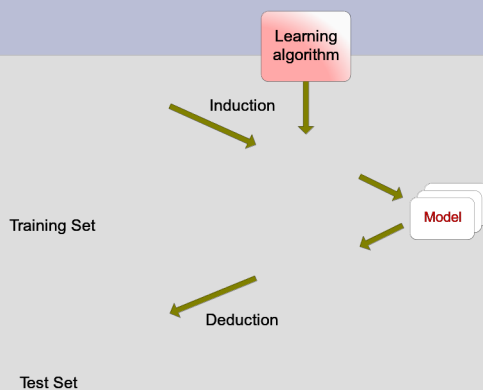# Classification: Basic Concepts, Decision Trees, and Model Evaluation

Chapter 4:
Introduction to Data Mining
by
Tan, Steinbach, Kumar

# Classification: Definition
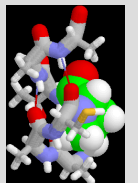
- Given a collection of records (*training set* )
  - Each record contains a set of *attributes*, one of the attributes is the *class*.
- Find a *model*  for class attribute as a function of the values of other attributes.
- Goal: previously unseen records should be assigned a class as accurately as 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
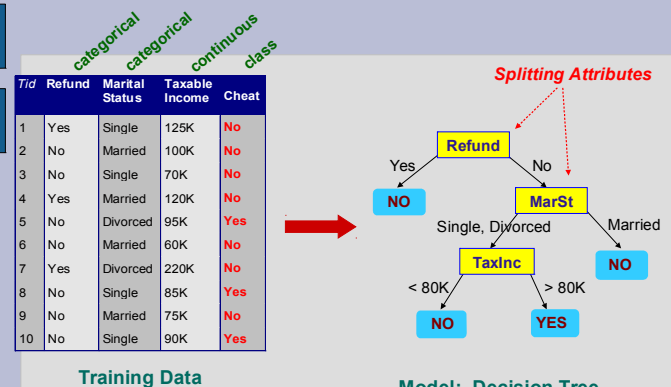


# Examples of Classification Task

- 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
- K-Nearest Neighbor
- 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 |

Training Data

Splitting Attributes

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

Model:  Decision Tree
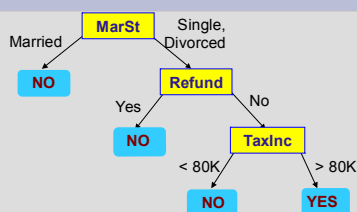
BETTER to exchange the 'Refund' with 'House Owner' (ref: scanned photocopies of book, chapter 4)
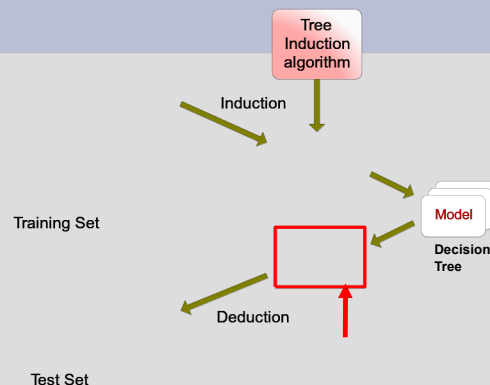
# Another Example of Decision Tree

| 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*

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

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

---

# Decision Tree Classification Task

Tree Induction algorithm

Training Set → Induction → Model (Decision Tree)

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 | ? |

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

---

# Apply Model to Test Data

**Test Data**

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

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

---

# Apply Model to Test Data

**Test Data**

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

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

---

# Apply Model to Test Data

**Test Data**

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

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

## Apply Model to Test Data

**Test Data**

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

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

## Apply Model to Test Data

**Test Data**

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

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

Assign Cheat to "No"

