# Data Mining -- Classification

Instructor: Jen-Wei Huang

Office: 92528 in the EE building
jwhuang@mail.ncku

# Classification

▸ Predicts categorical class labels (discrete or nominal)

▸ Classifies data (constructs a model) based on the training set and the values (class labels) in a classifying attribute and uses it in classifying new data

▸ Typical applications
  ◦ Credit approval
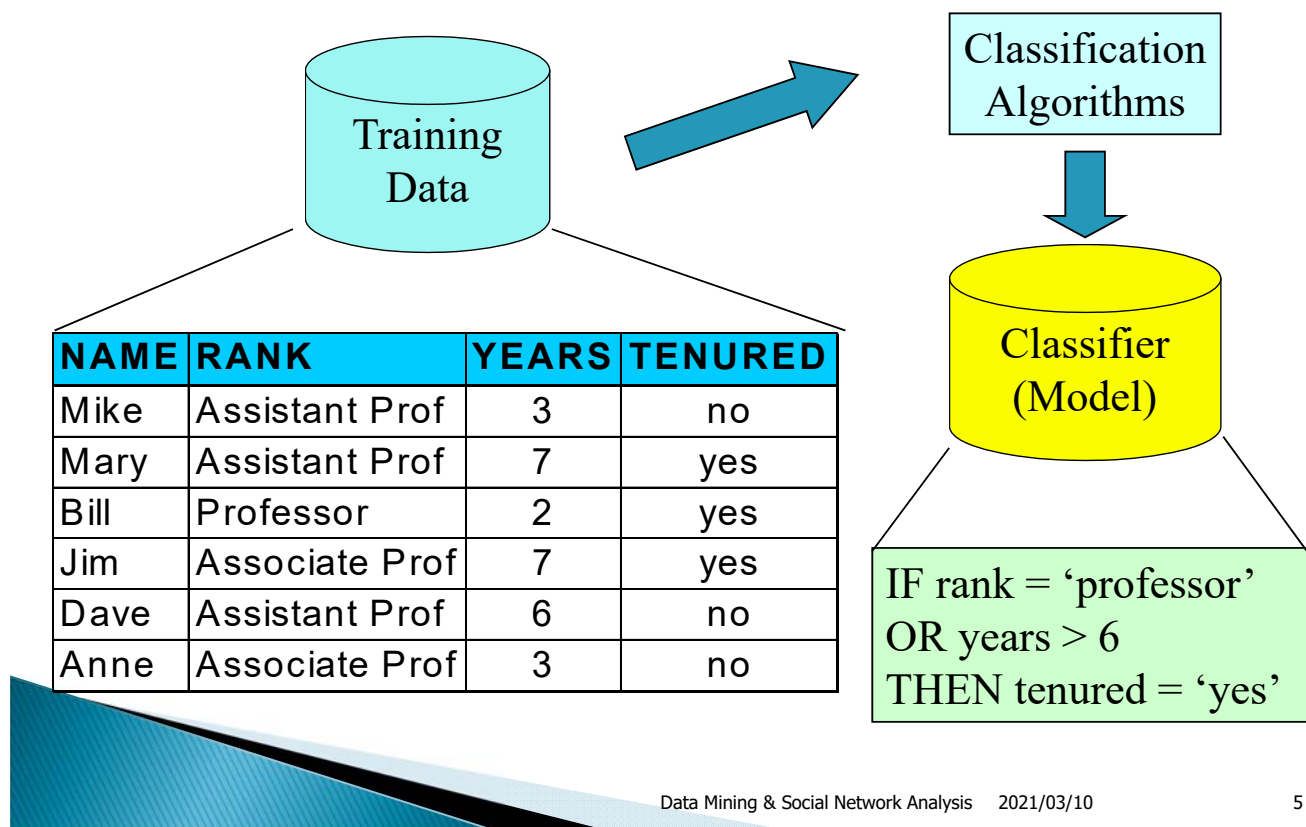  ◦ Target marketing
  ◦ Medical diagnosis
  ◦ Fraud detection

