



Classification and Prediction

——Basic Concepts——

徐华

清华大学 计算机系 智能技术与系统国家重点实验室

xuhua@tsinghua.edu.cn

1

Classification and Prediction



- ◉ **Basic Concepts**
- ◉ **Issues Regarding Classification and Prediction**
- ◉ **Decision Tree**
- ◉ **Bayesian Classification**
- ◉ **Neural Networks**
- ◉ **Support Vector Machine**
- ◉ **K-Nearest Neighbor**
- ◉ **Associative classification**
- ◉ **Classification Accuracy**

2



Classification vs. Prediction



- ◉ **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
- ◉ **Prediction:**
 - ◆ models continuous-valued functions, i.e., predicts unknown or missing values
- ◉ **Typical Applications**
 - ◆ Credit approval
 - ◆ Target marketing
 - ◆ Medical diagnosis
 - ◆ Fraud detection

3



Classification—A Two-Step Process

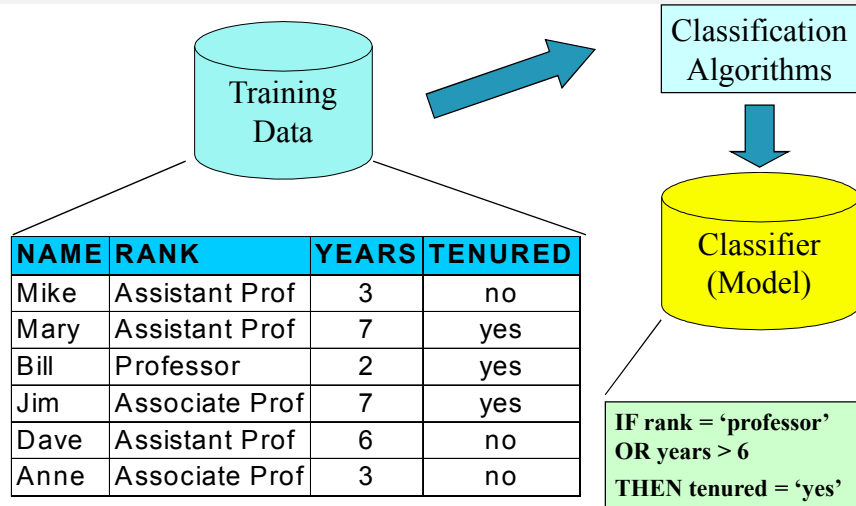


- ◉ **Model construction:** describing a set of predetermined classes
 - ◆ Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label attribute**
 - ◆ The set of tuples used for model construction is **training set**
 - ◆ The model is represented as classification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

4

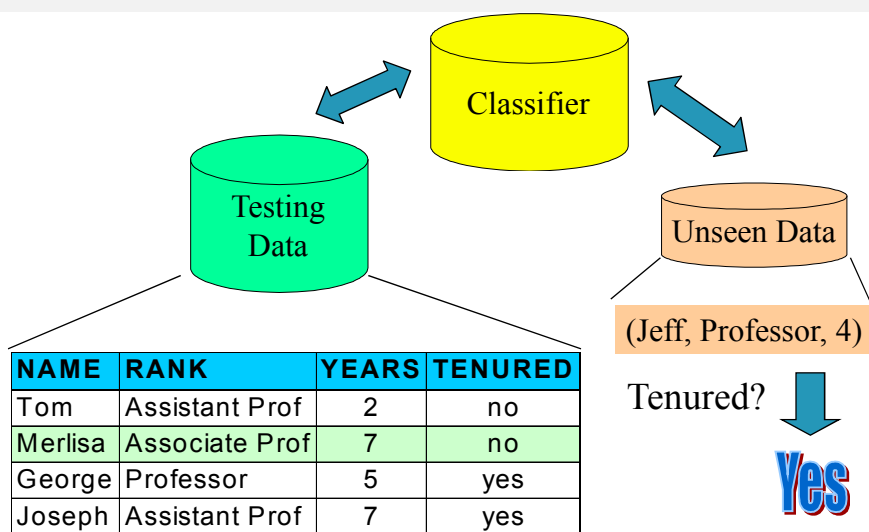


Classification Process (1): Model Construction



5

Classification Process (2): Use the Model in Prediction



6

Supervised vs. Unsupervised Learning



- ◉ **Supervised learning (classification)**
 - ◆ Supervision: The training data (observations, measurements, etc.) are accompanied by labels indicating the class of the observations
 - ◆ New data are classified based on the training set
- ◉ **Unsupervised learning (clustering)**
 - ◆ The class labels of training data is unknown
 - ◆ Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

7



Thanks !

8

