

Supervised Machine Learning - C

Road Map

- Basic concepts
- Regression
- Naïve Bayesian classification
- K-nearest neighbor
- Support vector machines
- **Decision tree induction**
- Ensemble methods: Bagging and Boosting
- Summary

Introduction

- Decision tree learning is one of the most widely used techniques for classification.
 - Its classification accuracy is competitive with other methods, and
 - It is very efficient.
- The classification model is a tree, called **decision tree**.

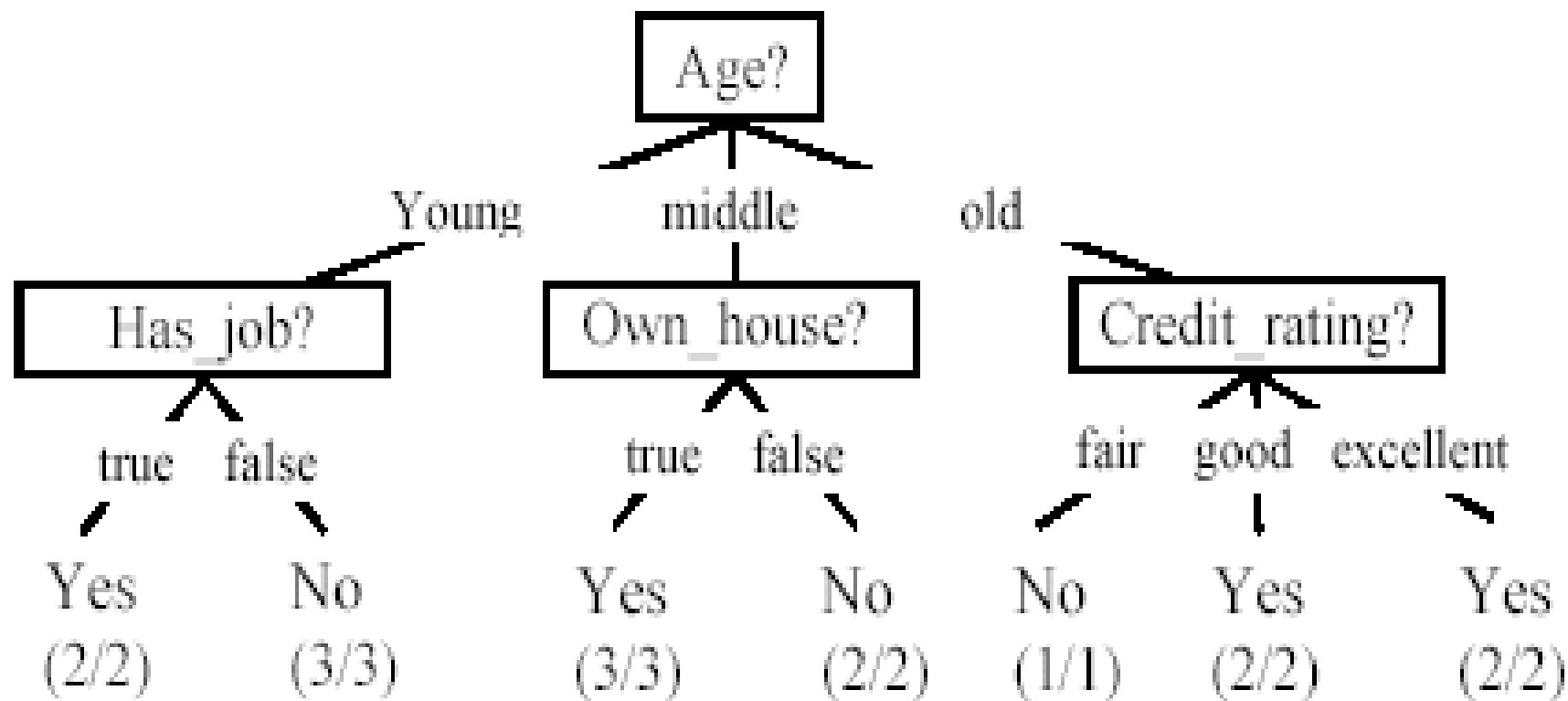
The loan data (reproduced)

Approved or not

ID	Age	Has_Job	Own_House	Credit_Rating	Class
1	young	false	false	fair	No
2	young	false	false	good	No
3	young	true	false	good	Yes
4	young	true	true	fair	Yes
5	young	false	false	fair	No
6	middle	false	false	fair	No
7	middle	false	false	good	No
8	middle	true	true	good	Yes
9	middle	false	true	excellent	Yes
10	middle	false	true	excellent	Yes
11	old	false	true	excellent	Yes
12	old	false	true	good	Yes
13	old	true	false	good	Yes
14	old	true	false	excellent	Yes
15	old	false	false	fair	No