# Two-Step Process

- **Model construction**: describing a set of predetermined classes
  - Each tuple/sample is assumed to belong to a predefined class, as determined by the class label attribute
  - The set of tuples used for model construction is the training set
  - The model is represented as classification rules, decision trees, or mathematical formulae
- **Model usage**: for classifying future or unknown objects
  - Testing set is the set of tuples used to estimate the accuracy of the model
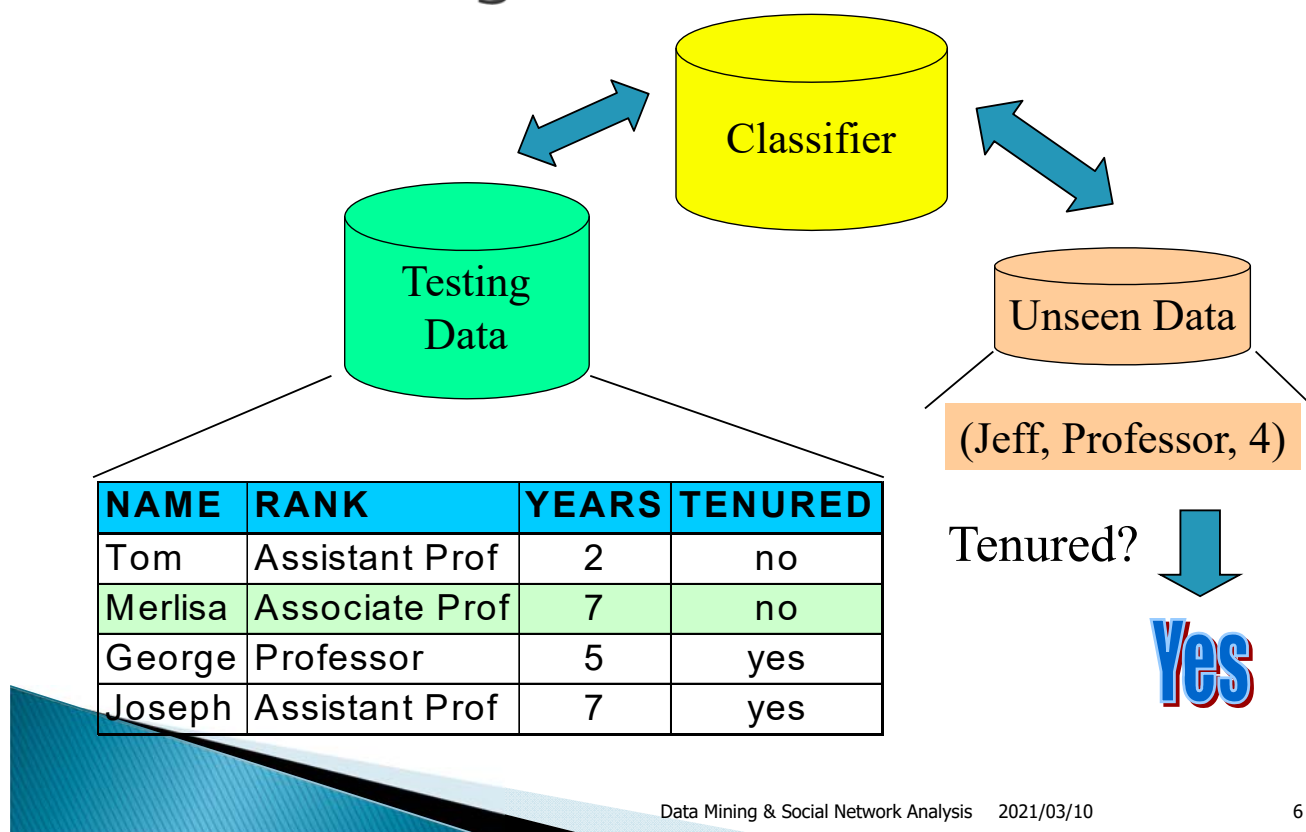  - If the accuracy is acceptable, use the model to classify data tuples whose class labels are not known

# Accuracy

- The known label of test sample is compared with the classified result from the model
- Accuracy rate is the percentage of test set samples that are correctly classified by the model
- Test set is independent of training set, otherwise over-fitting will occur

