# Data Mining

## Lecture 8

Ananya Jana
CS360

Fall 2024

# Classification

A classification technique (or classifier) is a systematic approach to building classification models from an input data set.

Examples include decision tree classifiers, rule-based classifiers, neural networks, support vector machines, and naïve Bayes classifiers.

Each technique employs a learning algorithm to identify a model that best fits the relationship between the attribute set and class label of the input data. The model generated by a learning algorithm should both fit the input data well and correctly predict the class labels of records it has never seen before. Therefore, **a key objective of the learning algorithm is to build models with good generalization capability**; i.e., models that accurately predict the class labels of previously unknown records.

# Classification

A general approach for solving classification problems: first, a training set consisting of records whose class labels are known must be provided. The training set is used to build a classification model, which is subsequently applied to the test set, which consists of records with unknown class labels.

# Classification

Evaluation of the performance of a classification model is based on the counts of test records correctly and incorrectly predicted by the model.

Confusion matrix for a 2-class problem.

| | | Predicted Class | |
|---|---|---|---|
| | | $Class = 1$ | $Class = 0$ |
| Actual | $Class = 1$ | $f_{11}$ | $f_{10}$ |
| Class | $Class = 0$ | $f_{01}$ | $f_{00}$ |

These counts are tabulated in a table known as a **confusion matrix**. The table depicts the confusion matrix for a binary classification problem. Each entry $f_{ij}$ in this table denotes the number of records from class i predicted to be of class j. For instance, $f_{01}$ is the number of records from class 0 incorrectly predicted as class 1. Based on the entries in the confusion matrix, the total number of correct predictions made by the model is $(f_{11} + f_{00})$ and the total number of incorrect predictions is $(f_{10} + f_{01})$.

# Classification

Although a confusion matrix provides the information needed to determine how well a classification model performs, summarizing this information with a single number would make it more convenient to compare the performance of different models. This can be done using a performance metric such as **accuracy**, which is defined as follows:

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} = \frac{f_{11} + f_{00}}{f_{11} + f_{10} + f_{01} + f_{00}}.$$

Equivalently, the performance of a model can be expressed in terms of its **error rate**, which is given by the following equation:

$$\text{Error rate} = \frac{\text{Number of wrong predictions}}{\text{Total number of predictions}} = \frac{f_{10} + f_{01}}{f_{11} + f_{10} + f_{01} + f_{00}}$$

# Decision Trees

We can solve a classification problem by asking a series of carefully crafted questions about the attributes of the test record. Each time we receive an answer, a follow-up question is asked until we reach a conclusion about the class label of the record. The series of questions and their possible answers can be organized in the form of a decision tree, which is a hierarchical structure consisting of nodes and directed edges.

# Classification

The decision tree has three types of nodes:

- A **root node** that has no incoming edges and zero or more outgoing edges.
- **Internal nodes**, each of which has exactly one incoming edge and two or more outgoing edges.
- **Leaf** or **terminal nodes**, each of which has exactly one incoming edge and no outgoing edges.

In a decision tree, each leaf node is assigned a class label. The non-terminal nodes, which include the root and other internal nodes, contain attribute test conditions to separate records that have different characteristics.
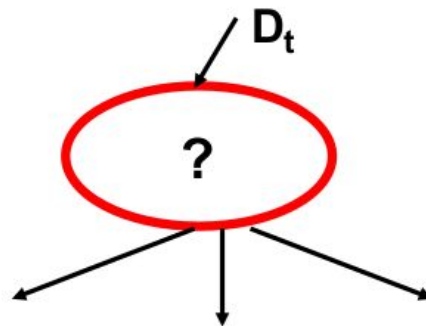
# Decision Tree Induction

- 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 $D_t$ 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.

| 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 |

$D_t$

?

