



Laboratorio 4: Clasificador Bayesiano

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

Tabla de contenidos

1. Introducción	1
1.1. Objetivos	1
2. Marco Teórico	2
2.1. Clasificador Bayesiano Ingenuo	2
2.2. Probabilidad a priori	2
2.3. Probabilidad a posteriori	2
3. Obtención del Clasificador	3
3.1. Pre Procesamiento	3
3.1.1. Eliminación de valores vacíos	3
3.1.2. Eliminación de variables	3
3.2. Clase a utilizar	4
3.3. Conjuntos de entrenamiento y de prueba	4
4. Análisis de los resultados	5
4.1. Comparación con experiencias anteriores	6
5. Conclusiones	8
6. Anexos	9
Bibliografía	18

Índice de figuras

Índice de cuadros

1. Matriz de Confusión del clasificador.	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 clasificadores bayesianos es una de las técnicas conocidas para lograr esto y permite establecer que la presencia o ausencia de una característica particular, no está relacionada con la presencia o ausencia de otra característica, es decir, un clasificador de bayes considera que cada característica contribuye de manera independiente a la probabilidad de que algo se cumpla.

Esta cuarta experiencia de laboratorio consiste en, mediante el proceso de minería de datos y a través de la aplicación de clasificadores bayesianos 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 clasificadores bayesianos como medida de obtención de conocimiento.
2. Corroborar la teoría vista en clases mediante la aplicación del clasificador bayesiano en el lenguaje estadístico R, concretamente en el paquete "*e1071*".
3. Obtener conocimiento de la base de datos "*Allhypo*" a través de esta técnica de minería de datos.

A continuación se define un marco teórico con la descripción del clasificador bayesiano, y la probabilidad apriori y aposteriori. Luego se presentan los resultados de la ejecución del script en R con la obtención de los clasificadores bayesianos para posteriormente realizar un análisis de los resultados y finalizar con las conclusiones respecto a esta experiencia.

2. Marco Teórico

Para contextualizar el ambiente de desarrollo relacionado a esta experiencia se disponen de las definiciones de los siguientes conceptos.

2.1. Clasificador Bayesiano Ingenuo

Este clasificador es de tipo probabilístico el cual es muy utilizado para resolver problemas de clasificación. Su funcionamiento está basado en el Teorema de Bayes, el cual corresponde al teorema que permite encontrar la probabilidad a posteriori de eventos suponiendo que se conoce la probabilidad a priori de este evento, la probabilidad condicional de una o mas condiciones dada el evento y la probabilidad a priori de las condiciones mencionadas anteriormente; como también esta basado en la suposición de que las variables predictoras son independientes, la cual es la razón de que este clasificador se considere ingenuo.

2.2. Probabilidad a priori

Esta corresponda a la probabilidad de ocurrencia de un evento sin condiciones adicionales. Si se ve de forma mas formal: La probabilidad a priori $p(c_{ie})$ será la probabilidad de que el sujeto se clasifique en la clase c_i , representado de la siguiente manera:

$$p(c_i) = \lim_{n \rightarrow +\infty} \frac{n_i}{n} \quad (1)$$

2.3. Probabilidad a posteriori

Corresponde a la probabilidad de ocurrencia de un evento dada ciertas condiciones. En palabras formales, se tiene que es la probabilidad de que un sujeto pertenezca a la clase c_i dado un valor y . Por ejemplo, para dos clases se tiene lo siguiente:

$$\sum_{i=1}^2 p(c_i/y) = 1 \quad (2)$$

Para obtener la probabilidad $p(c_i/y)$, se requiere conocer las relaciones de probabilidades condicionales.

3. Obtención del Clasificador

3.1. 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:

3.1.1. 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, para que así no entorpezca el estudio de esta base de datos.

3.1.2. 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 columnas correspondientes a la variables que son eliminadas son:

- Variable TBG: Corresponde a la proteína fijador de la Tiroxina (T4). 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 paciente de la muestra, por lo tanto se opta por eliminarla de la base de datos.
- Variables de verificación de medición (measured): Indican si el paciente se ha realizado alguna medición, siendo esta información no relevante para el estudio. Lo que si importa es el valor numérico de estas mediciones de las hormonas por litro de sangre en los pacientes. Entonces, se decide dejar fuera las siguientes variables de medición: *TSH measured*, *T3 measuted*, *TT4 measured*, *T4U measured*, *FTI measured* y *TBG measured*.
- Variable de fuente: Como también sucede con las variables anteriores, 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.

