# Supervised vs. Unsupervised Learning

- Supervised learning (**classification**)
  - Supervision: The training data instances and their attributes/features are accompanied by labels indicating the class of the instances.
  - **Predict labels** for testing data instances.
- Unsupervised learning (**clustering**)
  - The class **labels** of training data is **unknown**
  - Given a set of attributes, with the aim of establishing the existence of classes or clusters.

| Machine learning types | Data mining tasks |
| --- | --- |
| Supervised learning | **Classification**<br>Regression<br>… |
| Unsupervised learning | Clustering<br>Pattern/Association mining<br>… |

# Classification: Definition

- Given a collection of records (training set )
  - Each record is by characterized by a tuple $(x,y)$, where $x$ is the attribute set and $y$ is the class label
    - $x$: attribute, predictor, independent variable, input
    - $y$: class, response, dependent variable, output

- Task:
  - Learn a model that maps each attribute set $x$ into one of the predefined class labels $y$

# Examples of Classification Task

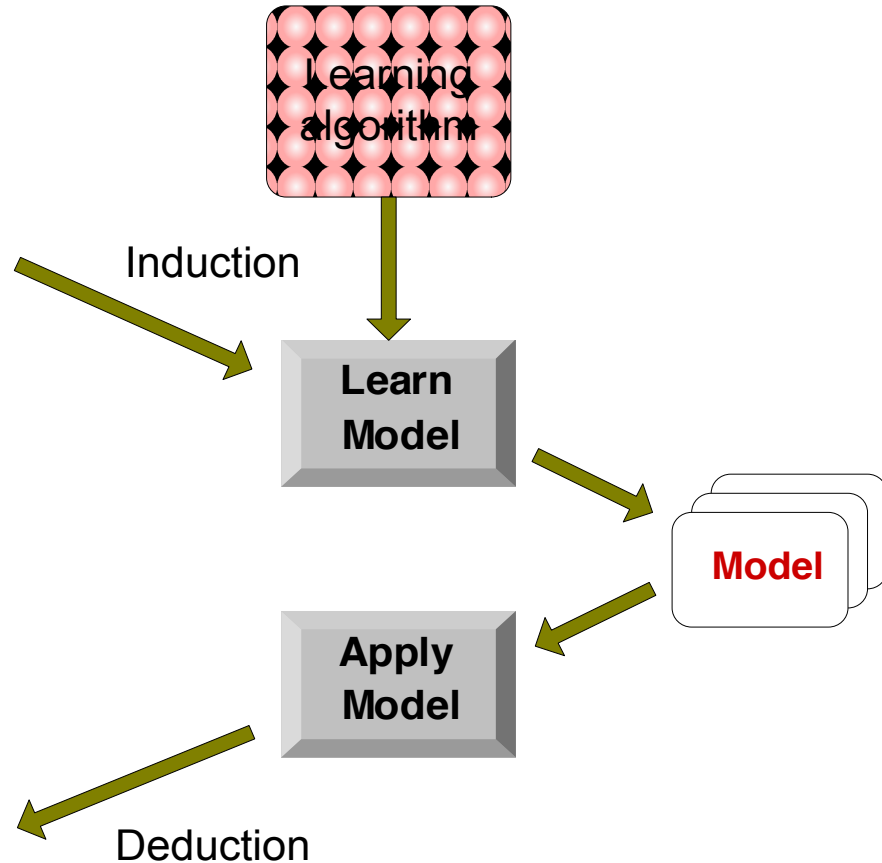| Task | Attribute set, $x$ | Class label, $y$ |
|---|---|---|
| Categorizing email messages | Features extracted from email message header and content | spam or non-spam |
| Identifying tumor cells | Features extracted from MRI scans | malignant or benign cells |
| Cataloging galaxies | Features extracted from telescope images | Elliptical, spiral, or irregular-shaped galaxies |

# General Approach for Building Classification Model

| 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**

| 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**

Learning algorithm

Induction

**Learn Model**

**Model**

**Apply Model**

Deduction

# Classification Techniques

- Base Classifiers
    - Decision Tree based Methods
    - Rule-based Methods
    - Nearest-neighbor
    - Neural Networks
    - Naïve Bayes and Bayesian Belief Networks
    - Support Vector Machines

