



جامعة الجلالة
GALALA UNIVERSITY

Artificial Intelligence

Chapter 4: Learning from Examples

Forms of Learning

- The technology of machine learning has become a standard part of software engineering.
- Any time you are building a software system, even if you don't think of it as an AI agent, components of the system can potentially be improved **with machine learning**.
 - For example, software to analyze images of galaxies with a machine-learned model.



Forms of Learning

- An agent is **learning** if it improves its performance after making observations about the world.
- When the **agent is a computer**,
 - we call it **machine learning**: a computer observes some data, builds a model based on the data, and uses it.
- Any component of an agent program can be improved by machine learning.



Types of learning

- In **supervised learning** the agent observes input-output pairs and learns a function that maps from input to output.
 - For example, the inputs could be camera images, each one is either “**bus**” or “**pedestrian**,” etc. An output like this is called a **label**.
- In **unsupervised learning** the agent learns patterns in the input without any explicit feedback.
 - The most common unsupervised learning task is **clustering**.
 - For example, when shown millions of images taken from the Internet, a computer vision system can identify a large cluster of similar images “cats.”



Types of learning

- In **reinforcement learning** the agent learns from a series of reinforcements: **rewards** and **punishments**.
 - For example, at the end of a chess game the agent is told that it has won (a reward) or lost (a punishment).



Supervised Learning

- Given a training set of example input–output pairs

$$(x_1, y_1), (x_2, y_2), \dots (x_N, y_N),$$

- where each pair was generated by an unknown function $y = f(x)$
discover a function that approximates the true function f
- We can say y is the **ground truth**
- We can evaluate that with a second sample of (x_i, y_i) pairs called a **test set**.



Prediction Problems: Classification vs. Numeric Prediction

- **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.
- **Numeric Prediction**
 - models continuous-valued functions, i.e., predicts unknown or missing values
- **Typical applications**
 - Credit/loan approval:
 - Medical diagnosis: if a tumor is cancerous or benign.
 - Web page categorization: which category it is.

