición (measured): Indican si el paciente se ha realizado alguna medición, siendo esta información no relevante. Se decide dejar fuera las variables: *TSH measured*, *T3 measuted*, *TT4 measured*, *T4U measured*, *FTI measured* y *TBG measured*.
- Variable de fuente: La variable *referral.source* no entrega información que sea pertinente para el análisis ya que solo indica la fuente de la cual se obtuvo los datos de sujeto en cuestión.

## 2.2. Categorización de los datos

Se transforman ciertos datos a variables nominales para que se tenga de esta manera una categorización de los datos que representan esas variables. Para las hormonas, se realiza lo mismo que en la experiencia pasada de *Reglas de Asociación*. Con respecto a la media de cada una, se va clasificando a la variable, es decir, si el valor de la variable es menor a la media se deja categorizada como 'f', caso contrario, se deja como 't'. Lo mismo se hace con la edad, si el valor es menor a la media se deja categorizada como 'f' y si es mayor, se deja como 't'.

## 2.3. Árbol de decisión

Se usa la biblioteca *C50* para la obtención del árbol cual retorna un objeto de clase C5.0 que puede ser el árbol de decisión, reglas, entre otros. Para la generación del árbol, es necesario definir su raíz, el cual debe ser el atributo adecuado que dé garantía de que los resultados que se entreguen sean lo mas cercano a la realidad. Por esta razón el atributo escogido es el *state*, el cual corresponde a la clase de Hipotiroidismo que presenta el sujeto.

Para generar el árbol, lo que se necesita es la base de datos en donde se encuentran la información necesaria y el atributo que se utilizará como raíz de este árbol. Entonces lo que se genera queda mostrado en la Figura 1:

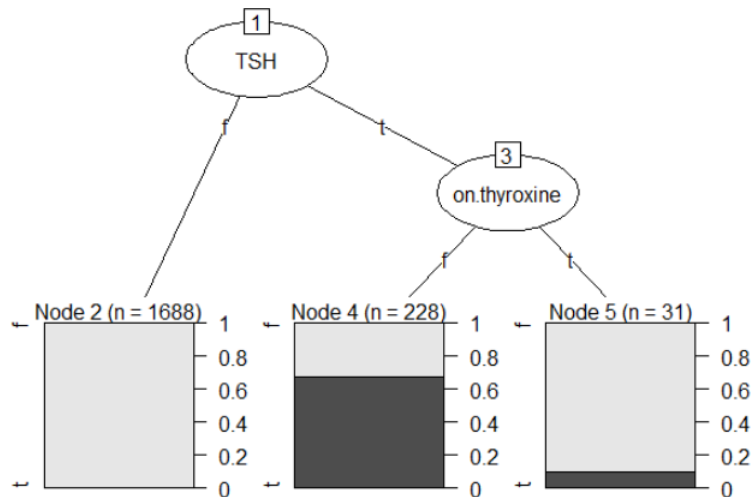


Figura 1: Árbol obtenido.

### 3. Análisis de los resultados

A continuación se realiza el análisis con respecto al árbol obtenido durante el desarrollo de esta experiencia. De acuerdo a lo obtenido en el modelo basado en los árboles de decisiones en el Capítulo 2, este árbol se puede expresar de acuerdo al modelo basado en *Reglas de Asociación* las cuales se encuentran descritas en el Cuadro 1 a continuación:

N° de Regla	Antecedentes	Consecuente
1	TSH=F	hypothyroid=F
2	TSH=T, on.thyroxine=F	hypothyroid=T
3	TSH=T, on.thyroxine=T	hypothyroid=F

Cuadro 1: Reglas obtenidas del árbol de decisión.

Del árbol de decisión correspondiente a la Figura 1 del Capítulo 2 y de las reglas presentadas en el Cuadro 1, se obtiene que la variable más importante resulta ser *TSH*, es decir, si el paciente presenta altos o bajos niveles de la hormona TSH (Tirotopina) considerándola como suficiente para determinar si un paciente no posee Hipotiroidismo, para los casos en los que la hormona TSH resulta ser alta, se incorpora la variable *on.thyroxine* considerada relevante en el árbol de decisión, la cual indica si un paciente se encuentra bajo tratamiento de reemplazo de Tiroxina, donde en los casos negativos, en su mayoría, el paciente presenta Hipotiroidismo y en los casos positivos, pocos de los pacientes presentan hipotiroidismo. Lo anterior se puede asociar a que dado el modelo de árboles de decisión de Quinlan se puede obtener la entropía de las variables que son importantes para la generación del árbol, para así saber que tan importantes son para el modelo. A continuación se muestra en el Cuadro 2, la entropía de las dos variables que se usaron para la generación del árbol:

Variable	Entropía
TSH	100 %
on.thyroxine	13,3 %

Cuadro 2: Entropía de las variables.

