# COSC 3337 : Data Science I

# N. Rizk

College of Natural and Applied Sciences

Department of Computer Science

University of Houston

Introduction _to_Classification

# 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
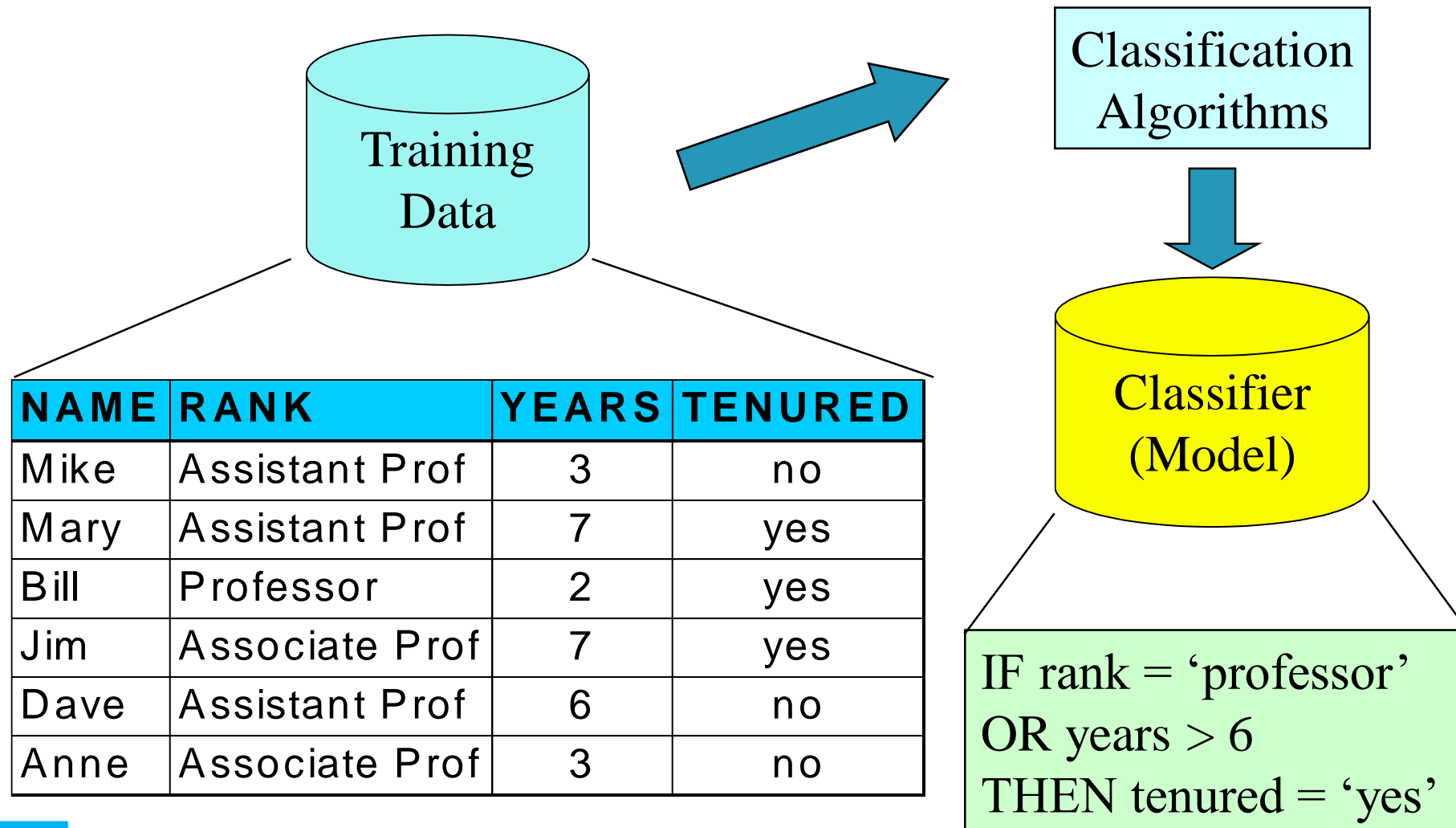
# Classification vs. Prediction

- 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

- Prediction
    - models continuous-valued functions, i.e., predicts unknown or missing values

- Typical applications
    - Credit approval
    - Target marketing
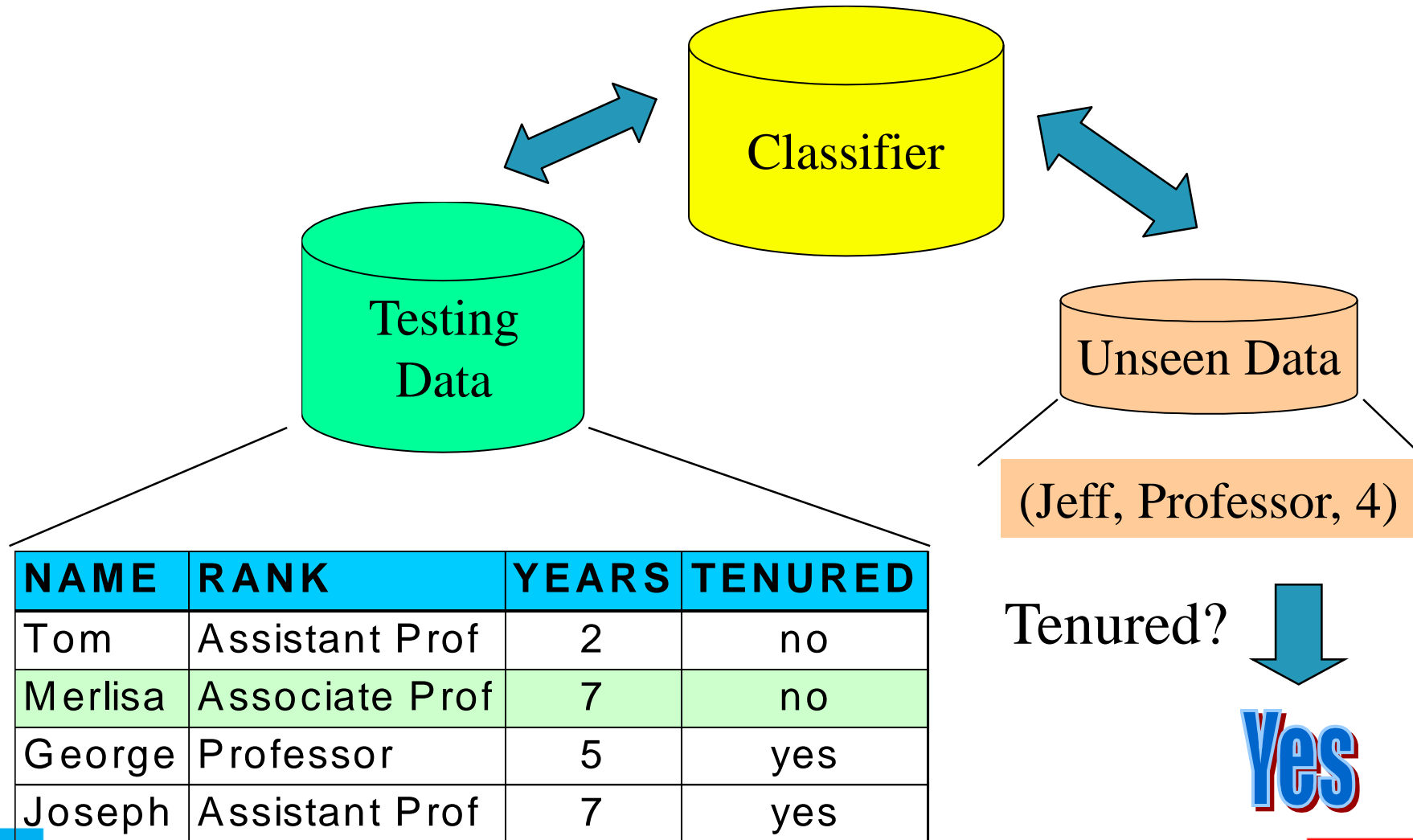    - Medical diagnosis
    - Fraud detection

