

Actividad 5

Leonardo Flores Torres

23 de diciembre de 2022

- Utilizando los siguientes conjuntos de entrenamiento del repositorio UCL (<https://archive.ics.uci.edu/ml/index.php>):

- Car evaluation
- Tic-tac-toe
- Qualitative bankruptcy

llevar a cabo las siguientes actividades:

1. Implementar en Lisp las mejoras al algoritmo ID3 que adoptaron en su implementación en Prolog. [30/100]
2. Con base en la validación cruzada, implementar una función de clasificación por votación en donde participen los árboles construidos en este proceso. Si se construyen diez árboles, los diez se utilizan para clasificar un nuevo caso cuya clase será la más votada por esos diez árboles ¿Cómo se compara la eficiencia de este enfoque contra el mejor árbol encontrado? [50/100]
3. Implementar una función que traduzca un árbol de decisión dado a un programa en Prolog equivalente, probar el clasificador en ese lenguaje. [20/100]

Solución:

Para poder utilizar el programa escrito en `sbcl` del algoritmo ID3 fue necesario modificar un par de archivos. Primero se agregó el archivo de `cl-id3-cross-validation` en el archivo para la definición del sistema `cl-id3.asd` para poder acceder a las funciones escritas ahí que tienen que ver con la validación cruzada:

```
1 (asdf:defsystem :cl-id3
2   :depends-on (:split-sequence)
3   :components ((:file "cl-id3-package")
4     (:file "cl-id3-algorithm"
5       :depends-on ("cl-id3-package"))
6     (:file "cl-id3-load"
7       :depends-on ("cl-id3-package"
8         "cl-id3-algorithm"))
9     (:file "cl-id3-classify"
```

```

10      :depends-on ("cl-id3-package"
11                  "cl-id3-algorithm"
12                  "cl-id3-load"))
13      ;;; Making cross validation available
14      (:file "cl-id3-cross-validation"
15          :depends-on ("cl-id3-package"
16                      "cl-id3-algorithm"
17                      "cl-id3-load"
18                      "cl-id3-classify"))))

```

También se agregó el archivo de validación cruzada en el archivo donde se hace la definición del paquete `cl-id3-package.lisp` :

```

1  (defpackage :cl-id3
2    (:use :cl :split-sequence)
3    (:export :load-file
4            :induce
5            :print-tree
6            :classify
7            :classify-new-instance
8            ;;; Making cross validation available
9            :cross-validation))

```

La mejora que se realizó en `prolog` fue tomar en cuenta la información individual para cada atributo y usarla como denominador en la división con la ganancia para así encontrar la razón de ganancia. Se escribió una función equivalente en el archivo `cl-id3-algorithm.lisp` para computar la información individual, como se muestra a continuación:

```

1  ;; Improvement to id3 by considering the information value of an attribute.
2  (defun split-information (examples attribute)
3    "It computes the split-information for an ATTRIBUTE in EXAMPLES"
4    (let ((parts (get-partition attribute examples))
5          (no-examples (count-if #'atom examples)))
6      (-
7        (apply #'+
8              (mapcar
9                #'(lambda (part)
10                   (let* ((size-part (count-if #'atom
11                                                (cdr part)))
12                         (proportion (if (eq size-part 0)
13                                         0
14                                         (* (/ size-part no-examples)
15                                              (log (/ size-part no-examples) 2))))))
16                  (* proportion 1)))
17        (cdr parts)))))