3.2. Clase a utilizar

Una vez realizado el pre-procesamiento, se procede a encontrar la Clase a utilizar de acuerdo a lo que se quiere buscar con este análisis utilizando Clasificador Bayesiano. Para esto, es importante plantear lo que se está buscando en relación a la información que se tiene de la base de datos. Entonces, la pregunta que se obtiene es:

¿Cuál es la probabilidad de que un sujeto padezca Hipotiroidismo a partir de información como la edad, sexo, resultado de sus exámenes hormonales, entre otros datos?

En base a esta pregunta, que una persona padezca de Hipotiroidismo no depende solo de un examen hormonal, por ejemplo, se tiene que esta condición depende de una serie de síntomas, que debido a la investigación que se realizó para la primera experiencia se pudo dar cuenta que los distintos atributos correspondientes a la base de datos *allhypo.data* de alguna manera se relacionan entre sí, como por ejemplo las medidas hormonales como la TSH y la TT4 o T4U. Entonces la conclusión que se obtiene es que con un solo atributo no se puede definir si un sujeto va a padecer o no de Hipotiroidismo. Si una persona tiene niveles altos o bajos de la hormona T3, no quiere decir que padecerá Hipotiroidismo, es necesario tener en cuenta otros factores para tener una certeza de que el diagnóstico será positivo. Dicho lo anterior, como atributo predictor se decide que la variable *State* es la mas indicada, la que muestra el cuadro de Hipotiroidismo para cada sujeto.

3.3. Conjuntos de entrenamiento y de prueba

Para poder llevar a cabo el modelo, se necesita un *Conjunto de entrenamiento* el cual se usa para la construcción y ajuste del modelo de acuerdo a información previa que se tiene. Y también es necesario de un *Conjunto de prueba (o test)* para luego aplicar el modelo. Para el *Conjunto de entrenamiento*, en vez de realizar algún método para obtener este conjunto, se decide ocupar la base de datos *allhypo.data*, lo mismo ocurre para el *Conjunto de prueba*, para evitar implementar un método y simplificar el trabajo, se decide trabajar con la base de datos *allhypo.test* asumiendo que los datos son representativos para los dos conjuntos. Ambas bases de datos utilizadas para estos conjuntos, se debe realizar el pre procesamiento para la limpieza de ellas.

4. Análisis de los resultados

Para observar los resultados del Clasificador, se tiene la Matriz de Confusión, en donde se va comparando la predicción de cada clase con las instancias de las clases reales, que en este caso corresponden las del conjunto de prueba.

A continuación en el Cuadro 1 se presenta la Matriz de confusión del Clasificador Bayesiano ingenuo. Las columnas corresponden a la clase real (conjunto de prueba) y las filas corresponden a las clasificación que se realiza.

Predicciones	Compensated Hypothyroid	Negative	Primary Hypothyroid
Compensated Hypothyroid	11	1	1
Negative	19	634	3
Primary Hypothyroid	0	2	25
Secondary Hypothyroid	0	0	0

Cuadro 1: Matriz de Confusión del clasificador.

Cabe destacar que en el Cuadro 1 en la columna correspondiente a las clases reales del conjunto de prueba no se aprecia el tipo de Hipotiroidismo denominado *Secondary Hypothyroid*, pero si en la que tiene ver con la clasificación, esto ocurre debido a que la base de datos que corresponde al conjunto de prueba, no se aprecian resultados de este tipo de Hipotiroidismo, solo se encuentran las otras 3 en el conjunto de prueba. Aunque esta situación no se ve a simple vista que no genera problema, se desconoce si realmente lo genera, debido a que está la posibilidad de que igual se realizara la clasificación aunque no se tenga el conjunto para realizarlo lo que provocaría un error en el clasificador. A partir de la Matriz de Confusión mostrada anteriormente, se pueden obtener las precisiones con respecto al rendimiento de este clasificador. Se obtuvieron 4 medidas distintas en donde se mide esta precisión, una correspondiente a una precisión general y las otras 3 en relación a los valores predictivos, la cuales son:

1. **Precisión General:** Ésta tiene una proporción de 0.9626437, correspondiendo a un 96.26 %.

2. **Valor Predictivo "Negative"**: Ésta tiene una proporción de 0.9664634, correspondiendo a un 96.64 %.
3. **Valor Predictivo "Primary Hypothyroid"**: Ésta tiene una proporción de 0.9259259, correspondiendo a un 92.59 %.
4. **Valor Predictivo "Compensated Hypothyroid"**: Ésta tiene una proporción de 0.8461538, correspondiendo a un 84.61 %.