Introduction _to_Classification

# 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 is training set
  - The model is represented as classification rules, decision trees, or mathematical formulae
- 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 over-fitting will occur
  - If the accuracy is acceptable, use the model to classify data tuples whose class labels are not known

# Process (1): Model Construction



Training Data

| 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 |

Classification Algorithms

Classifier (Model)

IF rank = 'professor'
OR years > 6
THEN tenured = 'yes'

N.Rizk (University of Houston)

Introduction _to_Classification

COSC 3337:DS 1

| NAME | RANK | YEARS | TENURED |
|------|------|-------|---------|
| Tom | Assistant Prof | 2 | no |
| Merlisa | Associate Prof | 7 | no |
| George | Professor | 5 | yes |
| Joseph | Assistant Prof | 7 | yes |

Classifier

Testing Data

Unseen Data

(Jeff, Professor, 4)

Tenured?

Yes

Introduction _to_Classification

COSC 3337:DS 1

# Issues: 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
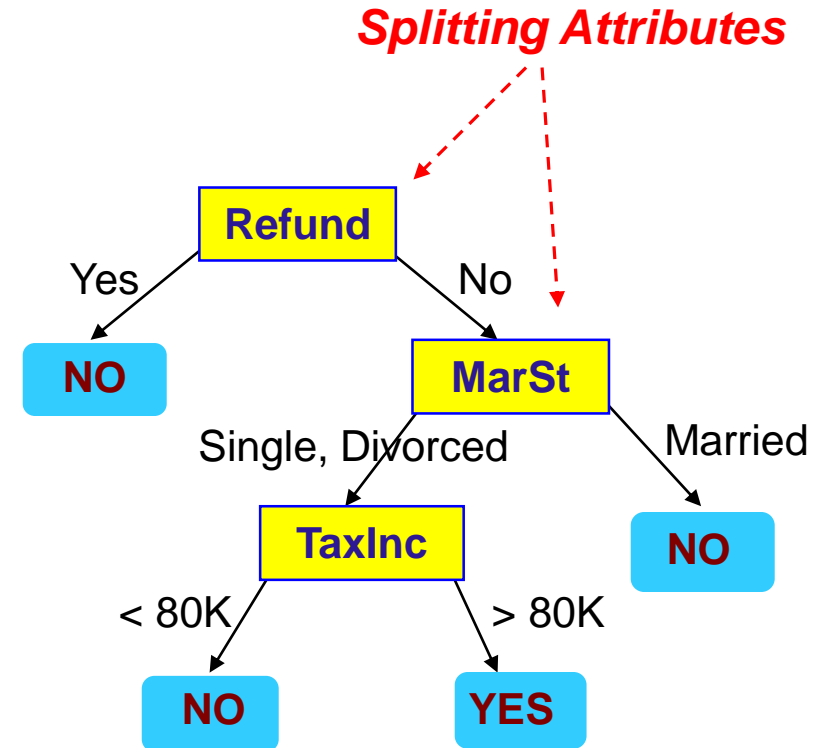  - Generalize and/or normalize data

Introduction _to_Classification

# Issues: Evaluating Classification Methods

- Accuracy
  - classifier accuracy: predicting class label
  - predictor accuracy: guessing value of predicted attributes
- Speed
  - time to construct the model (training time)
  - time to use the model (classification/prediction time)
- Robustness: handling noise and missing values
- Scalability: efficiency in disk-resident databases
- Interpretability
  - understanding and insight provided by the model
- Other measures, e.g., goodness of rules, such as decision tree size or compactness of classification rules

COSC 3337:DS 1

# Example of a Decision Tree

categorical

categorical

continuous

class

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

**Training Data**

*Splitting Attributes*

```
            Refund
          Yes /    \ No
            /        \
          NO        MarSt
                Single, Divorced / \ Married
                        /           \
                     TaxInc         NO
                  < 80K /  \ > 80K
                      /      \
                    NO       YES
```

**Model:  Decision Tree**

Introduction _to_Classification

# Another Example of Decision Tree



|     | categorical | categorical | continuous | class |
|-----|-------------|-------------|------------|-------|
| Tid | Refund | Marital Status | Taxable Income | Cheat |
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

**There could be more than one tree that fits the same data!**

Introduction _to_Classification

# Decision Tree Classification Task

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 1 | Yes | Large | 125K | No |
| 2 | No | Medium | 100K | No |
| 3 | No | Small | 70K | No |
| 4 | Yes | Medium | 120K | No |
| 5 | No | Large | 95K | Yes |
| 6 | No | Medium | 60K | No |
| 7 | Yes | Large | 220K | No |
| 8 | No | Small | 85K | Yes |
| 9 | No | Medium | 75K | No |
| 10 | No | Small | 90K | Yes |

**Training Set**

Tree Induction algorithm

Induction

**Learn Model**

**Model**

**Decision Tree**

**Apply Model**

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

**Test Set**

Deduction

# Apply Model to Test Data

**Test Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

Start from the root of tree.