```

Al haber implementado `split-information` se procedió a modificar la función `best-partition`, en el mismo archivo, para cambiar la manera en como se calcula a la mejor partición, en vez de usar solo la ganancia ahora será usando a la razón de ganancia.

```

18 (defun best-partition (attributes examples)
19   "It computes one of the best partitions induced by ATTRIBUTES over EXAMPLES"
20   (let* ((info-gains
21          (loop for attrib in attributes collect
22                (let ((ig (information-gain examples attrib))
23                      ;; Improvement using split information to consider the gain ratio
24                      ;; instead of just information gain.
25                      (si (split-information examples attrib))
26                      (p (get-partition attrib examples))))
27          (when *trace*
28            (format t "Partición inducida por el atributo ~s:~%~s~%"
29                    attrib p)
30            (format t "Razón de ganancia: ~s~%"
31                    ;; GAIN RATIO
32                    (/ ig (+ 1 si))))
33          (list (/ ig (+ 1 si)) p))))
34   (best (cadar (sort info-gains #'(lambda(x y) (> (car x) (car y))))))
35   (when *trace* (format t "Best partition: ~s~%-----~%" best))
36   best))

```

Además, se arregló el mismo problema que surgió anteriormente en la implementación en `prolog` en que si el contenido de información $IV(A_i)$ de un atributo era 0 entonces la razón de ganancia diverge. Se puede observar que se tomó la misma solución, sumar una unidad al contenido de información para cada atributo $IV'(A_i) = IV(A_i) + 1$.

Ahora que ya se agregaron las modificaciones mencionadas se procederá a encontrar la mejor partición para las 3 bases de datos. Los conjuntos de entrenamiento son lo suficientemente grandes como para que el poner los resultados en este reporte sea impracticable, por lo que se decidió incluir algunos de los resultados en archivos de texto individuales. Para calcular la mejor partición de cualquiera de estas bases de datos primero se tiene que hacer lo siguiente:

```

1 * (asdf:load-system :cl-id3)
2 * (in-package :cl-id3)
3 * (load-file "./path/to/myfile.arff")
4 * (best-partition (remove *target* *attributes*) *examples*)

```

Los resultados de las mejores particiones iniciales para cada uno de los archivos de las bases de datos se guardó en los siguientes archivos dentro del directorio `/outputfiles`,

- `/outputfiles/bestpartition_car.txt`,
- `/outputfiles/bestpartition_tictactoe.txt`,
- `/outputfiles/bestpartition_bankruptcy.txt`.

También se incluyó la información desglosada para cada partición posible, junto con sus respectivas razones de ganancia en

- `/outputfiles/bestpartition_car_trace.txt`,
- `/outputfiles/bestpartition_tictactoe_trace.txt`,

- `/outputfiles/bestpartition_bankruptcy_trace.txt` .

En realidad, la mejor partición inicial sin haber habilitado el `*trace*` es la que corresponde a la que tiene una razón de ganancia mayor. En los archivos con `*trace*` habilitado se pueden ver las razones de ganancia para cada uno y verificar que sí es el caso.

No es necesario computar las particiones iniciales, pero no está de más ver como es que las razones de ganancia que se calculan para cada suposición en cada base de datos. Se dice que no es necesario porque el algoritmo ID3 implementado en la función `id3` llama de manera recursiva a la función `best-partition`, y la función `induce` realiza la construcción del árbol producido en `id3`. Se utilizará la función `print-tree` para imprimir en consola una representación del árbol producido al ejecutar `induce`. Los árboles producidos se incluyen en los archivos:

- `/outputfiles/tree_car.txt` ,
- `/outputfiles/tree_tictactoe.txt` ,
- `/outputfiles/tree_bankruptcy.txt` .

El árbol más pequeño, el cual es razonable incluir directamente, es el de la base de datos del banco. Dicho árbol se muestra a continuación:

```

1  * (print-tree (induce))
2  CO
3  - P → NB
4  - N → B
5  - A
6      CR
7      - N
8          FF
9          - A → NB
10         - N → B
11     - P → NB
12     - A → NB
13  NIL

```

Para el segundo punto de esta actividad se modificó la función de validación cruzada `cross-validation` en el archivo `cl-id3-cross-validation.lisp` para no imprimir cada árbol generado dentro del `loop`. En casos de árboles pequeños tiene sentido pensar en ver los árboles que se puedan generar, pero cuando dichos árboles son muy grandes tienden a ocupar demasiado espacio en terminal lo que personalmente produce más confusión que claridad.

```

1  ;; Optional variable if the user doesn't want to print the tree
2  (defun cross-validation (k &optional (print_ t))
3    ;; Clean the tree list
4    (setq *trees* nil)
5    (let* ((long (length *examples*)))
6      (loop repeat k do
7        (let* ((training-data (folding (- long k) long))

