# EE 4146 Data Engineering and Learning Systems

## Lecture 10: Decision Tree
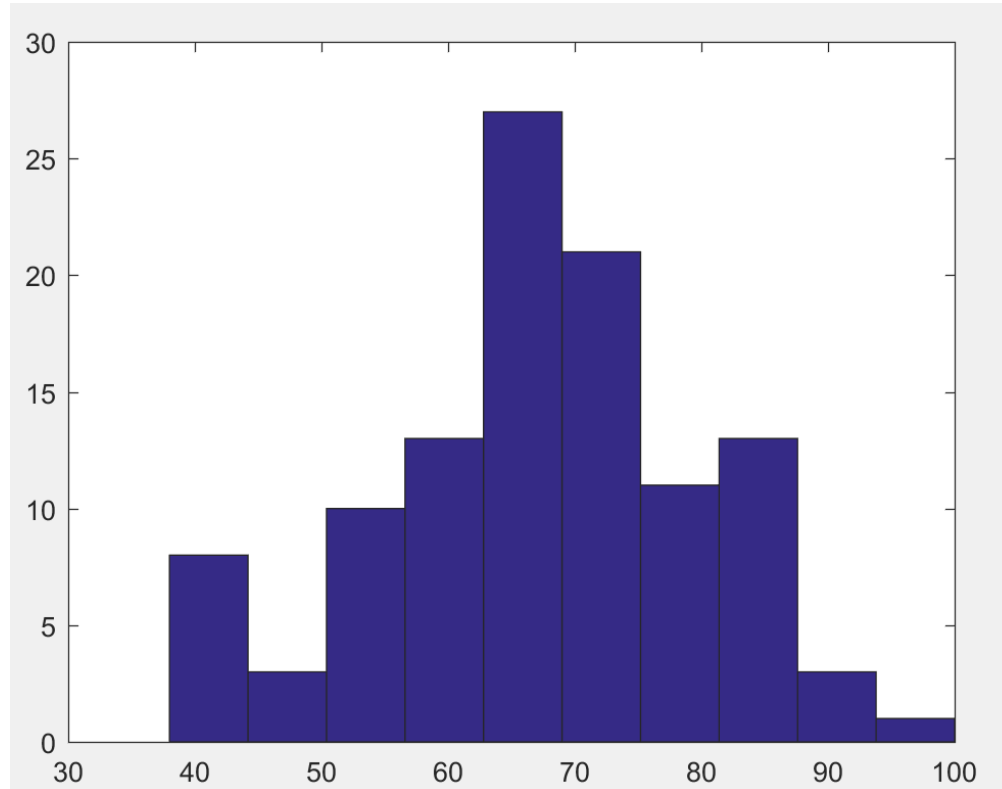
Semester A, 2021-2022

# Schedules

| Week | Date | Topics |
|------|------|--------|
| 1 | Sep. 1 | Introduction |
| 2 | Sep. 8 | Data exploration |
| 3 | Sep. 15 | Feature reduction and selection (HW1 out) |
| 4 | Sep. 22 | Mid-Autumn Festival |
| 5 | Sep. 29 | Clustering I: Kmeans based models (HW1 due in this weekend) |
| 6 | Oct. 6 | Clustering II: Hierarchical/density based/fuzzing clustering |
| 7 | Oct. 13 | Midterm (no tutorials this week) |
| 8 | Oct. 20 | Adverse Weather |
| 9 | Oct. 27 | Linear classifiers |
| 10 | Nov. 3 | Classification based on decision tree (Tutorial on project) (HW2 out in Friday) |
| 11 | Nov. 10 | Bayes based classifier (Tutorial on codes) (HW2 due in this weekend) |
| 12 | Nov. 17 | KNN and classifier ensemble |
| 13 | Nov. 24 | Deep learning based models (Quiz) |
| 14 | Make up video | Summary |

# Information about Midterm

- ## Mean=67.5273
  - You can find TA at G2325, AC1 in the open hour to get the examination paper Open hour: Monday 1:00-7:00 PM

# Project

- Group information has already been uploaded to 202109EE4146->Files->Project Related->EE4146 project group information

- Please contact me, if there are some mistakes in the grouping information

- TA will illustrate the project information, related to codes, topics, submission, etc. (6:00 PM-6:50 PM)

# Outline

- Basic introduction of decision tree
- Hunt's Algorithm
- Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
- Determine when to stop splitting

# Review: Classification definition

- Given a collection of records (*training set* )
  - Each record contains a set of attributes, one of the attributes is the class.
- Find a model for class attribute as a function of the values of other attributes.
- Goal: previously unseen records should be assigned a class as accurately as possible.
  - A test set is used to determine the accuracy of the  model. Usually, the given data set is divided into training  and test sets, with training set used to build the model  and test set used to validate it.
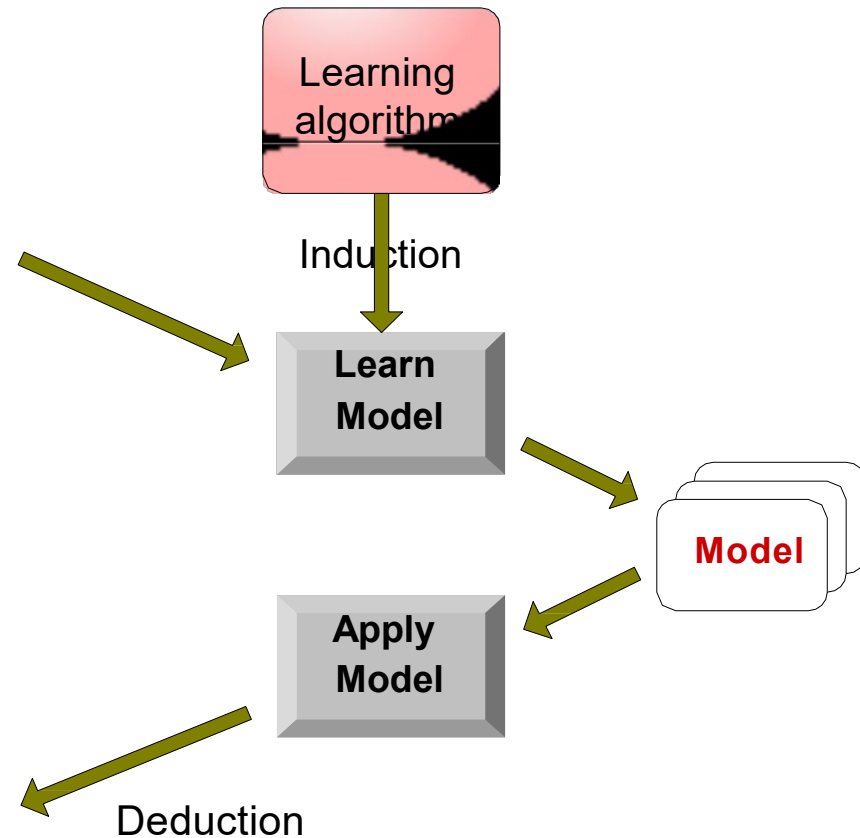
# Illustrating Classification Task

