

Content

- Introduction
- SPRINT Algorithm
- PUBLIC Algorithm
- Training Dataset
- Experimental Results
- Future Work

Advantage of Decision Tree

- Classification Techniques
- -Decision trees
- -Bayesian classification
- -Neural networks ...
- Advantage of Decision Tree
 -Easily comprehended by humans.
 - -Efficient and thus suitable for large training sets.
 - -Do not require additional information

Building a Decision Tree

• Building Phase

- The tree is built by recursively partitioning the dataset until each partition is "pure".
- Pruning Phase
- After fully grown, the tree is pruned to remove noises that may be particular only to the training datasets .

SPRINT and PUBLIC

- Using large training datasets, we can improve the accuracy of classification model
- Traditional decision tree algorithms(ID3 and C4.5) are established for small datasets
- SPRINT : Building Phase+Pruning Phase
- PUBLIC : Integrating the second "Pruning Phase" with the initial "Building Phase"



Attribute Lists

- One **attribute list** is for one attribute in the datasets
- Each record in the attribute list consists of an **attribute value**, a **class label** and the **record ID**.

rid	Age	Car Type	Risk High	
0	23	family		
1	17	sports	High	
2	43	sports	High	
3	68	family	Low	
4	32	truck	Low	
5	20	family	High	

Age	Class	nid
17	High	1
20	High	5
23	High	0
32	Low	4
43	High	2
68	Low	3





Pruning Principle

MDL(Minimum Description Length) principle states that the "**best**" tree is the one that can be encoded using the **fewest** number of bits.

- Cost of Encoding Tree
- -Encoding the structure of the tree
- -Encoding each split
 - -Encoding the classes of data records in each leaf

Encoding the Structure

Using a single bit to specify whether a node of the tree is an internal node(1) or leaf(0).



Encoding the Tree(Cont')

Let a set S contain n records each belonging to one of k classes, n_i being the number of records with class i, let " α " be the number of attributes and let v be the number of distinct values for the splitting attribute in records at the node.

 The cost of encoding each split C_{split}(N) is Log α +Log(v-1) or Log α+Log(2^v-2)

• The cost of encoding the classes of data records

 $C(S) = \sum n_i \log(n/n_i) + (k-1)\log(n/2) + \log(pi^{k/2}/\Gamma(k/2))$



Experimental Results										
Table 2 Execution Time (secs)										
Data Set	Car	letter	satimage	shuttle	market					
 SPRINT	45.58	3283.2	1471	457.78	21.14					
 PUBLIC	33.93	2786.57	1036.34	455.15	19.03					
Max Ratio	38%	18%	43%	0.60%	11%					

Future Work

- How often should PUBLIC to prune the partially built tree?
- How can we estimate the subtree cost more accurately?
- Using large training sets, we may get much better performance.