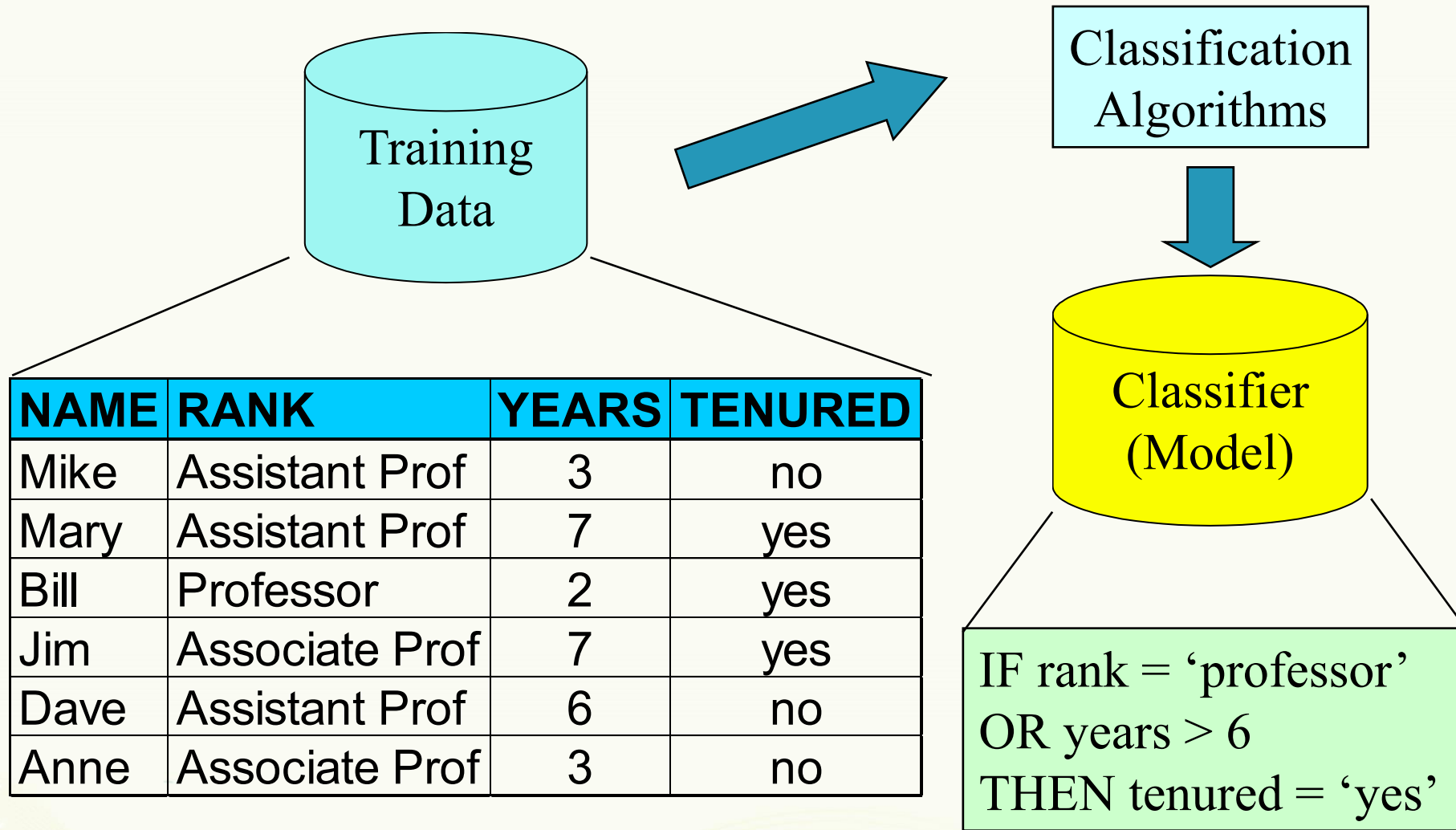


Classification—A Two-Step Process

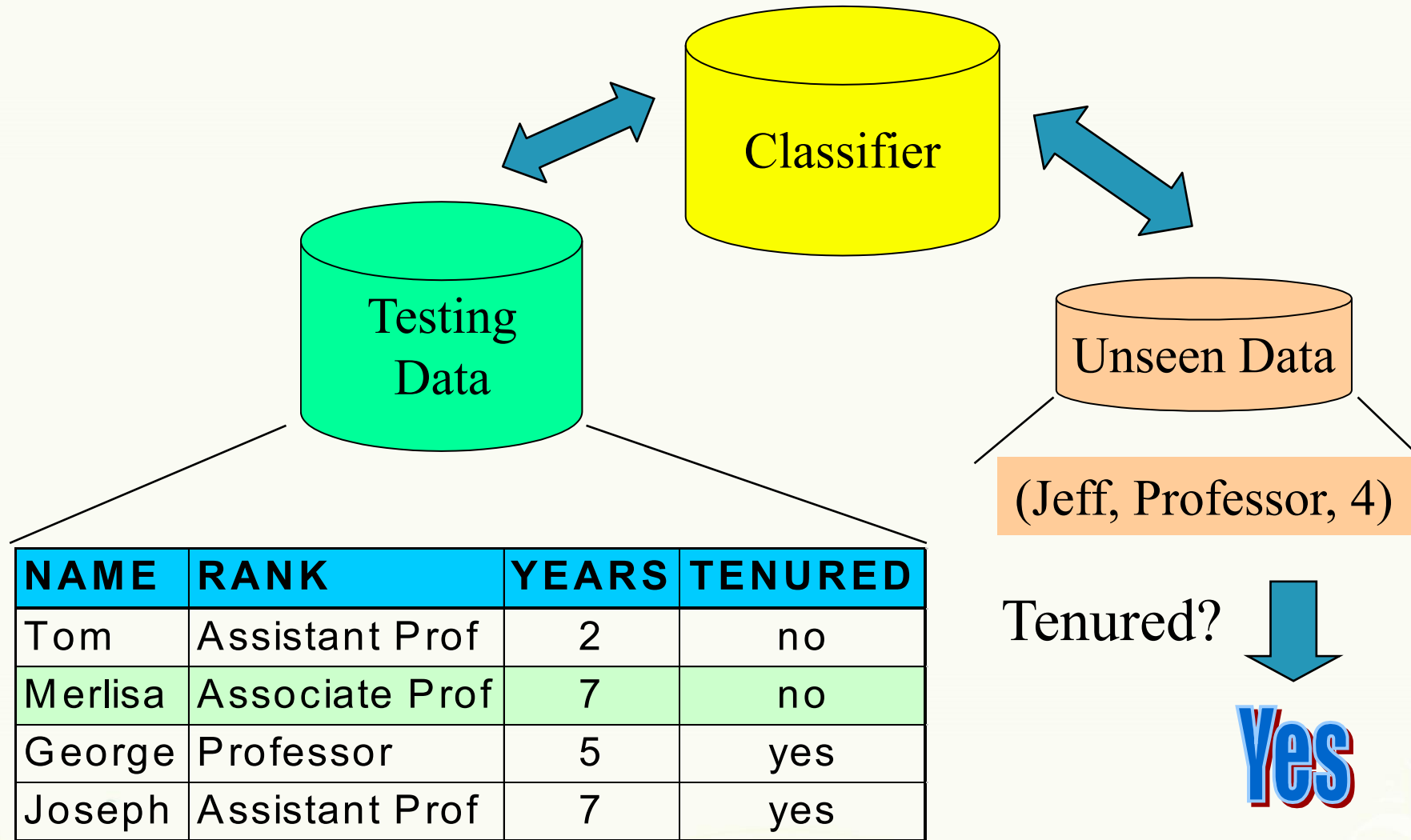
- **Model construction**: describing a set of predetermined classes
 - Each sample is assumed to belong to a predefined class, as determined by the **class label attribute**
 - The set of tuples used for model construction is **training set**
 - The model is represented as classification rules, decision trees, or mathematical formula
- **Model usage**: for classifying future or unknown objects
 - **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 overfitting)
 - If the accuracy is acceptable, use the model to **classify new data**
- Note: If *the test set* is used to select models, it is called **validation (test) set**