### 3.1. Comparación con Reglas de Asociación

En la tercera experiencia de laboratorio correspondiente a las reglas de asociación, se obtienen todas las reglas existentes en la base de datos de las cuales se obtiene un ranking con las 10 mejores a partir de medidas de calidad como el soporte, la confianza y el lift, estas reglas se presentan a continuación en el Cuadro 3:

N° de Regla	Antecedentes	Consecuente
1	on.thyroxine=F, TSH=T, TT4=F	hypothyroid=t
2	on.thyroxine=F, thyroid.surgery=F, TSH=T	hypothyroid=t
3	on.thyroxine=F, TSH=T, T3=F	hypothyroid=t
4	on.thyroxine=F, I131=F, TSH=T	hypothyroid=t
5	on.thyroxine=F, TSH=T, T4U=F	hypothyroid=t
6	on.thyroxine=F, TSH=T, FTI=F	hypothyroid=t
7	on.thyroxine=F, age=F, TSH=T	hypothyroid=t
8	on.thyroxine=F, psych=F, TSH=T	hypothyroid=t
9	on.thyroxine=F, on.antithyroid.med=F, TSH=T	hypothyroid=t
10	on.thyroxine=F, lithium=F, TSH=T	hypothyroid=t

Cuadro 3: Reglas obtenidas en la experiencia 3.

Como se observa en el Cuadro anterior, todas las reglas tienen como antecedente que el paciente posee altos niveles de la hormona TSH y no se encuentran bajo tratamiento de reemplazo de Tiroxina, además de otros factores, de los cuales tienen como consecuente, la presencia de Hipotiroidismo en el paciente. El hecho de que las 10 reglas posean como antecedente altos niveles de hormona TSH y que el paciente no se encuentre en tratamiento de reemplazo de Tiroxina, permite concluir que estas dos variables son gatillantes para determinar si el usuario posee hipotiroidismo, no así, las otras variables como TT4, T3, T4U, entre otras.

Estos resultados concuerdan con lo obtenido en el árbol de decisión, al observar el Cuadro 1 se tiene que si un paciente presenta altos niveles de hormona TSH y no se encuentra en tratamiento de reemplazo de Tiroxina tenderá a poseer Hipotiroidismo, lo cual encaja con



las 10 reglas obtenidas en la experiencia 3. Con respecto a los otros dos casos, si el paciente presenta altos niveles de hormona TSH pero se encuentra en tratamiento de reemplazo de tiroxina tenderá a no poseer Hipotiroidismo y de forma más tajante, si el paciente posee bajos niveles de hormona TSH no presentaría Hipotiroidismo.

Al obtener resultados similares a los de la experiencia 3, un árbol de decisión resulta ser una medida más eficiente para predecir la ocurrencia de un evento debido a que se construye en base a las variables que aportan una mayor cantidad de información al modelo, contrario a las reglas de asociación, donde se obtiene la totalidad de reglas a partir de una cantidad de antecedentes y de las cuales se debe realizar un filtrado posterior a modo de obtener las más relevantes en base a medidas de calidad.

