# Decision Tree Classifier

# 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

Induction

Tree Induction algorithm

Learn Model

Model

Decision Tree

Apply Model

Deduction

| 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

# Decision Tree Example

**Learn decision rules from a dataset**: Do we want to play tennis?

| Day | Outlook | Temp. | Humidity | Wind | P.Tennis |
|---|---|---|---|---|---|
| $d_1$ | Sunny | Hot | High | Weak | No |
| $d_2$ | Sunny | Hot | High | Strong | No |
| $d_3$ | Overcast | Hot | High | Weak | Yes |
| $d_4$ | Rain | Mild | High | Weak | Yes |
| $d_5$ | Rain | Cool | Normal | Weak | Yes |
| $d_6$ | Rain | Cool | Normal | Strong | No |
| $d_7$ | Overcast | Cool | Normal | Strong | Yes |
| $d_8$ | Sunny | Mild | High | Weak | No |
| $d_9$ | Sunny | Cool | Normal | Weak | Yes |
| $d_{10}$ | Rain | Mild | Normal | Weak | Yes |
| $d_{11}$ | Sunny | Mild | Normal | Strong | Yes |
| $d_{12}$ | Overcast | Mild | High | Strong | Yes |
| $d_{13}$ | Overcast | Hot | Normal | Weak | Yes |
| $d_{14}$ | Rain | Mild | High | Strong | No |

❏ 4 discrete-valued attributes (Outlook, Temperature, Humidity, Wind)

❏ Play tennis?:"Yes/No" classification problem

# Decision Tree Example

| Day | Outlook | Temp. | Humidity | Wind | P.Tennis |
|-----|---------|-------|----------|------|----------|
| $d_1$ | Sunny | Hot | High | Weak | No |
| $d_2$ | Sunny | Hot | High | Strong | No |
| $d_3$ | Overcast | Hot | High | Weak | Yes |
| $d_4$ | Rain | Mild | High | Weak | Yes |
| $d_5$ | Rain | Cool | Normal | Weak | Yes |
| $d_6$ | Rain | Cool | Normal | Strong | No |
| $d_7$ | Overcast | Cool | Normal | Strong | Yes |
| $d_8$ | Sunny | Mild | High | Weak | No |
| $d_9$ | Sunny | Cool | Normal | Weak | Yes |
| $d_{10}$ | Rain | Mild | Normal | Weak | Yes |
| $d_{11}$ | Sunny | Mild | Normal | Strong | Yes |
| $d_{12}$ | Overcast | Mild | High | Strong | Yes |
| $d_{13}$ | Overcast | Hot | Normal | Weak | Yes |
| $d_{14}$ | Rain | Mild | High | Strong | No |

# Example of a Decision Tree

| 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**

# 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 → **MarStat**

Single, Divorced → **TaxInc**

Married → **NO**

< 80K → **NO**

> 80K → **YES**

**Model:  Decision Tree**

# 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 |

MarSt

Married → NO

Single, Divorced → Refund

Refund:
- Yes → NO
- No → TaxInc
  - < 80K → NO
  - > 80K → YES

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

# Apply Model to Test Data

Start from the root of tree.