Con respecto a estas medidas, se puede decir que el clasificador tiene una alta precisión, siendo la más alta la del Valor Predictivo *Negative* con un error del solo del 3,35366 %. Si se realiza una comparación con los otros valores predictivos se tiene que la *Presición General* tiene un error del 3,73563 %, *Primary Hypothyroid* tiene un error del 7,40741 % y *Compensated Hypothyroid* tiene un error del 15,38462 %.

4.1. Comparación con experiencias anteriores

Con respecto al *Método de las K-medias* realmente es complicado realizar una comparación con esta experiencia que tiene que ver con *Clasificador Bayesiano* debido a que lo que se realiza en ambas difieren en su propósito y uso. En donde el objetivo de K-medias esta enfocado en realizar agrupamientos de datos de acuerdo a la distancia que hay entre ellos, en cambio el enfoque del Clasificador Bayesiano está en obtener probabilidades a priori y posteriori (Teorema de Bayes) para luego realizar la clasificación, y se le agrega el concepto de *ingenuo* que tiene que ver con la independencia entre las variables.

Los resultados obtenidos de las K-medias en donde se realizaron dos Clusters, en el primero, se pudo determinar que el factor mas notorio para la agrupación tenía que ver con el sexo del sujeto y para el segundo, se determinó que debido a la variación de las distintas variables correspondiente a las hormonas en las distintas agrupaciones, existe una alta correlación de estas variables. Los resultados correspondientes a esta experiencia, la clasificación se enfocó de acuerdo al cuadro del paciente.

Con respecto a lo realizado con las *Reglas de Asociación*, al igual que lo realizado en esta experiencia, ambas se enfocan en la clase o variable predicha, en donde se trata de

predecir el cuadro de un sujeto de acuerdo variables involucradas (en este caso, las que se encuentran en la base de datos utilizada para las experiencias). Lo que las diferencia, es que ambas utilizan metodologías distintas para realizar lo dicho anteriormente. Para las Reglas de Asociación, lo que se hace es obtener reglas que buscan relacionar distintas variables de tal manera que su asociación en conjunto, tenga como consecuencia la clase, en cambio, con el Clasificador Bayesiano ingenuo lo que se hace es, mediante a probabilidades, se busca la probabilidad a posteriori, que de acuerdo al Teorema de Bayes, se encuentra ésta suponiendo que se conoce la probabilidad a priori del evento.

Los resultados que muestran las Reglas de Asociación son las reglas que tienen el Hipotiroidismo como consecuente, sin especificar el tipo de Hipotiroidismo (Primary Hypothyroid, Secondary Hypothyroid o Compensated Hypothyroid), en cambio, los resultados del Clasificador son mas precisos de acuerdo a la clase se que eligió para la clasificación.

5. Conclusiones

Para obtener el clasificador bayesiano se realiza en primer lugar una limpieza a la base de datos, en la que se eliminan variables que no entregan información, tanto para el conjunto de datos de entrenamiento, como para el conjunto de prueba. Luego se utiliza el algoritmo naiveBayes para la obtención de las probabilidades dada la variable predictora que sería "State", correspondiente a la clasificación de hipotiroidismo del paciente, obtienen las predicciones con respecto al modelo obtenido con naiveBayes y se representan los resultados a través de una matriz de confusión.

Los objetivos de esta experiencia se cumplen, ya que al analizar los resultados se observa que el clasificador obtenido con respecto a la clase State posee altos valores predictivos y que la precisión general de éste es de 96,26 %.

Para el desarrollo de esta experiencia se presentan algunas dificultades, pero mayormente, la interpretación de los resultados obtenidos por el clasificador bayesiano.

