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

Lecture Notes for Chapter 4

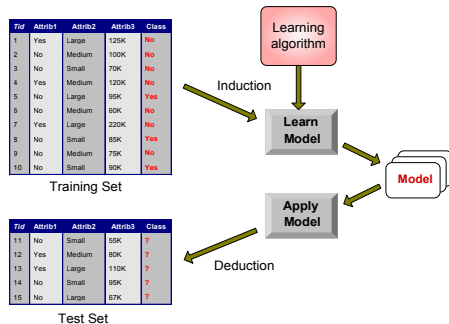Data Mining
by
Zhaonian Zou

## Outline

- Basic Concepts
- Decision Tree based Methods
- Model Evaluation
- Nearest Neighbor Methods
- Neural Networks
- Naïve Bayes and Bayesian Belief Networks
- Support Vector Machines

## 4.1. Basic Concepts
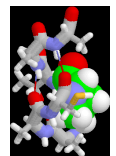
## Classification: Definition

- 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
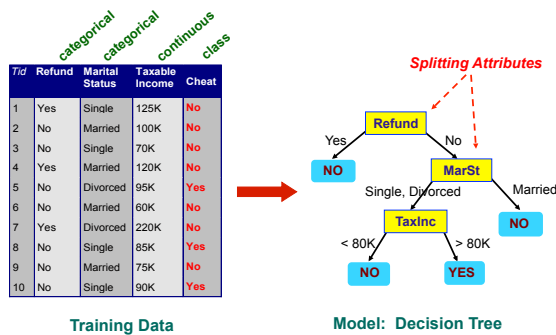
## Classification Techniques

- Decision Tree based Methods
- Rule-based Methods
- Memory based Reasoning
- Neural Networks
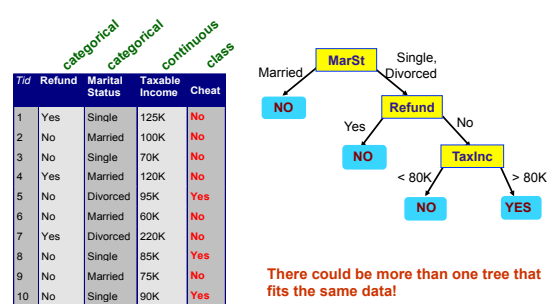- Naïve Bayes and Bayesian Belief Networks
- Support Vector Machines

---

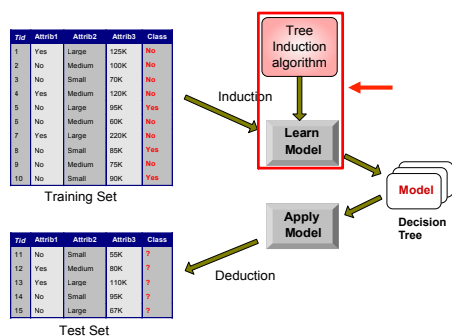## 4.2. Decision Tree based Methods
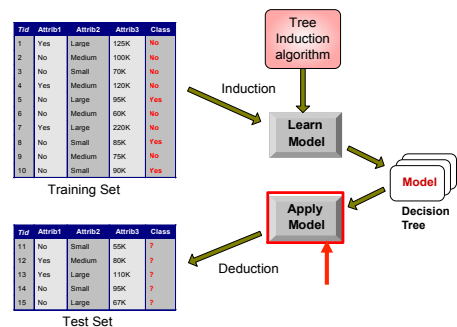
4.2.1. Basic Concepts

---

## 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 |

*Splitting Attributes*

**Training Data**          **Model: Decision Tree**

---

## 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 |

**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

Induction

Tree Induction algorithm

Learn Model

Apply 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

Deduction

---

## 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

Apply Model

Model → Decision Tree

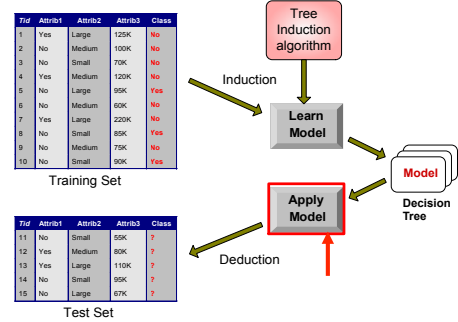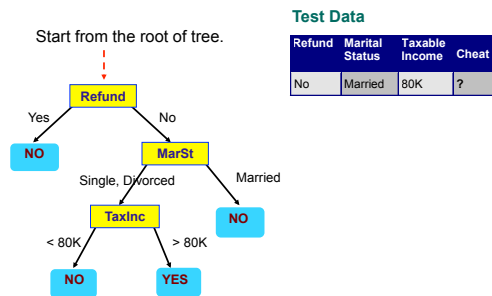| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

