

Databases are rich with hidden information that can be used as intelligent decision making.

Classification and prediction are two forms of data analysis.

Used to extract models describing important data classes or to predict future data trends.

Whereas classification predicts categorical (discrete, unordered) labels.

Prediction models continuous valued function.

For example, we can build a classification model to categorize bank loan applications as either safe or risky.

Predictive methods have been proposed by researchers to predict the expenditures in dollars of potential customers on computer equipment given their income and occupation.

Classification and prediction have numerous applications including fraud detection, target marketing, performance prediction, manufacturing, and medical diagnosis.

What is Classification?

A bank loan officer needs analysis of his data in order to learn which loan applicants are "safe" and which are "risky" for the bank.

A marketing manager at All Electronics needs data analysis to help guess whether a customer with a given profile will buy a new computer.

- In each of these examples, the data analysis task is classification.
- Where, a model or classifier is constructed to predict categorical labels such as "safe" or "sky" for the loan application data.
- "yes" or "no" for the marketing data.
- "treatment A", "treatment B", or "treatment C" for the medical data.

- Suppose that the marketing manager would like to predict how much a given customer will spend during a sale at AllElectronics?
- This data analysis task is an example of numeric prediction, where the model constructed predicts a continuous valued function, or ordered value; as opposed to a categorical label.
 - This model is a predictor.
 - Regression analysis is a statistical methodology that is most often used for numeric prediction.

How does classification work?

- Data classification is a two-step process.
- In the first step, a classifier is built describing a predetermined set of data classes or concepts.
- This is the learning step (or training phase).
- Where, a classification algorithm builds the classifier by analyzing or "learning from" a ^{training} set of made up of database tuples and their associated class labels.

(6) Learning: Training data analyze by classification
 Here no. class label attribute is loan decision, and learned model is represented in the form of classification rules. (2)

Training data

Classification algorithm

name age income loan-decision

Sandy Jones Young low risky

Bill Lee Young low risky

Caroline Middle-aged high safe

Rick Middle-aged low risky

Cherie Senior medium safe

Susan Senior low safe

Classification rule

If age = young THEN loan-decision = risky

If income = high THEN loan-decision = safe

If age = middle-aged AND income = low
THEN loan-decision = risky

fig. @

Classification rules

Test data

New data

name age Income loan-decision

Juan Bello Senior low safe

Sylvia Gort middle-aged low risky

Anne Middle-aged high safe

(John Henry, middle-aged, low

Loan decision?

risky

fig. b Classification: Test data are used to estimate the accuracy of the classification rules. If the accuracy is considered acceptable, the rules can be applied to the classification of new data tuples. 57

- A tuple X is represented by an n -dimensional attribute vector.
- $X = (x_1, x_2, \dots, x_n)$ depicting n measurements made on the tuple from n -database attributes respectively. A_1, A_2, \dots, A_n .
- Each tuple, X is assumed to belong to a pre-defined class as determined by another database attribute called the class Label attribute.
- The class Label attribute is discrete-valued and unordered.
- It is categorical in that each value serves as a category or class.
- The individual tuples making up the training set are referred as training tuples and selected from the database under analysis.
- In the context of classification, data tuples can be referred to as samples, examples, instances, data points, as objects.
- Because the class label of each training tuple is provided this step is also known as supervised learning (i.e. the learning of the classifier is "supervised," in that it is told to which class each training tuple belongs).

i) Data for Classification and Prediction

Several preprocessing steps may be applied to the data to help improve the accuracy, efficiency, and scalability of the classification or prediction process.