| 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

Learning algorithm

Induction

**Learn Model**

**Model**

**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

# Examples of Classification Task

- Predicting tumor cells as benign or malignant
- Classifying credit card transactions as legitimate or fraudulent
- Classifying emails as spams or normal emails
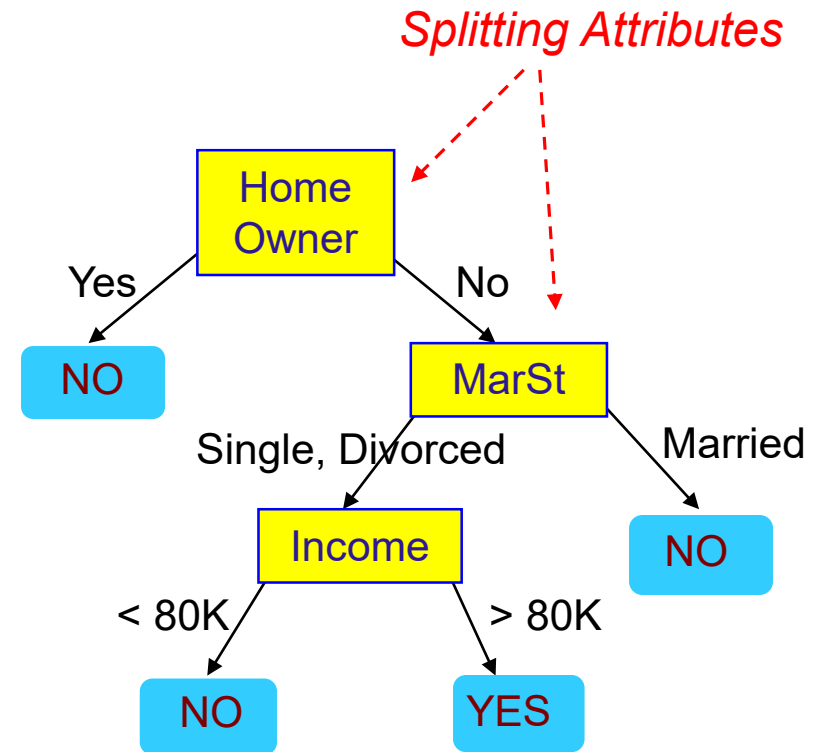- Categorizing news stories as finance, weather, entertainment, sports, etc

# Example of a Decision Tree

| ID | Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|----|-----------|---------------|---------------|---------------------|
| 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 |

categorical  categorical  continuous  class

Training Data

*Splitting Attributes*

Home Owner
Yes → NO
No → MarSt

MarSt
Single, Divorced → Income
Married → NO

Income
< 80K → NO
> 80K → YES

Model:  Decision Tree

# Another Example of Decision Tree

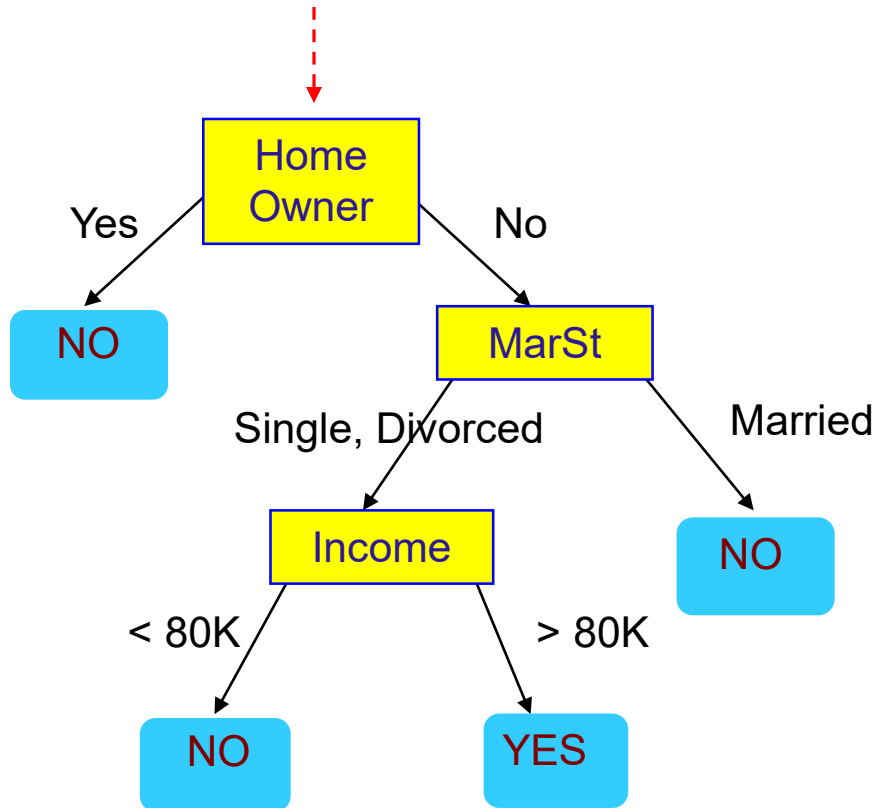|     | categorical | categorical | continuous | class |
|-----|-------------|-------------|------------|-------|
| **ID** | **Home Owner** | **Marital Status** | **Annual Income** | **Defaulted Borrower** |
| 1 | Yes | Single | 125K | **No** |
| 2 | No | Married | 100K | **No** |
| 3 | No | Single | 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!
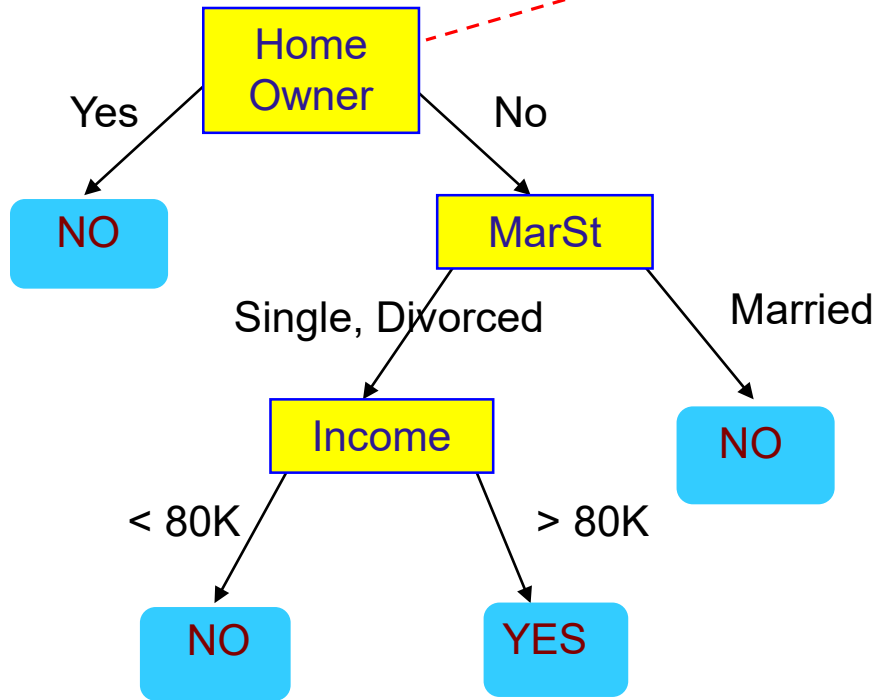
# Apply Model to Test Data

Test Data

| Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|------------|----------------|---------------|--------------------|
| No | Married | 80K | ? |

Start from the root of tree.