Process (1): Model Construction



Process (2): Using the Model in Prediction



Learning Decision Trees

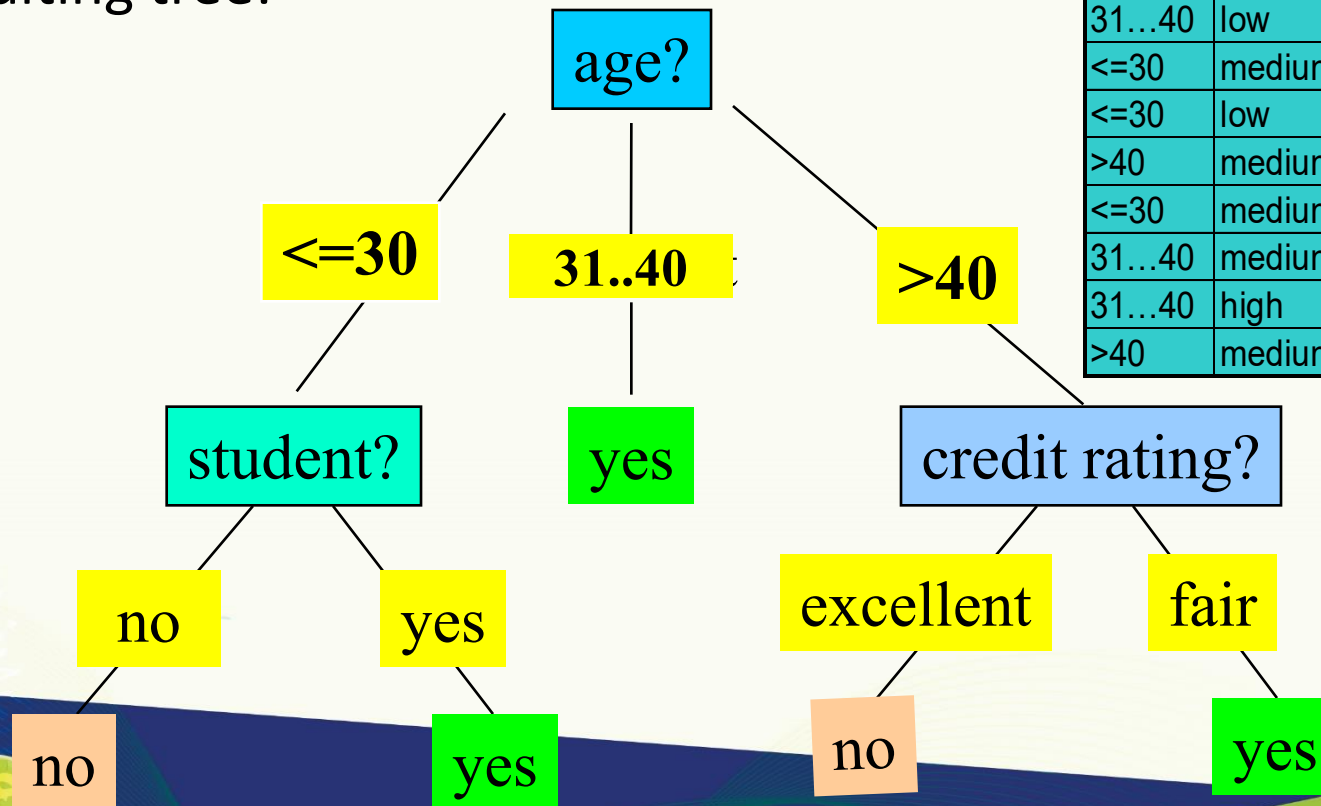
- A **decision tree** is a representation of a function that **maps a vector of attribute values** to a **single output value**—a “decision.”
- A decision tree reaches its decision by performing a sequence of tests, **starting at the root** and following the appropriate branch until a leaf is reached.



Decision Tree Induction: An Example

- ❑ Training data set: Buys_computer
- ❑ The data set follows an example of Quinlan's ID3 (Playing Tennis)
- ❑ Resulting tree:

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



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



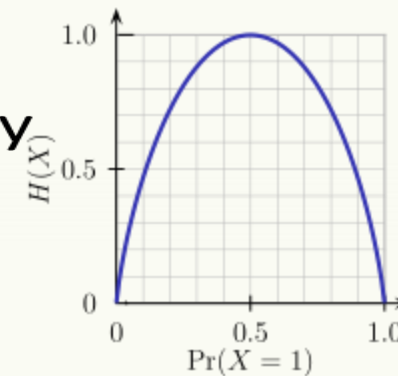
Brief Review of Entropy

- Entropy (Information Theory)

- A measure of uncertainty associated with a random variable
- Calculation: For a discrete random variable Y taking m distinct values $\{y_1, \dots, y_m\}$,
 - $H(Y) = -\sum_{i=1}^m p_i \log(p_i)$, where $p_i = P(Y = y_i)$
- Interpretation:
 - Higher entropy => higher uncertainty
 - Lower entropy => lower uncertainty

- Conditional Entropy

- $H(Y|X) = \sum_x p(x) H(Y|X = x)$



$m = 2$



Attribute Selection Measure: Information Gain (ID3/C4.5)

- Select the attribute with the highest information gain
- Let p_i be the probability that an arbitrary sample in D belongs to class C_i , estimated by $|C_{i,D}|/|D|$.

- Expected information (entropy) needed to classify a tuple in D :

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

- Information needed (after using A to split D into v partitions) to classify D :

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

- Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$



Attribute Selection: Information Gain

■ Class P: buys_computer = “yes”

■ Class N: buys_computer = “no”

$$Info(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

age	p _i	n _i	I(p _i , n _i)
<=30	2	3	0.971
31...40	4	0	0
>40	3	2	0.971

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

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

$\frac{5}{14} I(2,3)$ means “age <=30” has 5 out of 14 samples, with 2 yes’es and 3 no’s.

Hence

$$Gain(age) = Info(D) - Info_{age}(D) = 0.246$$

Similarly,

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit_rating) = 0.048$$



#Importing required libraries

```
import pandas as pd
import numpy as np
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
data = load_iris() print('Classes to predict: ', data.target_names)
#Extracting data attributes
X = data.data
### Extracting target/ class labels
y = data.target

print('Number of examples in the data:', X.shape[0])
#First four rows in the variable 'X'
X[:4]
```

