

# Classification and Regression

- Classification vs Prediction
  - Classification Process
- Decision Tree Model
  - Model
  - ID3 Algorithm
- More thoughts
- Regression Tree
- Take-home messages!

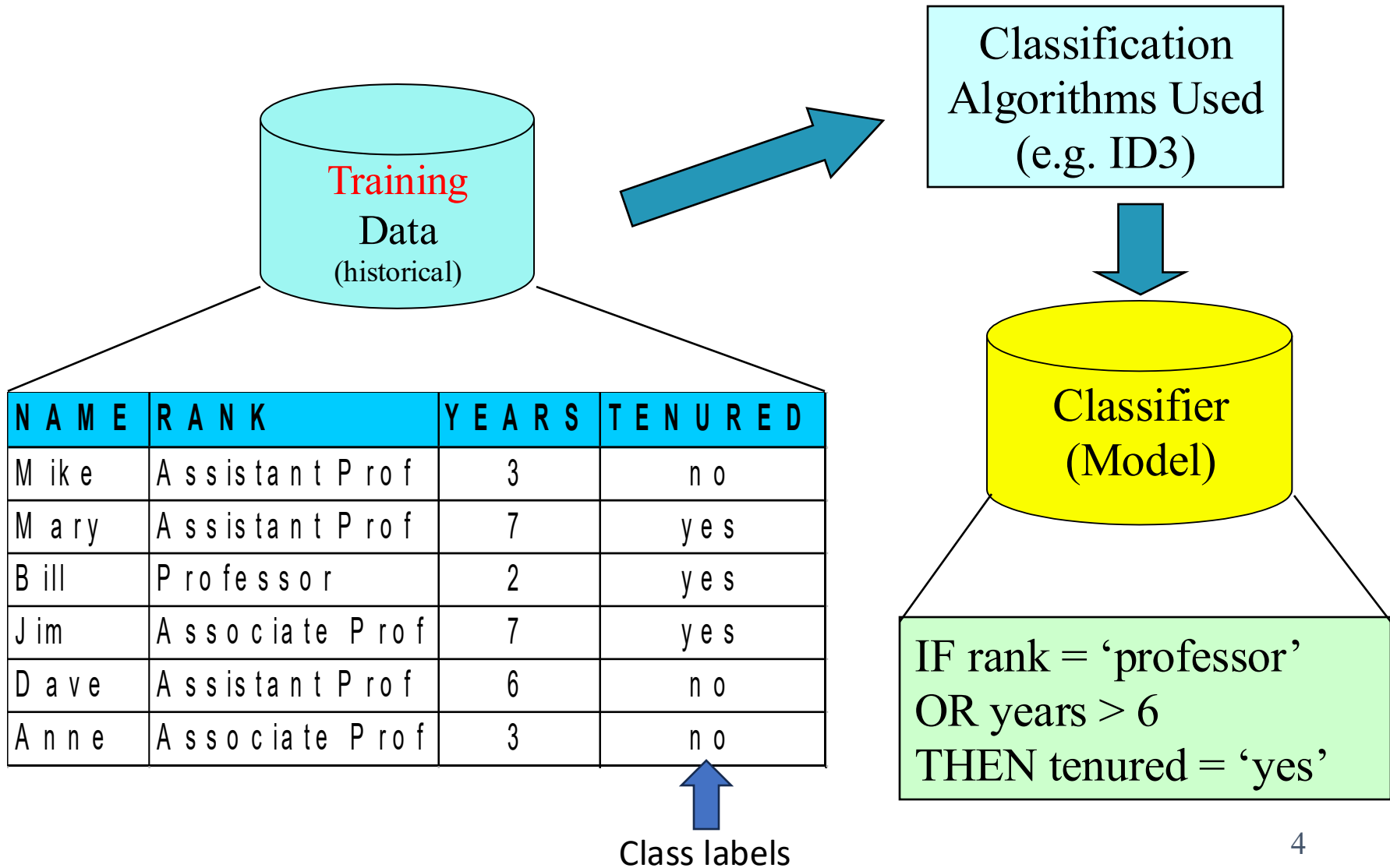
# Classification vs. Prediction

- **Classification:**
  - predicts 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
  - E.g. predicting the house price
  - Also referred as **regression**
- **Typical Applications**
  - credit approval
  - target marketing
  - medical diagnosis
  - Price prediction