# Apply Model to Test Data

Test Data

| Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|------------|----------------|---------------|--------------------|
| No | Married | 80K | ? |

Home Owner

Yes

No

NO

MarSt

Single, Divorced

Married

Income

NO

< 80K

> 80K

NO

YES

# Apply Model to Test Data

Test Data

| Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|---|---|---|---|
| No | Married | 80K | ? |

# Apply Model to Test Data

Test Data

| Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|---|---|---|---|
| No | Married | 80K | ? |

```
            Home
            Owner
       Yes  /    \  No
           /      \
         NO      MarSt
    Single, Divorced /  \  Married
                    /    \
                Income    NO
         < 80K  /   \  > 80K
               /     \
             NO      YES
```

# Apply Model to Test Data

Test Data

| Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|------------|----------------|---------------|--------------------|
| No | Married | 80K | ? |

Home Owner

Yes

No

NO

MarSt

Single, Divorced

Married
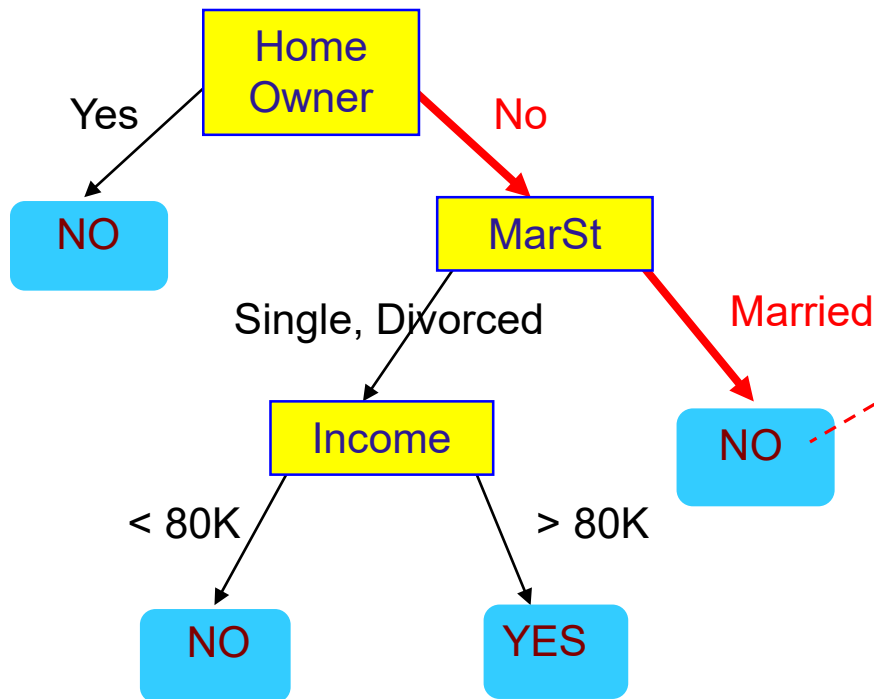
Income

NO

< 80K

> 80K

NO

YES

**How many tests do you need to reach a decision?**

# Apply Model to Test Data

Test Data

| Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|---|---|---|---|
| No | Married | 80K | ? |



Assign Defaulted to "No"

# 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

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 Induction

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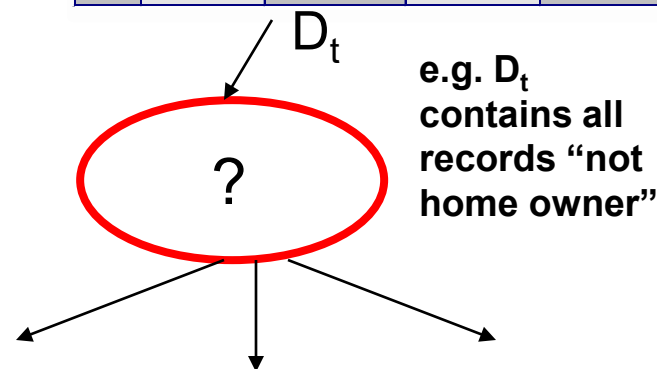  - SLIQ,SPRINT

# Outline

- Basic introduction of decision tree

- <span style="color:red">Hunt's Algorithm</span>

- Determine how to split the records

  - How to specify the attribute test condition?

  - How to determine the best split?

- Determine when to stop splitting

# 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 $D_t$ contains records that belong to more than one class, use an attribute to split the data into smaller subsets. Recursively apply the procedure to each subset.

| ID | Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|----|-----------|----------------|---------------|--------------------|
| 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 |

$D_t$

?

e.g. $D_t$ contains all records "not home owner"

# Hunt's Algorithm

| ID | Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|----|------------|----------------|---------------|--------------------|
| 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 |

Defaulted = No

(7,3)

(a)

# Hunt's Algorithm

Defaulted = No

(7,3)

(a)

Home Owner

Yes      No

Defaulted = No     Defaulted = No

(3,0)        (4,3)

(b)

| ID | Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|----|-----------|---------------|--------------|-------------------|
| 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 |

# Hunt's Algorithm



(a)



(b)



(c)

| ID | Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|---|---|---|---|---|
| 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 |

# Hunt's Algorithm

| ID | Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|----|-----------|----------------|---------------|--------------------|
| 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 |

(a) **Defaulted = No** (7,3)

(b) **Home Owner** — Yes: **Defaulted = No** (3,0); No: **Defaulted = No** (4,3)

(c) **Home Owner** — Yes: **Defaulted = No** (3,0); No: **Marital Status** — Single, Divorced: **Defaulted = Yes** (1,3); Married: **Defaulted = No** (3,0)