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

# Apply Model to Test Data

**Test Data**

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

**Refund**

Yes

No

**NO**

**MarSt**

Single, Divorced

Married

**TaxInc**

**NO**

< 80K

> 80K

**NO**

**YES**

# Apply Model to Test Data

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

# Apply Model to Test Data

# Apply Model to Test Data

**Test Data**

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

**Refund**

Yes → **NO**

No → **MarSt**

Single, Divorced → **TaxInc**

Married → **NO**

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

# 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"

# Decision Tree Induction

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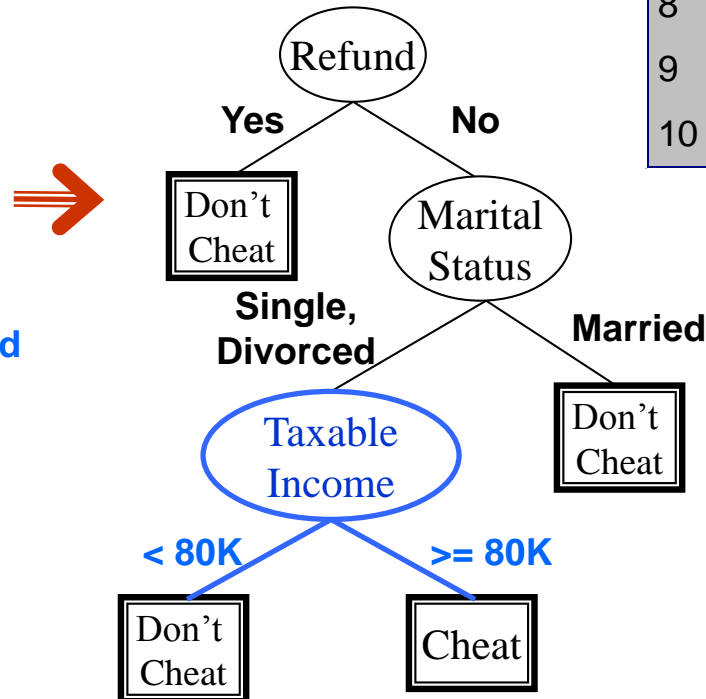  - CART
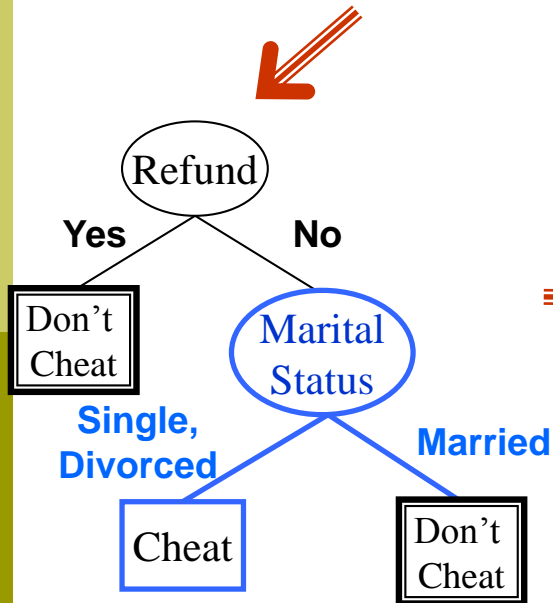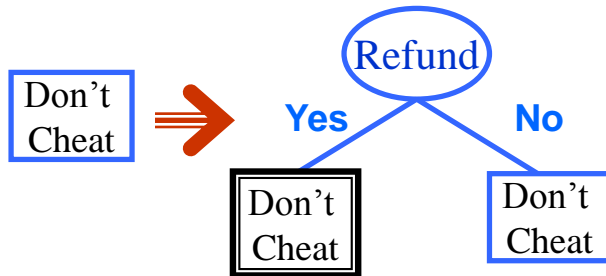  - ID3, C4.5
  - SLIQ,SPRINT

# 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 |

# 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?
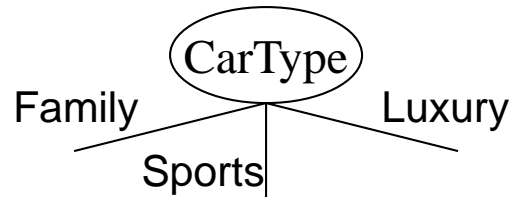  - Determine when to stop splitting

# How to Specify Test Condition?

- Depends on attribute types
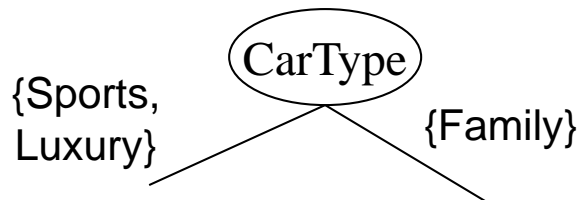  - Nominal
  - Ordinal
  - Continuous

- Depends on number of ways to split
  - 2-way split
  - Multi-way split
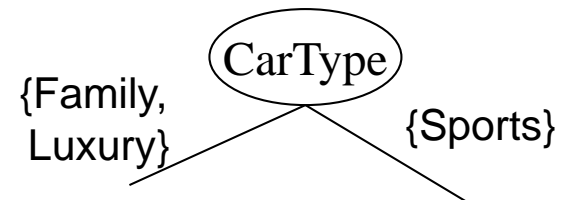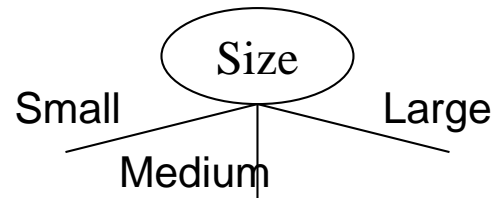
# Splitting Based on Nominal Attributes

- **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              OR              {Family,    CarType
Luxury}  /        \  {Family}                    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?
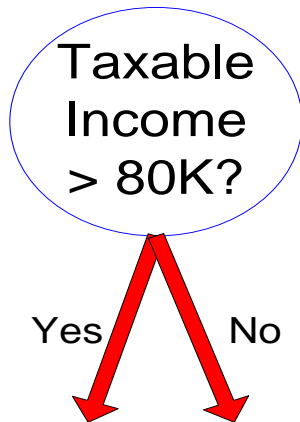
Size
{Small, Large}    {Medium}