6. Anexos

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

```
1 library("e1071")
2 # LIMPIEZA DE LA BASE DE DATOS DEL CONJUNTO DE ENTRENAMIENTO
3 data.entrenamiento <- read.csv("/home/nicole/allhypo.data",
4                                col.names = c("age", "sex", "on.thyroxine", "query.thyroxine",
5                                              "on.antithyroid.med",
6                                              "sick", "pregnant", "thyroid.surgery", "I131", "
7                                              query.hypothyroid",
8                                              "query.hyperthyroid", "lithium", "goitre", "
9                                              tumor", "hypopituitary",
10                                             "psych", "TSH.measured", "TSH", "T3.measured",
11                                             "T3", "TT4.measured",
12                                             "TT4", "T4U.measured", "T4U", "FTI.measured", "
13                                             FTI", "TBG.measured",
14                                             "TBG", "referral.source", "state"), header=FALSE
15                                , sep=",", stringsAsFactors=FALSE)
16 # Se limpia la variable de estado del paciente eliminando numeros
17 state <- data.entrenamiento$state
18 state <- strsplit(state, ".", fixed = TRUE)
19 state <- lapply(state, '[', 1)
20 state <- unlist(state)
21 data.entrenamiento$state <- state
22 # Se eliminan los pacientes que contengan algun NA
23 # Sex
24 error = data.entrenamiento$sex == '?'
25 data.entrenamiento = data.entrenamiento[!error,]
26 # Age
27 error = data.entrenamiento$age == '?'
28 data.entrenamiento = data.entrenamiento[!error,]
29 # On thyroxine
30 error = data.entrenamiento$on.thyroxine == '?'
31 data.entrenamiento = data.entrenamiento[!error,]
32 # Query on antithyroid thyroxine
```

```

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

```

```

63 error = data.entrenamiento$psych == '?'
64 data.entrenamiento = data.entrenamiento[!error,]
65 # TSH
66 error = data.entrenamiento$TSH == '?'
67 data.entrenamiento = data.entrenamiento[!error,]
68 # T3
69 error = data.entrenamiento$T3 == '?'
70 data.entrenamiento = data.entrenamiento[!error,]
71 # TT4
72 error = data.entrenamiento$TT4 == '?'
73 data.entrenamiento = data.entrenamiento[!error,]
74 # T4U
75 error = data.entrenamiento$T4U == '?'
76 data.entrenamiento = data.entrenamiento[!error,]
77 # FTI
78 error = data.entrenamiento$FTI == '?'
79 data.entrenamiento = data.entrenamiento[!error,]
80 # Se eliminan las variables que indican si se realizo la medicion de la
    hormona
81 data.entrenamiento$T3.measured <- NULL
82 data.entrenamiento$TSH.measured <- NULL
83 data.entrenamiento$TT4.measured <- NULL
84 data.entrenamiento$T4U.measured <- NULL
85 data.entrenamiento$FTI.measured <- NULL
86 data.entrenamiento$TBG.measured <- NULL
87 #Se elimina la variable TBG que no fue medida para ningun paciente
88 data.entrenamiento$TBG <- NULL
89 #Se elimina la variable referral.source la cual no es relevante
90 data.entrenamiento$referral.source <- NULL
91 #
92 #
93 # LIMPIEZA DE LA BASE DE DATOS DEL CONJUNTO DE PRUEBA
94 data.prueba <- read.csv("/home/nicole/allhypo.test",
95                          col.names = c("age", "sex", "on.thyroxine", "query.
    thyroxine", "on.antithyroid.med",

```

```

96         "sick", "pregnant", "thyroid.surgery", "
    I131", "query.hypothyroid",
97         "query.hyperthyroid", "lithium", "goitre"
    , "tumor", "hypopituitary",
98         "psych", "TSH.measured", "TSH", "T3.
    measured", "T3", "TT4.measured",
99         "TT4", "T4U.measured", "T4U", "FTI.
    measured", "FTI", "TBG.measured",
100         "TBG", "referral.source", "state"), header
    =FALSE, sep=",", stringsAsFactors=FALSE)
101 # Se limpia la variable de estado del paciente eliminando numeros
102 state <- data.prueba$state
103 state <- strsplit(state, ".", fixed = TRUE)
104 state <- lapply(state, '[', 1)
105 state <- unlist(state)
106 data.prueba$state <- state
107 # Se eliminan los pacientes que contengan algun NA
108 # Sex
109 error = data.prueba$sex == '?'
110 data.prueba = data.prueba[!error,]
111 # Age
112 error = data.prueba$age == '?'
113 data.prueba = data.prueba[!error,]
114 # On thyroxine
115 error = data.prueba$on.thyroxine == '?'
116 data.prueba = data.prueba[!error,]
117 # Query on antithyroid thyroxine
118 error = data.prueba$query.thyroxine == '?'
119 data.prueba = data.prueba[!error,]
120 # On antithyroid medication
121 error = data.prueba$on.antithyroid.med == '?'
122 data.prueba = data.prueba[!error,]
123 # Sick
124 error = data.prueba$sick == '?'
125 data.prueba = data.prueba[!error,]
126 # Pregnant