(d) **Home Owner** — Yes: **Defaulted = No** (3,0); No: **Marital Status** — Single, Divorced: **Annual Income** — < 80K: **Defaulted = No** (1,0); >= 80K: **Defaulted = Yes** (0,3); Married: **Defaulted = No** (3,0)
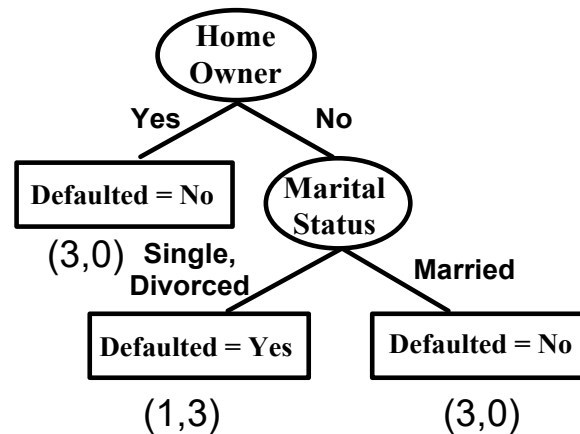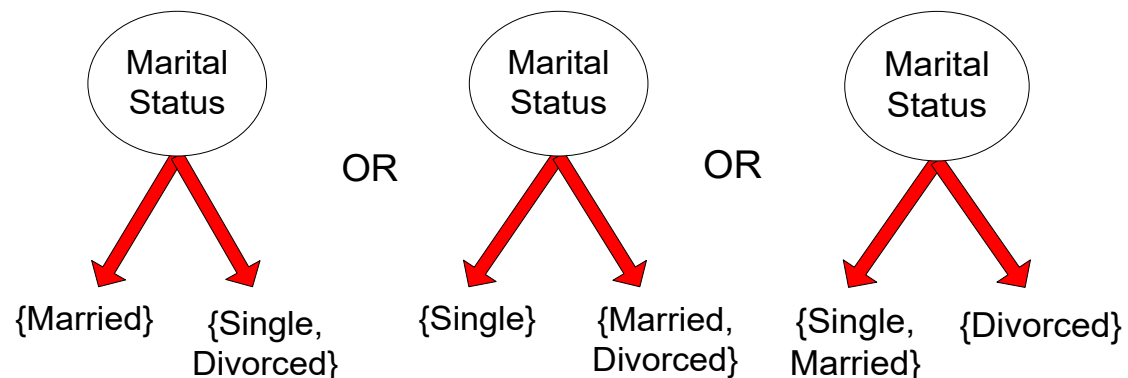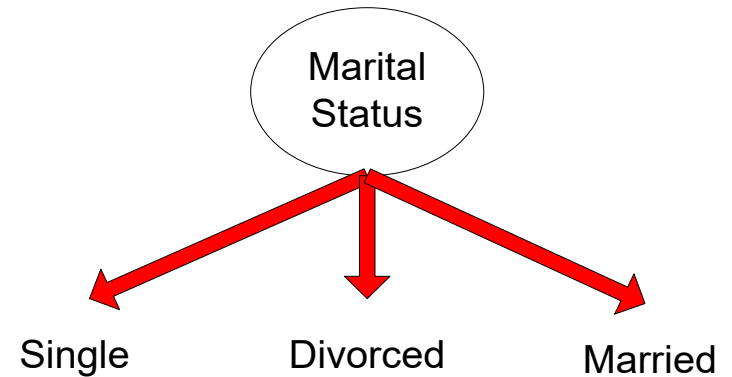
24

# 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

# Methods for Expressing Test Conditions

- Depends on attribute types
  - Nominal
  - Ordinal
  - Continuous

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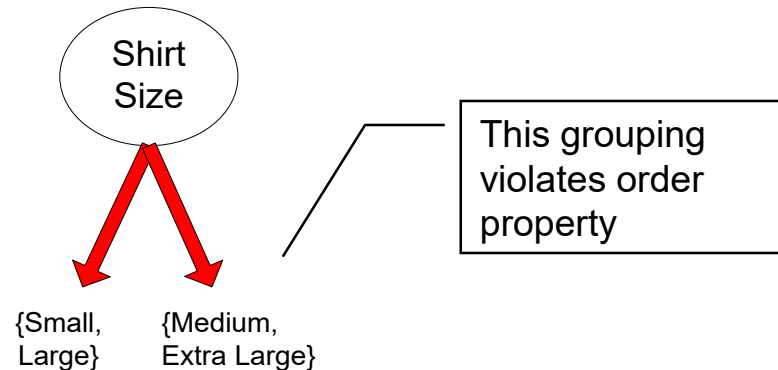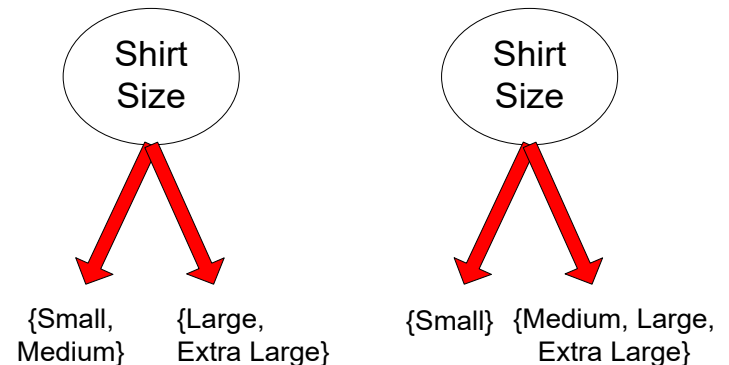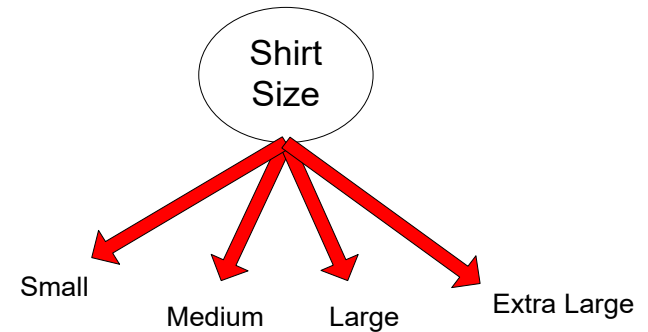  - 2-way split
  - Multi-way split

# Test Condition for Nominal Attributes

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

```
      ( Marital
        Status )
     /     |     \
  Single  Divorced  Married
```

- **Binary split:**
  - Divides values into two subsets

```
   ( Marital              ( Marital              ( Marital
     Status )      OR       Status )      OR       Status )
    /      \                /      \                /      \
{Married} {Single,      {Single} {Married,    {Single,  {Divorced}
          Divorced}              Divorced}     Married}
```

# Test Condition for Ordinal Attributes
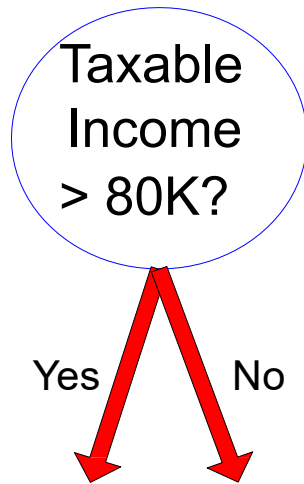
- **Multi-way split:**
  - Use as many partitions as distinct values

- **Binary split:**
  - Divides values into two subsets
  - Preserve order property among attribute values

Shirt Size

Small    Medium    Large    Extra Large

Shirt Size

{Small, Medium}    {Large, Extra Large}

Shirt Size

{Small}    {Medium, Large, Extra Large}

Shirt Size

{Small, Large}    {Medium, Extra Large}
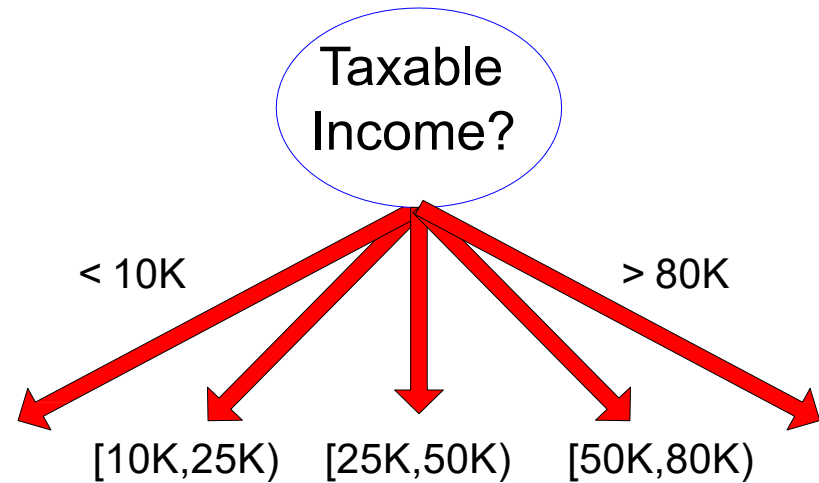
This grouping violates order property