Data Cleaning: This refers to the preprocessing of data in order to remove or reduce noise (by applying smoothing techniques) and the treatment of missing values (e.g. by replacing a missing value with the most commonly occurring value for that attribute, or with the most probable value based on statistics).
- This step can help reduce confusion during learning.

i) Relevance analysis: Many of the attributes in the data may be redundant. Correlation analysis can be used to identify whether any two given attributes are statistically related.
- Relevance analysis, in the form of correlation analysis and attribute subset selection, can be used to detect attributes that do not contribute to the classification or prediction task.
- Including such attributes may otherwise slow down, and possibly mislead, the learning step.

ii) Data transformation and reduction:
- The data may be transformed by normalization.
- When neural networks or methods involving distance measurements are used in learning step.
- Normalization involves scaling all values for a given attribute so that they fall within a small specified range, such as -1.0 to 1.0 or 0.0 to 1.0.

Classification by Decision Tree Induction \Rightarrow

- Decision tree induction is the learning of decision trees from class-labeled training tuples.

- A decision tree is a flow chart-like tree structure, where each internal node (nonleaf node) denotes a test on an attribute, each branch represents an outcome of the test.
- Each leaf (or terminal node) holds a class label.
- The topmost node in a tree is the root node.

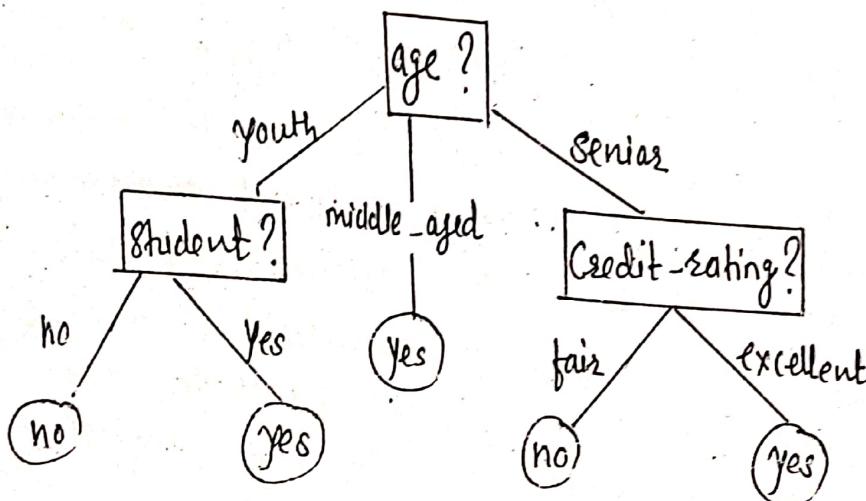


Fig. A decision tree for the concept `buys_computer`, indicating whether a customer at AllElectronics is likely to purchase a computer.

- Each leaf node represents a class (either `buys_computer = yes` or `buys_computer = no`).

- Figure represents the concept `buys_computer`, that is, it predicts whether a customer at AllElectronics is likely to purchase a computer.

- Internal nodes are denoted by rectangles, leaf nodes are denoted by ovals.

- Some decision trees algorithms produce only binary trees (where each internal node branches have exactly two other nodes).

(4)

The induction \Rightarrow

In far learning decision trees
during tree construction, attribute selection measures are used to
select the attribute that best partitions the tuples into
distinct classes.

Algorithm: generate_decision_tree \Rightarrow generate a decision tree from the
training tuples of data partition D.

Input:

- i). Data Partition, D, which is a set of training tuples and their associated class labels;
- attribute_list, the set of candidate attributes;
- Attribute_selection_method, a procedure to determine the splitting criteria that "best" partition the data tuples into individual class. Criterion consists of a splitting_attribute

Output: A decision tree.

Method:

- 1) Create a node N
- 2) If tuples in D are all of the same class, C then return it as a leaf node labeled with the class C;
- if attribute_list is empty then return it as a leaf node labeled with the majority class
- 3) D;
apply Attribute_selection_method (D, attribute_list) to find the "best" splitting_criterion;
- label node N with splitting_criterion;
- if splitting_attribute is discrete-valued and
multiple splits allowed then
attribute_list \leftarrow attribute_list - splitting_attribute //remove splitting_attribute

```

    for each outcome j of splitting-criterion
        // partition the tuples and grow subtrees for each partition
        let  $D_j$  be the set of data tuples in  $D$  satisfying outcome  $j$ ;
        // a partition
        if  $D_j$  is empty then
            attach a leaf labeled with the majority class in  $D$ 
            to node  $N$ ;
        else attach the node returned by generate-decision-tree( $D_j$ ,
            attribute-list) to node  $N$ ;
    end for
    return  $N$ .