# Model Construction

Training Data

Classification Algorithms

Classifier (Model)

| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Mike | Assistant Prof | 3 | no |
| Mary | Assistant Prof | 7 | yes |
| Bill | Professor | 2 | yes |
| Jim | Associate Prof | 7 | yes |
| Dave | Assistant Prof | 6 | no |
| Anne | Associate Prof | 3 | no |

IF rank = 'professor'
OR years > 6
THEN tenured = 'yes'

# Model Usage

Classifier

Testing Data

Unseen Data

(Jeff, Professor, 4)

| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Tom | Assistant Prof | 2 | no |
| Merlisa | Associate Prof | 7 | no |
| George | Professor | 5 | yes |
| Joseph | Assistant Prof | 7 | yes |

Tenured?

**Yes**

# Supervised Learning

▸ The training data (observations, measurements, etc.) are accompanied by labels indicating the class of the observations

▸ New data is classified based on the training set

▸ Compared to unsupervised learning (clustering):
  ▸ The class labels of training data is unknown

# Date Preparation

▸ Data cleaning
  ◦ Preprocess data in order to reduce noise and handle missing values
▸ Relevance analysis (feature selection)
  ◦ Remove the irrelevant or redundant attributes
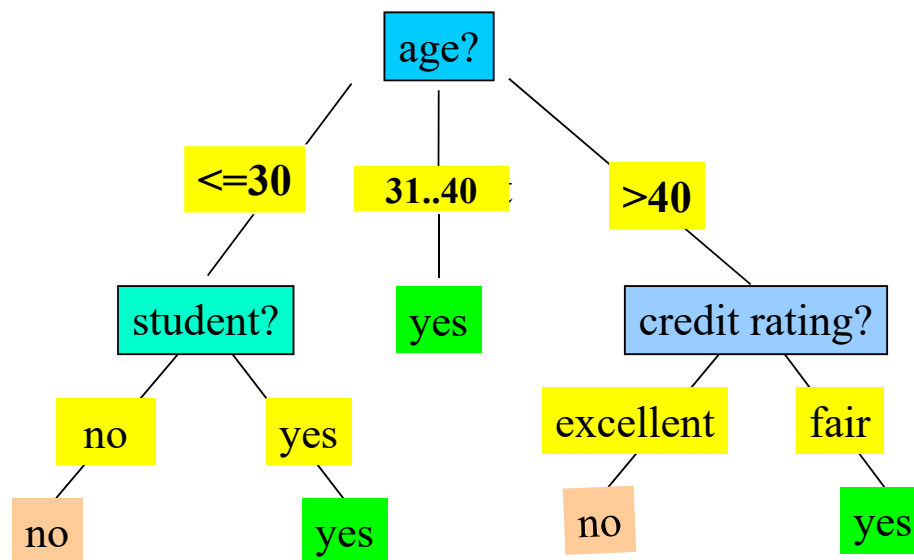▸ Data transformation
  ◦ Generalize and/or normalize data

# Evaluation

- Accuracy: the ratio of correctly predicted tuples in the testing set
- Speed
  - time to construct the model (training time)
  - time to use the model (classification time)
- Robustness: handling noise and missing values
- Scalability: efficiency in disk-resident databases
- Interpretability
  - understanding and insight provided by the model
- Other measures, e.g., goodness of rules, such as decision tree size or compactness of classification rules

# Example

| age | income | student | credit_rating | buys_computer |
|-----|--------|---------|---------------|---------------|
| <=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 |

# Decision Tree Model

```
                        age?
                      /   |   \
                 <=30   31..40   >40
                  /       |        \
            student?     yes    credit rating?
             /   \                /      \
           no    yes         excellent   fair
           |       \            |          \
          no       yes          no         yes
```

# Construction of Decision Tree

▸ Tree is constructed in a top-down recursive divide-and-conquer manner
▸ At start, all the training examples are at the root
▸ Attributes are categorical (if continuous-valued, they are discretized in advance)
▸ Examples are partitioned recursively based on selected attributes
▸ Test attributes are selected on the basis of a heuristic or statistical measure (e.g., information gain)