```

```

8      (test-data (difference training-data *examples*))
9      (tree (induce training-data)))
10     ;; Only report the tree if wanted, will do it by default.
11     (if print_
12         (report tree test-data))
13     ;; Save tree
14     (push tree *trees*))))))

```

Además, se declaró una variable global `*trees*` para guardar cada árbol producido durante la llamada de la función para tenerlos disponibles después durante el proceso de votación. Los resultados de llamar a la función `cross-validation` en cada una de las bases de datos se guardaron en los siguientes archivos

- `/outputfiles/crossvalidation_car.txt` ,
- `/outputfiles/crossvalidation_tictactoe.txt` ,
- `/outputfiles/crossvalidation_bankruptcy.txt` .

Se comenzó a realizar la validación cruzada con el archivo de los autos, pero la salida no se incluye aquí sino en un respectivo archivo:

```

1  * (asdf:load-system :cl-id3)
2
3  * (load-file "./databases/car.arff")
4
5  * (cross-validation 10)
6  ; output removed

```

La función `let-them-vote` toma un ejemplo como argumento, y clasifica el ejemplo respecto a cada árbol en `*trees*`, y guarda la clase más votada en una variable global `*main-class*`. Ahora a hacerlos votar, para esto se tomará el siguiente ejemplo `(high med 5more more big high acc)` ,

```

7  * (let-them-vote '(high med 5more more big high acc))
8
9  Classification for tree no. 1: ACC
10
11 Classification for tree no. 2: ACC
12
13 Classification for tree no. 3: ACC
14
15 Classification for tree no. 4: ACC
16
17 Classification for tree no. 5: ACC
18
19 Classification for tree no. 6: ACC
20
21 Classification for tree no. 7: ACC
22
23 Classification for tree no. 8: ACC

```

```

24
25 Classification for tree no. 9: ACC
26
27 Classification for tree no. 10: ACC
28
29 Class UNACC was detected 0 time(s).
30
31 Class ACC was detected 10 time(s).
32
33 Class GOOD was detected 0 time(s).
34
35 Class VGOOD was detected 0 time(s).
36
37 The most voted class was: ACC
38 NIL

```

Se obtienen 10 clasificaciones (una por árbol generado al llamar a la función de validación cruzada), el conteo de ocurrencias de las clases y la clase más votada, la cual en este caso es `acc`.

Ahora se repetirá lo mismo para la base de datos del banco, tomando como ejemplo `(P N N N N B)`:

```

1  * (asdf:load-system :cl-id3)
2  T
3
4  * (in-package :cl-id3)
5  #<PACKAGE "CL-ID3">
6
7  * (load-file "./databases/car.arff")
8  The ID3 setting has been reset.
9  Training set initialized after "./databases/car.arff".
10 NIL
11
12 * (load-file "./databases/bankruptcy.arff")
13 The ID3 setting has been reset.
14 Training set initialized after "./databases/bankruptcy.arff".
15 NIL
16
17 * (cross-validation 10)
18 CO
19 - N → B
20 - P → NB
21 - A
22   CR
23   - N
24     FF
25     - A → NB
26     - N → B
27   - P → NB
28   - A → NB
29