# Splitting Based on Continuous Attributes

- Different ways of handling
    - Discretization to form an ordinal categorical attribute
    - Static – discretize once at the beginning
    - Dynamic – repeat at each node for discretization

    - 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
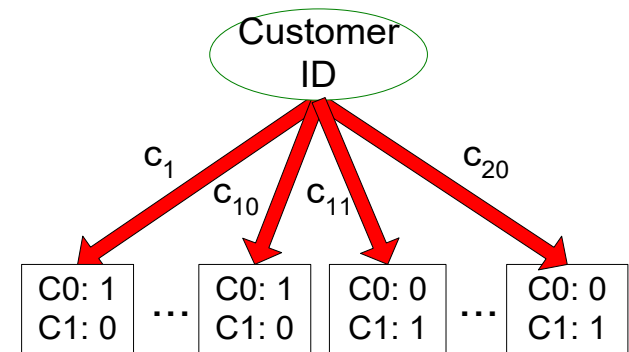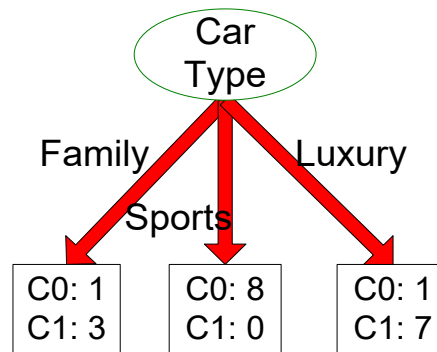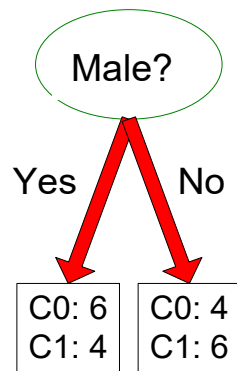


(i) Binary split

(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

| Customer Id | Gender | Car Type | Shirt Size | Class |
|---|---|---|---|---|
| 1 | M | Family | Small | C0 |
| 2 | M | Sports | Medium | C0 |
| 3 | M | Sports | Medium | C0 |
| 4 | M | Sports | Large | C0 |
| 5 | M | Sports | Extra Large | C0 |
| 6 | M | Sports | Extra Large | C0 |
| 7 | F | Sports | Small | C0 |
| 8 | F | Sports | Small | C0 |
| 9 | F | Sports | Medium | C0 |
| 10 | F | Luxury | Large | C0 |
| 11 | M | Family | Large | C1 |
| 12 | M | Family | Extra Large | C1 |
| 13 | M | Family | Medium | C1 |
| 14 | M | Luxury | Extra Large | C1 |
| 15 | F | Luxury | Small | C1 |
| 16 | F | Luxury | Small | C1 |
| 17 | F | Luxury | Medium | C1 |
| 18 | F | Luxury | Medium | C1 |
| 19 | F | Luxury | Medium | C1 |
| 20 | F | Luxury | Large | C1 |

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



Male?

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

Customer 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 purer class distribution are preferred

- Need a measure of node impurity:

| C0: 5 | | C0: 9 |
|-------|---|-------|
| C1: 5 | | C1: 1 |

High uncertainty                    Low uncertainty

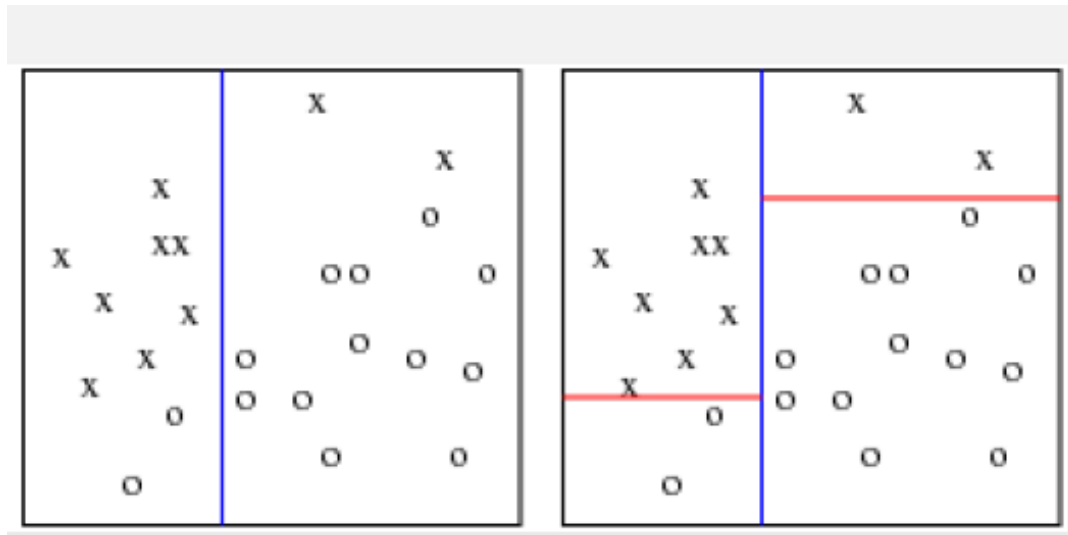High degree of impurity             Low degree of impurity

# Goodness of Split

- The goodness of split is measured by an impurity function defined for each node.

- Intuitively, we want each leaf node to be "pure", that is, one class dominates.

# Finding the Best Split

- Compute impurity measure (P) before splitting

- Compute impurity measure (M) after splitting

  - Compute impurity measure of each child node

  - M is the weighted impurity of children

- Choose the attribute test condition that produces the highest gain

$$\text{Gain} = P - M$$

or equivalently, lowest impurity measure after splitting (M)

# Finding the Best Split

Before Splitting:

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

$\longrightarrow$ P

N00: number of records with label 0 at node N0

A?

Yes — No

Node N1     Node N2

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

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

M11                M12

M1

B?

Yes — No

Node N3     Node N4

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

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

M21                M22

M2

Gain = P – M1    vs    P – M2

# Measures of Node Impurity

- Gini Index

$$GINI(t) = 1 - \sum_{j} [p(j \mid t)]^2$$   t is a node

- Entropy

$$Entropy(t) = -\sum_{j} p(j \mid t) \log p(j \mid t)$$

- Misclassification error

Entropy quantifies uncertainty

$$Error(t) = 1 - \max_{i} P(i \mid t)$$

# Measure of Impurity: GINI

- Gini Index for a given node t :

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

(NOTE: *p( j | t)* is the relative frequency of class j at node t).

- Maximum ($1 - 1/n_c$) when records are equally distributed among all classes, implying least information
- Minimum (0.0) when all records belong to one class, implying most information

# Measure of Impurity: GINI

- Gini Index for a given node t :

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

(NOTE: *p( j | t)* is the relative frequency of class j at node t).