# Construction of Decision Tree

▸ Conditions for stopping partitioning
  ◦ All samples for a given node belong to the same class
  ◦ There are no remaining attributes for further partitioning – majority voting is employed for classifying the leaf
  ◦ There are no samples left
▸ Attribute selection measeures:
  ◦ Information Gain (ID3/C4.5)
  ◦ Gini index (CART, IBM IntelligentMiner)
  ◦ Others: CHAID, C-SEP, G-statistics, MDL
  ◦ Multivariate splits

# Information Gain (ID3/C4.5)

■ Select the attribute with the highest information gain
■ Let $p_i$ be the probability that an arbitrary tuple in D belongs to class $C_i$, estimated by $|C_{i,D}|/|D|$
■ Expected information (entropy) needed to classify a tuple in D:

$$Info(D) = -\sum_{i=1}^{m} p_i \log_2(p_i)$$

■ Information needed (after using A to split D into v partitions) to classify D:

$$Info_A(D) = \sum_{j=1}^{v} \frac{|D_j|}{|D|} \times I(D_j)$$

■ Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

$$Info_{age}(D) = \frac{5}{14}I(2,3) + \frac{4}{14}I(4,0)$$

$$Info(D) = I(9,5) = -\frac{9}{14}\log_2(\frac{9}{14}) - \frac{5}{14}\log_2(\frac{5}{14}) = 0.940$$

$$+\frac{5}{14}I(3,2) = 0.694$$

| age | $p_i$ | $n_i$ | $I(p_i, n_i)$ |
|-----|-------|-------|---------------|
| <=30 | 2 | 3 | 0.971 |
| 31...40 | 4 | 0 | 0 |
| >40 | 3 | 2 | 0.971 |

$\frac{5}{14}I(2,3)$ means "age <=30" has 5 out of 14 samples, with 2 yes'es and 3 no's. Hence,

$$Gain(age) = Info(D) - Info_{age}(D) = 0.246$$

| age | income | student | credit_rating | buys_computer |
|-----|--------|---------|---------------|---------------|
| <=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 |

Similarly,

$$Gain(income) = 0.029$$
$$Gain(student) = 0.151$$
$$Gain(credit\_rating) = 0.048$$

# Continuous-Value Attributes

▸ Must determine the *best split point* for continuous-value attribute, A

  ◦ Sort the value A in increasing order

  ◦ Typically, the midpoint between each pair of adjacent values is considered as a possible *split point*

    • $(a_i + a_{i+1})/2$ is the midpoint between the values of $a_i$ and $a_{i+1}$

  ◦ The point with the *minimum expected information requirement* for A is selected as the split-point for A

# Gain Ratio (C4.5)

- Information gain measure is biased towards attributes with a large number of values
- C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

$$SplitInfo_A(D) = -\sum_{j=1}^{v} \frac{|D_j|}{|D|} \times \log_2(\frac{|D_j|}{|D|})$$

- GainRatio(A) = Gain(A)/SplitInfo(A)
  - Ex. $SplitInfo_A(D) = -\frac{4}{14} \times \log_2(\frac{4}{14}) - \frac{6}{14} \times \log_2(\frac{6}{14}) - \frac{4}{14} \times \log_2(\frac{4}{14}) = 0.926$
  - gain_ratio(income) = 0.029/0.926 = 0.031
- The attribute with the maximum gain ratio is selected as the splitting attribute

# Gini index (CART, IntelligentMiner)

- If a data set $D$ contains examples from $n$ classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^{n} p_j^2$$

  where $p_j$ is the relative frequency of class $j$ in $D$

- If a data set $D$ is split on A into two subsets $D_1$ and $D_2$, the $gini$ index $gini(D)$ is defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

- Reduction in Impurity:

$$\Delta gini(A) = gini(D) - gini_A(D)$$

# Gini index (CART, IntelligentMiner)

- All attributes are assumed continuous-valued
- The attribute provides the smallest $gini_{split}(D)$ or the largest reduction in impurity is chosen to split the node (*need to enumerate all the possible splitting points for each attribute*)
- May need other tools, e.g., clustering, to get the possible split values
- Can be modified for categorical attributes