Test Set

Deduction

# 4.2. Decision Tree based Methods

4.2.2. Decision Tree Deduction

---

# 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 → Learn Model ← Tree Induction algorithm

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

Deduction → Apply Model

---

# 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
No → MarSt
MarSt: Single, Divorced → TaxInc ; Married → NO
TaxInc: < 80K → NO ; > 80K → YES

---

# Apply Model to Test Data

**Test Data**

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

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

---

# Apply Model to Test Data

**Test Data**

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

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

---

# Apply Model to Test Data

**Test Data**

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

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

TaxInc — NO

< 80K / > 80K

NO / YES

---

## 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"

---

# 4.2. Decision Tree based Methods

4.2.3. Decision Tree Induction

---

## 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

| 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 → Model → Decision Tree

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 $D_t$ 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?
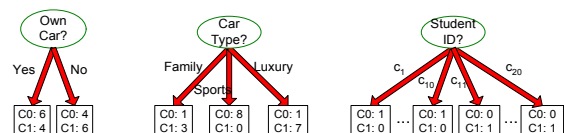    - How to determine the best split?
  - Determine when to stop splitting

---

## 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

---

## 4.2. Decision Tree based Methods

4.2.3. Decision Tree Induction

Issue 1: Specify Test Condition

---

## 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.

  Size
  Small — Medium — Large

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

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

- What about this split?   {Small, Large}   Size   {Medium}

## 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

Taxable Income > 80K?
Yes / No

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

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

---

### 4.2. Decision Tree based Methods

4.2.3. Decision Tree Induction

Issue 2: Determine the Best 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

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

Own Car?
Yes / No
C0: 6 / C1: 4      C0: 4 / C1: 6

Car Type?
Family / Sports / Luxury
C0: 1 / C1: 3      C0: 8 / C1: 0      C0: 1 / C1: 7

Student ID?
$c_1$ ... $c_{10}$ / $c_{11}$ ... $c_{20}$
C0: 1 / C1: 0   ...   C0: 1 / C1: 0      C0: 0 / C1: 1   ...   C0: 0 / 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 |

B?
Yes — No

Node N3        Node N4

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

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

M1        M2            M3        M4

M12                      M34

**Gain = M0 – M12 vs  M0 – M34**

---

## Measures of Node Impurity

- Gini Index

- Entropy

- Misclassification error

---

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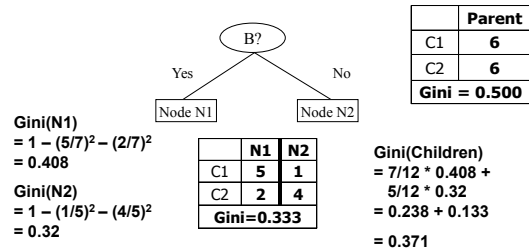- When a node p is split into k partitions (children), the quality of split is computed as,

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

where,  $n_i$ = number of records at child i,
  n = number of records at node p.

A?

a1  a2 ... ak

| Node N1 | Node N2 | Node Nk |

---

## 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
= 0.408

Gini(N2)
= 1 – (1/5)² – (4/5)²
= 0.32

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

Gini(Children)
= 7/12 * 0.408 +
  5/12 * 0.32
= 0.238 + 0.133

= 0.371

---

## 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 |
| 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

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

---

## Measures of Node Impurity

- Gini Index

- Entropy

- Misclassification error

## Measure of Impurity: Entropy

- Entropy at a given node t:

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

  (NOTE: $p(j\,|\,t)$ 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

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

| 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) $\log_2$ (5/6) = 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 Entropy

- 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 ID3 and C4.5
  – Disadvantage: Tends to prefer splits that result in large number of partitions, each being small but pure.

## Splitting Based on Entropy

- Gain Ratio:

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

  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
  – Designed to overcome the disadvantage of Information Gain

## Measures of Node Impurity

- Gini Index

- Entropy

- Misclassification error

## Measure of Impurity: Classification Error

- Classification error at a node t :

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

- Measures misclassification error made by a node.
    - 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

## Examples for Computing Error

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

| 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

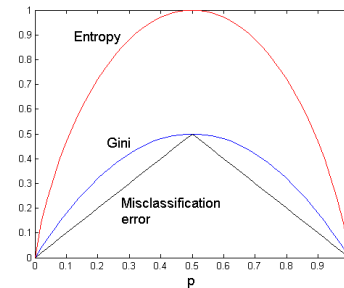Error = 1 – max (1/6, 5/6) = 1 – 5/6 = 1/6

| 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 Splitting Criteria

**For a 2-class problem:**



## Misclassification Error vs Gini



|  | Parent |
|----|--------|
| C1 | 7 |
| C2 | 3 |
| **Gini = 0.42** | |
| **Error = 0.3** | |

Gini(N1)
= 1 – (3/3)² – (0/3)²
= 0

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

Error(N1) = 0

Error(N2) = 3/7

|  | N1 | N2 |
|----|----|----|
| C1 | 3 | 4 |
| C2 | 0 | 3 |
| **Gini=0.361** | | |
| **Error=0.3** | | |