```

```

30 Instances classified correctly: 10
31 Instances classified incorrectly: 0
32
33 CO
34 - A
35     CR
36     - P → NB
37     - A → NB
38     - N
39     FF
40     - N → B
41     - A → NB
42 - N → B
43 - P → NB
44
45 Instances classified correctly: 10
46 Instances classified incorrectly: 0
47
48 CO
49 - P → NB
50 - N → B
51 - A
52     CR
53     - N
54     FF
55     - A → NB
56     - N → B
57     - A → NB
58     - P → NB
59
60 Instances classified correctly: 10
61 Instances classified incorrectly: 0
62
63 CO
64 - N → B
65 - P → NB
66 - A
67     CR
68     - N
69     FF
70     - A → NB
71     - N → B
72     - P → NB
73     - A → NB
74
75 Instances classified correctly: 10
76 Instances classified incorrectly: 0
77
78 CO
79 - N → B
80 - P → NB
81 - A
82     CR
83     - N

```

```

84         FF
85         - A → NB
86         - N → B
87     - A → NB
88     - P → NB
89
90 Instances classified correctly: 10
91 Instances classified incorrectly: 0
92
93 CO
94 - N → B
95 - A
96     CR
97     - N
98         FF
99         - A → NB
100        - N → B
101        - P → NB
102        - A → NB
103    - P → NB
104
105 Instances classified correctly: 10
106 Instances classified incorrectly: 0
107
108 CO
109 - A
110     CR
111     - A → NB
112     - P → NB
113     - N
114         FF
115         - A → NB
116         - N → B
117    - N → B
118    - P → NB
119
120 Instances classified correctly: 10
121 Instances classified incorrectly: 0
122
123 CO
124 - A
125     CR
126     - N
127         FF
128         - N → B
129         - A → NB
130        - P → NB
131        - A → NB
132    - P → NB
133    - N → B
134
135 Instances classified correctly: 10
136 Instances classified incorrectly: 0
137

```



```

138 CO
139 - A
140   CR
141   - N
142     FF
143     - A → NB
144     - N → B
145   - A → NB
146   - P → NB
147 - N → B
148 - P → NB
149
150 Instances classified correctly: 10
151 Instances classified incorrectly: 0
152
153 CO
154 - P → NB
155 - N → B
156 - A
157   CR
158   - N
159     FF
160     - A → NB
161     - N → B
162   - A → NB
163   - P → NB
164
165 Instances classified correctly: 10
166 Instances classified incorrectly: 0
167
168 NIL
169
170 * (let-them-vote '(P N N N N B))
171
172 Classification for tree no. 1: B
173
174 Classification for tree no. 2: B
175
176 Classification for tree no. 3: B
177
178 Classification for tree no. 4: B
179
180 Classification for tree no. 5: B
181
182 Classification for tree no. 6: B
183
184 Classification for tree no. 7: B
185
186 Classification for tree no. 8: B
187
188 Classification for tree no. 9: B
189
190 Classification for tree no. 10: B
191

```

```

192 Class B was detected 10 time(s).
193
194 Class NB was detected 0 time(s).
195
196 The most voted class was: B
197 NIL

```

El resultado de la validación cruzada y la votación se muestra solamente para el caso del banco por las mismas razones ya mencionadas anteriormente acerca de la longitud de las salidas en las otras bases de datos.

Finalmente, se repite lo mismo para el caso de la base de datos del juego de gato tomando como ejemplo `(x o o x o x b o x negative)`:

```

1  * (asdf:load-system :cl-id3)
2  T
3
4  * (in-package :cl-id3)
5  #<PACKAGE "CL-ID3">
6
7  * (load-file "./databases/bankruptcy.arff")
8  The ID3 setting has been reset.
9  Training set initialized after "./databases/bankruptcy.arff".
10 NIL
11
12 * (cross-validation 10)
13 ; output removed
14
15 * (let-them-vote '(x o o x o x b o x negative))
16
17 Classification for tree no. 1: NEGATIVE
18
19 Classification for tree no. 2: NEGATIVE
20
21 Classification for tree no. 3: NEGATIVE
22
23 Classification for tree no. 4: NEGATIVE
24
25 Classification for tree no. 5: NEGATIVE
26
27 Classification for tree no. 6: NEGATIVE
28
29 Classification for tree no. 7: NEGATIVE
30
31 Classification for tree no. 8: NEGATIVE
32
33 Classification for tree no. 9: NEGATIVE
34
35 Classification for tree no. 10: NEGATIVE
36
37 Class NEGATIVE was detected 10 time(s).
38

```

```

39 Class POSITIVE was detected 0 time(s).
40
41 The most voted class was: NEGATIVE
42 NIL

```

Para terminar la actividad hace falta implementar una función que convierta un árbol, como el mostrado para el caso del banco, a su equivalente en `prolog`. Para esto se añadió un archivo extra al conjunto de archivos que conforman el paquete `:cl-id3`, este archivo fue llamado `cl-id3-prolog-tree.lisp`.