#Output

```
Out: array([[5.1, 3.5, 1.4, 0.2],
           [4.9, 3. , 1.4, 0.2],
           [4.7, 3.2, 1.3, 0.2],
           [4.6, 3.1, 1.5, 0.2]])
```




```
#Using the train_test_split to create train and test sets.  
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state = 47, test_size = 0.25)  
#Importing the Decision tree classifier from the sklearn library.  
from sklearn.tree import DecisionTreeClassifier  
clf = DecisionTreeClassifier(criterion = 'entropy')
```

```
#Training the decision tree classifier.  
clf.fit(X_train, y_train)
```

```
#Output:
```

```
Out:DecisionTreeClassifier(class_weight=None, criterion='entropy', max_depth=None,  
    max_features=None, max_leaf_nodes=None,  
    min_impurity_decrease=0.0, min_impurity_split=None,  
    min_samples_leaf=1, min_samples_split=2,  
    min_weight_fraction_leaf=0.0, presort=False, random_state=None,  
    splitter='best')
```

```
#Predicting labels on the test set.  
y_pred = clf.predict(X_test)
```




```
#Importing the accuracy metric from sklearn.metrics library
```

```
from sklearn.metrics import accuracy_score
```

```
print('Accuracy Score on train data: ', accuracy_score(y_true=y_train, y_pred=clf.predict(X_train)))
```

```
print('Accuracy Score on test data: ', accuracy_score(y_true=y_test, y_pred=y_pred))
```

```
#Output:
```

```
Out: Accuracy Score on train data: 1.0
```

```
Accuracy Score on test data: 0.9473684210526315
```



Gain Ratio for Attribute Selection (C4.5)

- Information gain measure is **biased** towards attributes with a large number of values.
- C4.5 (a successor of ID3) **uses gain ratio** to overcome the problem (normalization to information gain).

$$SplitInfo_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

- $GainRatio(A) = Gain(A)/SplitInfo(A)$
- Ex.
$$SplitInfo_{income}(D) = -\frac{4}{14} \times \log_2\left(\frac{4}{14}\right) - \frac{6}{14} \times \log_2\left(\frac{6}{14}\right) - \frac{4}{14} \times \log_2\left(\frac{4}{14}\right) = 1.557.$$
 - $gain_ratio(income) = 0.029/1.557 = 0.019$
- The attribute with the **maximum gain ratio** is **selected** as the splitting **attribute**

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

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

- Calculate **Entropy** of Class attribute:

$$SplitInfo_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

$$Gain(A) = Info(D) - Info_A(D)$$

$$GainRatio(A) = \frac{Gain(A)}{SplitInfo(A)}$$

buys_computer

yes	no
9	5

$$Info(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = \underline{0.9403}$$

- Calculate **Gain Ratio** of all other attributes:

$$-\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

		Class		
		yes	no	
age	youth	2	3	5
	middle_aged	4	0	4
	senior	3	2	5
				14

$$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2)$$

$$= \frac{5}{14} \cdot 0.971 + \frac{4}{14} \cdot 0 + \frac{5}{14} \cdot 0.971 = 0.3467 + 0 + 0.3467 = \underline{0.6934}$$

$$Gain(age) = Info(D) - Info_{age}(D) = 0.9403 - 0.6934 = \underline{0.2469}$$

$$SplitInfo_{age}(D) = -\frac{5}{14} \log_2\left(\frac{5}{14}\right) - \frac{4}{14} \log_2\left(\frac{4}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = \underline{1.5774}$$

$$GainRatio(age) = \frac{Gain(A)}{SplitInfo(A)} = \frac{0.246}{1.5774} = \underline{0.1559}$$

		Class		
		yes	no	
income	low	3	1	4
	medium	4	2	6
	high	2	2	4
				14

$$Info_{income}(D) = \frac{4}{14} I(3,1) + \frac{6}{14} I(4,2) + \frac{4}{14} I(2,2)$$

$$= \frac{4}{14} \cdot 0.8113 + \frac{6}{14} \cdot 0.9183 + \frac{4}{14} \cdot 1 = 0.2318 + 0.3935 + 0.2857 = \underline{0.911}$$

$$Gain(income) = 0.9403 - 0.911 = \underline{0.0293}$$

$$SplitInfo_{income}(D) = -\frac{4}{14} \log_2\left(\frac{4}{14}\right) - \frac{6}{14} \log_2\left(\frac{6}{14}\right) - \frac{4}{14} \log_2\left(\frac{4}{14}\right) = \underline{1.5566}$$

$$GainRatio(income) = \frac{0.0293}{1.5566} = \underline{0.0188}$$

age	income	student	credit_rating	buys_computer
youth	high	no	fair	no
youth	high	no	excellent	no
middle_aged	high	no	fair	yes
senior	medium	no	fair	yes
senior	low	yes	fair	yes
senior	low	yes	excellent	no
middle_aged	low	yes	excellent	yes
youth	medium	no	fair	no
youth	low	yes	fair	yes
senior	medium	yes	fair	yes
youth	medium	yes	excellent	yes
middle_aged	medium	no	excellent	yes
middle_aged	high	yes	fair	yes
senior	medium	no	excellent	no

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

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

$$SplitInfo_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

$$Gain(A) = Info(D) - Info_A(D)$$

$$GainRatio(A) = \frac{Gain(A)}{SplitInfo(A)}$$

- Calculate **Entropy** of Class attribute:

buys_computer	
yes	no
9	5

$$Info(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = \mathbf{0.9403}$$

- Calculate **Gain Ratio** of all other attributes:

		Class		
		yes	no	
student	yes	6	1	7
	no	3	4	7
				14

$$Info_{student}(D) = \frac{7}{14} I(6,1) + \frac{7}{14} I(3,4)$$

$$= \frac{7}{14} * 0.5917 + \frac{7}{14} * 0.9852 = 0.2958 + 0.4926 = \mathbf{0.7884}$$

$$Gain(student) = 0.9403 - 0.7884 = \mathbf{0.1519}$$

$$SplitInfo_{student}(D) = -\frac{7}{14} * \log_2\left(\frac{7}{14}\right) - \frac{7}{14} * \log_2\left(\frac{7}{14}\right) = 1$$

$$GainRatio(student) = \frac{0.1519}{1} = \mathbf{0.1519}$$

		Class		
		yes	no	
credit_rating	fair	6	2	8
	excellent	3	3	6
				14

$$Info_{credit_rating}(D) = \frac{8}{14} I(6,2) + \frac{6}{14} I(3,3)$$

$$= \frac{8}{14} * 0.8113 + \frac{6}{14} * 1 = 0.4636 + 0.4286 = \mathbf{0.8922}$$

$$Gain(credit_rating) = 0.9403 - 0.8922 = \mathbf{0.0481}$$

$$SplitInfo_{credit_rating}(D) = -\frac{8}{14} * \log_2\left(\frac{8}{14}\right) - \frac{6}{14} * \log_2\left(\frac{6}{14}\right) = 0.9852$$

$$GainRatio(credit_rating) = \frac{0.0481}{0.9852} = \mathbf{0.0488}$$

age	income	student	credit_rating	buys_computer
youth	high	no	fair	no
youth	high	no	excellent	no
middle_aged	high	no	fair	yes
senior	medium	no	fair	yes
senior	low	yes	fair	yes
senior	low	yes	excellent	no
middle_aged	low	yes	excellent	yes
youth	medium	no	fair	no
youth	low	yes	fair	yes
senior	medium	yes	fair	yes
youth	medium	yes	excellent	yes
middle_aged	medium	no	excellent	yes
middle_aged	high	yes	fair	yes
senior	medium	no	excellent	no

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

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

- Calculate **Entropy** of Class attribute:

buys_computer	
yes	no
9	5

$$Info(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = \mathbf{0.9403}$$

$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \log_2\left(\frac{|D_j|}{|D|}\right)$$

$$Gain(A) = Info(D) - Info_A(D)$$

$$GainRatio(A) = \frac{Gain(A)}{SplitInfo(A)}$$

- Calculate **Gain Ratio** of all other attributes:

		Class		
		yes	no	
student	yes	6	1	7
	no	3	4	7
				14

$$Info_{student}(D) = \frac{7}{14} I(6,1) + \frac{7}{14} I(3,4)$$

$$= \frac{7}{14} * 0.5917 + \frac{7}{14} * 0.9852 = 0.2958 + 0.4926 = \mathbf{0.7884}$$

$$Gain(student) = 0.9403 - 0.7884 = \mathbf{0.1519}$$

$$SplitInfo_{student}(D) = -\frac{7}{14} \log_2\left(\frac{7}{14}\right) - \frac{7}{14} \log_2\left(\frac{7}{14}\right) = 1$$

$$GainRatio(student) = \frac{0.1519}{1} = \mathbf{0.1519}$$

		Class		
		yes	no	
credit_rating	fair	6	2	8
	excellent	3	3	6
				14

$$Info_{credit_rating}(D) = \frac{8}{14} I(6,2) + \frac{6}{14} I(3,3)$$

$$= \frac{8}{14} * 0.8113 + \frac{6}{14} * 1 = 0.4636 + 0.4286 = \mathbf{0.8922}$$

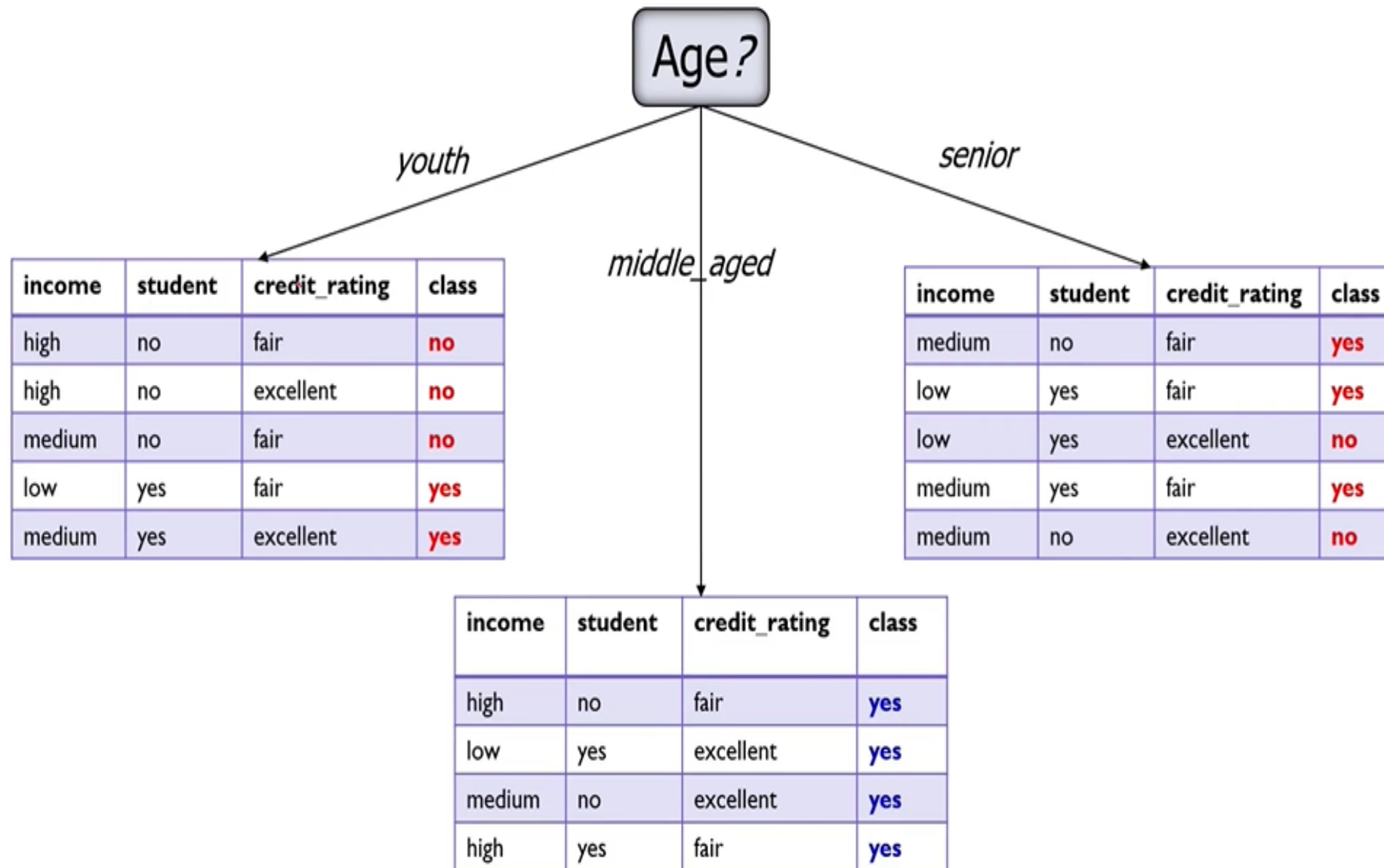
$$Gain(credit_rating) = 0.9403 - 0.8922 = \mathbf{0.0481}$$