```

```

127 error = data.prueba$pregnant == '?'
128 data.prueba = data.prueba[!error,]
129 # Thyroid surgery
130 error = data.prueba$thyroid.surgery == '?'
131 data.prueba = data.prueba[!error,]
132 # I131 treatment
133 error = data.prueba$I131 == '?'
134 data.prueba = data.prueba[!error,]
135 # Query hypothyroid
136 error = data.prueba$query.hypothyroid == '?'
137 data.prueba = data.prueba[!error,]
138 # Query hyperthyroid
139 error = data.prueba$query.hyperthyroid == '?'
140 data.prueba = data.prueba[!error,]
141 # Lithium
142 error = data.prueba$lithium == '?'
143 data.prueba = data.prueba[!error,]
144 # Goitre
145 error = data.prueba$goitre == '?'
146 data.prueba = data.prueba[!error,]
147 # Tumor
148 error = data.prueba$tumor == '?'
149 data.prueba = data.prueba[!error,]
150 # Hypopituitary
151 error = data.prueba$hypopituitary == '?'
152 data.prueba = data.prueba[!error,]
153 # Psych
154 error = data.prueba$psych == '?'
155 data.prueba = data.prueba[!error,]
156 # TSH
157 error = data.prueba$TSH == '?'
158 data.prueba = data.prueba[!error,]
159 # T3
160 error = data.prueba$T3 == '?'
161 data.prueba = data.prueba[!error,]
162 # TT4

```



```

163 error = data.prueba$TT4 == '?'
164 data.prueba = data.prueba[!error,]
165 # T4U
166 error = data.prueba$T4U == '?'
167 data.prueba = data.prueba[!error,]
168 # FTI
169 error = data.prueba$FTI == '?'
170 data.prueba = data.prueba[!error,]
171 # Se eliminan las variables que indican si se realizo la medicion de la
    hormona
172 data.prueba$T3.measured <- NULL
173 data.prueba$TSH.measured <- NULL
174 data.prueba$TT4.measured <- NULL
175 data.prueba$T4U.measured <- NULL
176 data.prueba$FTI.measured <- NULL
177 data.prueba$TBG.measured <- NULL
178 #Se elimina la variable TBG que no fue medida para ningun paciente
179 data.prueba$TBG <- NULL
180 #Se elimina la variable referral.source la cual no es relevante
181 data.prueba$referral.source <- NULL
182 #
183 #
184 # Una vez limpiada la base de datos, se transforman los datos correspondientes
185
186 # CONJUNTO DE ENTRENAMIENTO
187 # Variables continuas
188 data.entrenamiento$T3 <- as.numeric(data.entrenamiento$T3)
189 data.entrenamiento$age <- as.numeric(data.entrenamiento$age)
190 data.entrenamiento$TSH <- as.numeric(data.entrenamiento$TSH)
191 data.entrenamiento$TT4 <- as.numeric(data.entrenamiento$TT4)
192 data.entrenamiento$T4U <- as.numeric(data.entrenamiento$T4U)
193 data.entrenamiento$FTI <- as.numeric(data.entrenamiento$FTI)
194 # Variables binarias
195 data.entrenamiento$sex <- as.factor(data.entrenamiento$sex)
196 data.entrenamiento$sick <- as.factor(data.entrenamiento$sick)
197 data.entrenamiento$I131 <- as.factor(data.entrenamiento$I131)