La parte importante para realizar esto es generar un conjunto de ramas

```

1  ;; Convierte un arbol generado con el paquete en su respectivo conjunto de ramas.
2  (defun tree-to-branches (tree)
3    (if (leaf-p tree)
4        (list (list tree))
5        (mapcan (lambda (node)
6                  (mapcar (lambda (path)
7                          (cons (root tree) path))
8                          (tree-to-branches node)))
9                  (children tree))))
10
11  ;; Funcion auxiliar para la conversion de un arbol en lisp a clausulas en
12  ↪ prolog.
13  (defun prolog-tree-aux (tree filename)
14    (let ((branches (tree-to-branches tree)))
15      (dribble filename)
16      (loop for branch in branches
17            do (format t "branch(~a, [" (string-downcase (car (last branch))))
18            do (let ((attributes (butlast *attributes*)))
19                  (dotimes (n (length attributes))
20                    (if (member (nth n attributes) branch)
21                        (when t
22                          (downcase-attr (nth n attributes))
23                          (downcase-value (nth (+ (position (nth n attributes) branch)
24                                                         ↪ :test #'equal) 1) branch)))
24                    (print-empty))
25                    (if (eql n (- (length attributes) 1))
26                        (print-dot)
27                        (print-comma))))
28            do (print-newline))
29      (dribble)))
30
31  ;; Convierte un arbol en lisp a prolog. Guarda el arbol en el directorio
32  ↪ /home/usr
33  ;; por defecto.
34  (defun prolog-tree (ramas &optional (filename "~/prolog_tree.pl"))
35    (if (not (probe-file filename))
36        (prolog-tree-aux ramas filename)
37        (format t "File already exists.")))

```

Por ejemplo, las ramas del árbol que corresponde a la base de datos del banco se vería

como se muestra a continuación:

```

1  * (asdf:load-system :cl-id3)
2  T
3
4  * (in-package :cl-id3)
5  #<PACKAGE "CL-ID3">
6
7  * (load-file "./databases/bankruptcy.arff")
8  The ID3 setting has been reset.
9  Training set initialized after "./databases/bankruptcy.arff".
10 NIL
11
12 * (induce)
13 (CO (P NB) (N B) (A (CR (N (FF (A NB) (N B))) (P NB) (A NB))))
14
15 * (tree-to-branches (induce))
16 ((CO P NB) (CO N B) (CO A CR N FF A NB) (CO A CR N FF N B) (CO A CR P NB)
17  (CO A CR A NB))

```

Ya al tener todas las ramas solo es necesario convertir esto a su equivalente en `prolog` usando la función `prolog-tree` :

```

1  * (prolog-tree (induce) "~/prolog_tree_bank.pl")
2  branch(nb, [_ , _ , _ , _ , co/p, _]).
3  branch(b, [_ , _ , _ , _ , co/n, _]).
4  branch(nb, [_ , _ , ff/a, cr/n, co/a, _]).
5  branch(b, [_ , _ , ff/n, cr/n, co/a, _]).
6  branch(nb, [_ , _ , _ , cr/p, co/a, _]).
7  branch(nb, [_ , _ , _ , cr/a, co/a, _]).

```

No es necesario computar las ramas de antemano ya que no se guardan en ninguna variable global, éstas son computadas internamente en la función auxiliar `prolog-tree-aux` en una variable interna `branches` . La función `prolog-tree` envuelve a su contraparte auxiliar solamente para verificar que no existe ya un archivo en el camino especificado por el `usr` con el mismo nombre, si no existe entonces se procede a guardar el resultado.