# Example (Gini Index)

- Ex.  D has 9 tuples in buys_computer = "yes" and 5 in "no"

$$gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

- Suppose the attribute income partitions D into 10 in $D_1$: {low, medium} and 4 in $D_2$

$$gini_{income \in \{low, medium\}}(D) = \left(\frac{10}{14}\right) Gini(D_1) + \left(\frac{4}{14}\right) Gini(D_1)$$
$$= \frac{10}{14}(1 - (\frac{6}{10})^2 - (\frac{4}{10})^2) + \frac{4}{14}(1 - (\frac{1}{4})^2 - (\frac{3}{4})^2)$$
$$= 0.450$$
$$= Gini_{income \in \{high\}}(D)$$

    but $gini_{\{medium, high\}}$ is 0.30 and thus the best since it is the lowest

# Comparison

▸ The three measures return good results but
  ◦ Information gain:
    • biased towards multivalued attributes
  ◦ Gain ratio:
    • tends to prefer unbalanced splits in which one partition is much smaller than the others
  ◦ Gini index:
    • biased to multivalued attributes
    • has difficulty when # of classes is large
    • tends to favor tests that result in equal-sized partitions and purity in both partitions

# Other Attribute Selection Measures

▸ CHAID: a popular decision tree algorithm, measure based on $\chi^2$ test for independence
▸ C-SEP: performs better than info. gain and gini index in certain cases
▸ G-statistics: has a close approximation to $\chi^2$ distribution
▸ MDL (Minimal Description Length) principle (i.e., the simplest solution is preferred):
  ◦ The best tree as the one that requires the fewest # of bits to both (1) encode the tree, and (2) encode the exceptions to the tree

# Other Attribute Selection Measures

- Multivariate splits (partition based on multiple variable combinations)
  - CART: finds multivariate splits based on a linear comb. of attrs.
- Which attribute selection measure is the best?
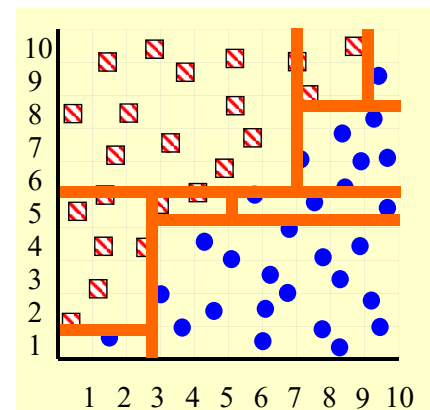  - Most give good results, none is significantly superior than others

---

# Large Databases

- Scalability: Classifying data sets with millions of examples and hundreds of attributes with reasonable speed
- Why decision tree induction in data mining?
  - relatively faster learning speed (than other classification methods)
  - convertible to simple and easy to understand classification rules
  - can use SQL queries for accessing databases
  - comparable classification accuracy with other methods
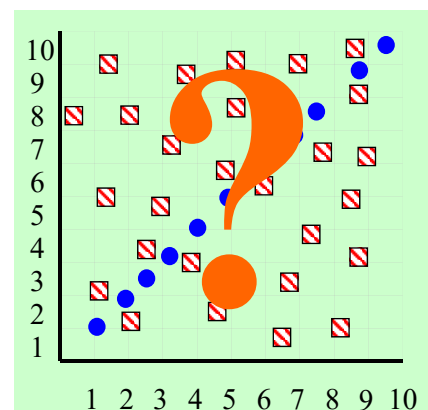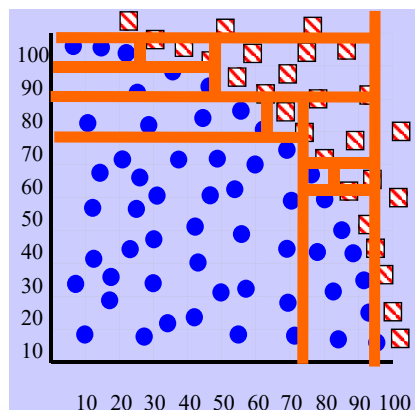
# Scalable Decision Tree Induction