```

→ fig. Basic algorithm for inducing a decision tree from training tuples.

- The algorithm is called with three parameters: D , attribute-list, Attribute_selection_method
- Initially D refers as a data partition, it is the complete set of training tuples and their associated class labels.
- The parameter attribute-list is a list of attributes describing the tuples.
- Attribute-selection-method specifies a heuristic procedure for selecting the attribute that "best" discriminates the given tuples according to class.
- The procedure employs an attribute selection measure, such as information gain or the gini index.

we can treat it as a single node, N , representing the training data.

If the types in D are all of the same class, then node N becomes leaf and is labeled with that class.

Otherwise, the algorithm calls Attribute-selection-method to determine the splitting criterion.

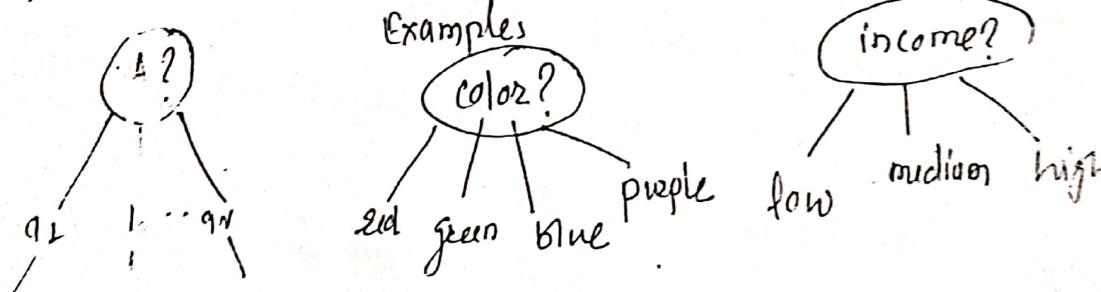
The splitting criterion tells us which branches to grow from node N with respect to the outcomes of the chosen test.

The node N is labeled with splitting criterion, which serves as a test at the node.

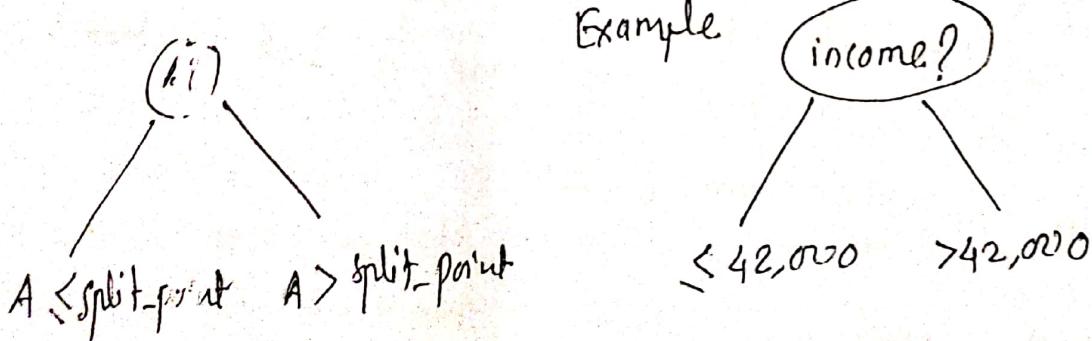
A branch is grown from node N for each of the outcomes of the splitting criterion.

Let A be the splitting attribute. A has v distinct values $\{a_1, a_2, \dots, a_v\}$, based on the training data.