Los árboles equivalentes en `prolog` de las otras dos bases de datos se incluyen en la misma carpeta correspondiente a los archivos de salida con los siguientes nombres:

- `/outputfiles/prolog_tree_car.pl` ,
- `/outputfiles/prolog_tree_tictactoe.pl` ,
- `/outputfiles/prolog_tree_bank.pl` .

Ahora se deben comparar los resultados de la representación del árbol en ambos lenguajes. Cargando el archivo donde se guardó el árbol equivalente de `prolog` de la base de datos del banco en `prolog` y probando con el siguiente ejemplo (corresponde a la línea 88 del archivo `bankruptcy.arff` cuyo valor de clase es `nb`),

- `sbcl : (a p a p a p) ,`
- `prolog : [ir/a, mr/p, ff/a, cr/p, co/a, op/p] .`

Después de haber cargado el archivo correspondiente en `sbcl` , el resultado es el siguiente:

```
1 * (setq *current-tree* (induce))
2 (CO (P NB) (N B) (A (CR (N (FF (A NB) (N B)))) (P NB) (A NB))))
3
4 * (classify-new-instance '(a p a p a p) *current-tree*)
5 NB
```

De manera similar, en `prolog` :

```
1 ?- [prolog_tree_bank].
2 true.
3
4 ?- branch(X, [ir/a, mr/p, ff/a, cr/p, co/a, op/p]).
5 X = nb.
```

Probando lo mismo ahora con la base de datos del juego de gato con el siguiente ejemplo (línea 520 del con valor de clase `positive`)

- `sbcl : (b x o o x b o x x) ,`
- `prolog : [top-left/b, top-middle/x, top-right/o, middle-left/o,`
`middle-middle/x, middle-right/b, bottom-left/o, bottom-middle/x,`
`bottom-right/x] .`

El resultado, después de haber cargado la base de datos del juego de gato en `sbcl` , y asignado el resultado de la inducción a la variable `*current-tree*` se obtiene lo siguiente:

```
1 * (classify-new-instance '(b x o o x b o x x) *current-tree*)
2 POSITIVE
```

Mientras que el resultado equivalente en `prolog` es:

```
1 ?- [prolog_tree_tictactoe].
2 true.
3
4 ?- branch(X, [top-left/b, top-middle/x, top-right/o, middle-left/o,
5 ↪ middle-middle/x, middle-right/b, bottom-left/o, bottom-middle/x,
6 ↪ bottom-right/x]).
7 X = positive.
```

Para concluir, se muestra lo equivalente para el archivo de la base de datos de los autos tomando el siguiente ejemplo (línea 1420 del con valor de clase `unacc`):

- `sbcl : (low high 2 2 med high) ,`
- `prolog : [buying/low, maint/high, doors/2, persons/2, lug_boot/med, safety/high] .`

El resultado para este último caso, después de haber cargado la base de datos de los autos en `sbcl`, y asignado el resultado de la inducción a la variable `*current-tree*` se obtiene lo siguiente:

```
1 * (classify-new-instance '(low high 2 2 med high) *current-tree*)
2 UNACC
```

Y su equivalente resultado en `prolog`:

```
1 ?- [prolog_tree_car].
2 true.
3
4 ?- branch(X, [buying/low, maint/high, doors/2, persons/2, lug_boot/med,
5 ↪ safety/high]).
X = unacc.
```

Referencias

- [1] Alejandro Guerra-Hernandez. Programación para la inteligencia artificial. <https://www.uv.mx/personal/aguerra/pia/>, 2022. Visitado: 2022-10-07.
- [2] Peter Seibel. *Practical common lisp*. Apress, 2006.
- [3] Dheeru Dua and Casey Graff. Uci machine learning repository. <http://archive.ics.uci.edu/ml>, 2017.