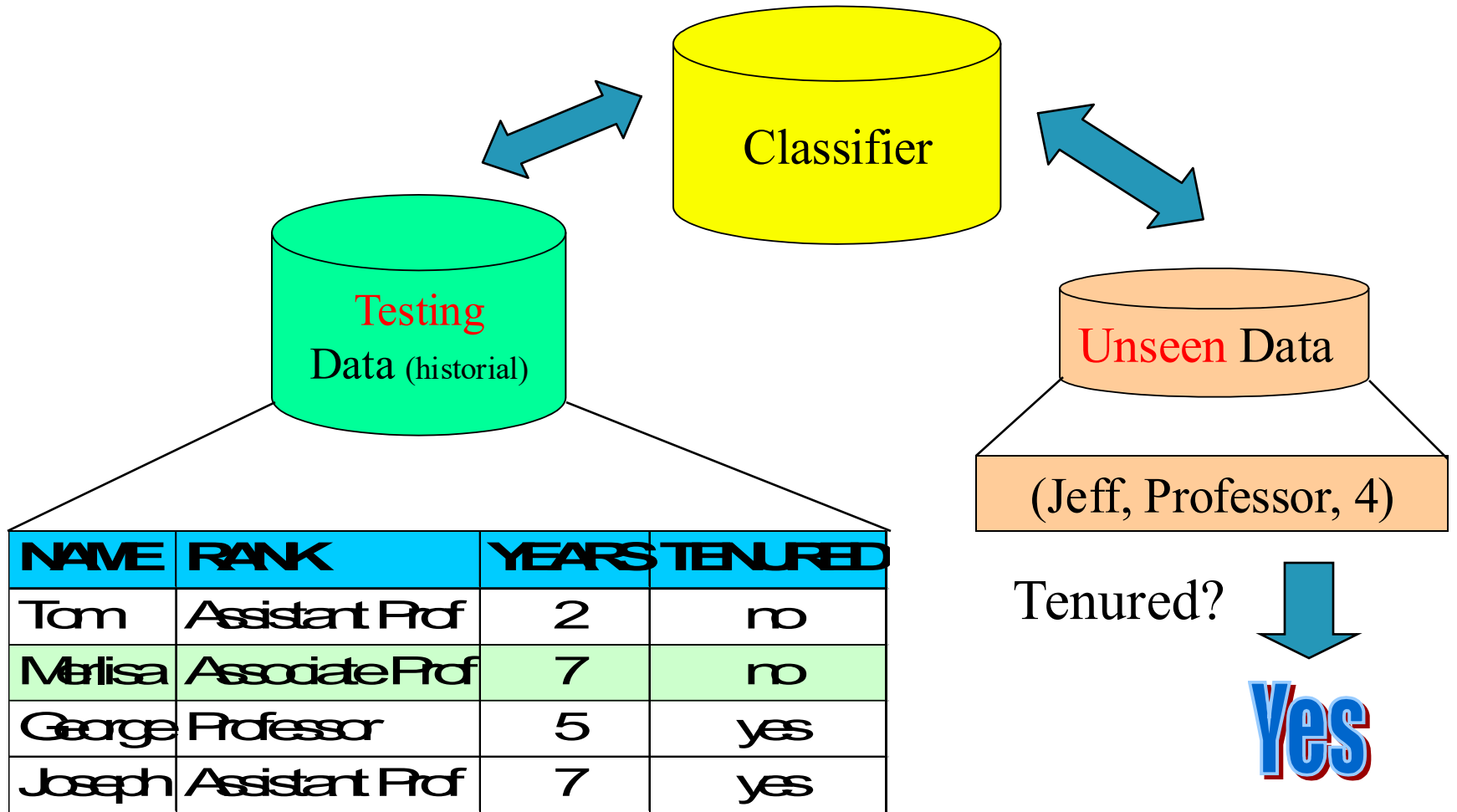
# Classification—A Two-Step Process

- ***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: **training set**
  - The model is represented as classification rules, decision trees, layered neural networks or mathematical formulae
- ***Model usage***: classifying future or unknown objects
  - The known label of test sample is compared with the classified result from the model
  - Accuracy is the percentage of test set samples that are correctly classified by the model
    - Need to estimate the accuracy of model
  - Test set is independent of training set, otherwise over-fitting will occur

