

PROGRAMACIÓN LÓGICA Y FUNCIONAL

Tarea 5. LISP en IA (ID3)

Angel García Báez
ZS24019400@estudiantes.uv.mx

Maestría en Inteligencia Artificial

IIIA Instituto de Investigaciones en Inteligencia Artificial
Universidad Veracruzana
Campus Sur, Calle Paseo Lote II, Sección 2a, No 112
Nuevo Xalapa, Xalapa, Ver., México 91097

16 de diciembre de 2024

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

- Car Evaluation
- Tic-Tac-Toe
- Zoo

llevar a cabo las siguientes actividades:

1. Implementar en Lisp las mejoras al algoritmo ID3 que adoptaron en su implementación en Prolog [65/100].
2. . Describa brevemente (máximo una página) las diferencias y similitudes del uso de la programación lógica (prolog) y funcional (lisp) para implementar ID3. [35/100]

1. Implementación de las mejoras en el ID3 (Split information y gain ratio)

Para realizar las mejoras al ID3 básico fue necesario retomar lo propuesto por Quinlan (1986) en su artículo Induction of Decision Trees con la finalidad de entender la estructura base del árbol así como de las oportunidades de mejora que puede tener la implementación inicial del mismo.

Después de revisar el artículo junto con las notas de clase de Guerra-Hernández (2024) y retomando la implementación hecha en prolog, se optó por mejorar el algoritmo base mediante el criterio la split information y el gain ratio. Estos criterios se aplican después de calcular la entropía y la ganancia de información de tal forma que se pueda penalizar a aquellos atributos que tengan muchas clases y que no estén aportando realmente a realizar la correcta clasificación de las clases objetivo.

Las fórmulas necesarias se expresan como siguen:

Entropía de la información por atributo:

$$\text{entropia}(S) = - \sum_{i=1}^k p_i \cdot \log_2(p_i)$$

Donde:

$$p_i = \frac{|S_i|}{|S|}$$

Ganancia de información:

$$\text{gain}(S, A) = \text{entropia}(S) - \sum \frac{|S_v|}{|S|} \cdot \text{entropia}(S_v)$$

División de la información:

$$\text{splitInformation}(S, A) = \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

Ratio de información:

$$\text{gainRatio}(S, A) = \frac{\text{gain}(S, A)}{\text{splitInformation}(S, A)}$$

Una vez identificada la mejora y recordados los pasos dados en prolog, se procede a buscar la manera de implementar esto en LISP.

Primero se identificaron las funciones del código que se encargan de calcular la entropía y la ganancia de información dentro del código.

A continuación se muestra la función de entropía dentro del código, la cual no fue cambiada dado que es necesaria para poder hacer el cálculo del gain ratio:

```
1  ;;; Entropia ;;;
2  (defun entropy (examples attrib)
3    "It computes the entropy of EXAMPLES with respect to an ATTRIB"
4    (let ((partition (get-partition attrib examples))
5          (number-of-examples (length examples)))
6      (apply #'(+
7              (mapcar #'(lambda(part)
8                          (let* ((size-part (count-if #'atom
9                                                         (cdr part)))
10                             (proportion
11                              (if (eq size-part 0) 0
12                                  (/ size-part
13                                     number-of-examples))))
13              (* -1.0 proportion (log proportion 2))))
14              (cdr partition)))))
15
```

A continuación se muestra la función de ganancia de información básica que de igual forma, no fue modificada dado que es necesaria:

```
1  ;;; Ganancia de la información ;;;
2  (defun information-gain (examples attribute)
3    "It computes information-gain for an ATTRIBUTE in EXAMPLES"
4    (let ((parts (get-partition attribute examples))
5          (no-examples (count-if #'atom examples)))
6      (- (entropy examples *target*)
7         (apply #'(+
8                 (mapcar
9                  #'(lambda(part)
10                      (let* ((size-part (count-if #'atom
11                                                         (cdr part)))
12                         (proportion (if (eq size-part 0) 0
13                                           (/ size-part
14                                              no-examples))))
13                  (* proportion (entropy (cdr part) *target*))))
14                 (cdr parts)))))
15
16
```

Aquí es donde se añaden las mejoras al código, comenzando por añadir el cálculo del split-information al código principal como sigue:

```
1  ;;; Split Information
2  (defun split-information (examples attribute)
3    "It computes the split information for an ATTRIBUTE in EXAMPLES"
4    (let ((partition (get-partition attribute examples))
5          (no-examples (length examples)))
6      (apply #'(+
7              (mapcar #'(lambda (part)
8                          (let* ((size-part (length (cdr part)))
9                                (proportion (if (zerop size-part) 0
11                                                (/ size-part no-examples))))
12                            (if (zerop proportion) 0
13                                (* -1.0 proportion (log proportion 2))))
14                          (cdr partition)))))
```

La siguiente mejora necesaria es añadir el cálculo del Gain Ratio como sigue:

```
1  ;;; Gain Ratio
2  (defun gain-ratio (examples attribute)
3    "It computes the gain ratio for an ATTRIBUTE in EXAMPLES"
4    (let ((info-gain (information-gain examples attribute))
5          (split-info (split-information examples attribute)))
6      (if (zerop split-info)
7          0
8          (/ info-gain split-info))))
```

Una vez que se tienen todas las funciones necesarias, es necesario modificar la forma en que el algoritmo determina la mejor rama del árbol para que se base en el gain ratio. Dicho código cambia la selección basada en la ganancia de información para que en su lugar, considere al atributo que maximiza el gain ratio en cada iteración. El resultado se presenta a continuación como modificación de la función base best-partition:

```

1  ;;; Mejor partición basada en Gain Ratio
2  (defun best-partition (attributes examples)
3    "Computar la mejor partición en base al gain ratio"
4    (let* ((gain-ratios
5            (loop for attrib in attributes collect
6                  (let ((gr (gain-ratio examples attrib))
7                        (p (get-partition attrib examples)))
8                    (when *trace*
9                      (format t "Partición inducida por el atributo ~s:~%~s~%"
10                             attrib p)
11                      (format t "Gain Ratio: ~s~%" gr))
12                    (list gr p))))
13          (best (cadar (sort gain-ratios #'(lambda (x y) (> (car x) (car y)))))))
14    (when *trace* (format t "Best partition: ~s~%-----~%" best))
15    best))

```

Referencias

- Guerra-Hernández, D. A. (2024). Programación Lógica y Funcional - Programación Lógica.
- Quinlan, J. R. (1986). Induction of decision trees. *Machine learning*, 1:81–106.