Refund
- Yes → NO
- No → MarSt
  - Single, Divorced → TaxInc
    - < 80K → NO
    - > 80K → YES
  - Married → NO

Introduction _to_Classification

COSC 3337:DS 1

# Apply Model to Test Data

**Test Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

N.Rizk (University of Houston)

COSC 3337:DS 1

Introduction _to_Classification

# Apply Model to Test Data

**Test Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

Introduction __to_Classification

# Apply Model to Test Data

**Test Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

N.Rizk (University of Houston)

Introduction _to_Classification

COSC 3337:DS 1

# Apply Model to Test Data

**Test Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|---------------|----------------|-------|
| No | Married | 80K | ? |

COSC 3337:DS 1

Introduction _to_Classification

# Apply Model to Test Data

**Test Data**

| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No | Married | 80K | ? |

**Refund**

Yes → **NO**

No → **MarSt**

MarSt:
- Single, Divorced → **TaxInc**
- Married → **NO**

TaxInc:
- < 80K → **NO**
- > 80K → **YES**

Assign Cheat to "No"

COSC 3337:DS 1

Introduction _to_Classification

# Decision Tree Classification Task

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 1 | Yes | Large | 125K | No |
| 2 | No | Medium | 100K | No |
| 3 | No | Small | 70K | No |
| 4 | Yes | Medium | 120K | No |
| 5 | No | Large | 95K | Yes |
| 6 | No | Medium | 60K | No |
| 7 | Yes | Large | 220K | No |
| 8 | No | Small | 85K | Yes |
| 9 | No | Medium | 75K | No |
| 10 | No | Small | 90K | Yes |

**Training Set**

Tree Induction algorithm

Induction

**Learn Model**

**Model**

**Apply Model**

**Decision Tree**

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

**Test Set**

Deduction

COSC 3337:DS 1

Introduction _to_Classification

# 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

# Decision Tree Induction

- Many Algorithms:
    - Hunt's Algorithm (one of the earliest)
    - CART
    - ID3, C4.5
    - SLIQ,SPRINT

# General Structure of Hunt's Algorithm

- Let $D_t$ be the set of training records that reach a node t

- General Procedure:
  - If $D_t$ contains records that belong the same class $y_t$, then t is a leaf node labeled as $y_t$
  - If $D_t$ is an empty set, then t is a leaf node labeled by the default class, $y_d$
  - If $D_t$ contains records that belong to more than one class, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

$D_t$

?

# Hunt's Algorithm

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|---------------|---------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

Introduction _to_Classification

COSC 3337:DS 1

# Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.

- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting

# Tree Induction

- # Greedy strategy.
  - ## Split the records based on an attribute test that optimizes certain criterion.

- # Issues
  - ## Determine how to split the records
    - ### How to specify the attribute test condition?
    - ### How to determine the best split?
  - ## Determine when to stop splitting

# How to Specify Test Condition?

- Depends on attribute types
  - Nominal
  - Ordinal
  - Continuous

- Depends on number of ways to split
  - 2-way split
  - Multi-way split

# Splitting Based on Nominal Attributes

- **Multi-way split:** Use as many partitions as distinct values.

```
              CarType
Family      /    |    \      Luxury
                Sports
```

- **Binary split:** Divides values into two subsets.
  Need to find optimal partitioning.

```
{Sports,   CarType                          {Family,   CarType
Luxury}   /       \  {Family}      OR       Luxury}   /       \  {Sports}
```

# Splitting Based on Ordinal Attributes

- **Multi-way split:** Use as many partitions as distinct values.

Size
- Small
- Medium
- Large

- **Binary split:** Divides values into two subsets. Need to find optimal partitioning.

Size — {Small, Medium} / {Large}     OR     Size — {Medium, Large} / {Small}

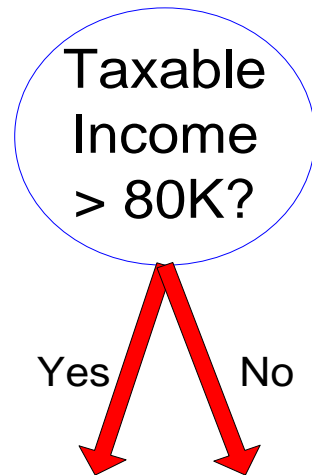- **What about this split?**

Size — {Small, Large} / {Medium}

# Splitting Based on Continuous Attributes
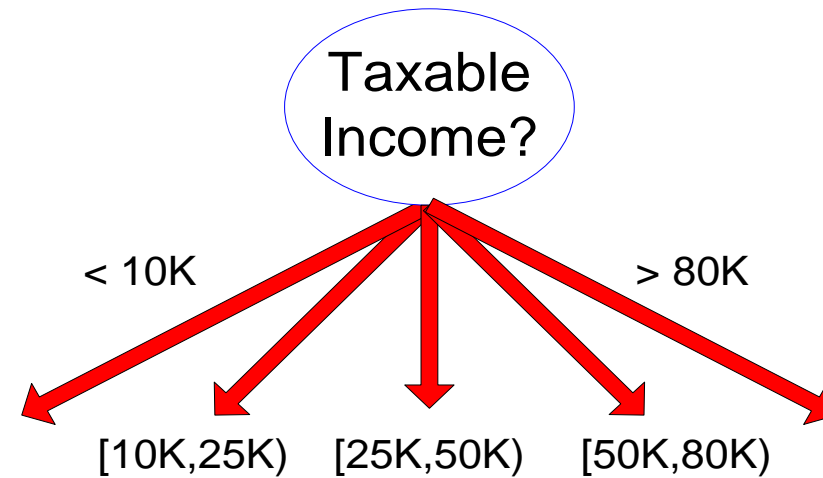
- Different ways of handling
  - Discretization to form an ordinal categorical attribute
    - Static – discretize once at the beginning
    - Dynamic – ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.

  - Binary Decision: (A < v) or (A $\geq$ v)
    - consider all possible splits and finds the best cut
    - can be more compute intensive

# Splitting Based on Continuous Attributes

Taxable Income > 80K?

Yes          No

(i) Binary split

Taxable Income?