- Ensemble Classifiers
    - Boosting, Bagging, Random Forests

# Rule-Based Classifier

- Classify records by using a collection of "if…then…" rules

- Rule:    (*Condition*) $\rightarrow$ *y*
  - where
    - *Condition* is a conjunctions of attributes
    - *y* is the class label
  - *LHS*: rule antecedent or condition
  - *RHS*: rule consequent
  - Examples of classification rules:
    - (Blood Type=Warm) $\wedge$ (Lay Eggs=Yes) $\rightarrow$ Birds
    - (Taxable Income < 50K) $\wedge$ (Refund=Yes) $\rightarrow$ Evade=No

# Rule-based Classifier (Example)

| Name | Body Temperature | Skin Cover | Gives Birth | Aquatic Creature | Aerial Creature | Has Legs | Hiber-nates | Class Label |
|------|------------------|------------|-------------|------------------|-----------------|----------|-------------|-------------|
| human | warm-blooded | hair | yes | no | no | yes | no | Mammals |
| python | cold-blooded | scales | no | no | no | no | yes | Reptiles |
| salmon | cold-blooded | scales | no | yes | no | no | no | Fishes |
| whale | warm-blooded | hair | yes | yes | no | no | no | Mammals |
| frog | cold-blooded | none | no | semi | no | yes | yes | Amphibians |
| komodo dragon | cold-blooded | scales | no | no | no | yes | no | Reptiles |
| bat | warm-blooded | hair | yes | no | yes | yes | yes | Mammals |
| pigeon | warm-blooded | feathers | no | no | yes | yes | no | Birds |
| cat | warm-blooded | fur | yes | no | no | yes | no | Mammals |
| guppy | cold-blooded | scales | yes | yes | no | no | no | Fishes |
| alligator | cold-blooded | scales | no | semi | no | yes | no | Reptiles |
| penguin | warm-blooded | feathers | no | semi | no | yes | no | Birds |
| porcupine | warm-blooded | quills | yes | no | no | yes | yes | Mammals |
| eel | cold-blooded | scales | no | yes | no | no | no | Fishes |
| salamander | cold-blooded | none | no | semi | no | yes | yes | Amphibians |

$r_1$: (Gives Birth = no) $\land$ (Aerial Creature = yes) $\longrightarrow$ Birds
$r_2$: (Gives Birth = no) $\land$ (Aquatic Creature = yes) $\longrightarrow$ Fishes
$r_3$: (Gives Birth = yes) $\land$ (Body Temperature = warm-blooded) $\longrightarrow$ Mammals
$r_4$: (Gives Birth = no) $\land$ (Aerial Creature = no) $\longrightarrow$ Reptiles
$r_5$: (Aquatic Creature = semi) $\longrightarrow$ Amphibians

# Application of Rule-Based Classifier

- A rule *r* <span style="color:red">covers</span> an instance **x** if the attributes of the instance satisfy the condition of the rule

$r_1$:  (Gives Birth = no) $\wedge$ (Aerial Creature = yes) $\longrightarrow$ Birds
$r_2$:  (Gives Birth = no) $\wedge$ (Aquatic Creature = yes) $\longrightarrow$ Fishes
$r_3$:  (Gives Birth = yes) $\wedge$ (Body Temperature = warm-blooded) $\longrightarrow$ Mammals
$r_4$:  (Gives Birth = no) $\wedge$ (Aerial Creature = no) $\longrightarrow$ Reptiles
$r_5$:  (Aquatic Creature = semi) $\longrightarrow$ Amphibians

| Name | Body Temperature | Skin Cover | Gives Birth | Aquatic Creature | Aerial Creature | Has Legs | Hiber-nates |
|------|------------------|------------|-------------|------------------|-----------------|----------|-------------|
| hawk | warm-blooded | feather | no | no | yes | yes | no |
| grizzly bear | warm-blooded | fur | yes | no | no | yes | yes |

The rule r1 covers a hawk => Bird

The rule r3 covers the grizzly bear => Mammal

# Rule Coverage and Accuracy

- Coverage of a rule:
  - Fraction of records that satisfy the antecedent of a rule
- Accuracy of a rule:
  - Fraction of records that satisfy both the antecedent and consequent of a rule

| Tid | Refund | Marital Status | Taxable Income | Class |
|-----|--------|----------------|----------------|-------|
| 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** |