# Classification Process: Model Construction



# Classification Process: Model Usage



# Different Types of Data in Classification

- **Seen Data** = Historical Data
- **Unseen Data** = Current and Future Data
- **Seen Data:**
  - Training Data – for constructing the model
  - Testing Data – for validating the model (It is unseen during the model construction process)
- **Unseen Data:**
  - Application (actual usage) of the constructed model

\*\*\* We will later come back to data for validation \*\*\*

# 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 regarding classification and prediction: 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 regarding classification and prediction: 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 provided by the model
- Goodness of rules
  - decision tree size
  - compactness of classification rules

# Decision Tree

**A Data Analytics Perspective + Feature Engineering**

# Classification by Decision Tree Induction

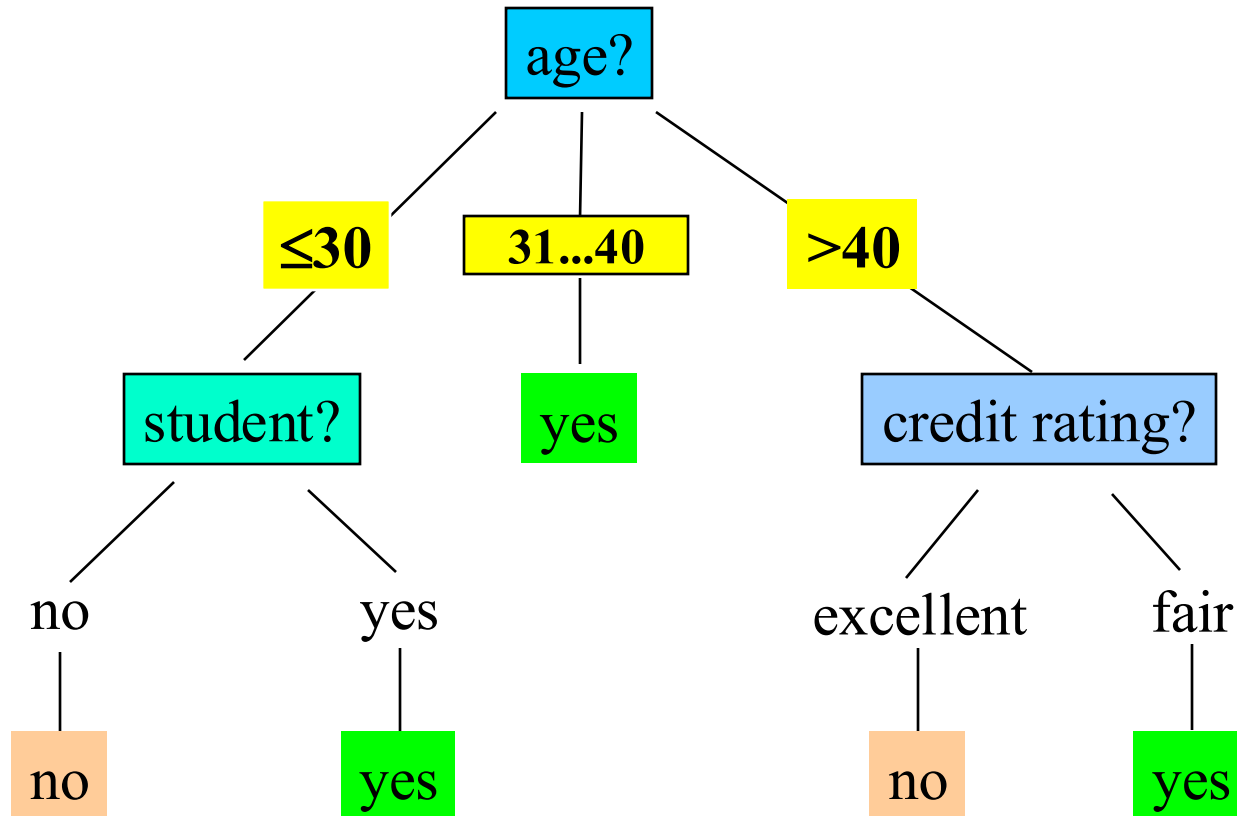
- Decision tree Structure
  - A flow-chart-like tree structure
  - Internal node denotes a test on an attribute
  - Branch represents an outcome of the test
  - Leaf nodes represent **class labels** or **class distribution**
- Decision tree generation consists of two phases
  - Tree construction
    - 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:

Following 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*”



# Extracting Classification Rules (Knowledge) from Decision Trees

- Represent the knowledge in the form of **IF-THEN** rules
- One rule is created for each path from the root to a leaf
- Each attribute-value pair along a path forms a conjunction
- The leaf node holds the class prediction
- Rules are easier for humans to understand
- Example

IF *age* = " $\leq 30$ " AND *student* = "no" THEN *buys\_computer* = "no"

IF *age* = " $\leq 30$ " AND *student* = "yes" THEN *buys\_computer* = "yes"

IF *age* = "31...40" THEN *buys\_computer* = "yes"

IF *age* = " $>40$ " AND *credit\_rating* = "excellent" THEN *buys\_computer* = "no"

IF *age* = " $>40$ " AND *credit\_rating* = "fair" THEN *buys\_computer* = "yes"

# Algorithm for Decision Tree Construction

## ⦿ 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

# Which attribute to choose?

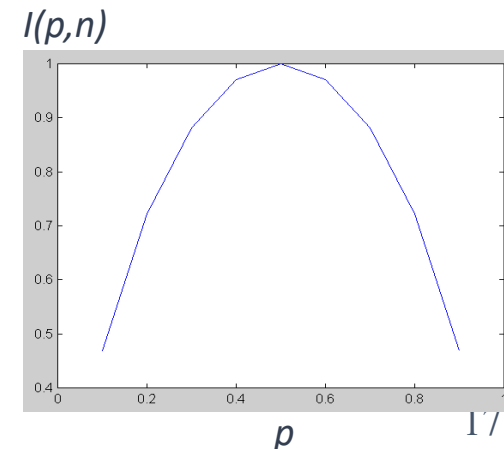
## Attribute Selection Measure

- **Information gain** (ID3/C4.5)
  - All attributes are assumed to be categorical
  - Can be modified for continuous-valued attributes
- **Gini index** (IBM IntelligentMiner)
  - All attributes are assumed continuous-valued
  - Assume there exist several possible split values for each attribute
  - May need other tools, such as clustering, to get the possible split values
  - Can be modified for categorical attributes

# Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- Assume there are two classes,  $P$  and  $N$ 
  - Let the set of examples  $S$  contain  $p$  elements of class  $P$  and  $n$  elements of class  $N$
  - The amount of information, needed to decide if an arbitrary example in  $S$  belongs to  $P$  or  $N$  is defined as

$$I(p, n) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$



# Information Gain in Decision Tree Construction

- Assume that using attribute  $A$  a set  $S$  will be partitioned into sets  $\{S_1, S_2, \dots, S_v\}$ 
  - If  $S_i$  contains  $p_i$  examples of  $P$  and  $n_i$  examples of  $N$ , the **entropy**, or the expected information needed to classify objects in all subtrees  $S_i$  is

$$E(A) = \sum_{i=1}^v \frac{p_i + n_i}{p + n} I(p_i, n_i)$$

- The encoding information that would be gained by branching on  $A$

$$Gain(A) = I(p, n) - E(A)$$

# Attribute Selection by Information Gain Computation

- Class P: buys\_computer = “yes”
- Class N: buys\_computer = “no”
- $I(p, n) = I(9, 5) = 0.940$
- Compute the entropy for *age*:

age	$p_i$	$n_i$	$I(p_i, n_i)$
$\leq 30$	2	3	0.971
30...40	4	0	0
$> 40$	3	2	0.971

$$E(\text{age}) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

Hence,

$$Gain(\text{age}) = I(p, n) - E(\text{age}) = 0.246$$

Similarly,

$$Gain(\text{income}) = 0.029$$

$$Gain(\text{student}) = 0.151$$

$$Gain(\text{credit\_rating}) = 0.048$$

Thus, we should select “age” as the root node of the decision tree.