## 4. Anexos

A continuación se presenta el script en R utilizado para el desarrollo de esta experiencia.

```
1 library(C50)
2
3 data <- read.csv("C:\\Users\\diego\\Dropbox\\Usach\\ANLISIS DE DATOS\\Lab\\
  allhypo.data",
4               col.names = c("age", "sex", "on.thyroxine", "query.thyroxine",
  "on.antithyroid.med",
5               "sick", "pregnant", "thyroid.surgery", "I131", "
  query.hypothyroid",
6               "query.hyperthyroid", "lithium", "goitre", "
  tumor", "hypopituitary",
7               "psych", "TSH.measured", "TSH", "T3.measured",
  "T3", "TT4.measured",
8               "TT4", "T4U.measured", "T4U", "FTI.measured", "
  FTI", "TBG.measured",
9               "TBG", "referral.source", "state"), header=FALSE
  , sep="," , stringsAsFactors=FALSE)
10
11 # LIMPIEZA DE LA BASE DE DATOS
12
13 # Se limpia la variable de estado del paciente eliminando numeros
14 state <- data$state
15 state <- strsplit(state, ".", fixed = TRUE)
16 state <- lapply(state, '[', 1)
17 state <- unlist(state)
18 data$state <- state
19
20 # Se eliminan los pacientes que contengan algun NA
21 # Sex
22 error = data$sex == '?'
23 data = data[!error,]
24 # Age
25 error = data$age == '?'
```

```

26 data = data[!error,]
27 # On thyroxine
28 error = data$on.thyroxine == '?'
29 data = data[!error,]
30 # Query on antithyroid thyroxine
31 error = data$query.thyroxine == '?'
32 data = data[!error,]
33 # On antithyroid medication
34 error = data$on.antithyroid.med == '?'
35 data = data[!error,]
36 # Sick
37 error = data$sick == '?'
38 data = data[!error,]
39 # Pregnant
40 error = data$pregnant == '?'
41 data = data[!error,]
42 # Thyroid surgery
43 error = data$thyroid.surgery == '?'
44 data = data[!error,]
45 # I131 treatment
46 error = data$I131 == '?'
47 data = data[!error,]
48 # Query hypothyroid
49 error = data$query.hypothyroid == '?'
50 data = data[!error,]
51 # Query hyperthyroid
52 error = data$query.hyperthyroid == '?'
53 data = data[!error,]
54 # Lithium
55 error = data$lithium == '?'
56 data = data[!error,]
57 # Goitre
58 error = data$goitre == '?'
59 data = data[!error,]
60 # Tumor
61 error = data$tumor == '?'

```

```

62 data = data[!error,]
63 # Hypopituitary
64 error = data$hypopituitary == '?'
65 data = data[!error,]
66 # Psych
67 error = data$psych == '?'
68 data = data[!error,]
69 # TSH
70 error = data$TSH == '?'
71 data = data[!error,]
72 # T3
73 error = data$T3 == '?'
74 data = data[!error,]
75 # TT4
76 error = data$TT4 == '?'
77 data = data[!error,]
78 # T4U
79 error = data$T4U == '?'
80 data = data[!error,]
81 # FTI
82 error = data$FTI == '?'
83 data = data[!error,]
84
85 # Se eliminan las variables que indican si se realizo la medicion de la
    hormona
86 data$T3.measured <- NULL
87 data$TSH.measured <- NULL
88 data$TT4.measured <- NULL
89 data$T4U.measured <- NULL
90 data$FTI.measured <- NULL
91 data$TBG.measured <- NULL
92 #Se elimina la variable TBG que no fue medida para ningun paciente
93 data$TBG <- NULL
94 #Se elimina la variable referral.source
95 data$referral.source <- NULL
96

```

```

97 #
98
99 # Una vez limpiada la base de datos, se transforman los datos correspondientes
100
101 # Variables continuas
102 data$T3 <- as.numeric(data$T3)
103 data$age <- as.numeric(data$age)
104 data$TSH <- as.numeric(data$TSH)
105 data$TT4 <- as.numeric(data$TT4)
106 data$T4U <- as.numeric(data$T4U)
107 data$FTI <- as.numeric(data$FTI)
108
109 # Variables binarias
110 data$sex <- as.factor(data$sex)
111 data$sick <- as.factor(data$sick)
112 data$I131 <- as.factor(data$I131)
113 data$state <- as.factor(data$state)
114 data$tumor <- as.factor(data$tumor)
115 data$psych <- as.factor(data$psych)
116 data$goitre <- as.factor(data$goitre)
117 data$lithium <- as.factor(data$lithium)
118 data$pregnant <- as.factor(data$pregnant)
119 data$on.thyroxine <- as.factor(data$on.thyroxine)
120 data$hypopituitary <- as.factor(data$hypopituitary)
121 data$query.thyroxine <- as.factor(data$query.thyroxine)
122 data$thyroid.surgery <- as.factor(data$thyroid.surgery)
123 data$query.hypothyroid <- as.factor(data$query.hypothyroid)
124 data$query.hyperthyroid <- as.factor(data$query.hyperthyroid)
125 data$on.antithyroid.med <- as.factor(data$on.antithyroid.med)
126
127 #Se obtiene la media de los datos continuos para establecerla de manear
    binaria
128 mediaAge <- mean(data$age)
129 mediaT3 <- mean(data$T3)