**(Status=Single) $\rightarrow$ No**

**Coverage = 40%,  Accuracy = 50%**

# How does Rule-based Classifier Work?

$r_1$:  (Gives Birth = no) ∧ (Aerial Creature = yes) ⟶ Birds
$r_2$:  (Gives Birth = no) ∧ (Aquatic Creature = yes) ⟶ Fishes
$r_3$:  (Gives Birth = yes) ∧ (Body Temperature = warm-blooded) ⟶ Mammals
$r_4$:  (Gives Birth = no) ∧ (Aerial Creature = no) ⟶ Reptiles
$r_5$:  (Aquatic Creature = semi) ⟶ Amphibians

| Name | Body Temperature | Skin Cover | Gives Birth | Aquatic Creature | Aerial Creature | Has Legs | Hiber-nates |
|---|---|---|---|---|---|---|---|
| lemur | warm-blooded | fur | yes | no | no | yes | yes |
| turtle | cold-blooded | scales | no | semi | no | yes | no |
| dogfish shark | cold-blooded | scales | yes | yes | no | no | no |

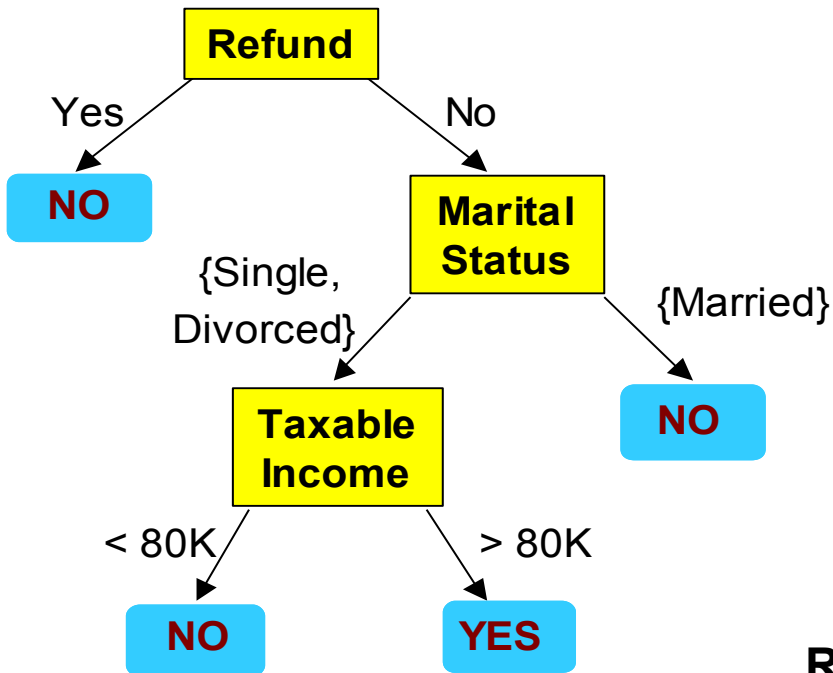A lemur triggers rule r3, so it is classified as a mammal

A turtle triggers both r4 and r5

A dogfish shark triggers none of the rules

# Characteristics of Rule-Based Classifier

- Mutually exclusive rules
    - Classifier contains mutually exclusive rules if the rules are independent of each other
    - Every record is covered by at most one rule

- Exhaustive rules
    - Classifier has exhaustive coverage if it accounts for every possible combination of attribute values
    - Each record is covered by at least one rule

