# UNIT III - CLASSIFICATION

# 18CS54 – DATA MINING

# Outline

- Recap
- Classification
  - *Definition*
  - *Illustrating Classification Task*
  - *Classification Techniques*
- Decision Tree
  - *Decision Tree Induction*
  - *Hunt's Algorithm*
    - *Measure of Node Impurity*
      - GINI
      - Entropy
      - Misclassification Error
    - *CART , SLIQ, SPRINT*
    - *C4.5*
- Rule based Classifiers
- Nearest Neighbor classifiers

# LECTURE 14

Dr.Vani V

# Recap & Moving forward

- Unit –I Data Mining: Introduction, KDD Process, Challenges, Data Mining Tasks, Data Mining Trends and Applications.

- Unit –II Data, Types of Data, Data Pre-processing, Measures of Similarity And Dissimilarity

- **Unit –III Classification: Basics, General Approach to Solve Classification Problem, Decision Tree Induction, Rule Based Classifiers, Nearest Neighbor Classifiers.**

# 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: <u>previously unseen</u> 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.*
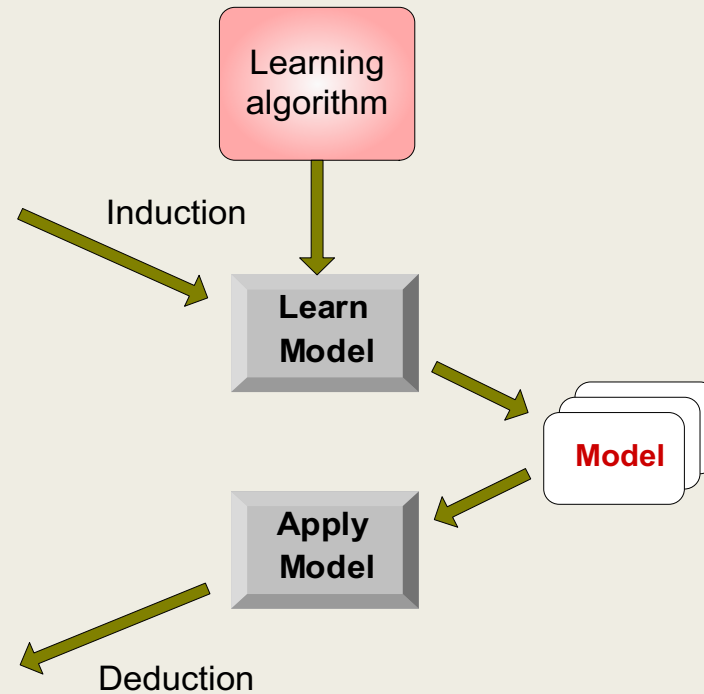
# Illustrating Classification Task

| 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**

| 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

**Apply Model**

Deduction

# Examples of Classification Task



■ Predicting tumor cells as benign or malignant

■ Classifying credit card transactions
as legitimate or fraudulent

■ Classifying secondary structures of protein
as alpha-helix, beta-sheet, or random
coil



■ Categorizing news stories as finance,
weather, entertainment, sports, etc

# Classification Techniques

- Decision Tree based Methods

- Rule-based Methods

- Memory based reasoning

- Neural Networks

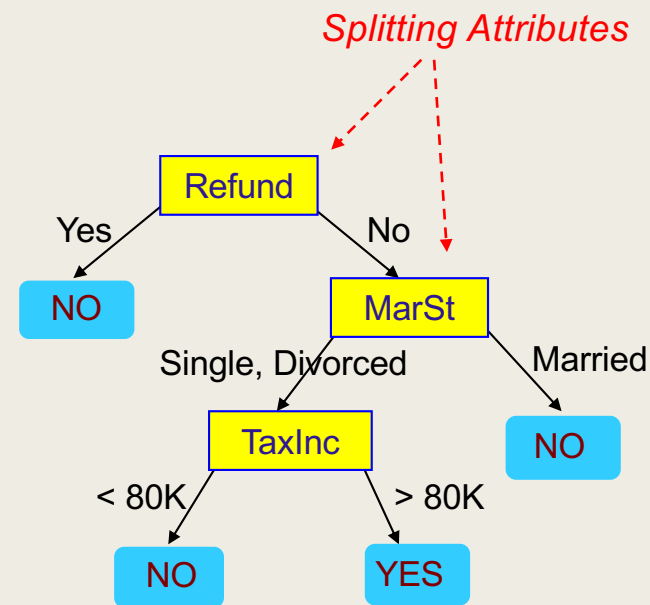- Naïve Bayes and Bayesian Belief Networks