$\triangleright A$ is discrete-valued :-



$\triangleright A$ is continuous-valued :-



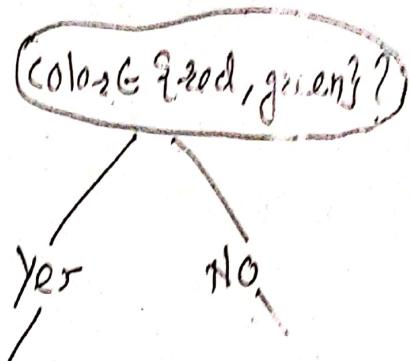
- A is discrete-valued and binary tree must be produced
- the test at node N is of the form,

$A \in S_A ?$

- S_A is the splitting subset for A , returned by Attribute selection method as part of the splitting criterion.



Example :-



- The algorithm uses the same process recursively to form a decision tree for the tuples at each resulting partition D_j of D .
- The recursive partitioning stops only when any one of the following terminating conditions is true:-
- All of the tuples in partition D (represented at node N) belongs to the same class, or
- There are no remaining attributes on which the tuples may be further partitioned.
- There are no tuples for a given branch, that is, a partition D_j is empty. In this case, a leaf is created with the majority class in D .

Attribute Selection Measures :-

An attribute selection measure is a heuristic for selecting the splitting criterion that "best" separates a given data partition, D , into labeled training tuples into individual classes.

If we were to split D into smaller partitions according to the outcomes of the splitting criterion, ideally each partition would be pure.

i.e. all of the tuples that fall into a given partition would belong to the same class).

- Attribute selection measures are also known as splitting rules because they determine how the tuples at given node are to be split.

- The attribute selection measure provides ranking.

- The attribute having the best score for the measure is chosen as the splitting attribute for the given tuples.

This selection describes three popular attribute selection measures

- i) Information gain
- ii) Gini ratio
- iii) Gini index

i) Information gain \Rightarrow

- This measure is based on pioneering work by Claude Shannon on information theory, which studied the value as "information content" of messages.

- Let node N represent or hold the tuples of partition D .

- The attribute with the highest information gain is chosen as the splitting attribute for node N .

- The expected information needed to classify a tuple in D is given by,

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

- where p_i is the probability that an arbitrary tuple in D belongs to class C_i .

- and is estimated by $\frac{|C_i, D|}{|D|}$

- A log function to the base 2 is used, because the information is encoded in bits.

- $\text{Info}(D)$ is just the average amount of information needed to identify the class label of a tuple in D .

- $\text{Info}(D)$ is also known as the entropy of D .

- Now suppose we were to partition the tuples in D on some attribute A having n distinct values $\{q_1, q_2, \dots, q_n\}$, as observed from the training data.

- If A is discrete-valued, these values correspond directly to the n outcomes of a test on A .

- Attribute A can be used to split D into n partitions or subsets.

- $\{D_1, D_2, \dots, D_n\}$.

- Where D_j contains those tuples in D that have outcome q_j of A .

- These partitions would correspond to the branches going from node N .

- How much more information would we still need after the partitioning in order to arrive at an exact classification?

- This amount is measured by,

$$\text{Info}_A(D) = \sum_{j=1}^n \frac{|D_j|}{|D|} \times \text{info}(D_j).$$

- The term $\frac{|D_j|}{|D|}$ acts as the weight of the j^{th} partition.

is the expected information required to classify a tuple from
and on partitioning by A.

The smaller the expected information (still) required, the greater
the purity of the partition.

Information gain is defined as the difference between the original
information requirement and the new requirement. That is;

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

The attribute A with the highest information gain, ($\text{Gain}(A)$), is
chosen as the splitting attribute at node N.

Example: class-labeled training tuples from the ALLElectronics (automobiles)
database.

RFID	age	income	student	credit_rating	class: buys_computer
1	young	high	no	fair	no
2	young	high	no	excellent	no
3	middle-aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	excellent	no
6	senior	low	yes	excellent	yes
7	middle-aged	low	yes	fair	no
8	young	medium	no	fair	yes
9	young	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	young	medium	yes	excellent	yes
12	middle-aged	medium	no	excellent	yes
13	middle-aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

- In this example, each attribute is discrete-valued.
- The class label attribute, buys_computer, has two distinct values, namely, {yes, no}.