# From Decision Trees To Rules

**Refund**

Yes         No

**NO**

**Marital Status**

{Single, Divorced}       {Married}

**Taxable Income**

**NO**

< 80K     > 80K

**NO**     **YES**

---

## Classification Rules

(Refund=Yes) ==> No

(Refund=No, Marital Status={Single,Divorced}, Taxable Income<80K) ==> No

(Refund=No, Marital Status={Single,Divorced}, Taxable Income>80K) ==> Yes

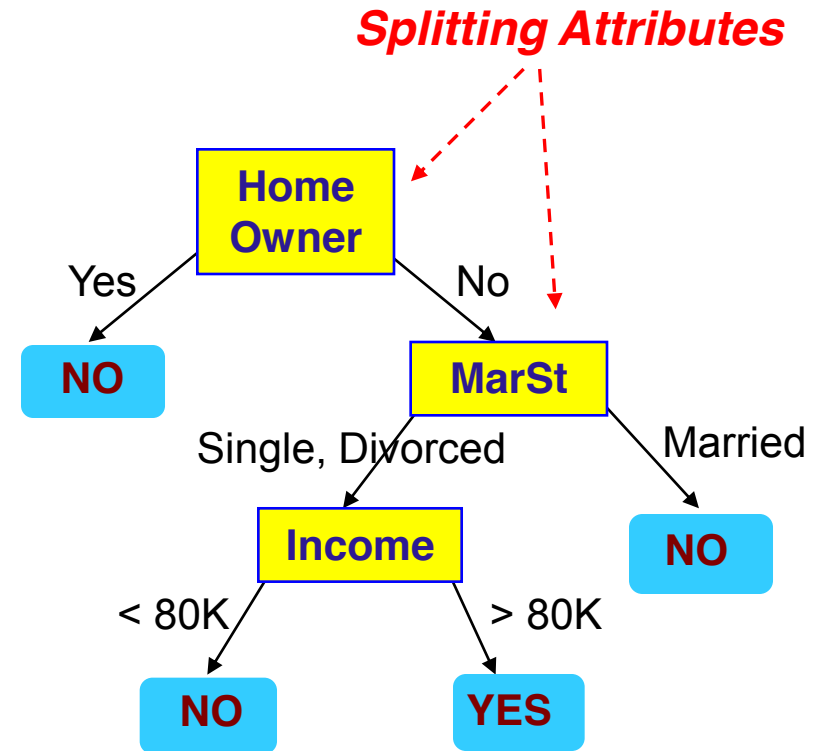(Refund=No, Marital Status={Married}) ==> No

**Rules are mutually exclusive and exhaustive**

**Rule set contains as much information as the tree**

# 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 |

**Training Data**

*Splitting Attributes*

**Model:  Decision Tree**

# Another Example of 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



**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

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

**Home Owner**

Yes → **NO**

No → **MarSt**

Single, Divorced → **Income**

Married → **NO**

< 80K → **NO**

> 80K → **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

MarSt

Single, Divorced → Income

Married → NO

Income

< 80K → NO

> 80K → 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 → **Income**

Married → **NO**

< 80K → **NO**

> 80K → **YES**

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

| 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

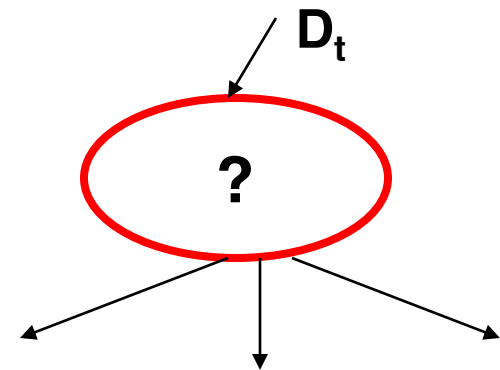Deduction

# Decision Tree Induction

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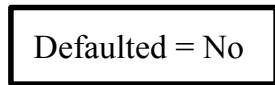    - ID3, C4.5
    - SLIQ,SPRINT

# 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.
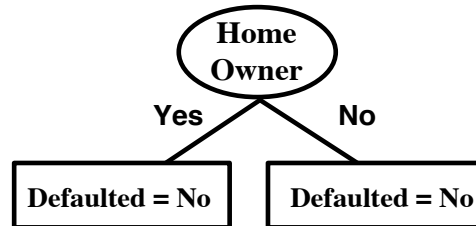
| ID | Home Owner | Marital Status | Annual Income | Defaulted Borrower |
|----|-----------|----------------|---------------|-------------------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | 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

| 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
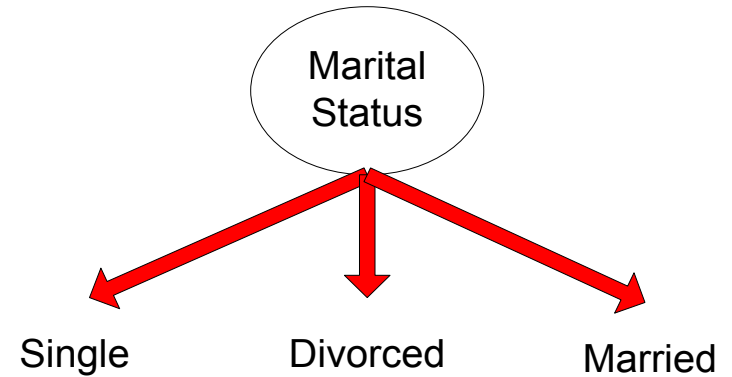  - Early termination