< 10K          > 80K

[10K,25K)   [25K,50K)   [50K,80K)
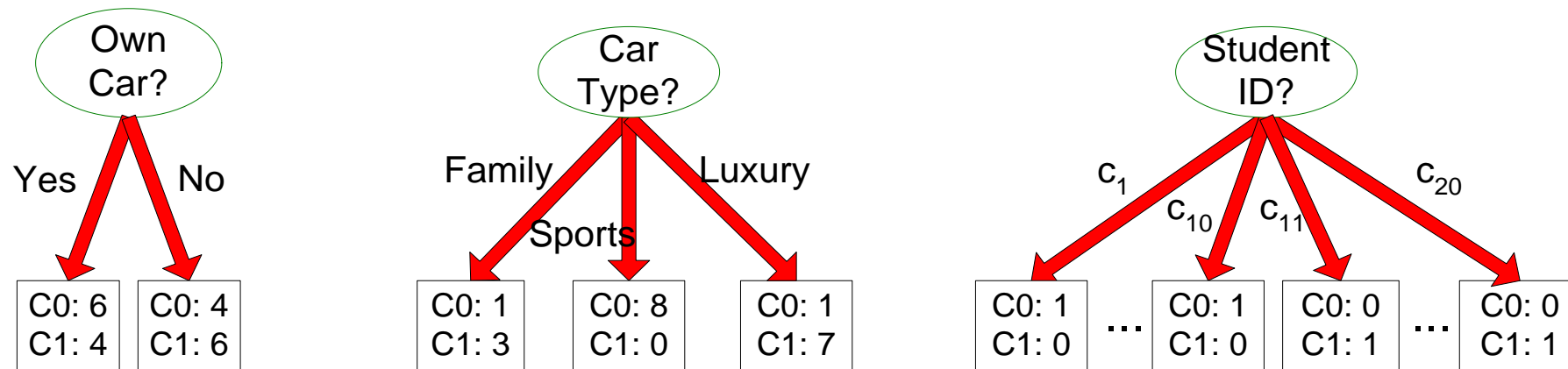
(ii) Multi-way split

# Tree Induction

- ## Greedy strategy.
  - ### Split the records based on an attribute test that optimizes certain criterion.

- ## Issues
  - ### Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - ### Determine when to stop splitting

# How to determine the Best Split

**Before Splitting: 10 records of class 0,
10 records of class 1**

Own Car?

Yes          No

| C0: 6 | C0: 4 |
|-------|-------|
| C1: 4 | C1: 6 |

Car Type?

Family          Luxury

Sports

| C0: 1 | C0: 8 | C0: 1 |
|-------|-------|-------|
| C1: 3 | C1: 0 | C1: 7 |

Student ID?

$c_1$          $c_{20}$

$c_{10}$   $c_{11}$

| C0: 1 |     | C0: 1 | C0: 0 |     | C0: 0 |
|-------|-----|-------|-------|-----|-------|
| C1: 0 | ... | C1: 0 | C1: 1 | ... | C1: 1 |

**Which test condition is the best?**

# How to determine the Best Split

- Greedy approach:
  - Nodes with homogeneous class distribution are preferred
- Need a measure of node impurity:

<table>
<tr><td>C0: 5<br>C1: 5</td><td>C0: 9<br>C1: 1</td></tr>
</table>

**Non-homogeneous,**

**High degree of impurity**

**Homogeneous,**
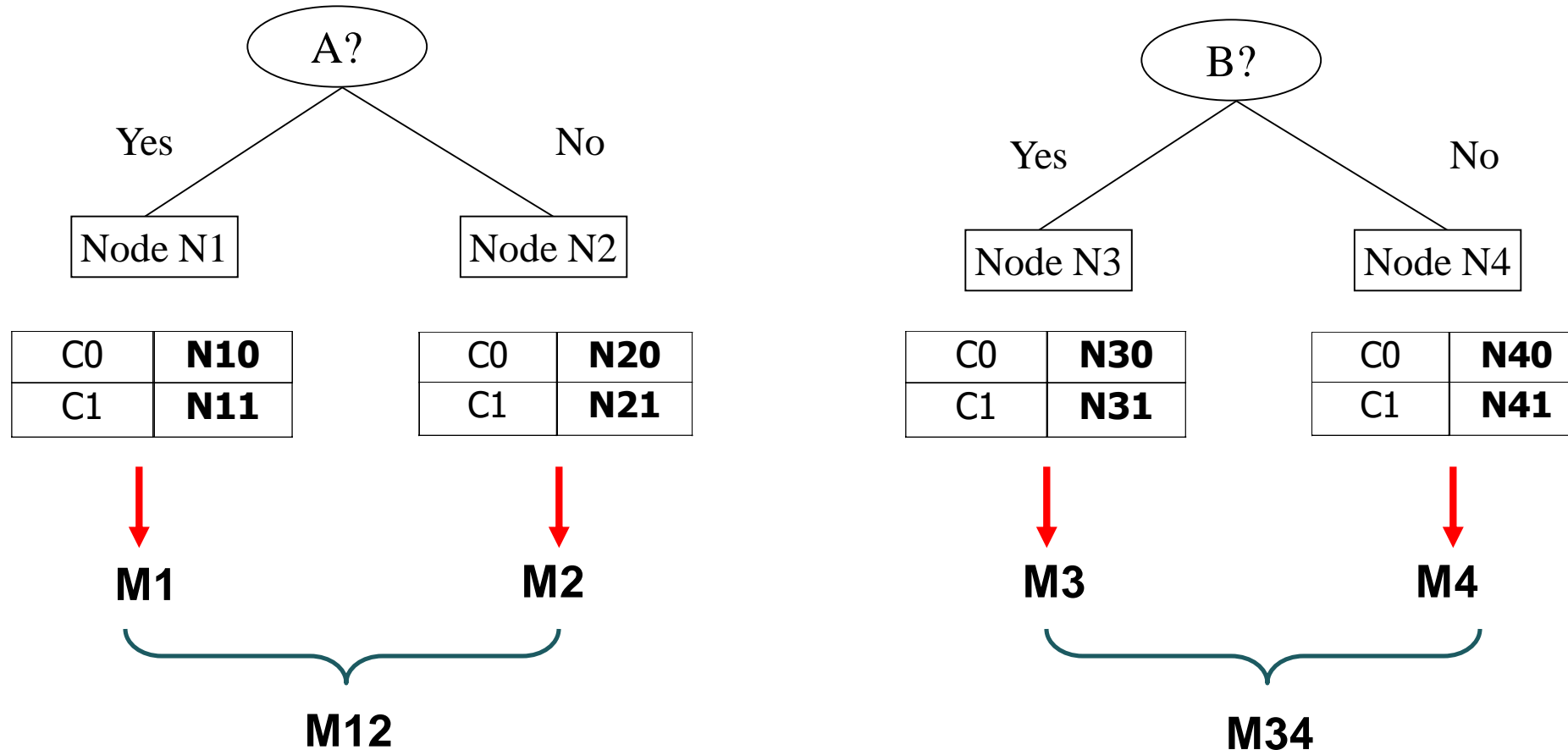
**Low degree of impurity**

# Measures of Node Impurity

- Gini Index

- Entropy

- Misclassification error

# How to Find the Best Split

**Before Splitting:**

| C0 | **N00** |
|----|---------|
| C1 | **N01** |

$\longrightarrow$ **M0**

**A?**

Yes — No

Node N1 — Node N2

| C0 | **N10** |
|----|---------|
| C1 | **N11** |

| C0 | **N20** |
|----|---------|
| C1 | **N21** |

**M1** — **M2**

**M12**

**B?**

Yes — No

Node N3 — Node N4

| C0 | **N30** |
|----|---------|
| C1 | **N31** |

| C0 | **N40** |
|----|---------|
| C1 | **N41** |

**M3** — **M4**

**M34**

**Gain = M0 – M12 vs M0 – M34**

# Examples

- Using Information Gain

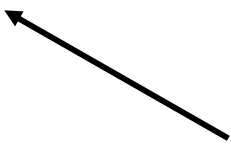# Information Gain in a Nutshell

$$InformationGain(A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} \cdot Entropy(S_v)$$

$$Entropy = \sum_{d \in Decisions} - p(d) * \log(p(d))$$

*typically yes/no*

# Playing Tennis

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

Introduction _to_Classification

# Choosing an Attribute

- We want to split our decision tree on one of the attributes

- There are four attributes to choose from:
  - Outlook
  - Temperature
  - Humidity
  - Wind

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

# How to Choose an Attribute

- Want to calculated the information gain of each attribute

- Let us start with Outlook

- What is Entropy(S)?

- -5/14*log2(5/14) – 9/14*log2(9/14)

= Entropy(5/14,9/14) = 0.9403

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

# Outlook Continued

- The expected conditional entropy is:
  5/14 * Entropy(3/5,2/5) +
  4/14 * Entropy(1,0) +
  5/14 * Entropy(3/5,2/5) = 0.6935

- So IG(Outlook) = 0.9403 – 0.6935 = 0.2468

# Temperature

- Now let us look at the attribute Temperature

- The expected conditional entropy is:
  4/14 * Entropy(2/4,2/4) +
  6/14 * Entropy(4/6,2/6) +
  4/14 * Entropy(3/4,1/4) = 0.9111

- So IG(Temperature) = 0.9403 – 0.9111 = 0.0292

# Humidity

- Now let us look at attribute Humidity

- What is the expected conditional entropy?

- 7/14 * Entropy(4/7,3/7) +
  7/14 * Entropy(6/7,1/7) = 0.7885

- So IG(Humidity) = 0.9403 – 0.7885
  = 0.1518

# Wind

- What is the information gain for wind?

- Expected conditional entropy:
  8/14 * Entropy(6/8,2/8) +
  6/14 * Entropy(3/6,3/6) = 0.8922

- IG(Wind) = 0.9403 – 0.8922 = 0.048

# Information Gains

- Outlook            0.2468

- Temperature       0.0292

- Humidity          0.1518

- Wind              0.0481

- We choose Outlook since it has the highest information gain

# Decision Tree So Far



Outlook
— Sunny    Overcast    Rain

- Now must decide what to do when Outlook is:
  - Sunny
  - Overcast
  - Rain

Introduction _to_Classification

# Sunny Branch

- Examples to classify:
- *Temperature, Humidity, Wind, Tennis*
- Hot, High, Weak, no
- Hot, High, Strong, no
- Mild, High, Weak, no
- Cool, Normal, Weak, yes
- Mild, Normal, Strong, yes

Introduction _to_Classification

# Splitting Sunny on Temperature

- What is the Entropy of Sunny?
  - Entropy(2/5,3/5) = 0.9710

- How about the expected utility?
  - 2/5 * Entropy(1,0) +
    2/5 * Entropy(1/2,1/2) +
    1/5 * Entropy(1,0) = 0.4000

- IG(Temperature) = 0.9710 – 0.4000
  = 0.5710

Introduction _to_Classification

# Splitting Sunny on Humidity

- The expected conditional entropy is

- 3/5 * Entropy(1,0) +
  2/5 * Entropy(1,0) = 0

- IG(Humidity) = 0.9710 – 0 = 0.9710

Introduction _to_Classification

# Considering Wind?

- Do we need to consider wind as an attribute?

- No – it is not possible to do any better than an expected entropy of 0; i.e. humidity must maximize the information gain

Introduction _to_Classification
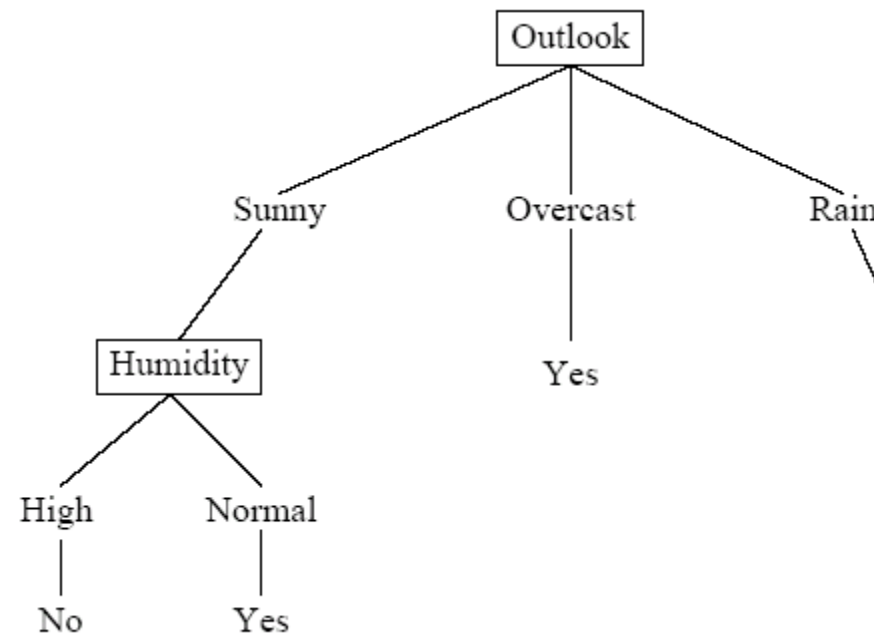
# New Tree