A decision tree from the loan data

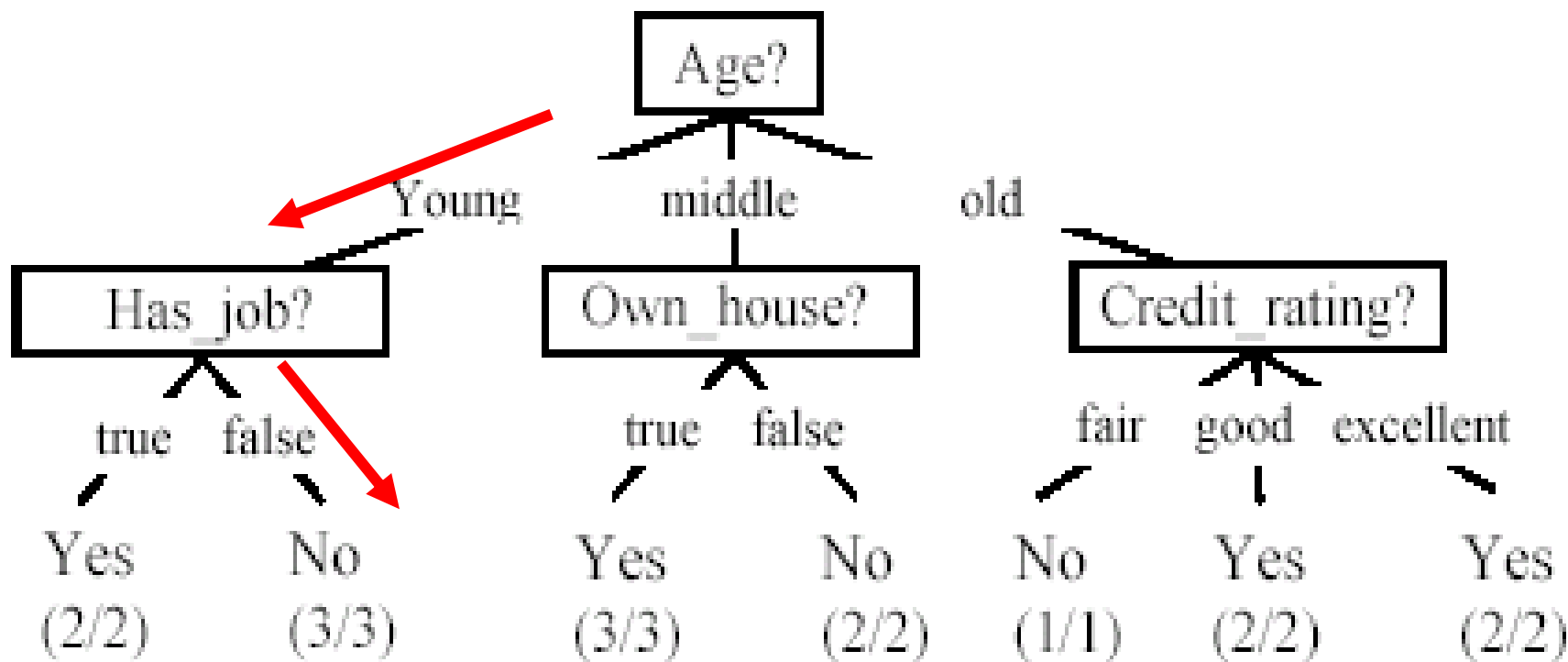
- **Decision nodes** and **leaf nodes (classes)**



Use the decision tree

Age	Has_Job	Own_house	Credit-Rating	Class
young	false	false	good	?

No



Is the decision tree unique?

- **No**. Here is a simpler tree.
- We want **smaller tree** and **accurate tree**.
 - Easy to understand and perform better.
- Finding the best tree is difficult.

From a decision tree to a set of rules

- A decision tree can be converted to a set of rules
- Each path from the root to a leaf is a rule.

Own_house = true \rightarrow Class = Yes

Own_house = false, Has_job = true \rightarrow Class = Yes

Own_house = false, Has_job = false \rightarrow Class = No

Algorithm for decision tree learning

- Basic algorithm (a greedy **divide-and-conquer** algorithm)
 - Assume attributes are categorical now
 - Tree is constructed in a **top-down manner**
 - At start, all the training examples are at the root
 - Examples are partitioned based on selected attributes
 - Attributes are selected on the basis of some functions
- Conditions for stopping partitioning
 - All examples for a given node belong to the same class
 - There are no remaining attributes for further partitioning – majority class is the leaf
 - There are no examples left

Choose an attribute to partition data

- The *key* to building a decision tree - which attribute to choose in order to branch.
- The objective is to reduce uncertainty in data as much as possible.
 - A subset of data is *pure* if all instances belong to the same class.

Information theory

- **Information theory** provides a mathematical basis for measuring the information content.
- To understand the notion of information, think about it as providing the answer to a question, for example, whether a coin will come up heads.
 - If one already has a good guess about the answer, then the actual answer is less informative.
 - If one already knows that the coin is rigged so that it will come with heads with probability 0.99, then a message (advanced information) about the actual outcome of a flip is worth less than it would be for a honest coin (50-50).