## Decision Tree Classification Task

Tree Induction algorithm

Induction

Training Set

Model
**Decision Tree**
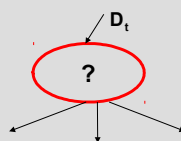
Deduction

Test Set

## Decision Tree Induction

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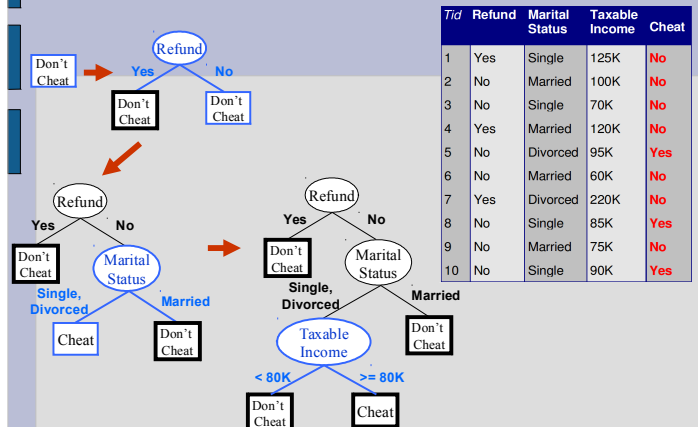  - CART
  - ID3, C4.5
  - SLIQ, SPRINT

## 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.

| 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 |

Don't Cheat →

**Refund**
- Yes → Don't Cheat
- No → Don't Cheat

**Refund**
- Yes → Don't Cheat
- No → **Marital Status**
  - Single, Divorced → Cheat
  - Married → Don't Cheat

**Refund**
- Yes → Don't Cheat
- No → **Marital Status**
  - Single, Divorced → **Taxable Income**
    - < 80K → Don't Cheat
    - >= 80K → Cheat
  - Married → Don't Cheat

## Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.

- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting

## How to Specify Test Condition?

- Depends on attribute types
  - Nominal
  - Ordinal
  - Continuous

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

## Splitting Based on Nominal Attributes

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

CarType — Family, Luxury, Sports

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

CarType {Sports, Luxury} {Family}   OR   CarType {Family, Luxury} {Sports}

## Splitting Based on Ordinal Attributes

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

Size — Small, Medium, Large

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

Size {Small, Medium} {Large}   OR   Size {Medium, Large} {Small}

- What about this split?

Size {Small, Large} {Medium}

## Splitting Based on Continuous Attributes

Taxable Income > 80K?
Yes   No

Taxable Income?
< 10K   [10K,25K]   [25K,50K]   [50K,80K]   > 80K

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

## How to determine the Best Split

- Greedy approach:
  - Nodes with homogeneous class distribution are preferred
- Need a measure of node impurity:

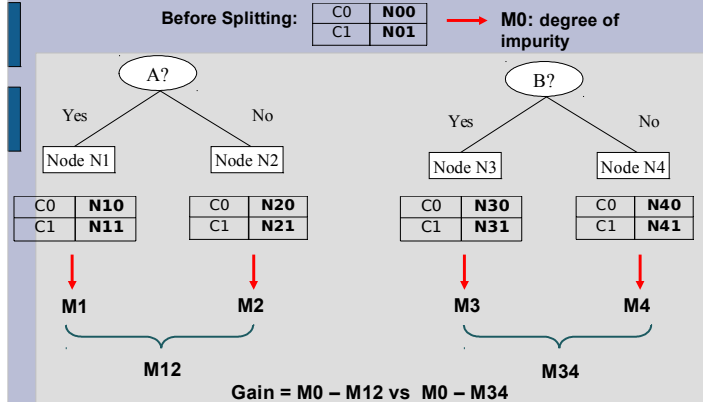C0: 5
C1: 5

C0: 9
C1: 1

Non-homogeneous,
High degree of impurity

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: degree of impurity

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
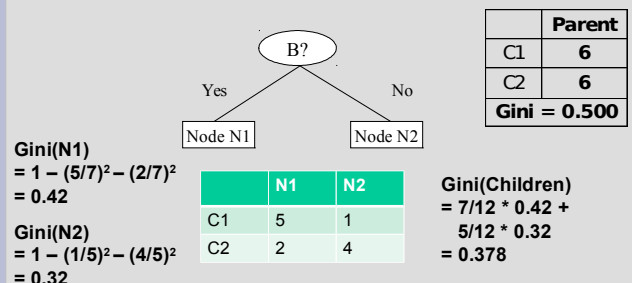
## Splitting Based on GINI

- 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.

B?

Yes    No

Node N1    Node N2

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

Gini(N1)
= 1 – (5/7)² – (2/7)²
= 0.42

Gini(N2)
= 1 – (1/5)² – (4/5)²
= 0.32

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

Gini(Children)
= 7/12 * 0.42 +
    5/12 * 0.32
= 0.378

## 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

## Splitting Based on INFO...

● 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)
– Disadvantage: Tends to prefer splits that result in large number of partitions, each being small but pure.

## 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