# Methods for Expressing Test Conditions

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

- Depends on number of ways to split
  - 2-way split
  - Multi-way split
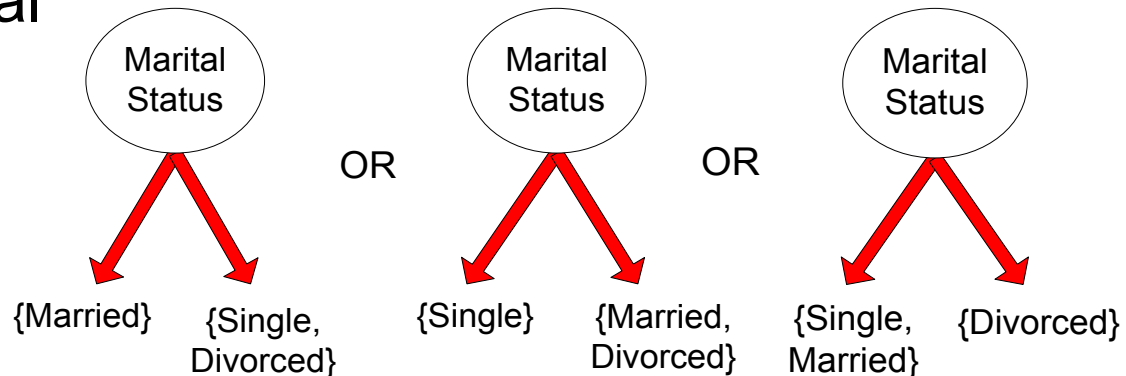
# Test Condition for Nominal Attributes

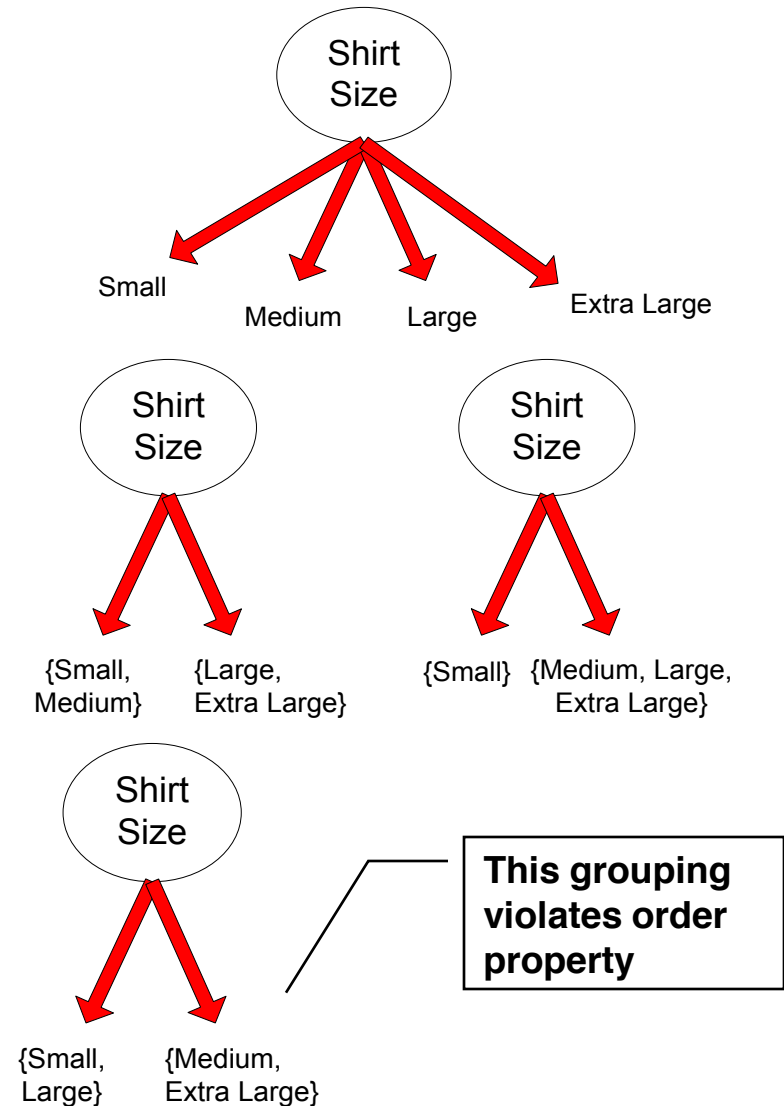- ### Multi-way split:
  - Use as many partitions as distinct values.

```
          Marital
          Status

    Single    Divorced    Married
```