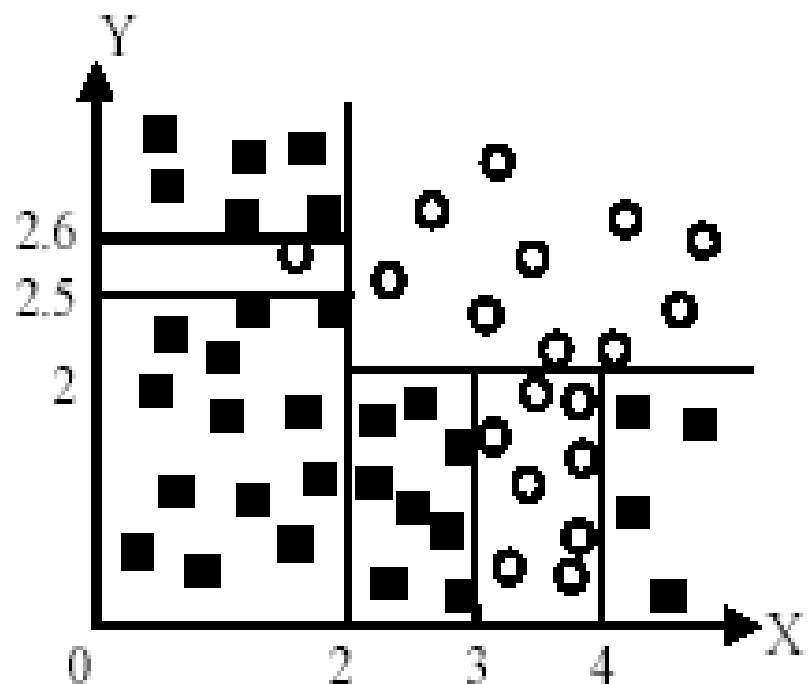
Information theory (cont ...)

- For a fair (honest) coin, you have no information, and you are willing to pay more (say in terms of \$) for advanced information - less you know, the more valuable the information.
- **Information theory** uses this same intuition, but instead of measuring the value for information in dollars, it measures information contents in **bits**.
- One bit of information is enough to answer a yes/no question about which one has no idea, such as the flip of a fair coin

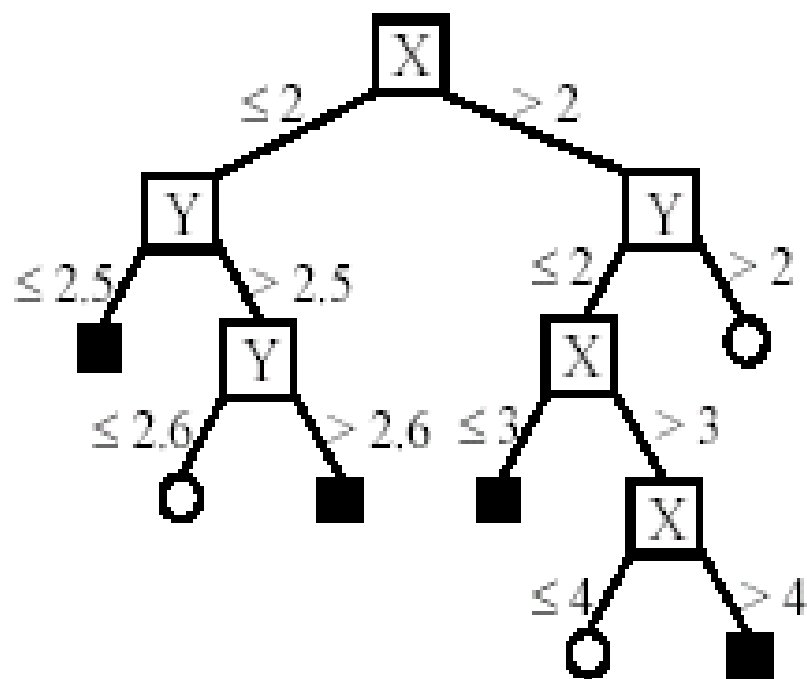
Handling continuous attributes

- Handle continuous attribute by splitting into two intervals (can be more) at each node.
- How to find the best threshold to divide?
 - Sort all the values of an continuous attribute in increasing order $\{v_1, v_2, \dots, v_r\}$,
 - There is one possible threshold between two adjacent values v_i and v_{i+1} . Try all possible thresholds and find the one that maximizes information gain.

An example in a continuous space



(A) A partition of the data space



(B). The decision tree

Avoid overfitting in classification

- **Overfitting**: A tree may overfit the training data
 - Good accuracy on training data but poor on test data
 - Symptoms: tree too deep and too many branches, some may reflect anomalies due to noise or outliers
- Two approaches to avoid overfitting
 - **Pre-pruning**: Halt tree construction early
 - Difficult to decide because we do not know what may happen subsequently if we keep growing the tree.
 - **Post-pruning**: Remove branches or sub-trees from a “fully grown” tree.
 - This method is commonly used.