



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

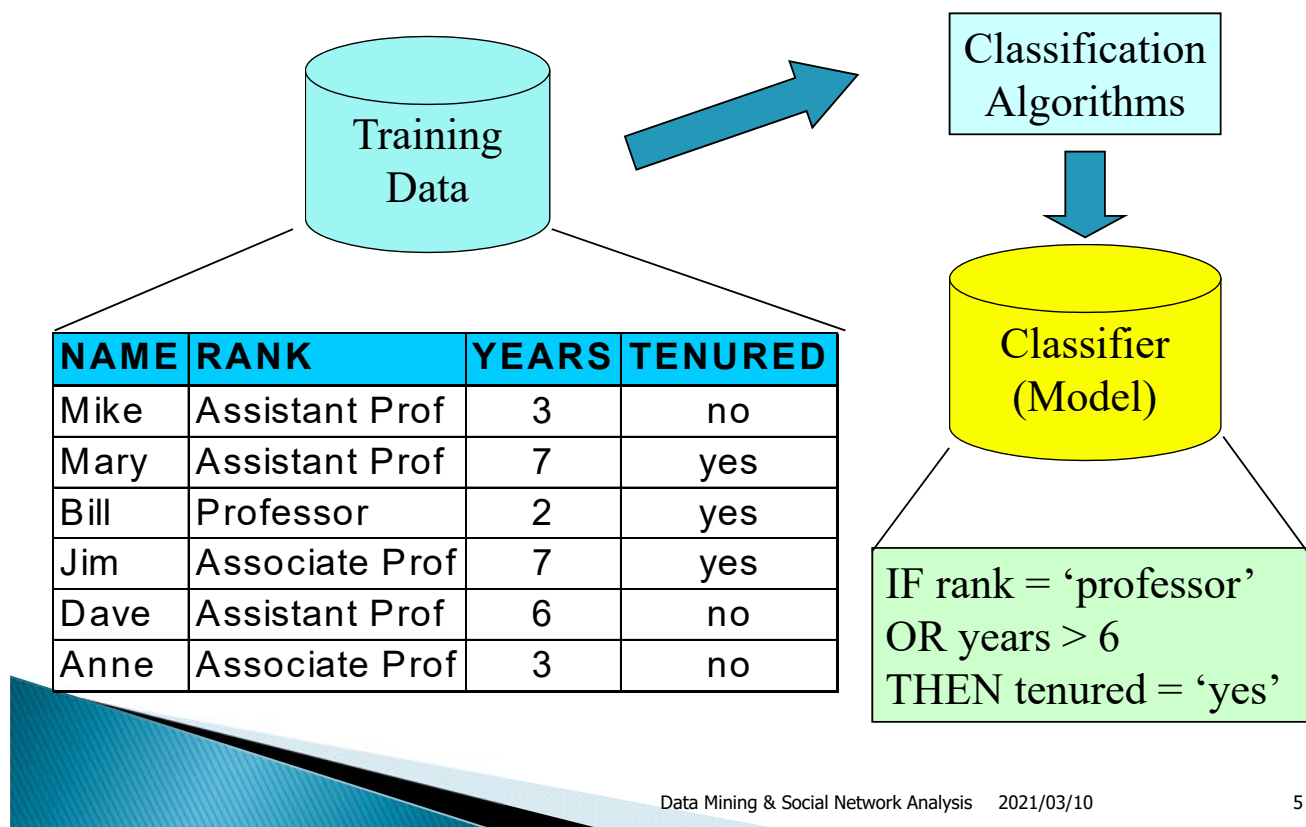
Two-Step Process

- ▶ **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

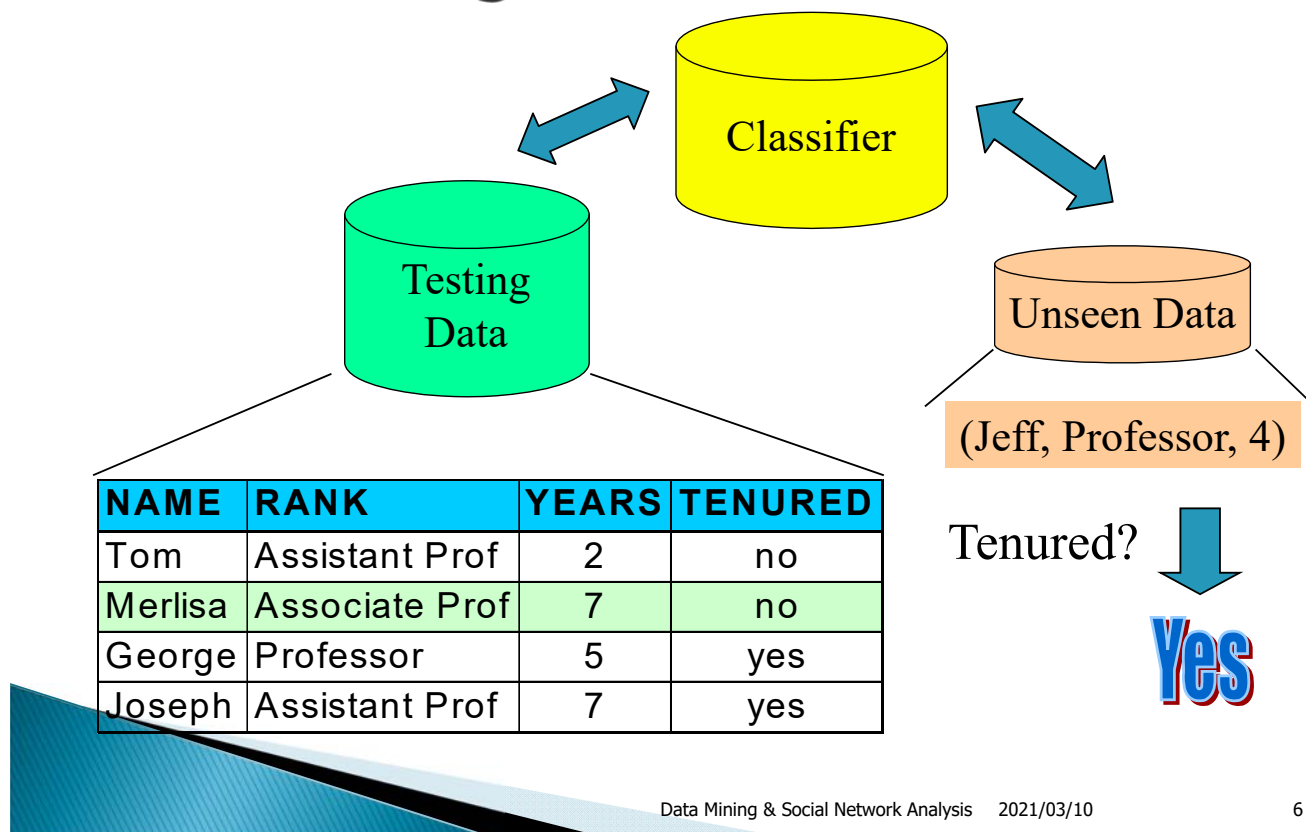
Accuracy

- ▶ 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



Model Usage