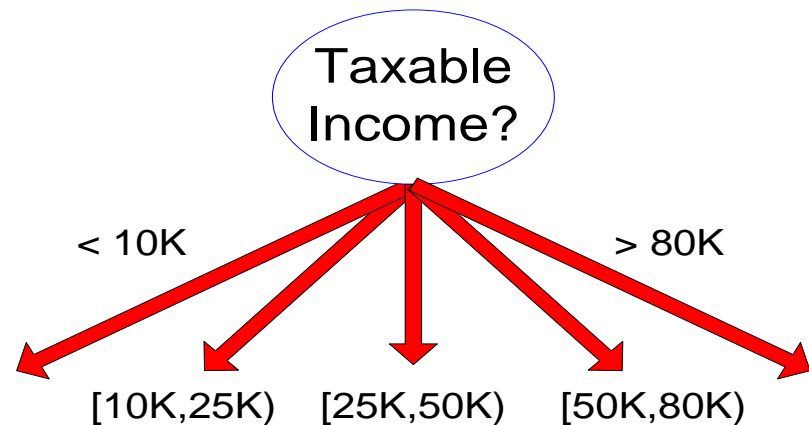# Splitting Based on Continuous Attributes

- **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

Taxable Income?

< 10K    > 80K

[10K,25K)    [25K,50K)    [50K,80K)

(i) Binary split

(ii) Multi-way split
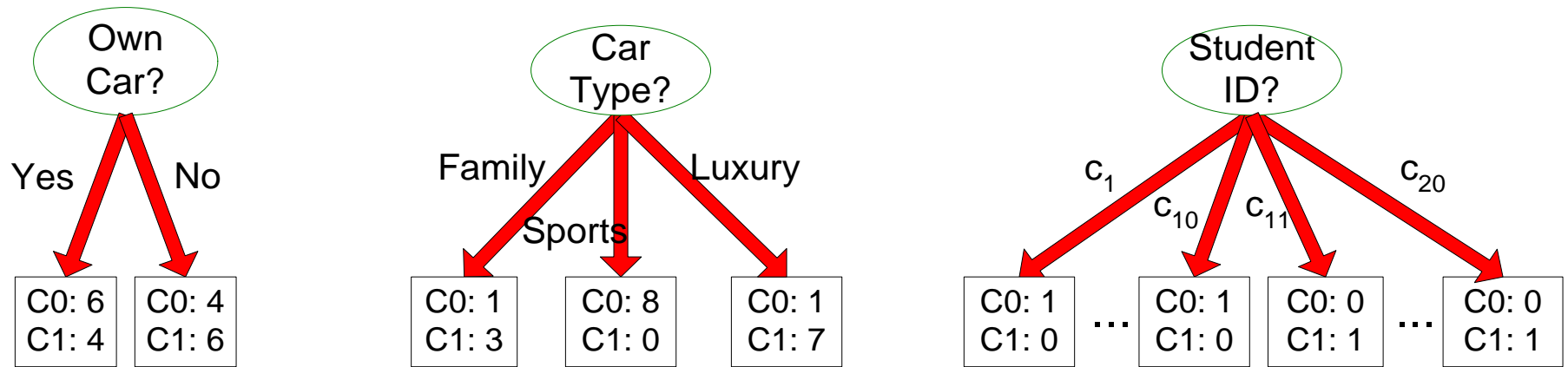
# Tree Induction

- **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

- Nodes with <span style="color:red">homogeneous</span> class distribution are preferred
- Need a measure of node impurity:

C0: 9
C1: 1

C0: 5
C1: 5

**Homogeneous,**

**Low degree of impurity**

**Non-homogeneous,**

**High degree of impurity**

# Measures of Node Impurity

- Gini Index

- Entropy

- Misclassification error

# Examples for computing GINI

$$GINI(t) = 1 - \sum_j [p(j \mid t)]^2$$

| C1 | 0 |
|----|---|
| C2 | 6 |

$P(C1) = 0/6 = 0 \qquad P(C2) = 6/6 = 1$

Gini $= 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$

| C1 | 1 |
|----|---|
| C2 | 5 |

$P(C1) = 1/6 \qquad P(C2) = 5/6$

Gini $= 1 - (1/6)^2 - (5/6)^2 = 0.278$

| C1 | 2 |
|----|---|
| C2 | 4 |

$P(C1) = 2/6 \qquad P(C2) = 4/6$

Gini $= 1 - (2/6)^2 - (4/6)^2 = 0.444$

# Decision Tree Based Classification

- Advantages:
    - Inexpensive to construct
    - Extremely fast at classifying unknown records
    - Easy to interpret for small-sized trees
    - Accuracy is comparable to other classification techniques for many simple data sets