# Avoid Overfitting in Classification

- The generated tree may overfit the training data
  - Too many branches, some may reflect anomalies due to noise or outliers
  - Result is in poor accuracy for unseen samples
- Two approaches to avoid overfitting
  - Prepruning: Halt tree construction early—do not split a node if this would result in the goodness measure falling below a threshold
    - Difficult to choose an appropriate threshold
  - Postpruning: Remove branches from a “fully grown” tree—get a sequence of progressively pruned trees
    - Use a set of data different from the training data to decide which is the “best pruned tree”

# Thinking time ...

The model is simple enough to understand many machine learning and data analytics concepts.

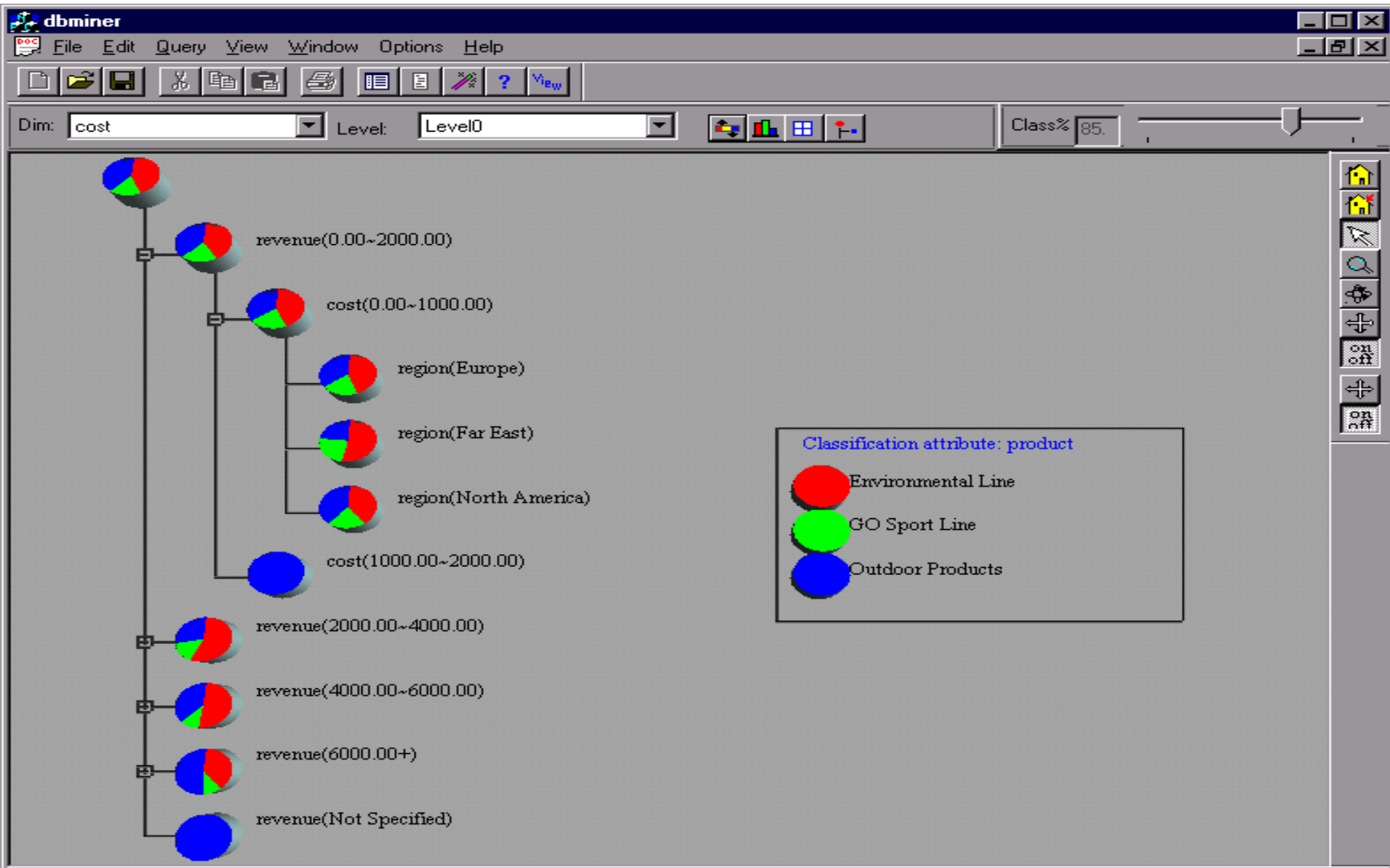
Q1. Can you add one more record to the `buys_computer` dataset so that the classification accuracy on the training set is less than 100%?