- **SLIQ** (EDBT'96 — Mehta *et al.*)
  - ◦ Builds an index for each attribute and only class list and the current attribute list reside in memory
- **SPRINT** (VLDB'96 — J. Shafer et al.)
  - ◦ Constructs an attribute list data structure
- **PUBLIC** (VLDB'98 — Rastogi *et al.*)
  - ◦ Integrates tree splitting and tree pruning: stop growing the tree earlier
- **RainForest** (VLDB'98 — Gehrke *et al.*)
  - ◦ Builds an AVC-list (attribute, value, class label)
- **BOAT** (PODS'99 — Gehrke *et al.*)
  - ◦ Uses bootstrapping to create several small samples

# Correlated Features

1) Deep Bushy Tree
2) Useless
3) Deep Bushy Tree

The Decision Tree has a hard time with correlated attributes

# Overfitting Problem [1]

- A model describes random error or noise instead of the underlying relationship.
- Overfitting occurs when a model is excessively complex, such as having too many degrees of freedom, in relation to the amount of data available.
- Sometimes overfitting depends on the conformability of the model structure with the data shape, and the magnitude of model error compared to the expected level of noise or error in the data.

# Overfitting Problem

- An induced tree may overfit the training data
  - Too many branches, some may reflect anomalies due to noise or outliers
  - Poor accuracy for unseen samples
- Two approaches to avoid overfitting
  - Prepruning: Halt tree construction early—do not split a node if this would result in the goodness measure falling below a threshold
    - Difficult to choose an appropriate threshold
  - Postpruning: Remove branches from a "fully grown" tree—get a sequence of progressively pruned trees

# Deep Learning

- https://dotblogs.com.tw/allanyiin/2016/03/12/222215

# References

- [1] http://en.wikipedia.org/wiki/Overfitting
- C. Apte and S. Weiss. **Data mining with decision trees and decision rules**. Future Generation Computer Systems, 13, 1997.
- R. O. Duda, P. E. Hart, and D. G. Stork. **Pattern Classification**, 2ed. John Wiley and Sons, 2001
- U. M. Fayyad. **Branching on attribute values in decision tree generation**. AAAI'94.
- M. Kamber, L. Winstone,  W. Gong,  S. Cheng, and J. Han. **Generalization and decision tree induction: Efficient classification in data mining**. RIDE'97.

# References

- T.-S. Lim, W.-Y. Loh, and Y.-S. Shih. **A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms.** Machine Learning, 2000.

- S. K. Murthy, **Automatic Construction of Decision Trees from Data: A Multi-Disciplinary Survey**, Data Mining and Knowledge Discovery 2(4): 345-389, 1998

- J. R. Quinlan. **Induction of decision trees.** *Machine Learning*, 1:81-106, 1986.

- J. R. Quinlan. **C4.5: Programs for Machine Learning**. Morgan Kaufmann, 1993.

- J. R. Quinlan. **Bagging, boosting, and c4.5**. AAAI'96

# References

- R. Rastogi and K. Shim. **Public: A decision tree classifier that integrates building and pruning**. VLDB'98.

- J. Shafer, R. Agrawal, and M. Mehta. **SPRINT : A scalable parallel classifier for data mining**. VLDB'96.

- X. Yin and J. Han. **CPAR: Classification based on predictive association rules**. SDM'03

- J. Gehrke, R. Ramakrishnan, and V. Ganti. **Rainforest: A framework for fast decision tree construction of large datasets**. VLDB'98.

- J. Gehrke, V. Gant, R. Ramakrishnan, and W.-Y. Loh, **BOAT -- Optimistic Decision Tree Construction**. SIGMOD'99*.

- Y. Freund and R. E. Schapire. **A decision-theoretic generalization of on-line learning and an application to boosting**. J. Computer and System Sciences, 1997.

# References

- Slides from Prof. J.-W. Han, UIUC
- Slides from Prof. M.-S. Chen, NTU
- Slides from Prof. W.-Z. Peng, NCTU