```

```

198 data.entrenamiento$state <- as.factor(data.entrenamiento$state)
199 data.entrenamiento$tumor <- as.factor(data.entrenamiento$tumor)
200 data.entrenamiento$psych <- as.factor(data.entrenamiento$psych)
201 data.entrenamiento$goitre <- as.factor(data.entrenamiento$goitre)
202 data.entrenamiento$lithium <- as.factor(data.entrenamiento$lithium)
203 data.entrenamiento$pregnant <- as.factor(data.entrenamiento$pregnant)
204 data.entrenamiento$on.thyroxine <- as.factor(data.entrenamiento$on.thyroxine)
205 data.entrenamiento$hypopituitary <- as.factor(data.entrenamiento$hypopituitary
    )
206 data.entrenamiento$query.thyroxine <- as.factor(data.entrenamiento$query.
    thyroxine)
207 data.entrenamiento$thyroid.surgery <- as.factor(data.entrenamiento$thyroid.
    surgery)
208 data.entrenamiento$query.hypothyroid <- as.factor(data.entrenamiento$query.
    hypothyroid)
209 data.entrenamiento$query.hyperthyroid <- as.factor(data.entrenamiento$query.
    hyperthyroid)
210 data.entrenamiento$on.antithyroid.med <- as.factor(data.entrenamiento$on.
    antithyroid.med)
211
212 #CONJUNTO DE PRUEBA
213 # Variables continuas
214 data.prueba$T3 <- as.numeric(data.prueba$T3)
215 data.prueba$age <- as.numeric(data.prueba$age)
216 data.prueba$TSH <- as.numeric(data.prueba$TSH)
217 data.prueba$TT4 <- as.numeric(data.prueba$TT4)
218 data.prueba$T4U <- as.numeric(data.prueba$T4U)
219 data.prueba$FTI <- as.numeric(data.prueba$FTI)
220 # Variables binarias
221 data.prueba$sex <- as.factor(data.prueba$sex)
222 data.prueba$sick <- as.factor(data.prueba$sick)
223 data.prueba$I131 <- as.factor(data.prueba$I131)
224 data.prueba$state <- as.factor(data.prueba$state)
225 data.prueba$tumor <- as.factor(data.prueba$tumor)
226 data.prueba$psych <- as.factor(data.prueba$psych)
227 data.prueba$goitre <- as.factor(data.prueba$goitre)

```

```

228 data.prueba$lithium <- as.factor(data.prueba$lithium)
229 data.prueba$pregnant <- as.factor(data.prueba$pregnant)
230 data.prueba$on.thyroxine <- as.factor(data.prueba$on.thyroxine)
231 data.prueba$hypopituitary <- as.factor(data.prueba$hypopituitary)
232 data.prueba$query.thyroxine <- as.factor(data.prueba$query.thyroxine)
233 data.prueba$thyroid.surgery <- as.factor(data.prueba$thyroid.surgery)
234 data.prueba$query.hypothyroid <- as.factor(data.prueba$query.hypothyroid)
235 data.prueba$query.hyperthyroid <- as.factor(data.prueba$query.hyperthyroid)
236 data.prueba$on.antithyroid.med <- as.factor(data.prueba$on.antithyroid.med)
237 #
238 #
239 # SE PROCEDE A REALIZAR LO NECESARIO PARA EL CLASIFICADOR BAYESIANO INGENUO
240
241 # Se realizan las probabilidades dada la variable predictora que es la State (
      clasificaci n del Hipotiroidismo)
242 modelo <- naiveBayes(state ~., data = data.entrenamiento)
243 # Se realizan las predicciones con respecto al modelo anterior del conjunto de
      Prueba.
244 predicciones <- predict(object = modelo, newdata=data.prueba, type = "class")
245 # Se obtiene la matriz de confusion con los resultados anteriores y la
      variable predictora
246 matriz.de.confusion <- table(predicciones, data.prueba$state)
247 # A partir de la matriz de confusion se calculan las distintas precisiones con
      respecto a la clasificacion del Hipotiroidismo
248 # (correspondiente a la variable "State")
249 precision.en.general <- sum(diag(matriz.de.confusion)) / sum(matriz.de.
      confusion)
250 # Precision para Primary Hypothyroid (Hipotiroidismo Primario)
251 precision.primary <- sum(matriz.de.confusion[3,3]) / sum(matriz.de.confusion
      [3,])
252 # Precision para Negative
253 precision.negative <- sum(matriz.de.confusion[2,2]) / sum(matriz.de.confusion
      [2,])
254 # Presicion para Compensated Hypothyroid (Hipotiroidismo compensado o
      subclinico)

```

```
255 precision.compensated <- sum(matriz.de.confusion[1,1]) / sum(matriz.de.  
    confusion[1,])
```

Bibliografía

Chacón, M. (2017). Clasificación bayesiana. [Online] <http://www.udesantiagovirtual.cl/moodle2/mod/resource/view.php?id=156965>.

Richter-Walsh, S. (2017). Clasificación naive bayes en r (parte 2). [Online] <https://www.r-bloggers.com/naive-bayes-classification-in-r-part-2>.

Sucar, L. E. (2012). Sesión 6: Clasificadores bayesianos. [Online] <https://ccc.inaoep.mx/~esucar/Clases-mgp/pgm06-clasif-2012.pdf>.

Vega, J. B. M. (2018). Naïve bayes con r para clasificación de texto. [Online] https://rstudio-pubs-static.s3.amazonaws.com/378731_524d5e2c353c4dbebc8a23e4235a383a.html#ajustando-naive-bayes.