Q2. When a given **training** dataset cannot be perfectly classified by DT, what can be done to boost up its performance?

# Classification in Large Databases

- Classification—a classical problem extensively studied by statisticians and machine learning researchers
- Scalability: Classifying data sets with **millions of examples** and **hundreds of attributes** with reasonable speed
- Why decision tree induction in data mining?
  - relatively faster learning speed (than other classification methods)
  - convertible to simple and easy to understand classification rules
  - can use SQL queries for accessing databases
  - comparable classification accuracy with other methods

# Presentation of Classification Results



How to generalize decision tree for regression task?

## Regression Trees or (CART)

Classification and Regression Tree

# Regression Trees for Prediction/Regression

## Similarity

- Used with continuous outcome variable
- Procedure similar to classification tree
- Many splits attempted, choose the one that minimizes impurity

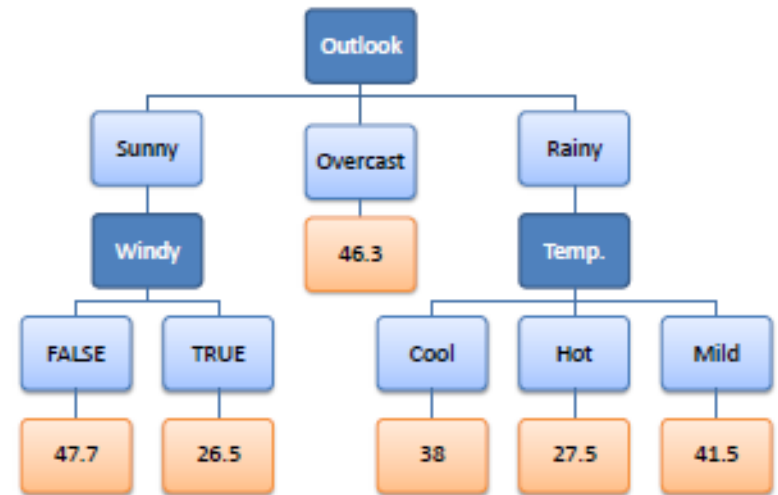
## Difference

- Prediction is computed as the **average** of numerical target variable of involved branches/rules (cf. majority voting in decision trees)
- Impurity measured by **sum of squared deviations** from leaf mean
- Performance measured by RMSE (root mean squared error)

# Regression Tree: Decision Tree for Regression

## Play Tennis Regression Problem

Predictors				Target
Outlook	Temperature	Humidity	Windy	Hours Played
rainy	hot	high	false	25
rainy	hot	high	true	30
overcast	hot	high	false	46
sunny	mild	high	false	45
sunny	cool	normal	false	52
sunny	cool	normal	true	23
overcast	cool	normal	true	43
rainy	mild	high	false	35
rainy	cool	normal	false	38
sunny	mild	normal	false	46
rainy	mild	normal	true	48
overcast	mild	high	true	52
overcast	hot	normal	false	44
sunny	mild	high	true	30



Outlook	Temperature	Humidity	Windy	Class
sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

## Play Tennis Classification Problem

# How to construct the Regression Tree?

## From information gain to deviation reduction

Hours Played
25
30
46
45
52
23
43
35
38
46
48
52
44
30



$$\text{Count} = n = 14$$

$$\text{Average} = \bar{x} = \frac{\sum x}{n} = 39.8$$

$$\text{Standard Deviation} = S = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} = 9.32$$

$$\text{Coefficient of Variation} = CV = \frac{S}{\bar{x}} * 100\% = 23\%$$

- Standard deviation (S) is for tree building (branching)
- Coefficient of deviation (CV) is used to determine branching termination.
- Average (Avg) is the value in the leaf nodes