- Support Vector Machines

# 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 |

*categorical  categorical  continuous  class*

Training Data

*Splitting Attributes*

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

NO          Refund

Yes          No

NO          TaxInc

< 80K          > 80K

NO          YES

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

# 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

Start from the root of tree.

Test Data

| 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 | ? |

# Apply Model to Test Data

Test Data

| 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

Test Data

| 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

Assign Cheat 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

Induction

**Tree Induction algorithm**

**Learn Model**

**Model**

Decision Tree

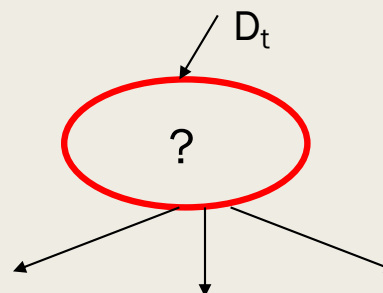| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

Test Set

**Apply Model**

Deduction

# Decision Tree Induction

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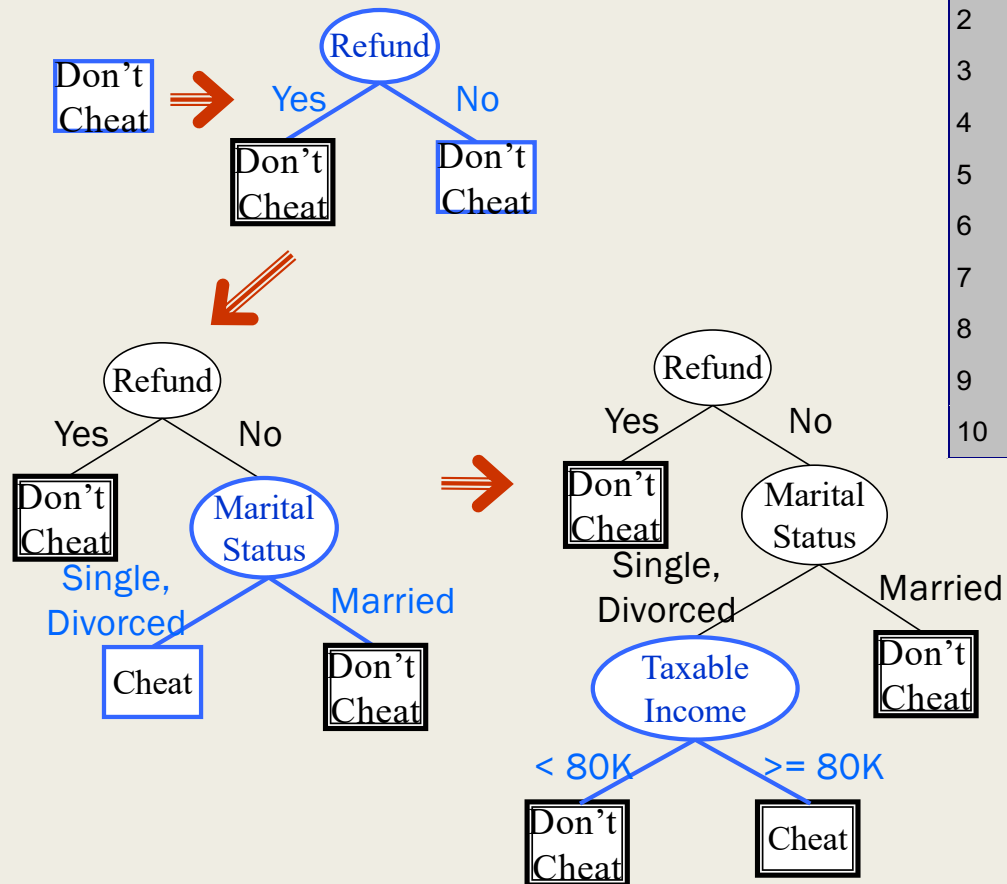  - *CART*
  - *ID3, C4.5*
  - *SLIQ,SPRINT*