Introduction _to_Classification

# What if it is Overcast?

- All examples indicate yes

- So there is no need to further split on an attribute

- The information gain for any attribute would have to be 0

- Just write yes at this node

Introduction _to_Classification

# New Tree

# What about Rain?

- Let us consider attribute temperature
- First, what is the entropy of the data?
  - Entropy(3/5,2/5) = 0.9710
- Second, what is the expected conditional entropy?
  - 3/5 * Entropy(2/3,1/3) + 2/5 * Entropy(1/2,1/2) = 0.9510
- IG(Temperature) = 0.9710 – 0.9510 = 0.020

Introduction _to_Classification

# Or perhaps humidity?
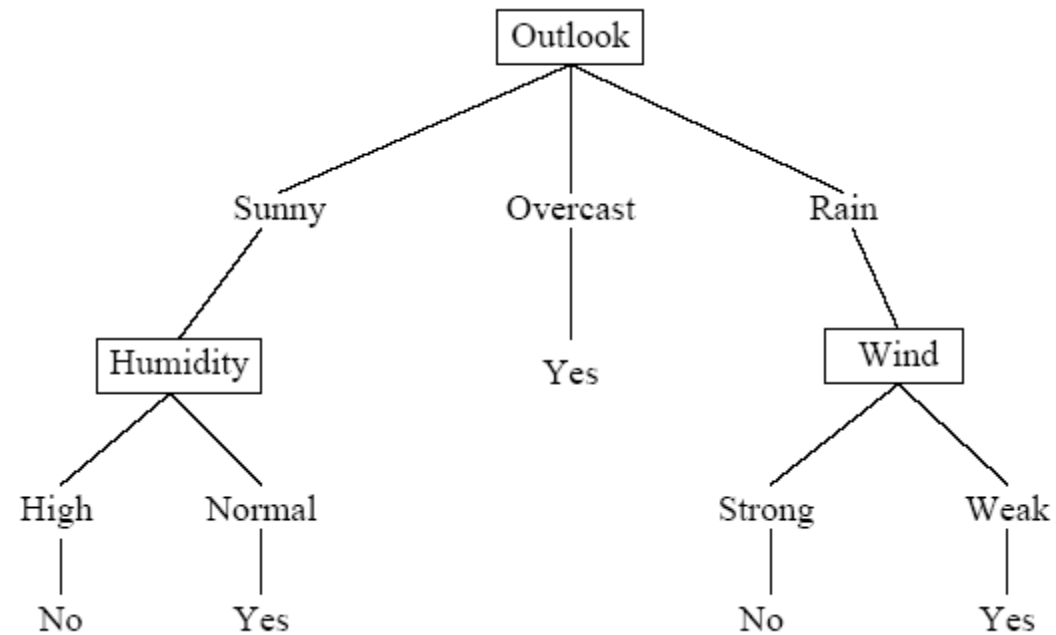
- What is the expected conditional entropy?
  - 3/5 * Entropy(2/3,1/3) + 2/5 * Entropy(1/2,1/2) = 0.9510 (the same)
- IG(Humidity) = 0.9710 – 0.9510 = 0.020 (again, the same)

Introduction _to_Classification

# Now consider wind
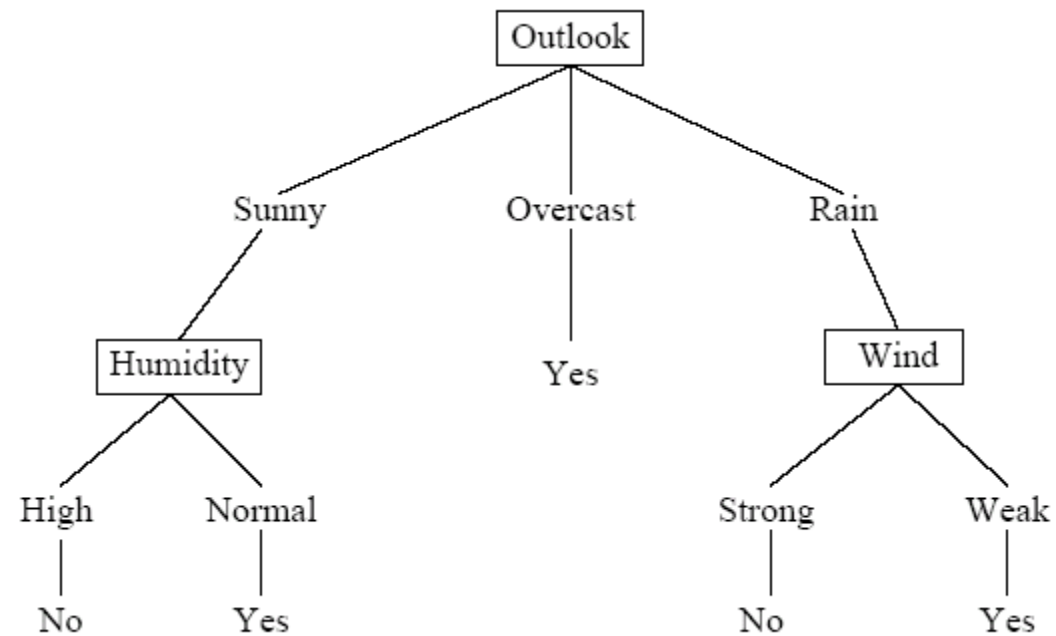
- Expected conditional entropy:
  - 3/5*Entropy(1,0) + 2/5*Entropy(1,0) = 0
- IG(Wind) = 0.9710 – 0 = 0.9710
- Thus, we split on Wind

Introduction _to_Classification

COSC 3337:DS 1

# Split Further?

Introduction _to_Classification

# Final Tree

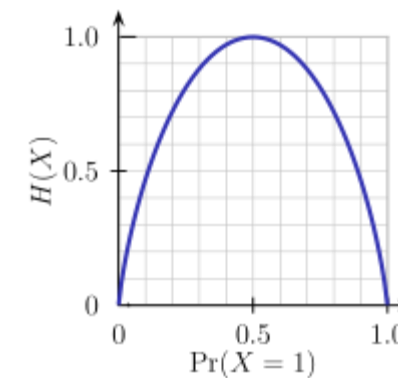Introduction _to_Classification

# Brief Review of Entropy

## Formula for Entropy :

$$H(X) = \Sigma\ p(x) * log2(p(x))$$

H(X) is Entropy, measure of uncertainty, associated with random variable X

P(x) is probability of occurrence of outcome x of variable X

Log(P(x)) is information encoded in outcome x of variable X. Based On Shannon's Information Theory