  - For 2-class problem (p, 1 − p):
  - GINI = 1 − p² − (1 − p)² = 2p (1-p)

| C1 | 0 |
|---|---|
| C2 | 6 |
| **Gini=0.000** ||

| C1 | 1 |
|---|---|
| C2 | 5 |
| **Gini=0.278** ||

| C1 | 2 |
|---|---|
| C2 | 4 |
| **Gini=0.444** ||

| C1 | 3 |
|---|---|
| C2 | 3 |
| **Gini=0.500** ||

# Computing Gini Index of a Single Node

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

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

P(C1) = 0/6 = 0     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        P(C2) = 5/6

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

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

P(C1) = 2/6        P(C2) = 4/6

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

# Computing Gini Index for a Collection of Nodes

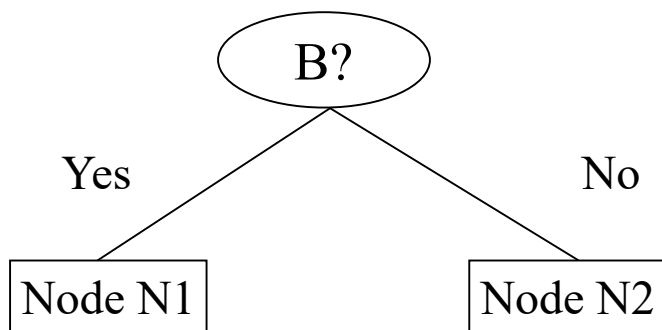- When a node p is split into k partitions (children)

$$GINI_{split} = \sum_{i=1}^{k} \frac{n_i}{n} GINI(i)$$

where,    $n_i$ = number of records at child i,

          n = number of records at parent node p.

- Choose the attribute that minimizes weighted average Gini index of the children

- Gini index is used in decision tree algorithms such as CART, SLIQ, SPRINT

# Binary Attributes: Computing GINI Index

- Splits into two partitions
- Effect of Weighing partitions:
  - Larger and Purer Partitions are sought for.



| | Parent |
|---|---|
| C1 | 7 |
| C2 | 5 |
| **Gini = 0.486** | |

Gini(P)
$= 1 - (5/12)^2 - (7/12)^2$
$= 0.486$

| | N1 | N2 |
|---|---|---|
| C1 | 5 | 2 |
| C2 | 1 | 4 |
| **Gini=0.361** | | |

Gini(N1)
$= 1 - (5/6)^2 - (1/6)^2$
$= 0.278$

Gini(N2)
$= 1 - (2/6)^2 - (4/6)^2$
$= 0.444$

Weighted Gini of N1 N2
$= 6/12 * 0.278 +$
  $6/12 * 0.444$
$= 0.361$

Gain = 0.486 – 0.361 = 0.125

# Binary Attributes: Computing GINI Index

- Splits into two partitions
- Effect of Weighing partitions:
  - Larger and Purer Partitions are sought for



|      | Parent |
|------|--------|
| C1   | 6      |
| C2   | 6      |
| **Gini = 0.500** | |

Gini(P)
$= 1 - (6/12)^2 - (6/12)^2$
$= 0.5$

|      | N1 | N2 |
|------|----|----|
| C1   | 5  | 1  |
| C2   | 2  | 4  |

Gini(N1)
$= 1 - (5/7)^2 - (2/7)^2$
$= 0.408$

Gini(N2)
$= 1 - (1/5)^2 - (4/5)^2$
$= 0.32$

Gini(Children)
$= 7/12 * 0.408 +$
$5/12 * 0.32$
$= 0.371$

Gain = 0.5 – 0.371 = 0.129

43

# Categorical Attributes: Computing Gini Index

- For each distinct value, gather counts for each class in the dataset
- Use the count matrix to make decisions

Multi-way split

| CarType | | |
|---|---|---|
| **Family** | **Sports** | **Luxury** |
| 1 | 8 | 1 |
| 3 | 0 | 7 |
| **Gini** | **0.1625** | |

(C1, C2 row labels)

| | CarType | | |
|---|---|---|---|
| | **Family** | **Sports** | **Luxury** |
| **C1** | 1 | 8 | 1 |
| **C2** | 3 | 0 | 7 |
| **Gini** | **0.1625** | | |

Two-way split
(find best partition of values)

Case 1

| | CarType | |
|---|---|---|
| | **{Sports, Luxury}** | **{Family}** |
| **C1** | 9 | 1 |
| **C2** | 7 | 3 |
| **Gini** | **0.468** | |

Case 2

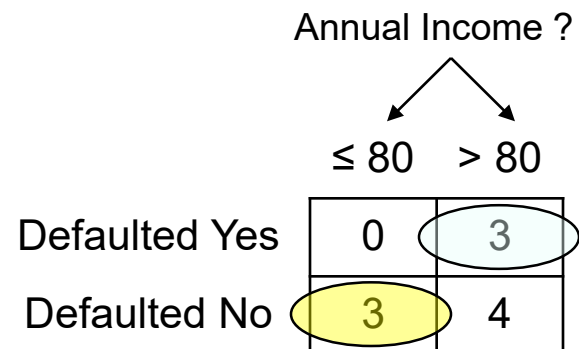| | CarType | |
|---|---|---|
| | **{Sports}** | **{Family, Luxury}** |
| **C1** | 8 | 2 |
| **C2** | 0 | 10 |
| **Gini** | **0.167** | |

16/20*(1-(7/16)^2-(9/16)^2)+4/20*(1-(1/4)^2-(3/4)^2)

8/20*(1-(0/8)^2-(8/8)^2)+12/20*(1-(2/12)^2-(10/12)^2)

Which of these is the best?

44

# Continuous Attributes: Computing Gini Index

- **Use Binary Decisions based on one value**

- **Each splitting value has a count matrix associated with it**
  - Class counts in each of the partitions, A < v and A $\geq$ v

- **Simple method to choose best v**
  - For each v, scan the database to gather count matrix and compute its Gini index
  - Computationally Inefficient! Repetition of work.

| ID | Home Owner | Marital Status | Annual Income | Defaulted |
|----|------------|----------------|---------------|-----------|
| 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 |

Annual Income ?

≤ 80      > 80

|  | ≤ 80 | > 80 |
|--------------|------|------|
| Defaulted Yes | 0 | 3 |
| Defaulted No | 3 | 4 |

# Continuous Attributes: Computing Gini Index...

- For efficient computation: for each attribute,
    - Sort the attribute on values
    - Linearly scan these values, each time updating the count matrix and computing gini index
    - Choose the split position that has the least gini index

| Cheat | No | No | No | Yes | Yes | Yes | No | No | No | No |
|-------|----|----|----|-----|-----|-----|----|----|----|----|
| **Annual Income** | | | | | | | | | | |
| Sorted Values → | 60 | 70 | 75 | 85 | 90 | 95 | 100 | 120 | 125 | 220 |

# Continuous Attributes: Computing Gini Index...

- For efficient computation: for each attribute,
    - Sort the attribute on values
    - Linearly scan these values, each time updating the count matrix and computing gini index
    - Choose the split position that has the least gini index

Sorted Values →

Split Positions →

| Cheat | No | No | No | Yes | Yes | Yes | No | No | No | No |
|---|---|---|---|---|---|---|---|---|---|---|
| **Annual Income** | | | | | | | | | | |
| | 60 | 70 | 75 | 85 | 90 | 95 | 100 | 120 | 125 | 220 |
| | 55 | 65 | 72 | 80 | 87 | 92 | 97 | 110 | 122 | 172 | 230 |
| | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > |

# Continuous Attributes: Computing Gini Index...

- For efficient computation: for each attribute,
  - Sort the attribute on values
  - Linearly scan these values, each time updating the count matrix and computing gini index
  - Choose the split position that has the least gini index

| Cheat | No | No | No | Yes | Yes | Yes | No | No | No | No |
|---|---|---|---|---|---|---|---|---|---|---|

Annual Income

Sorted Values →

| 60 | 70 | 75 | 85 | 90 | 95 | 100 | 120 | 125 | 220 |
|---|---|---|---|---|---|---|---|---|---|

Split Positions →

| 55 | | 65 | | 72 | | 80 | | 87 | | 92 | | 97 | | 110 | | 122 | | 172 | | 230 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > |
| **Yes** | | | | | | 0 | 3 | | | | | | | | | | | | | | |
| **No** | | | | | | 3 | 4 | | | | | | | | | | | | | | |
| **Gini** | | | | | | 0.343 | | | | | | | | | | | | | | | |

Gini(N1) = 1 – (0/3)$^2$ – (3/3)$^2$ = 0
Gini(N2) = 1 – (3/7)$^2$ – (4/7)$^2$ = 0.4898

Gini(Children) = 3/10*0 + 7/10 * 0.4898
= 0.343

48

# Continuous Attributes: Computing Gini Index…

- For efficient computation: for each attribute,
  - Sort the attribute on values
  - Linearly scan these values, each time updating the count matrix and computing gini index
  - Choose the split position that has the least gini index

| Cheat | No | No | No | Yes | Yes | Yes | No | No | No | No |
|---|---|---|---|---|---|---|---|---|---|---|

**Annual Income**

Sorted Values →

| 60 | 70 | 75 | 85 | 90 | 95 | 100 | 120 | 125 | 220 |
|---|---|---|---|---|---|---|---|---|---|

Split Positions →

| | 55 | | 65 | | 72 | | 80 | | 87 | | 92 | | 97 | | 110 | | 122 | | 172 | | 230 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > |
| **Yes** | | | | | | | | 0 | 3 | 1 | 2 | | | | | | | | | | | |
| **No** | | | | | | | | 3 | 4 | 3 | 4 | | | | | | | | | | | |
| **Gini** | | | | | | | | 0.343 | | 0.417 | | | | | | | | | | | | |

# Continuous Attributes: Computing Gini Index...

- For efficient computation: for each attribute,
  - Sort the attribute on values
  - Linearly scan these values, each time updating the count matrix and computing gini index
  - Choose the split position that has the least gini index

| Cheat | | No | | No | | No | | Yes | | Yes | | Yes | | No | | No | | No | | No | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Annual Income** | | | | | | | | | | | | | | | | | | | | |
| Sorted Values → | | 60 | | 70 | | 75 | | 85 | | 90 | | 95 | | 100 | | 120 | | 125 | | 220 | |
| Split Positions → | | 55 | | 65 | | 72 | | 80 | | 87 | | 92 | | 97 | | 110 | | 122 | | 172 | | 230 |
| | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > | <= | > |
| Yes | 0 | 3 | 0 | 3 | 0 | 3 | 0 | 3 | 1 | 2 | 2 | 1 | 3 | 0 | 3 | 0 | 3 | 0 | 3 | 0 | 3 | 0 |
| No | 0 | 7 | 1 | 6 | 2 | 5 | 3 | 4 | 3 | 4 | 3 | 4 | 3 | 4 | 4 | 3 | 5 | 2 | 6 | 1 | 7 | 0 |
| Gini | 0.420 | | 0.400 | | 0.375 | | 0.343 | | 0.417 | | 0.400 | | *0.300* | | 0.343 | | 0.375 | | 0.400 | | 0.420 | |

# Measure of Impurity: Entropy

- Entropy at a given node t:

$$Entropy(t) = -\sum_j p(j \mid t) \log p(j \mid t)$$

- (NOTE: p(j|t) is the relative frequency of class j at node t).

- Measures purity of a node
  - Maximum ($\log n_c$) when records are equally distributed among all classes implying least information
  - Minimum (0.0) when all records belong to one class, implying most information
  - Entropy based computations are quite similar to the GINI index computations

# Computing Entropy of a Single Node

$$Entropy(t) = -\sum_j p(j \mid t) \log_2 p(j \mid t)$$

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

P(C1) = 0/6 = 0    P(C2) = 6/6 = 1

Entropy = – 0 log 0 – 1 log 1 = – 0 – 0 = 0

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

P(C1) = 1/6        P(C2) = 5/6

Entropy = – (1/6) $\log_2$ (1/6) – (5/6) $\log_2$ (1/6) = 0.65

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

P(C1) = 2/6        P(C2) = 4/6

Entropy = – (2/6) $\log_2$ (2/6) – (4/6) $\log_2$ (4/6) = 0.92

# Computing Information Gain After Splitting

- Information Gain:

$$GAIN_{split} = Entropy(p) - \left( \sum_{i=1}^{k} \frac{n_i}{n} Entropy(i) \right)$$

Parent Node, p is split into k partitions; $n_i$ is number of records in partition i

- Measures reduction in entropy achieved because of the split. Choose the split that achieves most reduction (maximizes GAIN)

- Used in the ID3 and C4.5 decision tree algorithms

# Problem with large number of partitions

■ Node impurity measures tend to prefer splits that result in large number of partitions, each being small but pure

**Gender**

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 |

**Customer ID**

$c_1$    $c_{10}$   $c_{11}$     $c_{20}$

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

■ Customer ID has highest information gain because entropy for all the children is zero

# Gain Ratio

- Gain Ratio:

$$GainRATIO_{split} = \frac{GAIN_{Split}}{SplitINFO}$$

$$SplitINFO = -\sum_{i=1}^{k} \frac{n_i}{n} \log \frac{n_i}{n}$$

Parent Node, p is split into k partitions, $n_i$ is the number of records in partition i

- Adjusts Information Gain by the entropy of the partitioning (SplitINFO).
  - Higher entropy partitioning (large number of small partitions) is penalized!
- Used in C4.5 algorithm
- Designed to overcome the disadvantage of Information Gain

# Gain Ratio

■ Gain Ratio

$$GainRATIO_{split} = \frac{GAIN_{Split}}{SplitINFO}$$

$$SplitINFO = -\sum_{i=1}^{k} \frac{n_i}{n} \log \frac{n_i}{n}$$

Parent Node, p is split into k partitions $n_i$ is the number of records in partition i

| Case 1 | CarType | | |
|---|---|---|---|
| | Family | Sports | Luxury |
| C1 | 1 | 8 | 1 |
| C2 | 3 | 0 | 7 |
| Gini | 0.1625 | | |

SplitINFO = 1.52

-[(4/20)*log2(4/20)+(8/20)*log2(8/20)+(8/20)*log2(8/20)]

| Case 2 | CarType | |
|---|---|---|
| | {Sports, Luxury} | {Family} |
| C1 | 9 | 1 |
| C2 | 7 | 3 |
| Gini | 0.4688 | |

