



Decision Tree Construction in Large Datasets

Classification in Large Databases

- ❑ Why is decision tree induction popular?
 - ❑ Relatively fast learning speed
 - ❑ Convertible to simple and easy to understand classification rules
 - ❑ Easy to be adapted to database system implementations (e.g., using SQL)
 - ❑ Comparable classification accuracy with other methods
- ❑ Scalability: Classifying data sets with millions of examples and hundreds of attributes with reasonable speed
- ❑ **RainForest** (VLDB'98 - Gehrke, Ramakrishnan & Ganti)
 - ❑ Builds an AVC-list (attribute, value, class label)

RainForest: A Scalable Classification Framework

- The criteria that determine the quality of the tree can be computed separately
 - Builds an AVC-list: **AVC (Attribute, Value, Class_label)**
- **AVC-set** (of an attribute X)
 - Projection of training dataset onto the attribute X and class label, where counts of individual class label are aggregated

- **AVC-group** (of a node n)

- Set of AVC-sets of all predictor attributes at the node n

age	income	student	credit_rating	comp
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

The Training Data

AVC-set on Age

Age	Buy_Computer	
	yes	no
<=30	2	3
31..40	4	0
>40	3	2

AVC-set on Income

income	Buy_Computer	
	yes	no
high	2	2
medium	4	2
low	3	1

AVC-set on Student

student	Buy_Computer	
	yes	no
yes	6	1
no	3	4

AVC-set on Credit_Rating

Credit rating	Buy_Computer	
	yes	no
fair	6	2
excellent	3	3

Its AVC Sets