## From information gain to “deviation” reduction

- Standard deviation for **two** attributes

$$S(T, X) = \sum_{c \in X} P(c)S(c)$$

		Hours Played (StDev)	Count
Outlook	Overcast	3.49	4
	Rainy	7.78	5
	Sunny	10.87	5
			14



$$\begin{aligned} S(\text{Hours}, \text{Outlook}) &= P(\text{Overcast}) * S(\text{Overcast}) + P(\text{Rainy}) * S(\text{Rainy}) + P(\text{Sunny}) * S(\text{Sunny}) \\ &= (4/14) * 3.49 + (5/14) * 7.78 + (5/14) * 10.87 \\ &= 7.66 \end{aligned}$$

## Standard Deviation Reduction (SDR)

- Calculate the standard deviation for each branch

		Hours Played (StDev)
Outlook	Overcast	3.49
	Rainy	7.78
	Sunny	10.87
		SDR=1.66

		Hours Played (StDev)
Temp.	Cool	10.51
	Hot	8.95
	Mild	7.65
		SDR=0.17

		Hours Played (StDev)
Humidity	High	9.36
	Normal	8.37
		SDR=0.28

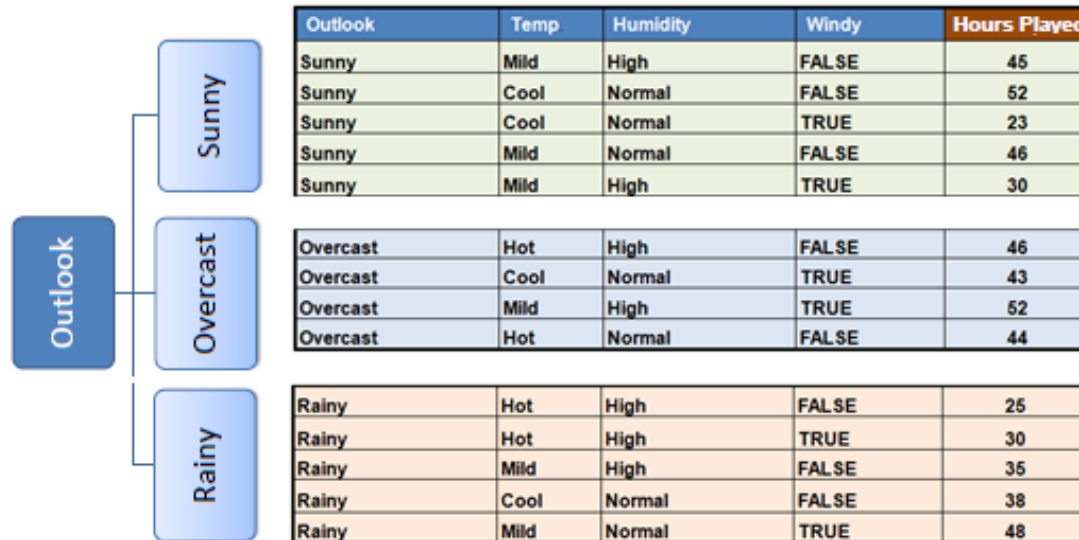
		Hours Played (StDev)
Windy	False	7.87
	True	10.59
		SDR=0.29

$$SDR(T, X) = S(T) - S(T, X)$$

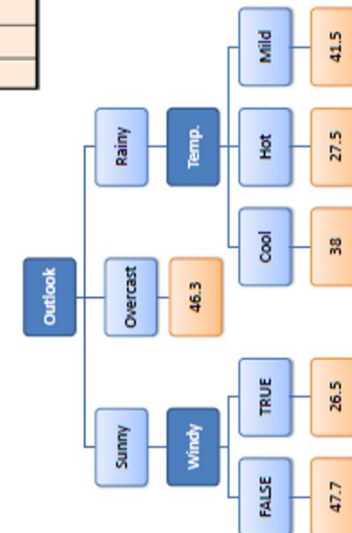
$$\begin{aligned} \text{SDR}(\text{Hours}, \text{Outlook}) &= S(\text{Hours}) - S(\text{Hours}, \text{Outlook}) \\ &= 9.32 - 7.66 = 1.66 \end{aligned}$$

# Constructing Regression Tree Based on SDR

- Then, select the attribute with the largest SDR.