# Hunt's Algorithm

Defaulted = No

**(7,3)**

(a)

| 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)

Defaulted = No
(7,3)

(b)

Home Owner
  Yes → Defaulted = No (3,0)
  No → Defaulted = No (4,3)

| 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)

(b)

(c)

(d)

# Design Issues of Decision Tree Induction

- How should training records be split?
  - Method for specifying test condition
    - depending on attribute types
  - Measure for evaluating the goodness of a test condition

- How should the splitting procedure stop?
  - Stop splitting if all the records belong to the same class or have identical attribute values
  - Early termination
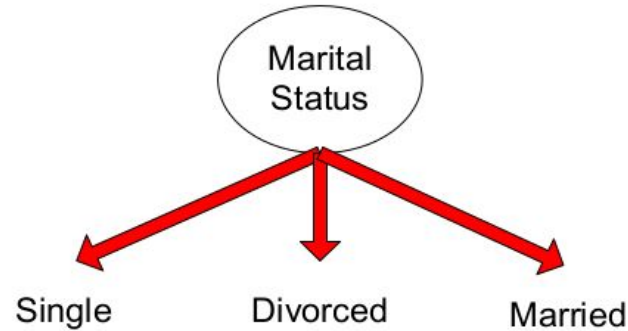
# Methods for specifying test conditions

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

- Depends on number of ways to split
    - 2-way split
    - Multi-way split
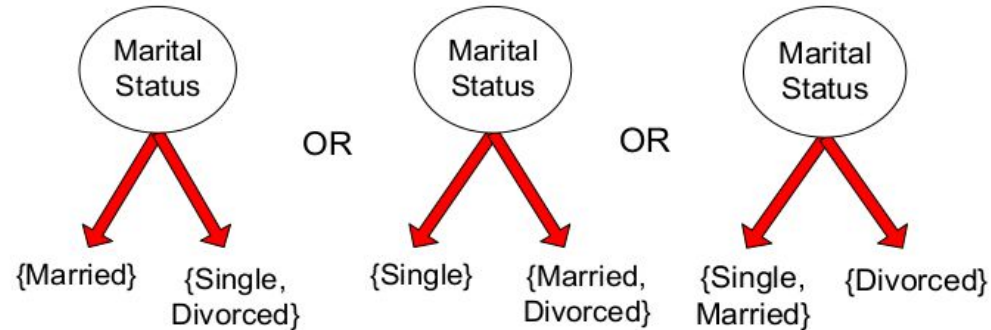
# Test condition for Nominal Attribute

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



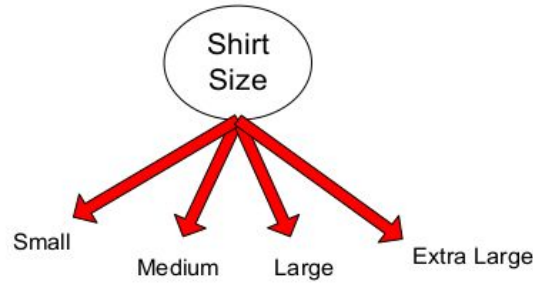Marital Status → Single, Divorced, Married

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



Marital Status: {Married} {Single, Divorced}  OR  Marital Status: {Single} {Married, Divorced}  OR  Marital Status: {Single, Married} {Divorced}
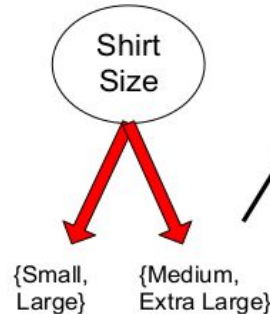
# Test condition for Ordinal Attribute

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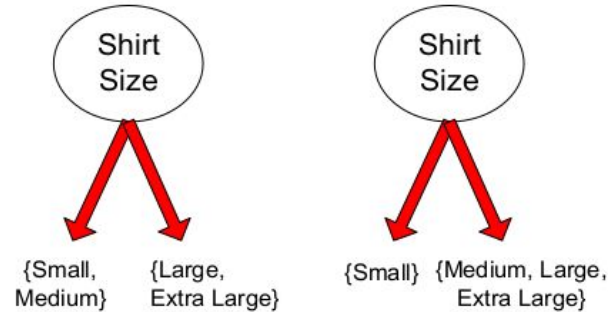  - Use as many partitions as distinct values

