



Data Mining: Classification

Classification vs. Prediction



■ Classification:

- estimates categorical class labels
- 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
- treatment effectiveness analysis

Classification Scenario



A model or classifier is constructed to predict the categorical labels.

- A bank loan officer wants to analyze the data in order to know which customer (loan applicant) are risky or which are safe.
 - Conclusion: risky or safe
- A marketing manager at a company needs to analyze a customer with a given profile, who will buy a new computer.
 - Conclusion: yes or no

Prediction Scenario



- Sometimes marketing manager needs to predict how much a given customer will spend during a sale at his company.
- Here, we have to predict a numeric value. Therefore the data analysis task is an example of numeric prediction.
- In this case, a model or a predictor will be constructed that predicts a continuous-valued-function or ordered value.
- Regression analysis is a statistical methodology that is popularly used for numeric prediction.

Classification—A Two-Step Process

1. Building Classifier (Model construction):
 - describing a set of predetermined classes

1. Using Classifier (Model usage):
 - for classifying future or unknown objects

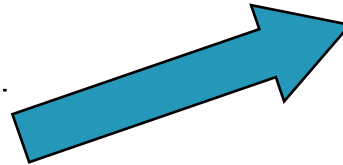
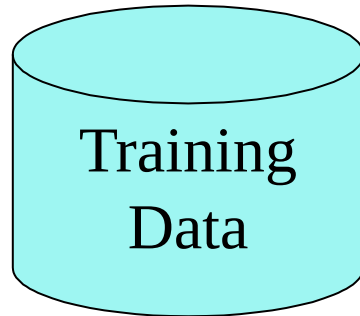
Classification—A Two-Step Process

- Building Classifier (Model construction):
 - Learning phase
 - Classification algorithm builds classifier
 - 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: training set
 - The model is represented as classification rules, decision trees, or mathematical formula

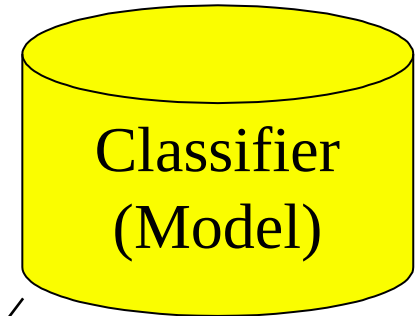
Classification—A Two-Step Process

- Using Classifier (Model usage) :
 - The classifier is used for classification.
 - The test data is used to estimate the accuracy of classification rules.
 - The classification rules can be applied to the new data tuples if the accuracy is considered acceptable.
 - 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