# 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 to 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?
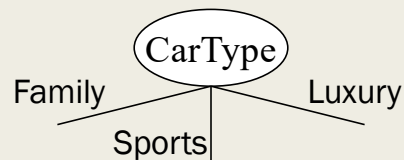  - *Determine when to stop splitting*

# How to Specify Test Condition?

- Depends on attribute types
  - *Nominal*
  - *Ordinal*
  - *Continuous*

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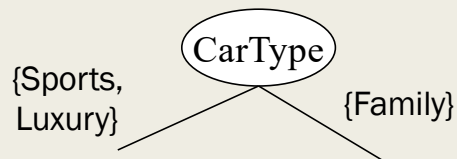  - *2-way split*
  - *Multi-way split*

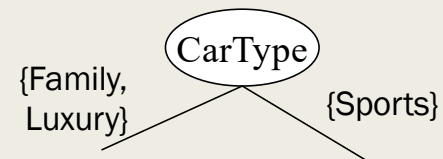# Splitting Based on Nominal Attributes

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



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

# Splitting Based on Ordinal Attributes

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



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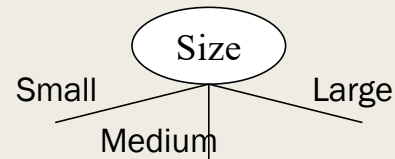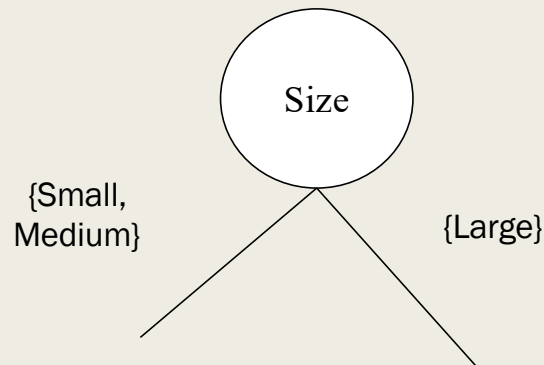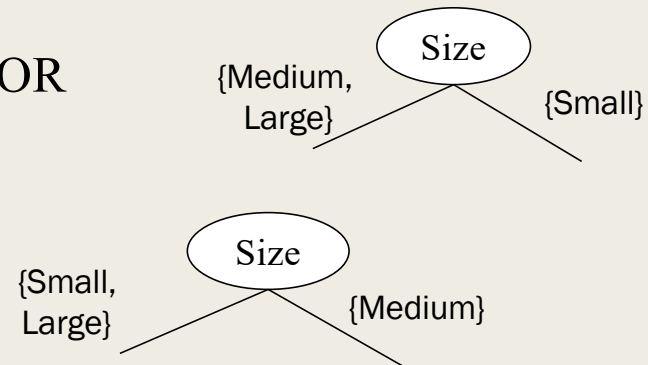  Need to find optimal partitioning.
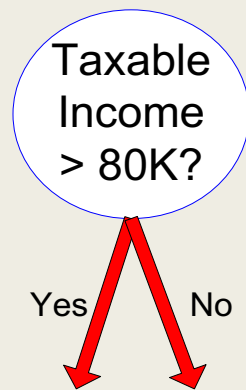
- What about this split?