- **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

# Test condition for Continuous Attribute



(i) Binary split        (ii) Multi-way split

# Test condition for Continuous Attribute

- Different ways of handling
  - Discretization to form an ordinal categorical attribute

    Ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.
    - ◆ Static – discretize once at the beginning
    - ◆ Dynamic – repeat at each node

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

# How to determine the best split
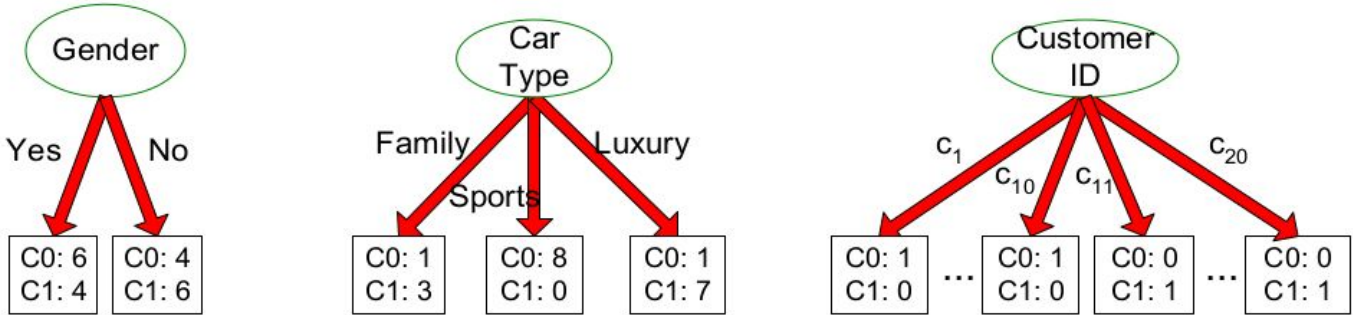
| Customer Id | Gender | Car Type | Shirt Size | Class |
|---|---|---|---|---|
| 1 | M | Family | Small | C0 |
| 2 | M | Sports | Medium | C0 |
| 3 | M | Sports | Medium | C0 |
| 4 | M | Sports | Large | C0 |
| 5 | M | Sports | Extra Large | C0 |
| 6 | M | Sports | Extra Large | C0 |
| 7 | F | Sports | 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**

**Gender**

Yes — No

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

**Car Type**

Family — Sports — Luxury

| 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 |

**Which test condition is the best?**

# How to determine the best split

- Greedy approach:
  - Nodes with <span style="color:red">purer</span> class distribution are preferred

- Need a measure of node impurity:

| C0: 5 |
|---|
| C1: 5 |

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

**High degree of impurity**          **Low degree of impurity**

# Measures of Node Impurity

- Gini Index

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

- Entropy

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

- Misclassification error

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

# Finding the best split

1. Compute impurity measure (P) before splitting
2. Compute impurity measure (M) after splitting
   - Compute impurity measure of each child node
   - M is the weighted impurity of children
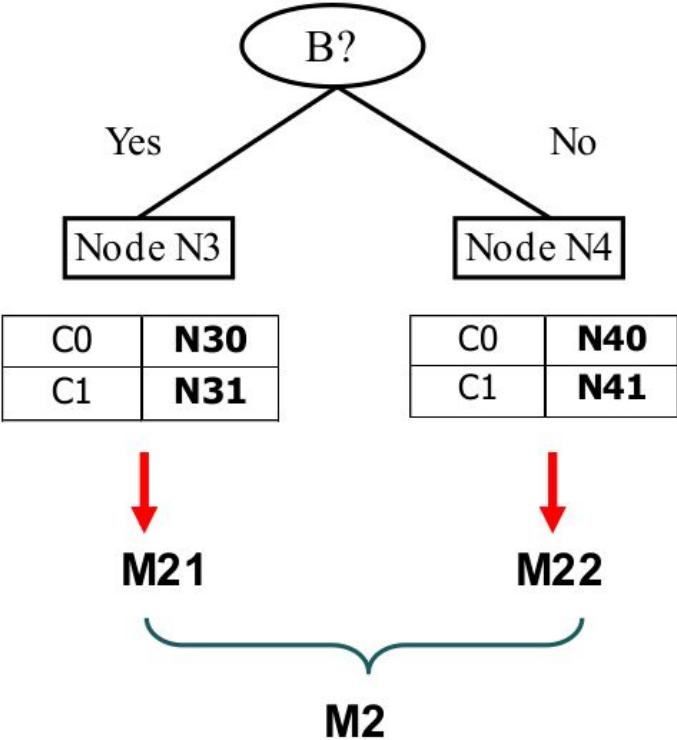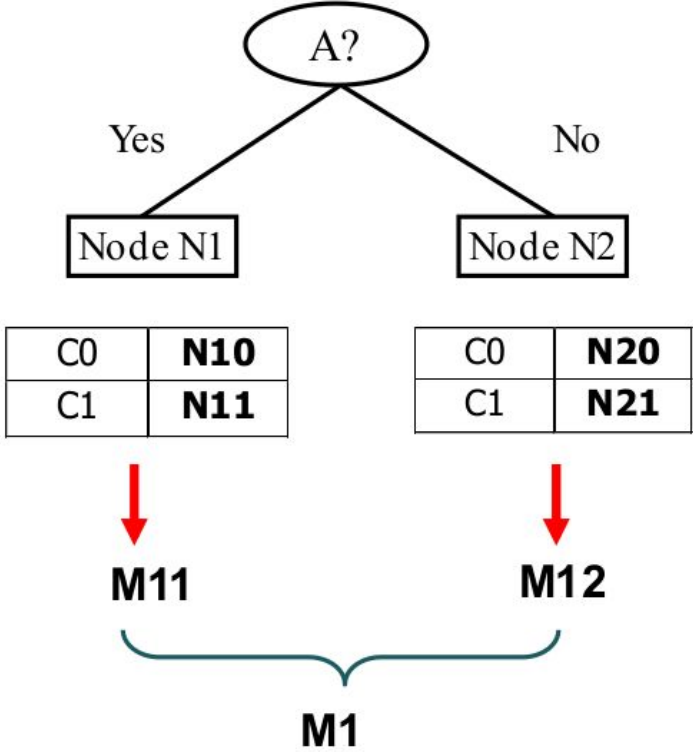3. Choose the attribute test condition that produces the highest gain

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

or equivalently, lowest impurity measure after splitting (M)

# Finding the best split



**Before Splitting:**

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

→ **P**

**A?**

Yes — Node N1

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

↓

**M11**

No — Node N2

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

↓

**M12**

**M1**

**B?**

Yes — Node N3

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

↓

**M21**

No — Node N4

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

↓

**M22**

**M2**

Gain = P – M1    vs    P – M2

# Measure of impurity: Entropy

- Entropy at a given node t:

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

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

- ◆ 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

# 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

– Choose the split that achieves most reduction (maximizes GAIN)

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

# Measure of Impurity: Gini

- Gini Index for a given node t :

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

(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

- Minimum (0.0) when all records belong to one class,

# 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).

- For 2-class problem (p, 1 − p):
  - GINI = $1 - p^2 - (1 - p)^2 = 2p\,(1\text{-}p)$

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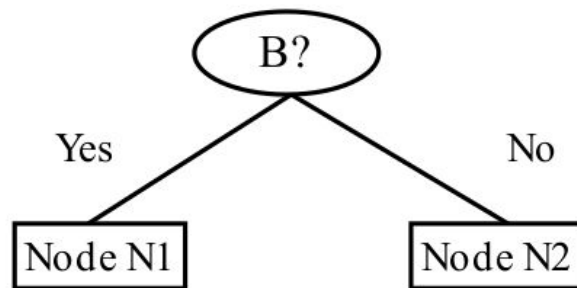
$$GINI_{split} = \sum_{i=1}^{k} \frac{n_i}{n} GINI(i)$$

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

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

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

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

Gain = 0.486 – 0.361 = 0.125

# Task

1. Compute the initial Gini index.
   Next, compute GINI index after you split based on the attribute
   i) Gender,
   ii) Car Type,
   iii) Shirt Size

2. Do the same exercise as 1 with entropy as the measure of impurity.

| 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 |