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 tuple 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:

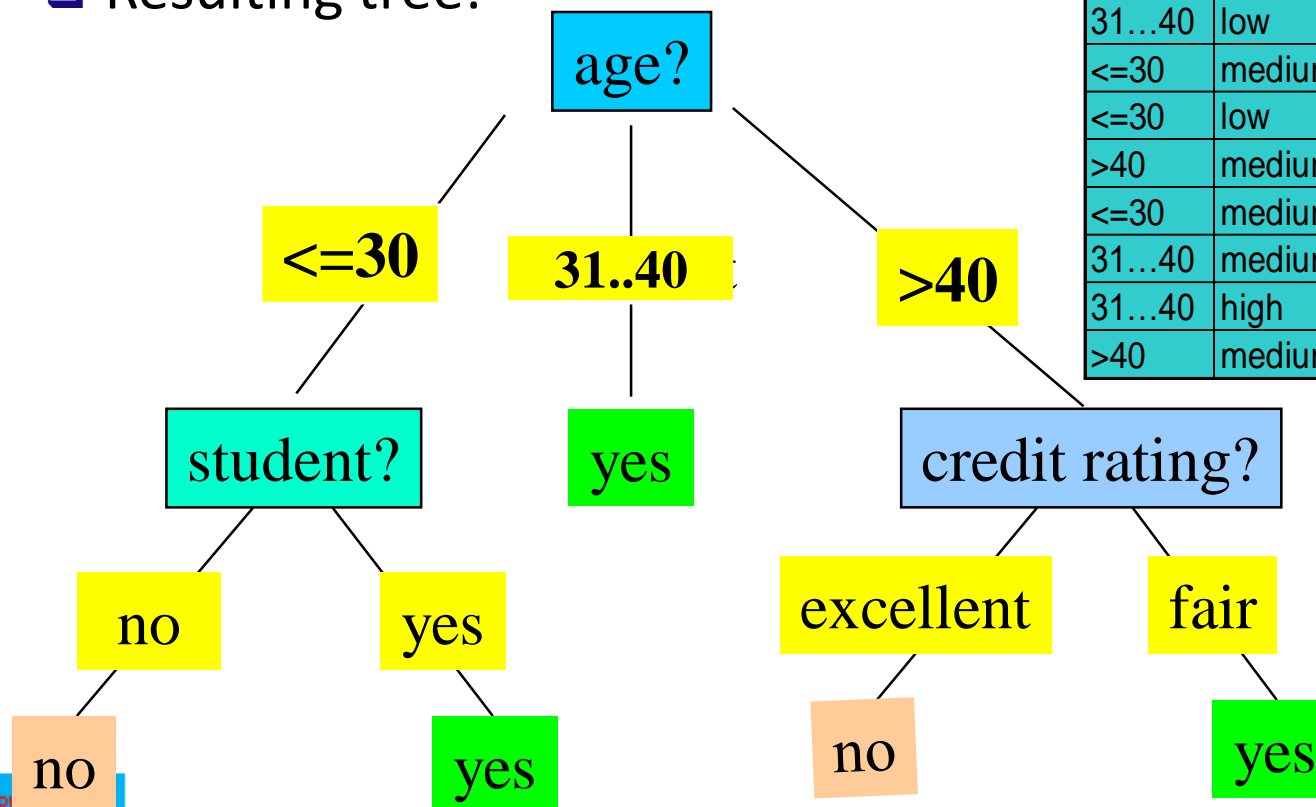$$Info_A(D) = \sum_{j=1}^{v} \frac{|D_j|}{|D|} \times Info(D_j)$$

- Information gained by branching on attribute A

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

# Decision Tree Induction: Another 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 |

age?

**<=30**      **31..40**      **>40**

student?      yes      credit rating?

no      yes      excellent      fair

no      yes      no      yes

Introduction _to_Classification

# Attribute Selection: Information Gain

■ Class P: buys_computer = "yes"
■ Class N: buys_computer = "no"

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

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

$$+ \frac{5}{14}I(3,2) = 0.694$$

| 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 |

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

| 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 |

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

Similarly,

$$Gain(income) = 0.029$$
$$Gain(student) = 0.151$$
$$Gain(credit\_rating) = 0.048$$

# Computing Information-Gain for Continuous-Value Attributes

- Let attribute A be a continuous-valued attribute

- Must determine the *best split point* for A

  - Sort the value A in increasing order

  - Typically, the midpoint between each pair of adjacent values is considered as a possible *split point*

    - $(a_i + a_{i+1})/2$ is the midpoint between the values of $a_i$ and $a_{i+1}$

  - The point with the *minimum expected information requirement* for A is selected as the split-point for A

- Split:

  - D1 is the set of tuples in D satisfying A ≤ split-point, and D2 is the set of tuples in D satisfying A > split-point

# 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)

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

  - GainRatio(A) = Gain(A)/SplitInfo(A)

- Ex. $SplitInfo_A(D) = -\frac{4}{14} \times \log_2(\frac{4}{14}) - \frac{6}{14} \times \log_2(\frac{6}{14}) - \frac{4}{14} \times \log_2(\frac{4}{14}) = 0.926$
  - gain_ratio(income) = 0.029/0.926 = 0.031

- The attribute with the maximum gain ratio is selected as the splitting attribute

# Gini index (CART, IBM IntelligentMiner)

- If a data set $D$ contains examples from $n$ classes, gini index, $gini(D)$ is defined as

$$gini(D) = 1 - \sum_{j=1}^{n} p_j^2$$

  where $p_j$ is the relative frequency of class $j$ in $D$

- If a data set $D$ is split on A into two subsets $D_1$ and $D_2$, the *gini* index $gini(D)$ is defined as

- Reduction in Impurity:

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

- 

$$\Delta gini(A) = gini(D) - gini_A(D)$$

- The attribute provides the smallest $gini_{split}(D)$ (or the largest reduction in impurity) is chosen to split the node (*need to enumerate all the possible splitting points for each attribute*)