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

Data 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

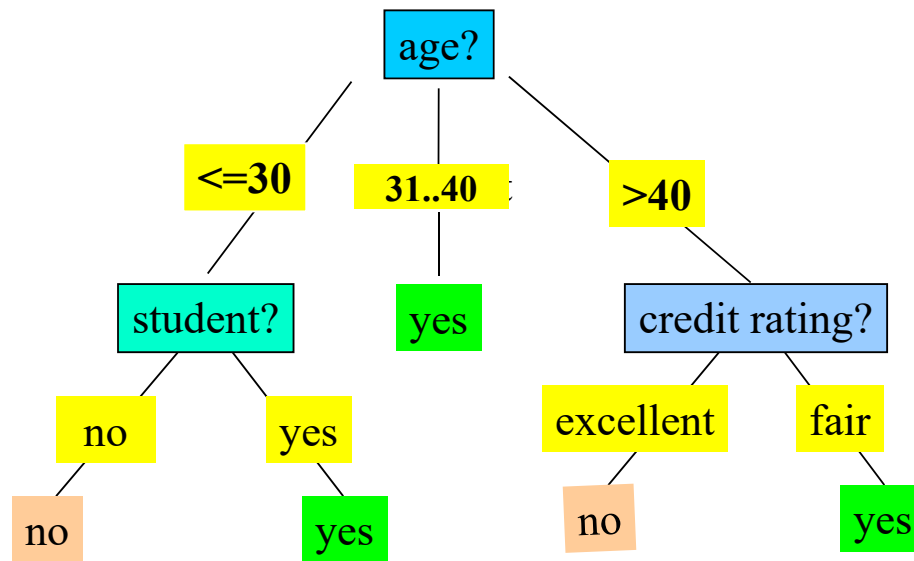
Evaluation

- ▶ 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



References

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