

## Árbol de clasificación

Dado un conjunto  $D$  de datos etiquetados el problema es construir un árbol que permita mediante preguntas sobre los atributos permitir asignar clase a un patrón, dato no etiquetado. La idea es conseguir lo anterior con el mínimo número de preguntas, con la mayor generalidad posible.

Para lograr lo anterior en cada etapa de la construcción del árbol se elige para la pregunta al atributo que maximiza la ganancia de información dada por

$G(H(D), a) = H(D) - \sum P_i H(D_i)$ ,  $D = \cup D_i$ ,  $D_i$  forman una partición de  $D$  preguntando por los valores del atributo  $a$ ,

$P_i = \frac{|D_i|}{|D|}$  y  $H(D) = - \sum_{k \in C} P_k \log(P_k)$ ,  $k$  índice de clase.

$H(D)$  es la entropía de los datos de acuerdo a la clase a que pertenecen, cuando todos son de la misma clase la entropía es cero es decir están totalmente en orden de acuerdo a las clases. De acuerdo a la anterior Euristicas el algoritmo de entrenamiento funciona como sigue:

create-Branch (Datos):

    chechar si cada dato es de la misma clase:

        Si es el caso regresa la etiqueta de clase

    Else

        Encontrar el mejor atributo para partir los datos

        Partir los datos mediante la pregunta

        crear nodo correspondiente

        Para cada partición  $D_i$

            create-Branch ( $D_i$ ) y conectar el resultado

        regresa el nodo creado