SplitINFO = 0.72

-[(16/20)*log2(16/20)+(4/20)*log2(4/20)]

| Case 3 | CarType | |
|---|---|---|
| | {Sports} | {Family, Luxury} |
| C1 | 8 | 2 |
| C2 | 0 | 10 |
| Gini | 0.1667 | |

SplitINFO = 0.97

-[(8/20)*log2(8/20)+(12/20)*log2(12/20)]

# Measure of Impurity: Classification Error

- Classification error at a node t :

$$Error(t) = 1 - \max_i P(i \mid t)$$

- Measures misclassification error made by a node.
  - Maximum ($1 - 1/n_c$) when records are equally distributed among all classes, implying least information
  - Minimum (0) when all records belong to one class, implying most information

# Computing Error of a Single Node

$$Error(t) = 1 - \max_i P(i \mid t)$$

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

P(C1) = 0/6 = 0    P(C2) = 6/6 = 1

Error = 1 – max (0, 1) = 1 – 1 = 0

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

P(C1) = 1/6      P(C2) = 5/6

Error = 1 – max (1/6, 5/6) = 1 – 5/6 = 1/6

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

P(C1) = 2/6      P(C2) = 4/6

Error = 1 – max (2/6, 4/6) = 1 – 4/6 = 1/3

# Comparison among Impurity Measures

For a 2-class problem:

# Misclassification Error vs Gini Index

A?

Yes — Node N1

No — Node N2

|  | Parent |
|---|---|
| C1 | 7 |
| C2 | 3 |
| Gini = 0.42 | |

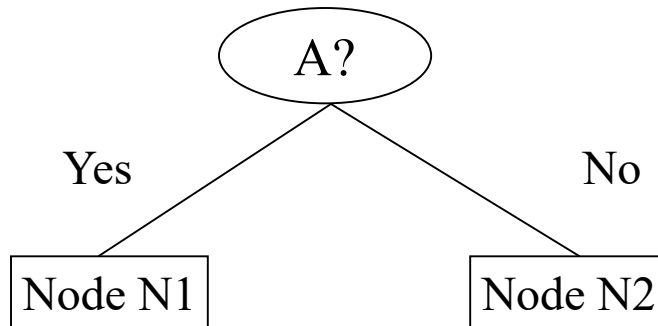|  | N1 | N2 |
|---|---|---|
| C1 | 3 | 4 |
| C2 | 0 | 3 |
| Gini=0.342 | | |

Error for children
N1: 1-max (3/3, 0/3) = 0
N2:  1-max(4/7, 3/7) = 3/7
Weighted sum:
3/10 * 0 + 7/10 * 3/7 = 0.3

Error for parent nodes
1-max(7/10,3/10)=0.3

# Misclassification Error vs Gini Index



A?

Yes — Node N1

No — Node N2

| | Parent |
|---|---|
| C1 | 7 |
| C2 | 3 |
| **Gini = 0.42** | |

| | N1 | N2 |
|---|---|---|
| C1 | 3 | 4 |
| C2 | 0 | 3 |
| **Gini=0.342** | | |

Gini(N1)
$= 1 - (3/3)^2 - (0/3)^2$
$= 0$

Gini(N2)
$= 1 - (4/7)^2 - (3/7)^2$
$= 0.489$

Gini(Children)
$= 3/10 * 0$
$+ 7/10 * 0.489$
$= 0.342$

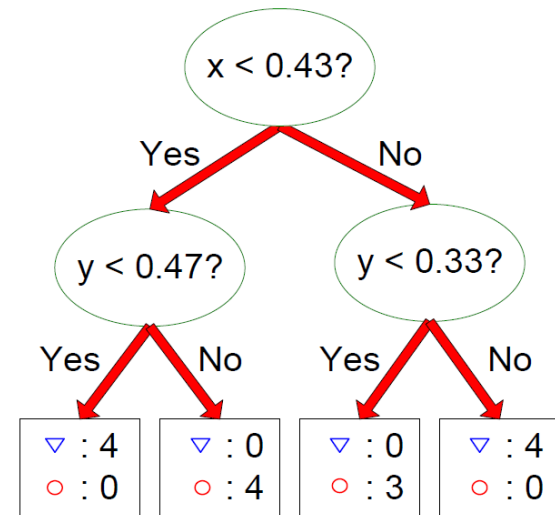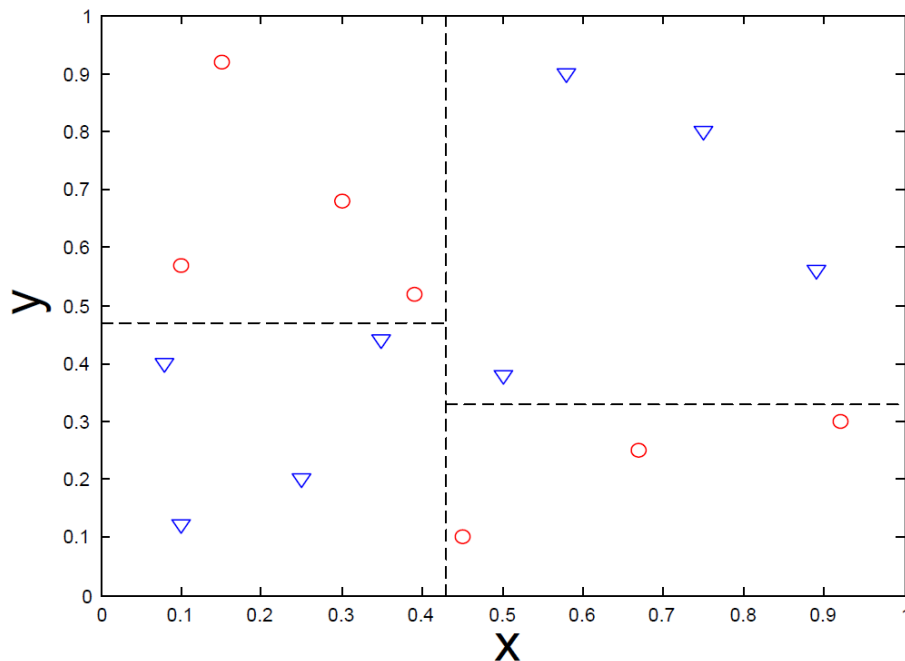Gini improves but error remains the same!!

# Outline

- Basic introduction of decision tree

- Hunt's Algorithm

- Determine how to split the records

    - How to specify the attribute test condition?

    - How to determine the best split?

- Determine when to stop splitting

# Determine when to stop splitting

- Stop expanding a node when all the records belong to the same class

- Stop expanding a node when all the records have similar attribute values

# Decision Boundary

- Border line between two neighboring regions of different classes is known as decision boundary

- Decision boundary is parallel to axes because test condition involves a single attribute at-a-time

# Decision Tree Based Classification

- Advantages:
  - Inexpensive to construct
  - Extremely fast at classifying unknown records
  - Easy to interpret for small-sized trees
  - Robust to noise
  - Can easily handle redundant or irrelevant attributes (unless the attributes are interacting)

- Disadvantages:
  - Space of possible decision trees is exponentially large. Greedy approaches are often unable to find the best tree.
  - Does not take into account interactions between attributes
  - Each decision boundary involves only a single attribute