Classification Process (1): Model Construction



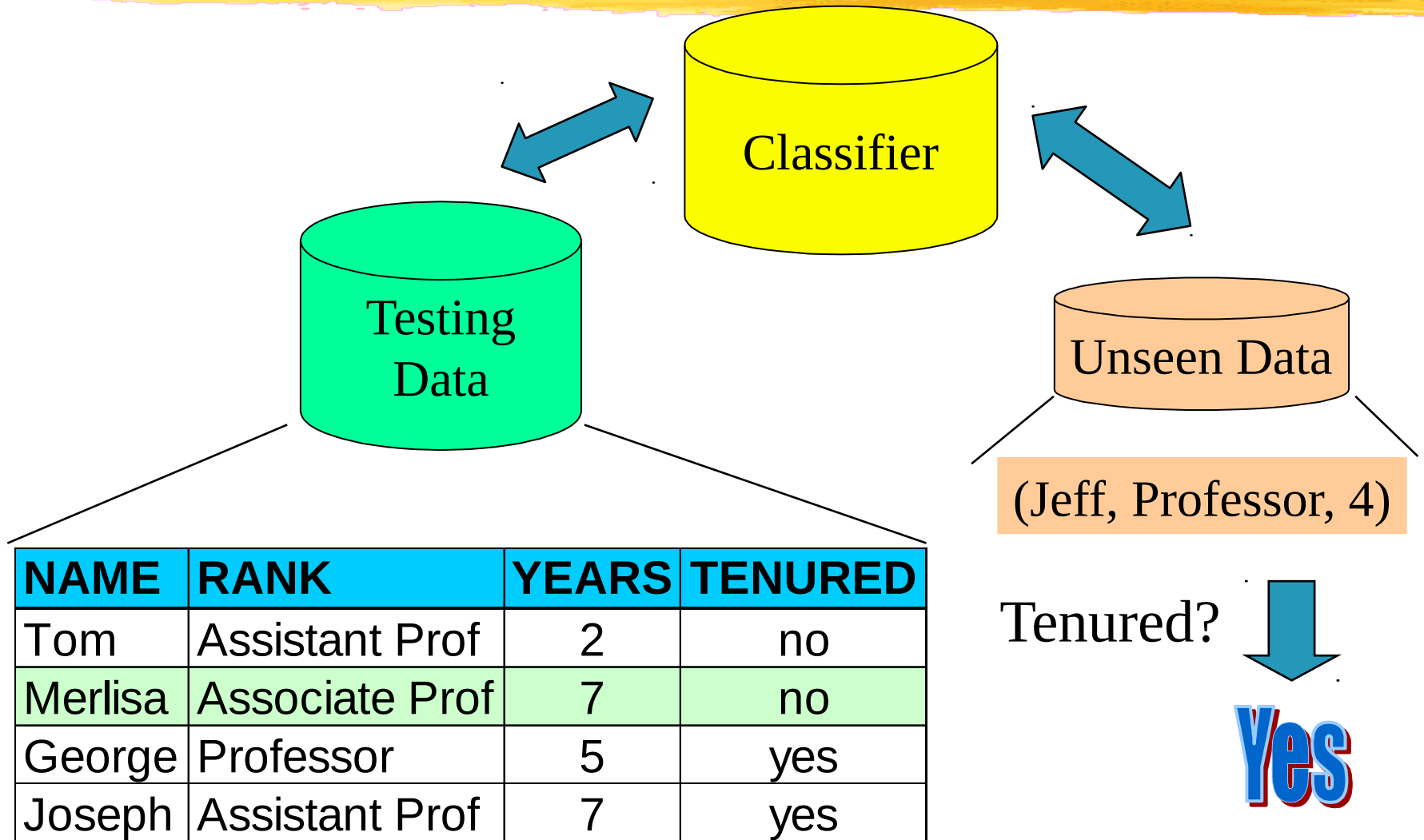
Classification
Algorithms



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'

Classification Process (2): Use the Model in Prediction



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 is 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

Issues (1): 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

Issues (2): Evaluating Classification Methods



- Predictive accuracy
- Speed and scalability
 - time to construct the model
 - time to use the model
- Robustness
 - handling noise and missing values
- Scalability
 - efficiency in disk-resident databases
- Interpretability:
 - understanding and insight proved by the model
- Goodness of rules
 - decision tree size
 - compactness of classification rules

Classification by Decision Tree Induction

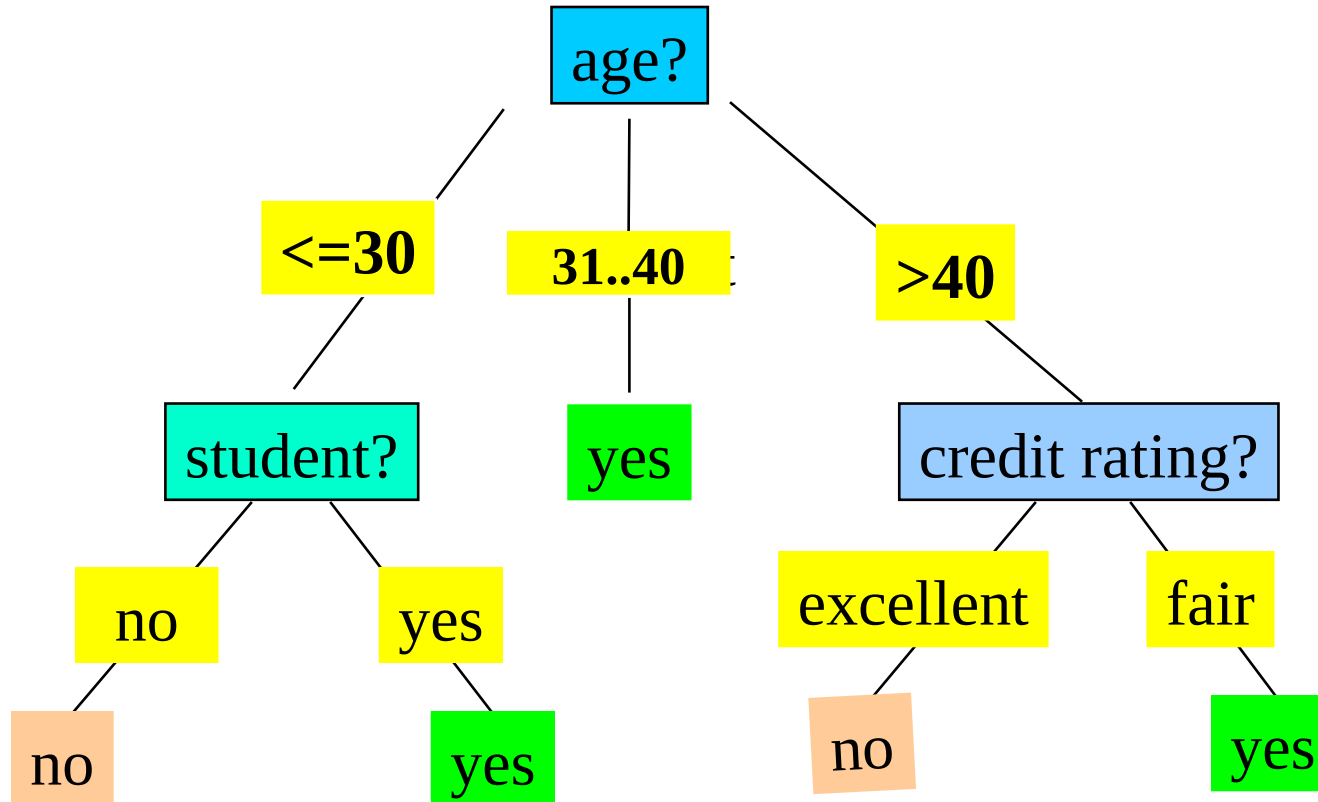
- Decision tree
 - A flow-chart-like tree structure
 - Internal node denotes a test on an attribute
 - Branch represents an outcome of the test
 - Leaf nodes represent class labels or class distribution
- Decision tree generation consists of two phases
 - Tree construction
 - At start, all the training examples are at the root
 - Partition examples recursively based on selected attributes
 - Tree pruning
 - Identify and remove branches that reflect noise or outliers
- Use of decision tree: Classifying an unknown sample
 - Test the attribute values of the sample against the decision tree

Training Dataset

This follows an example from Quinlan's ID3

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

Output: A Decision Tree for “*buys_computer*”



Algorithm for Decision Tree Induction

- Basic algorithm (a greedy algorithm)
 - 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)
- 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