- Continue to split and apply the stopping criteria. Ending up with

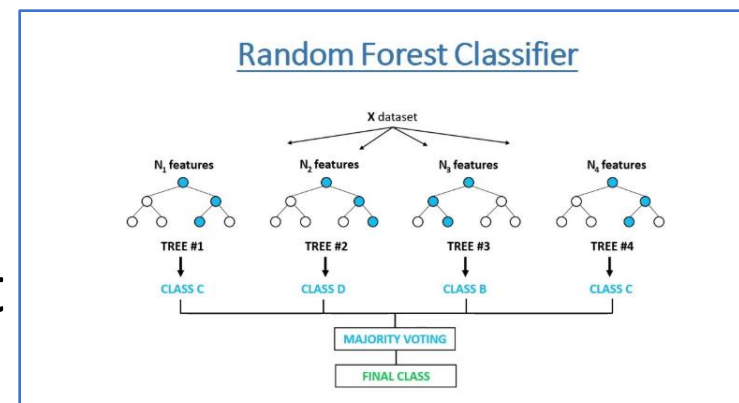


# Take-home messages

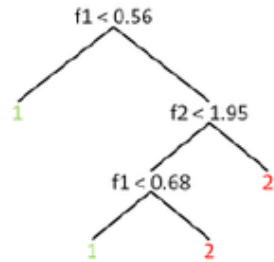
- Decision tree is a highly popular model in both data analytics and machine learning.
- The idea is very simple with tons of extensions (which involves many PhDs' hard works).
- Decision Trees or CART are an easily understandable and transparent method for predicting or classifying new records.
- If the feature engineering work is good, a simple model like DT can attain very good performance. As a data scientist/analyst, one may need to do a better **feature engineering work**. As a machine learning scientist/engineer, you may need to do a better modeling job (e.g. [Random Forest](#) and XGBoost).

# An Extension

- Decision Tree -> Random Forest



tree 1

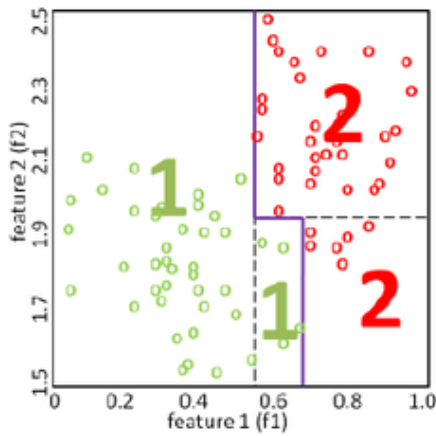


tree 500

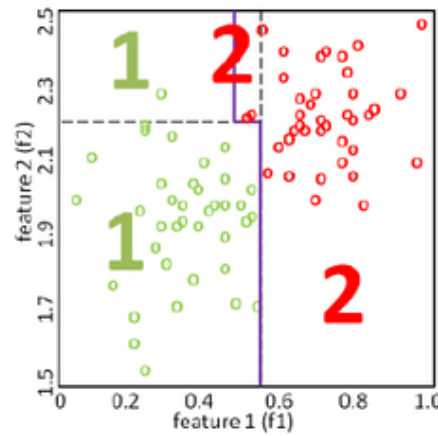


...

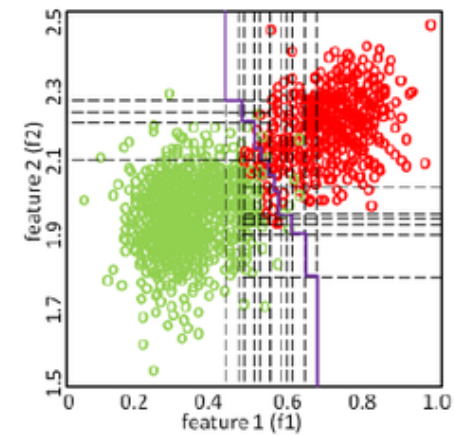
forest



Data subset 1



Data subset 500



Full dataset