$$SplitInfo_{credit_rating}(D) = -\frac{8}{14} \log_2\left(\frac{8}{14}\right) - \frac{6}{14} \log_2\left(\frac{6}{14}\right) = 0.9852$$

$$GainRatio(credit_rating) = \frac{0.0481}{0.9852} = \mathbf{0.0488}$$

age	income	student	credit_rating	buys_computer
youth	high	no	fair	no
youth	high	no	excellent	no
middle_aged	high	no	fair	yes
senior	medium	no	fair	yes
senior	low	yes	fair	yes
senior	low	yes	excellent	no
middle_aged	low	yes	excellent	yes
youth	medium	no	fair	no
youth	low	yes	fair	yes
senior	medium	yes	fair	yes
youth	medium	yes	excellent	yes
middle_aged	medium	no	excellent	yes
middle_aged	high	yes	fair	yes
senior	medium	no	excellent	no

- As, the Gain Ratio of "age" is highest,
- So "**age**" is the best attribute & becomes the root node of the decision tree.



- For Left subtree: Calculate **Entropy** of Class attribute:

buys_computer	
yes	no
2	3

$$Info(D) = I(2,3) = -\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) = \underline{0.971} \checkmark$$

- Calculate **Gain Ratio** of all other attributes:

		Class		
		yes	no	
income	low	1	0	1
	medium	1	1	2
	high	0	2	2
				5

$$Info_{income}(D) = \frac{1}{5} I(1,0) + \frac{2}{5} I(1,1) + \frac{2}{5} I(0,2)$$

$$= \frac{1}{5} * 0 + \frac{2}{5} * 1 + \frac{2}{5} * 0 = 0 + 0.4 + 0 = \underline{0.4} \checkmark$$

$$Gain(income) = 0.971 - 0.4 = \underline{0.571} \checkmark$$

$$SplitInfo_{income}(D) = -\frac{1}{5} * \log_2 \left(\frac{1}{5} \right) - \frac{2}{5} * \log_2 \left(\frac{2}{5} \right) - \frac{2}{5} * \log_2 \left(\frac{2}{5} \right) = \underline{1.5219} \checkmark$$

$$GainRatio(income) = \frac{0.571}{1.5219} = \underline{0.3751} \checkmark$$

		Class		
		yes	no	
credit_rating	fair	1	2	3
	excellent	1	1	2
				5

$$Info_{credit_rating}(D) = \frac{3}{5} I(1,2) + \frac{2}{5} I(1,1)$$

$$= \frac{3}{5} * 0.9183 + \frac{2}{5} * 1 = 0.3443 + 0.4 = \underline{0.7443}$$

$$Gain(credit_rating) = 0.971 - 0.7443 = \underline{0.2267}$$

$$SplitInfo_{credit_rating}(D) = -\frac{3}{5} * \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} * \log_2 \left(\frac{2}{5} \right) = \underline{0.9709}$$

$$GainRatio(credit_rating) = \frac{0.2267}{0.9709} = \underline{0.2335} \checkmark$$