- Therefore, there are two distinct classes that is $m = 2$.

- let class C_1 correspond to yes

class C_2 correspond to no

- A (root) node N is created for the tuples in D .

- we must compute the information gain of each attribute to find the splitting criterion for these tuples.

- To compute the expected information needed to classify tuple

In D , we use

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

yes = 9 times
no = 5 times

here $m=2$, $p_i = \frac{|C_i|}{|D|}$ total tuples = 14

$$\text{info}(D) = - \frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right) = 0.940 \text{ bit}$$

Next we need to compute the expected information requirement for each attribute.

Let's start with the attribute age.

we have,

$$\text{info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{info}(D_j)$$

yes	no
-----	----

Here $v=3$ i.e.

young (5 times)	2 times	3 times
middle-aged (4 times)	4	0
senior (5 times)	3	2

$$\begin{aligned}
 \text{Info}_{\text{age}}(i) &= \frac{5}{14} \times \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) \\
 &\quad + \frac{4}{14} \times \left(-\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} \right) \\
 &\quad + \frac{5}{14} \times \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) \\
 &= 0.694 \text{ bits}
 \end{aligned}$$

Hence, the gain in information from such a partitioning would be,

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D) = 0.940 - 0.694 = 0.246 \text{ bits}$$

Similarly, we can compute $\text{Gain}(\text{income}) = 0.029 \text{ bits}$,
 $\text{Gain}(\text{student}) = 0.151 \text{ bits}$ and
 $\text{Gain}(\text{credit_sahig}) = 0.048 \text{ bits}$.

⇒ Because age has the highest information gain among the attributes, it is selected as the splitting attribute.

⇒ Node N. is labeled with age, and branches are made for each of the attributes values.

age?

<p>Youth</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th colspan="4">income student credit-rating class</th> </tr> </thead> <tbody> <tr><td>high</td><td>no</td><td>fair</td><td>no</td></tr> <tr><td>high</td><td>no</td><td>excellent</td><td>no</td></tr> <tr><td>medium</td><td>no</td><td>fair</td><td>no</td></tr> <tr><td>low</td><td>yes</td><td>fair</td><td>yes</td></tr> <tr><td>medium</td><td>yes</td><td>excellent</td><td>yes</td></tr> </tbody> </table>	income student credit-rating class				high	no	fair	no	high	no	excellent	no	medium	no	fair	no	low	yes	fair	yes	medium	yes	excellent	yes	<p>Senior</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th colspan="4">income student credit-rating class</th> </tr> </thead> <tbody> <tr><td>medium</td><td>no</td><td>fair</td><td>yes</td></tr> <tr><td>low</td><td>yes</td><td>fair</td><td>yes</td></tr> <tr><td>low</td><td>yes</td><td>excellent</td><td>no</td></tr> <tr><td>medium</td><td>yes</td><td>fair</td><td>yes</td></tr> <tr><td>medium</td><td>no</td><td>excellent</td><td>no</td></tr> </tbody> </table>	income student credit-rating class				medium	no	fair	yes	low	yes	fair	yes	low	yes	excellent	no	medium	yes	fair	yes	medium	no	excellent	no
income student credit-rating class																																																	
high	no	fair	no																																														
high	no	excellent	no																																														
medium	no	fair	no																																														
low	yes	fair	yes																																														
medium	yes	excellent	yes																																														
income student credit-rating class																																																	
medium	no	fair	yes																																														
low	yes	fair	yes																																														
low	yes	excellent	no																																														
medium	yes	fair	yes																																														
medium	no	excellent	no																																														
<table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th colspan="4">income student credit-rating class</th> </tr> </thead> <tbody> <tr><td>high</td><td>no</td><td>fair</td><td>yes</td></tr> <tr><td>low</td><td>yes</td><td>excellent</td><td>yes</td></tr> <tr><td>medium</td><td>no</td><td>excellent</td><td>yes</td></tr> <tr><td>high</td><td>yes</td><td>fair</td><td>yes</td></tr> </tbody> </table>		income student credit-rating class				high	no	fair	yes	low	yes	excellent	yes	medium	no	excellent	yes	high	yes	fair	yes																												
income student credit-rating class																																																	
high	no	fair	yes																																														
low	yes	excellent	yes																																														
medium	no	excellent	yes																																														
high	yes	fair	yes																																														