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

```
   Marital              Marital              Marital
   Status      OR       Status      OR       Status

{Married}  {Single,  {Single}  {Married,  {Single,  {Divorced}
           Divorced}            Divorced}  Married}
```
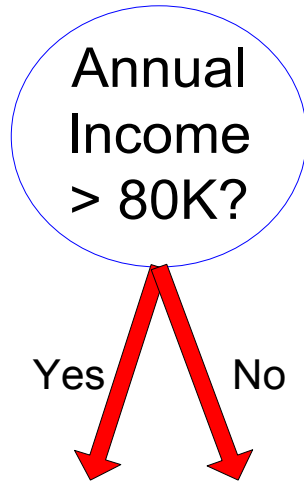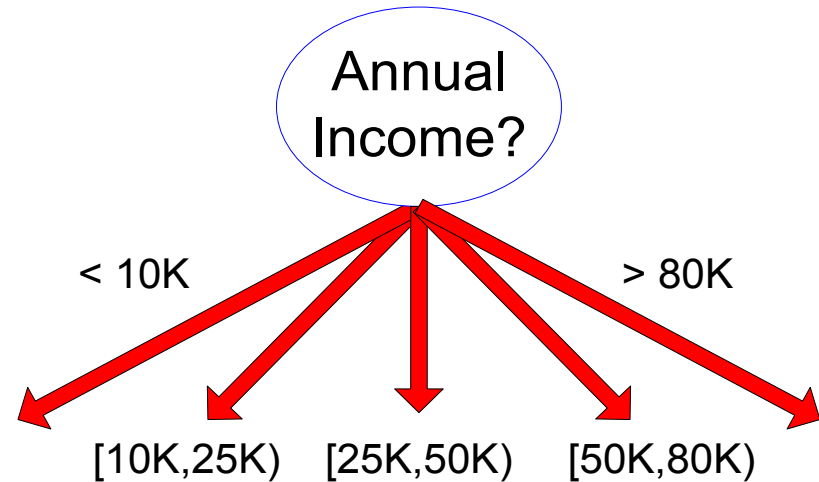
# Test Condition for Ordinal Attributes

- ● Multi-way split:
  - – Use as many partitions as distinct values

- ● Binary split:
  - – Divides values into two subsets
  - – Need to find optimal partitioning
  - – Preserve the 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 Attributes
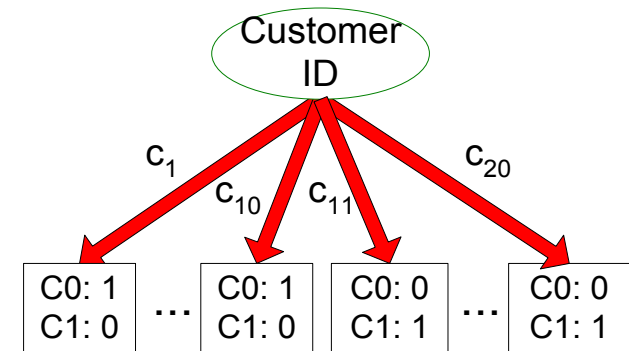


(i) Binary split

(ii) Multi-way 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

# 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**



**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 |

**High degree of impurity**

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

**Low degree of impurity**

# Measures of Node Impurity

- Gini Index

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

- Entropy

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

- Misclassification error

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

# 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
   - Compute the average impurity of the children (M)

3. Choose the attribute test condition that produces the highest gain

$$\textbf{Gain = P – M}$$

or equivalently, lowest impurity measure after splitting (M)