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

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

$$Gain(A) = Info(D) - Info_A(D)$$

$$GainRatio(A) = \frac{Gain(A)}{SplitInfo(A)}$$

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

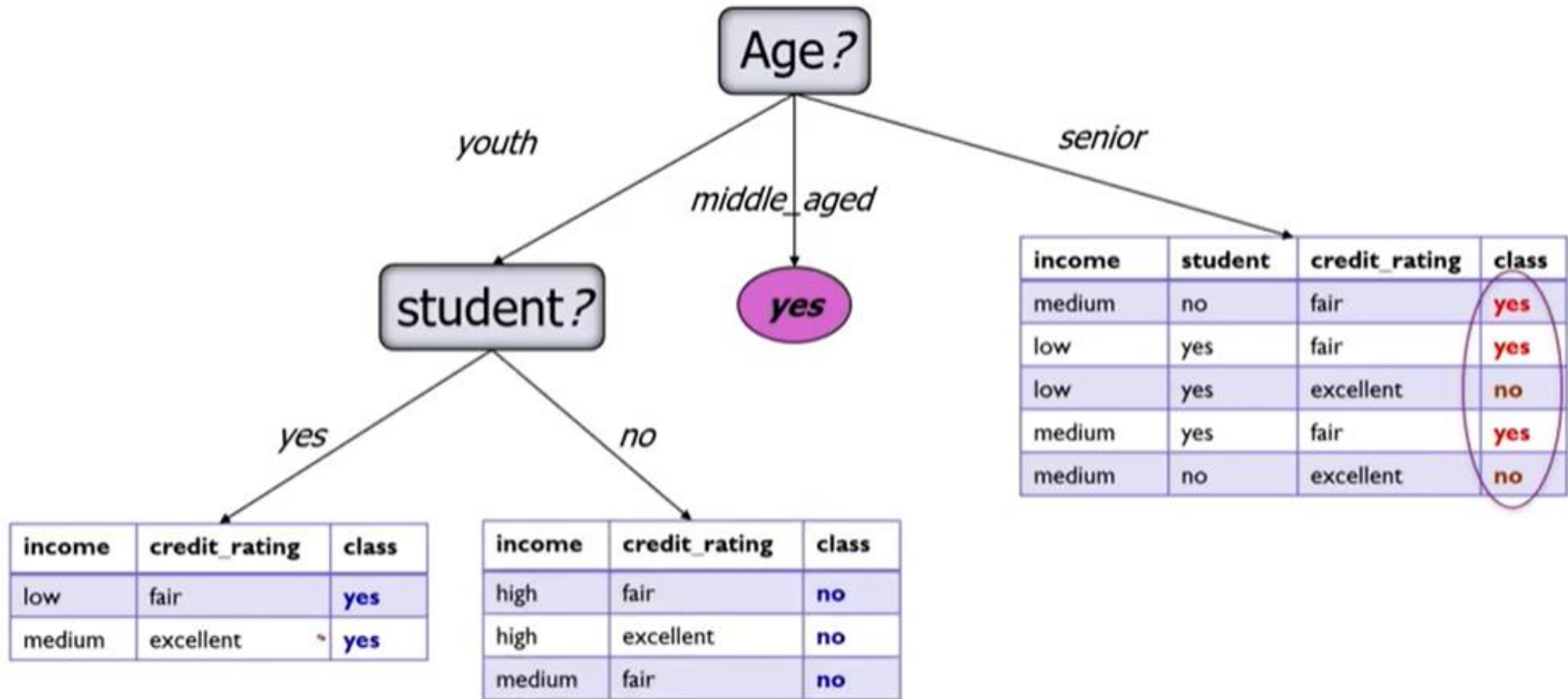
		Class		
		yes	no	
student	yes	2	0	2
	no	0	3	3
				5

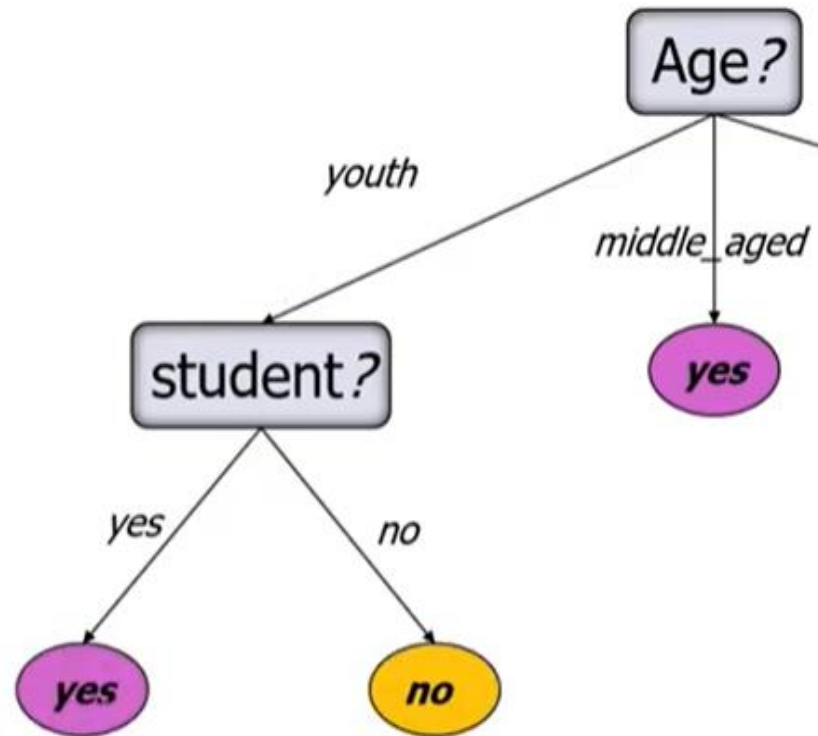
$$Info_{student}(D) = \frac{2}{5} I(2,0) + \frac{3}{5} I(0,3) = \frac{2}{5} * 0 + \frac{3}{5} * 0 = \underline{0}$$

$$Gain(student) = 0.971 - 0 = \underline{0.971}$$

$$SplitInfo_{student}(D) = -\frac{2}{5} * \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} * \log_2 \left(\frac{3}{5} \right) = \underline{0.9709}$$

$$GainRatio(student) = \frac{0.971}{0.9709} = \underline{1} \checkmark$$





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

Gain Ratio [C4.5] - Example

- For Right subtree: Calculate **Entropy** of Class attribute:

buys_computer	
yes	no
3	2

$$Info(D) = I(3,2) = -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right) = 0.971 \checkmark$$

- Calculate **Gain Ratio** of all other attributes:

		Class		
		yes	no	
income	low	1	1	2
	medium	2	1	3
	high	0	0	0
				5

$$Info_{income}(D) = \frac{2}{5} I(1,1) + \frac{3}{5} I(2,1)$$

$$= \frac{2}{5} * 1 + \frac{3}{5} * 0.9183 = 0.4 + 0.551 = 0.951$$

$$Gain(income) = 0.971 - 0.951 = 0.02$$

$$SplitInfo_{income}(D) = -\frac{2}{5} * \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} * \log_2 \left(\frac{3}{5} \right) = 0.9709$$

$$GainRatio(income) = \frac{0.02}{0.9709} = 0.0205 \checkmark$$

		Class		
		yes	no	
credit_rating	fair	3	0	3
	excellent	0	2	2
				5

$$Info_{credit_rating}(D) = \frac{3}{5} I(3,0) + \frac{2}{5} I(0,2) = \frac{3}{5} * 0 + \frac{2}{5} * 0 = 0$$

$$Gain(credit_rating) = 0.971 - 0 = 0.971$$

$$SplitInfo_{credit_rating}(D) = -\frac{3}{5} * \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} * \log_2 \left(\frac{2}{5} \right) = 0.9709$$

$$GainRatio(credit_rating) = \frac{0.971}{0.9709} = 1$$

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

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

$$Gain(A) = Info(D) - Info_A(D)$$

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

		Class		
		yes	no	
student	yes	2	1	3
	no	1	1	2
				5

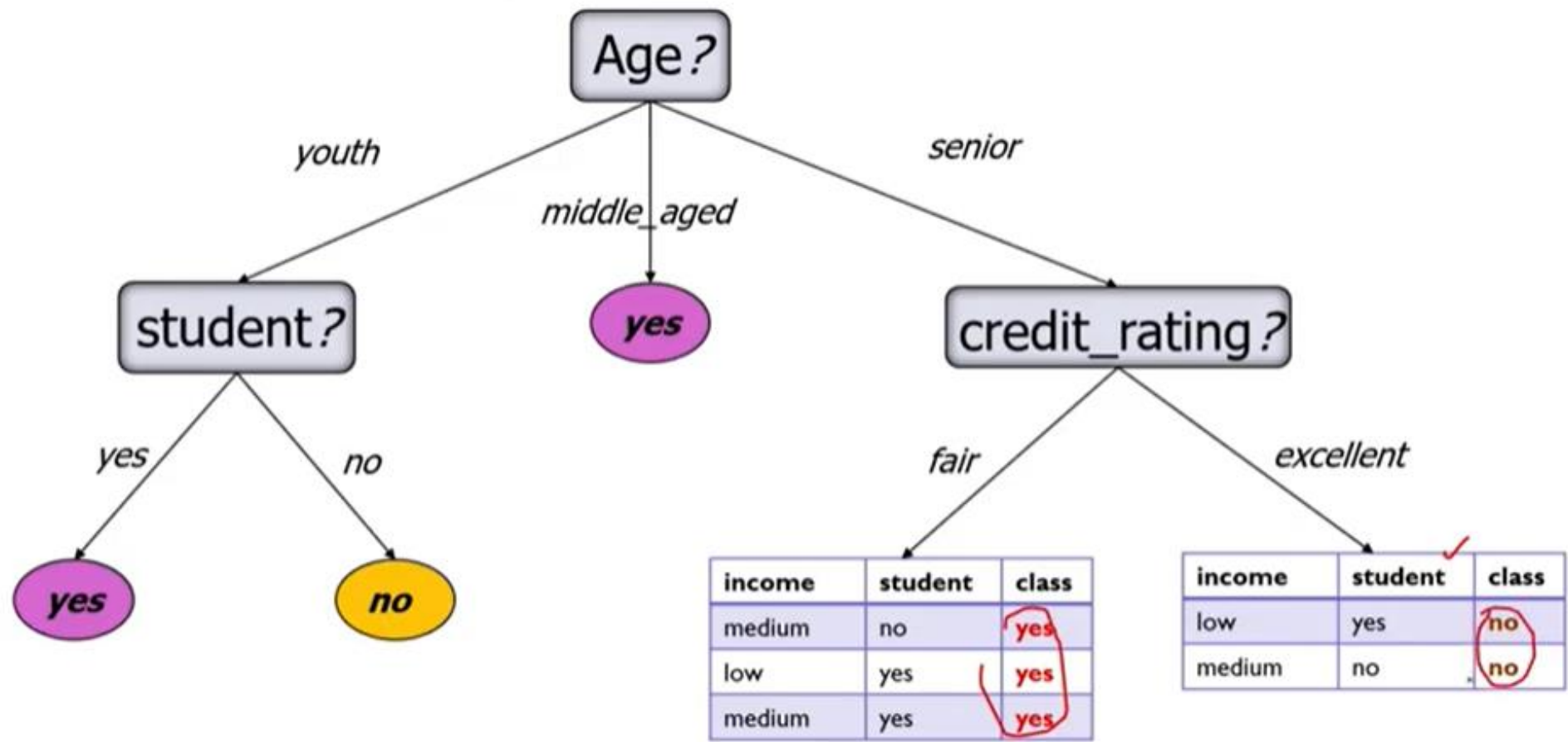
$$Info_{student}(D) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1)$$

$$= \frac{3}{5} * 0.9183 + \frac{2}{5} * 1 = 0.551 + 0.4 = 0.951$$

$$Gain(student) = 0.971 - .951 = 0.02$$

$$SplitInfo_{student}(D) = -\frac{3}{5} * \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} * \log_2 \left(\frac{2}{5} \right) = 0.9709$$

$$GainRatio(student) = \frac{0.02}{0.9709} = 0.0205$$



What is the decision for

- $X=[\text{age, income, student, credit}]=[15,\text{low},\text{no},\text{excellent}]$
- $X=[\text{age, income, student, credit}]=[40,\text{low},\text{no},\text{excellent}]$



Comparing Attribute Selection Measures

- The two measures, in general, return good results but
 - **Information gain:**
 - biased towards multivalued attributes
 - **Gain ratio:**
 - tends to prefer unbalanced splits in which one partition is much smaller than the others