```

```

130 mediaT4U <- mean(data$T4U)
131 mediaTT4 <- mean(data$TT4)
132 mediaTSH <- mean(data$TSH)
133 mediaFTI <- mean(data$FTI)
134
135 #Se declaran los arreglos con las nuevas variables binarias
136 newAge <- integer(length(data[[1]]))
137 newT3 <- integer(length(data[[1]]))
138 newT4U <- integer(length(data[[1]]))
139 newTT4 <- integer(length(data[[1]]))
140 newTSH <- integer(length(data[[1]]))
141 newFTI <- integer(length(data[[1]]))
142
143 #Se transforman los datos continuos a binarios
144 for(i in 1:length(data[[1]])){
145   #si la edad del paciente es mayor que la media, entonces es un paciente de
   edad avanzada
146   if(data$age[i] >= mediaAge){newAge[i] <- "t"}else{newAge[i] <- "f"}
147   #si la cantidad de hormona presente en el paciente es mayor que la media se
   considera como que posee altos niveles de esta
148   if(data$T3[i] >= mediaT3){newT3[i] <- "t"}else{newT3[i] <- "f"}
149   if(data$TSH[i] >= mediaTSH){newTSH[i] <- "t"}else{newTSH[i] <- "f"}
150   if(data$TT4[i] >= mediaTT4){newTT4[i] <- "t"}else{newTT4[i] <- "f"}
151   if(data$T4U[i] >= mediaT4U){newT4U[i] <- "t"}else{newT4U[i] <- "f"}
152   if(data$FTI[i] >= mediaFTI){newFTI[i] <- "t"}else{newFTI[i] <- "f"}
153 }
154 #Se binariza la columna state, aquellos valores falsos son lo que poseen
   negative y todos los dems son verdaderos, es decir, presentan
   hipotiroidismo
155 data$state <- ifelse(data$state %in% c("primary hypothyroid", "secondary
   hypothyroid", "compensated hypothyroid"), "t", "f")
156 #Se reemplazan las columnas correspondientes
157 data$state <- as.factor(data$state)
158 data$age <- as.factor(newAge)
159 data$T3 <- as.factor(newT3)
160 data$TT4 <- as.factor(newTT4)

```

```

161 data$T4U <- as.factor(newT4U)
162 data$TSH <- as.factor(newTSH)
163 data$FTI <- as.factor(newFTI)
164
165 #Se obtienen los datos que se utilizan para construir el arbol,
166 #es decir, todas las columnas menos state(la que define si el paciente
    presenta hipotiroidismo)
167 data_tree <- subset(data, select = -state)
168
169 #Se modela el arbol con la funcin C5.0 con los datos para el arbol y la
    variable data$state
170 model <- C5.0(data_tree, data$state)
171
172 #Se plotea el arbol
173 tree <- plot(model)

```

# Bibliografía

- [Bomarito] Bomarito, D. M. J. Hipotiroidismo: síntomas y tratamiento. [Online] <https://www.hospitalaleman.org.ar/mujeres/hipotiroidismo-sintomas-y-tratamiento/>.
- [2] Caparrini, F. S. (2017). Aprendizaje inductivo: Árboles de decisión. [Online] <http://www.cs.us.es/~fsancho/?e=104>.
- [3] Chacón, M. (2015). Árboles de decisión. [Online] [http://www.udesantiagovirtual.cl/moodle2/pluginfile.php?file=%2F115779%2Fmod\\_resource%2Fcontent%2F0%2FCapitulo%20VII%20%C3%81rboles%20de%20Decisi%C3%B3n.pdf](http://www.udesantiagovirtual.cl/moodle2/pluginfile.php?file=%2F115779%2Fmod_resource%2Fcontent%2F0%2FCapitulo%20VII%20%C3%81rboles%20de%20Decisi%C3%B3n.pdf).
- [Hurtado] Hurtado, C. Árboles de decisión (i). [Online] [https://www.u-cursos.cl/ingenieria/2007/1/CC52A/1/material\\_docente/bajar?id\\_material=119128](https://www.u-cursos.cl/ingenieria/2007/1/CC52A/1/material_docente/bajar?id_material=119128).