# Finding the Best Split

**Before Splitting:**

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

⟶ **P**

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**

# 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 | |

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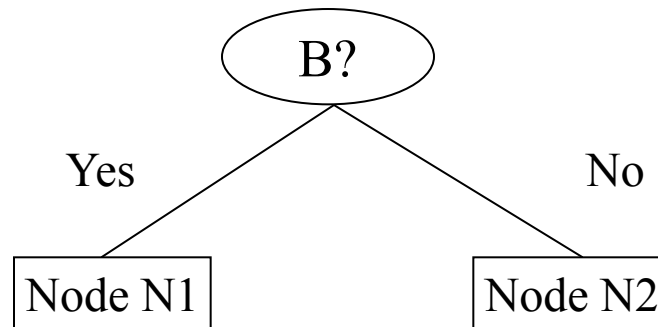
$$GINI_{split} = \sum_{i=1}^{k} \frac{n_i}{n} GINI(i)$$

  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 | 6 |
| C2 | 6 |
| **Gini = 0.500** | |

B?

Yes ⟶ Node N1

No ⟶ Node N2

**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** | | |

**Gini(Children)**
= 6/12 * 0.278 +
  6/12 * 0.444
= 0.361

# 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 | | Two-way split (find best partition of values) |
|---|---|---|

Multi-way split

| | CarType | | |
|---|---|---|---|
| | **Family** | **Sports** | **Luxury** |
| **C1** | 1 | 8 | 1 |
| **C2** | 3 | 0 | 7 |
| **Gini** | **0.163** | | |

Two-way split
(find best partition of values)

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

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

# 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 $\geq$ 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 > 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

| 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: $p(j \mid t)$ 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

   – Entropy based computations are quite similar to the GINI index computations

# Computing Entropy of a Single Node

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

| 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:

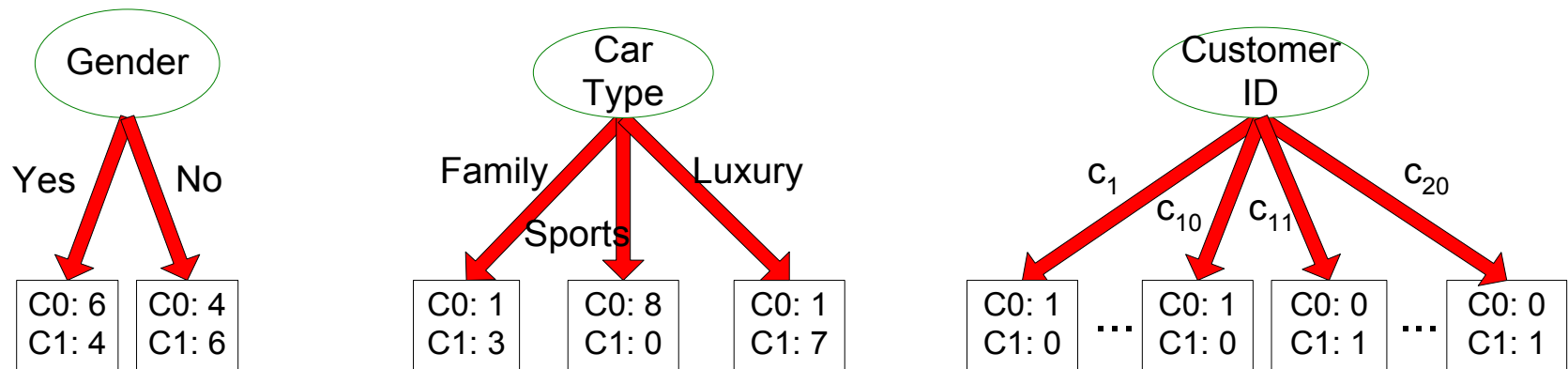$$GAIN_{split} = Entropy(p) - \left( \sum_{i=1}^{k} \frac{n_i}{n} Entropy(i) \right)$$

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

# Problems with Information Gain

- Info Gain tends to prefer splits that result in large number of partitions, each being small but pure



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

# Gain Ratio

- Gain Ratio:

$$GainRATIO_{split} = \frac{GAIN_{Split}}{SplitINFO} \qquad SplitINFO = -\sum_{i=1}^{k} \frac{n_i}{n} \log \frac{n_i}{n}$$

  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

# Measure of Impurity: Classification Error

- Classification error at a node t :

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

  - Maximum $(1 - 1/n_c)$ when records are equally distributed among all classes, implying least interesting information
  - Minimum (0) when all records belong to one class, implying most interesting 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