Introduction _to_Classification

# Gini index (CART, IBM IntelligentMiner)

- Ex.  D has 9 tuples in buys_computer = "yes" and 5 in "no"

$$gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

- Suppose the attribute income partitions D into 10 in $D_1$: {low, medium} and 4 in $D_2$

$$gini_{income \in \{low, medium\}}(D) = \left(\frac{10}{14}\right)Gini(D_1) + \left(\frac{4}{14}\right)Gini(D_1)$$

$$= \frac{10}{14}\left(1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2\right) + \frac{4}{14}\left(1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2\right)$$

$$= 0.450$$

$$= Gini_{income \in \{high\}}(D)$$

but gini$_{\{medium, high\}}$ is 0.30 and thus the best since it is the lowest

- All attributes are assumed continuous-valued

- May need other tools, e.g., clustering, to get the possible split values

- Can be modified for categorical attributes

# Comparing Attribute Selection Measures

- The three 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

  - Gini index:

    - biased to multivalued attributes

    - has difficulty when # of classes is large

    - tends to favor tests that result in equal-sized partitions and purity in both partitions
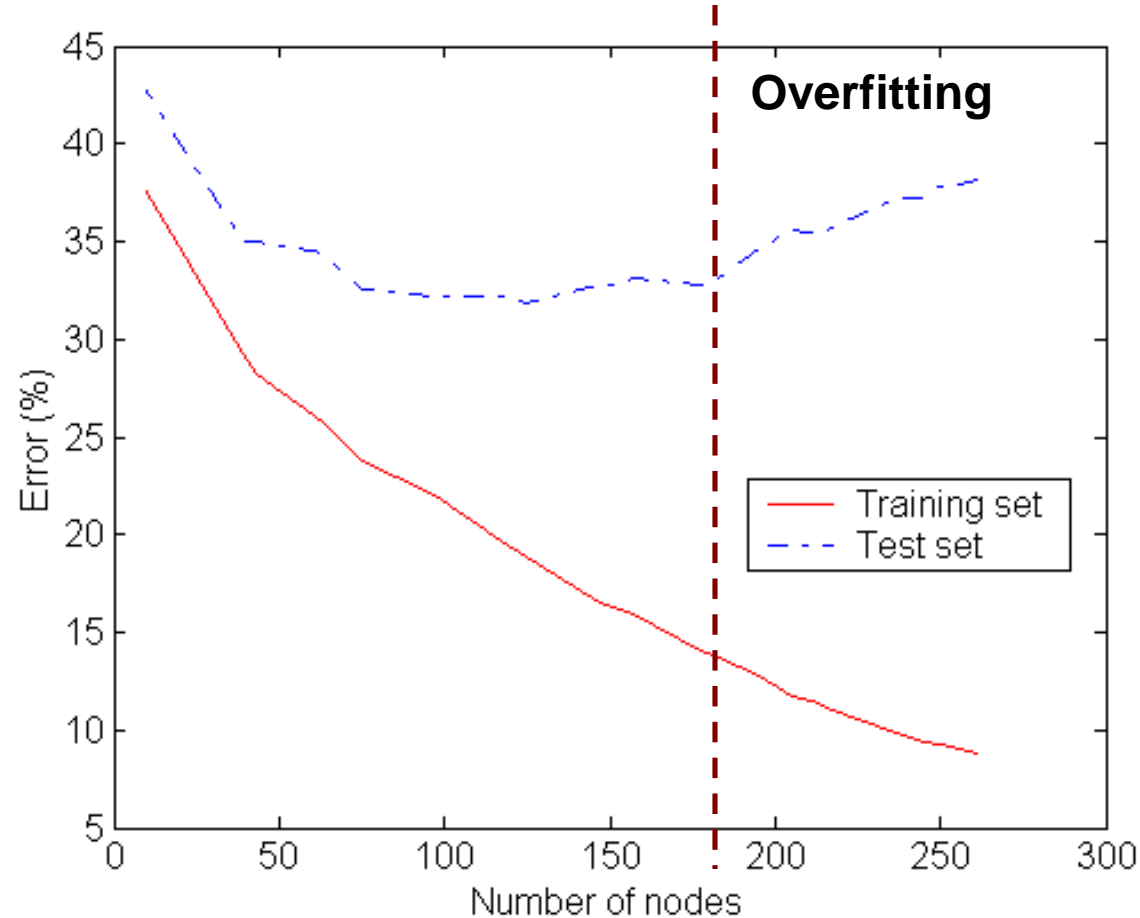
# Other Attribute Selection Measures

- CHAID: a popular decision tree algorithm, measure based on $\chi^2$ test for independence

- C-SEP: performs better than info. gain and gini index in certain cases

- G-statistics: has a close approximation to $\chi^2$ distribution

- MDL (Minimal Description Length) principle (i.e., the simplest solution is preferred):
    - The best tree as the one that requires the fewest # of bits to both (1) encode the tree, and (2) encode the exceptions to the tree

- Multivariate splits (partition based on multiple variable combinations)
    - CART: finds multivariate splits based on a linear comb. of attrs.

- Which attribute selection measure is the best?
    - Most give good results, none is significantly superior than others

Introduction _to_Classification

# Overfitting and Tree Pruning

- Overfitting:  An induced tree may overfit the training data

  - Too many branches, some may reflect anomalies due to noise or outliers

  - 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"

# Underfitting and Overfitting



**Underfitting**: when model is too simple, both training and test errors are large

# Overfitting in Classification

- Overfitting:  An induced tree may overfit the training data
    - Too many branches, some may reflect anomalies due to noise or outliers
    - Poor accuracy for unseen samples

Introduction _to_Classification

# Enhancements to Basic Decision Tree Induction

- Allow for continuous-valued attributes
  - Dynamically define new discrete-valued attributes that partition the continuous attribute value into a discrete set of intervals
- Handle missing attribute values
  - Assign the most common value of the attribute
  - Assign probability to each of the possible values
- Attribute construction
  - Create new attributes based on existing ones that are sparsely represented
  - This reduces fragmentation, repetition, and replication

Introduction _to_Classification

COSC 3337:DS 1

# 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

N.Rizk (University of Houston)

Introduction_to_Classification

COSC 3337:DS 1

# Scalable Decision Tree Induction Methods

- SLIQ (EDBT'96 — Mehta et al.)
  - Builds an index for each attribute and only class list and the current attribute list reside in memory

- SPRINT (VLDB'96 — J. Shafer et al.)
  - Constructs an attribute list data structure

- PUBLIC (VLDB'98 — Rastogi & Shim)
  - Integrates tree splitting and tree pruning: stop growing the tree earlier

- RainForest (VLDB'98 — Gehrke, Ramakrishnan & Ganti)
  - Builds an AVC-list (attribute, value, class label)

- BOAT (PODS'99 — Gehrke, Ganti, Ramakrishnan & Loh)
  - Uses bootstrapping to create several small samples

Introduction _to_Classification