# 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



(i) Binary split

(ii) Multi-way split

# LECTURE 15

Dr.Vani V

# 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?

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

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:

| C0: 5 |
| C1: 5 |

Non-homogeneous,

High degree of impurity

| C0: 9 |
| C1: 1 |

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** |

→ 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

# Measure of Impurity: GINI

■ Gini Index for a given node t :

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

(NOTE: $p(j \mid t)$ 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 interesting information*

– *Minimum (0.0) when all records belong to one class, implying most interesting information*

| 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** | |

# Examples for computing GINI

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

| | |
|----|---|
| 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

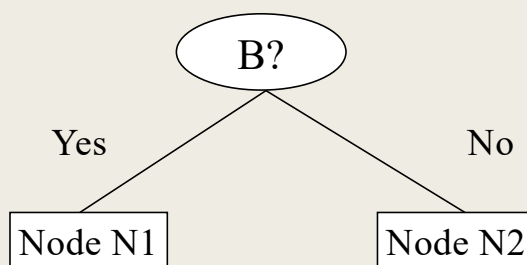Gini = 1 – (2/6)$^2$ – (4/6)$^2$ = 0.444

# Splitting Based on GINI

- Used in CART, SLIQ, SPRINT.

- When a node p is split into k partitions (children), the quality of split is computed as,

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

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

          n  = number of records at node p.

# 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** | |

B?

Yes          No

Node N1      Node N2

Gini(N1)
= $1 - (5/6)^2 - (2/6)^2$
= 0.194

|  | N1 | N2 |
|---|---|---|
| C1 | 5 | 1 |
| C2 | 2 | 4 |
| **Gini=0.333** | | |

Gini(N2)
= $1 - (1/6)^2 - (4/6)^2$
= 0.528

Gini(Children)
= 7/12 * 0.194 +
   5/12 * 0.528
= 0.333

# Categorical Attributes: Computing Gini Index

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

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

Multi-way split

| | CarType | | |
|---|---|---|---|
| | Family | Sports | Luxury |
| C1 | 1 | 2 | 1 |
| C2 | 4 | 1 | 1 |
| Gini | 0.393 | | |

Two-way split
(find best partition of values)

| | CarType | |
|---|---|---|
| | {Sports, Luxury} | {Family} |
| C1 | 3 | 1 |
| C2 | 2 | 4 |
| Gini | 0.400 | |

| | CarType | |
|---|---|---|
| | {Sports} | {Family, Luxury} |
| C1 | 2 | 2 |
| C2 | 1 | 5 |
| Gini | 0.419 | |

# Continuous Attributes: Computing Gini Index

- Use Binary Decisions based on one value
- Several Choices for the splitting value
  - *Number of possible splitting values = Number of distinct values*
- Each splitting value has a count matrix associated with it
  - *Class counts in each of the partitions, A < v and A ≥ 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.*

| 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 |

Taxable Income > 80K?

Yes     No

# 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 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **Taxable Income** | | | | | | | | | | | | | | | | | | | |
| | 60 | | 70 | | 75 | | 85 | | 90 | | 95 | | 100 | | 120 | | 125 | | 220 | |
| | 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 |
| 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 | |

## Alternative Splitting Criteria based on INFO

- Entropy at a given node t:

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

(NOTE: $p(j \mid t)$ is the relative frequency of class j at node t).

- *Measures homogeneity 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 similar to the GINI index computations*

# Examples for computing Entropy

$$Entropy(t) = -\sum_{j} p(j \mid t) \log_2 p(j \mid t)$$

| 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

# References

TEXTBOOKS :
1. Pang-Ning Tan, Vipin Kumar, Michael Steinbach: **Introduction to Data Mining**, Pearson, 2012.
2. Jiawei Han and Micheline Kamber: **Data Mining - Concepts and Techniques**, 3rd Edition, MorganKaufmann Publisher, 2014