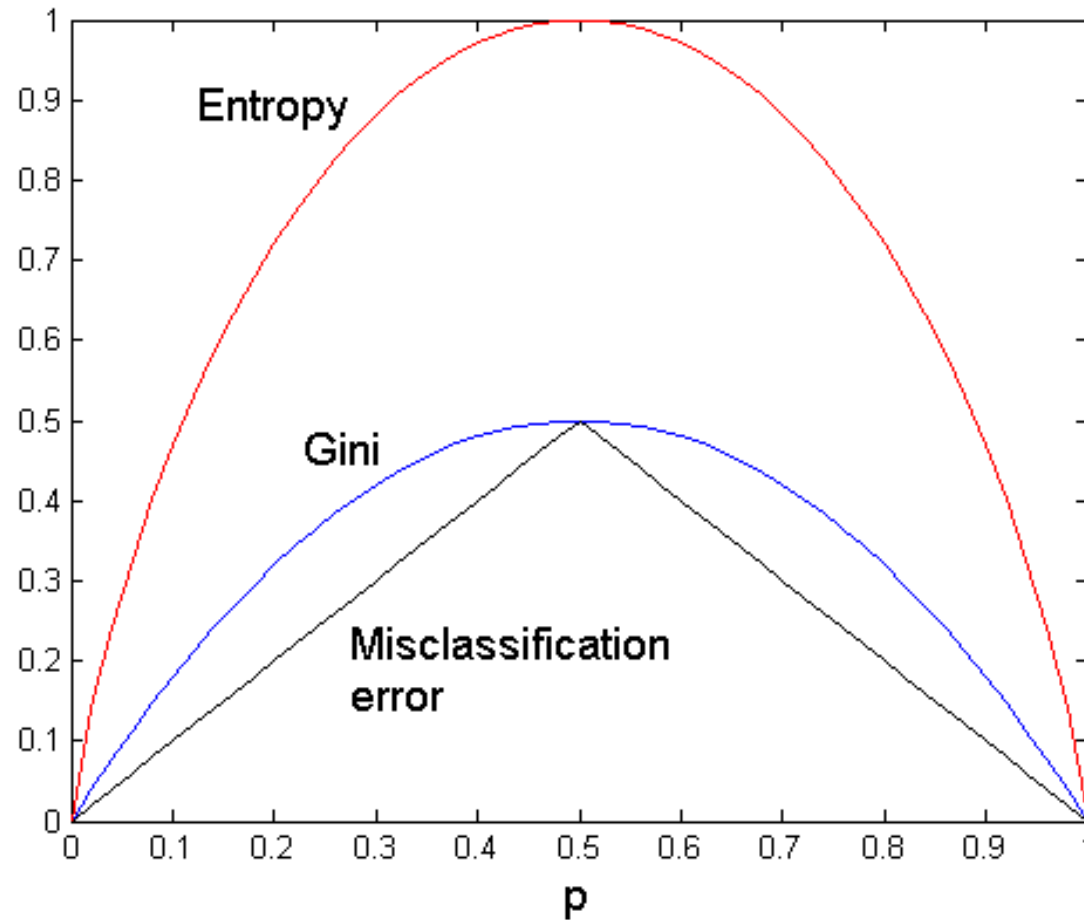Error = 1 – max (1/6, 5/6) = 1 – 5/6 = 1/6

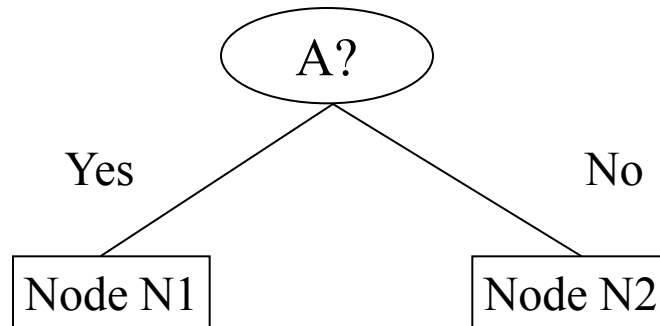| C1 | 2 |
|----|---|
| C2 | 4 |

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



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

**Gini(N1)**
**= 1 − (3/3)² − (0/3)²**
**= 0**

**Gini(N2)**
**= 1 − (4/7)² − (3/7)²**
**= 0.489**

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

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

**Gini improves but error remains the same!!**

# 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