→ Final tree

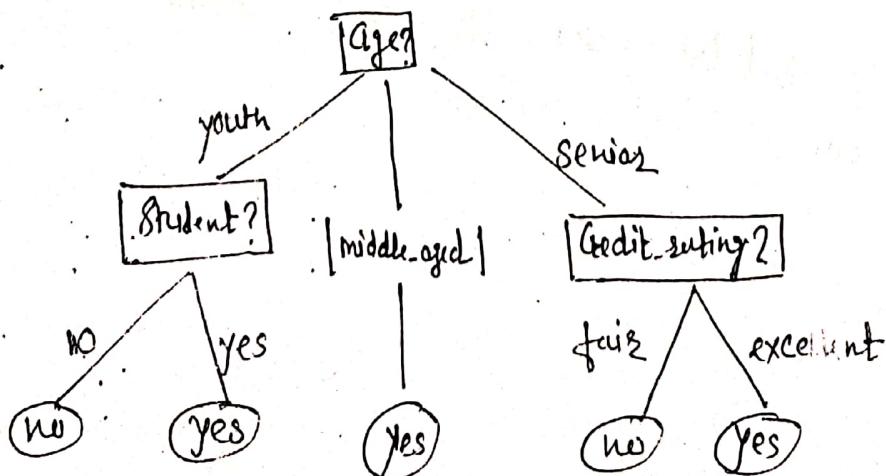


Fig: A decision tree for the concept buys-computer

1. the process of grouping the data into classes or, clusters, in which objects within a cluster have high similarity in comparison to another but are very dissimilar to objects in other clusters.

What is cluster Analysis?

The process of grouping a set of physical or abstract objects into classes of similar objects is called clustering.

A cluster is a collection of data objects that are similar to one another within the same cluster and are dissimilar to the objects in other clusters.

- A cluster of data objects can be treated collectively as one group and so may be considered as a form of data compression.

- Cluster analysis is an important human activity.

- Cluster analysis has been widely used in numerous applications including market research, pattern recognition, data analysis and image processing.

- In business, clustering can help marketers discover distinct groups in their customer bases and characterize customers groups based on purchasing patterns.

- Clustering is also called data segmentation. In some applications, because clustering partitions large data sets into groups according to their similarity.

- Clustering can also be used for outlier detection, where outliers (values that are "far away" from any cluster) may be more interesting than common cases.

- Applications of outlier detection include the detection of credit card fraud and the monitoring of criminal activities in electronic commerce.

Observing is a form of learning by observation, testing, trial, error, and success.

clustering is a challenging field of research in which its potential applications pose their own special requirements.

The following are typical requirements of clustering in data mining.

- scalability! - Highly scalable clustering algorithms are needed.

- Ability to deal with different types of attributes.

- Discovery of clusters with arbitrary shape

- Minimal requirements for domain knowledge to determine input parameters.

Ability to deal with noisy data

- Incremental clustering and insensitivity to the order of input records.

-- **High dimensionality:** A database or a data warehouse can contain several dimensions as attributes. Many clustering algorithms are good at handling low-dimensional data, involving two to three dimensions. Finding clusters of data points in high dimensional space is challenging.

constraint-based clustering! Real-world applications may need to perform clustering under various kinds of constraints.

interpretability and usability: - Users expect clustering results to be interpretable, comprehensible, and usable. It is important to study how an application goal may influence the selection of clustering features and methods.