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

Error(Children)

=3/10 * 0 + 7/10 * 3/7 = 0.3

**Gini improves !!**

## 4.2. Decision Tree based Methods

4.2.3. Decision Tree Induction

Issue 3: Stopping Criteria

## 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

## Stopping Criteria for Tree Induction

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

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

- Early termination (to be discussed later)

## 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

---

## 4.2. Decision Tree based Methods

4.2.4. Model Overfitting

---

## Classification Errors

|  | PREDICTED CLASS | |
|---|---|---|
|  |  | Class=Yes | Class=No |
| ACTUAL CLASS | Class=Yes | a (TP) | b (FN) |
|  | Class=No | c (FP) | d (TN) |

● Most widely-used metric:

$$\text{Accuracy} = \frac{a+d}{a+b+c+d} = \frac{TP+TN}{TP+TN+FP+FN}$$

$$Error = \frac{b+c}{a+b+c+d} = \frac{FP+FN}{TP+TN+FP+FN}$$

---

## Underfitting and Overfitting (Example)



**500 circular and 500 triangular data points.**

**Circular points:**

$0.5 \leq \text{sqrt}(x_1^2+x_2^2) \leq 1$

**Triangular points:**

$\text{sqrt}(x_1^2+x_2^2) < 0.5$ or
$\text{sqrt}(x_1^2+x_2^2) > 1$

---

## Underfitting and Overfitting



**Underfitting**: when model is too simple, both training and test errors are large

---

## Overfitting due to Noise



**Decision boundary is distorted by noise point**

## Overfitting due to Insufficient Examples



**Lack of data points in the lower half of the diagram makes it difficult to predict correctly the class labels of that region**

**- Insufficient number of training records in the region causes the decision tree to predict the test examples using other training records that are irrelevant to the classification task**

## Notes on Overfitting

- Overfitting results in decision trees that are more complex than necessary

- Training error no longer provides a good estimate of how well the tree will perform on previously unseen records

- Need new ways for estimating errors

---

# 4.2. Decision Tree based Methods

4.2.4. Model Overfitting

Generalization Errors

## Estimating Generalization Errors

- Training errors: error on training ($\Sigma$ e(t) )
- Generalization errors: error on testing ($\Sigma$ e' (t))
- Methods for estimating generalization errors:
  - Optimistic approach: e' (t) = e(t)
  - Pessimistic approach:
    - For each leaf node: e' (t) = (e(t)+0.5)
    - Total errors: e' (T) = e(T) + N × 0.5 (N: number of leaf nodes)
    - For a tree with 30 leaf nodes and 10 errors on training (out of 1000 instances):
      Training error = 10/1000 = 1%
      Generalization error = (10 + 30×0.5)/1000 = 2.5%
  - Reduced error pruning (REP):
    - uses validation data set to estimate generalization error

## Occam's Razor

- Given two models of similar generalization errors, one should prefer the simpler model over the more complex model

- For complex models, there is a greater chance that it was fitted accidentally by errors in data

- Therefore, one should include model complexity when evaluating a model

## Minimum Description Length (MDL)



- Cost(Model,Data) = Cost(Data|Model) + Cost(Model)
  - Cost is the number of bits needed for encoding.
  - Search for the least costly model.
- Cost(Data|Model) encodes the misclassification errors.
- Cost(Model) uses node encoding (number of children) plus splitting condition encoding.

## 4.2. Decision Tree based Methods

4.2.4. Model Overfitting

Decision Tree Pruning

---

## How to Address Overfitting

- Pre-Pruning (Early Stopping Rule)
  - Stop the algorithm before it becomes a fully-grown tree
  - Typical stopping conditions for a node:
    - Stop if all instances belong to the same class
    - Stop if all the attribute values are the same
  - More restrictive conditions:
    - Stop if number of instances is less than some user-specified threshold
    - Stop if class distribution of instances are independent of the available features (e.g., using $\chi^2$ test)
    - Stop if expanding the current node does not improve impurity measures (e.g., Gini or information gain).

---

## How to Address Overfitting…

- Post-pruning
  - Grow decision tree to its entirety
  - Trim the nodes of the decision tree in a bottom-up fashion
  - If generalization error improves after trimming, replace sub-tree by a leaf node.
  - Class label of leaf node is determined from majority class of instances in the sub-tree
  - Can use MDL for post-pruning

---

## Example of Post-Pruning

| Class = Yes | 20 |
| --- | --- |
| Class = No | 10 |
| Error = 10/30 | |

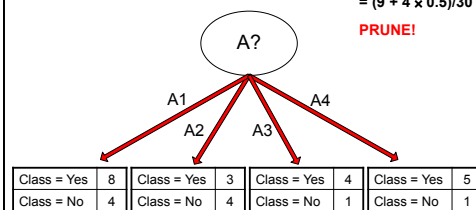**Training Error (Before splitting) = 10/30**
**Pessimistic error = (10 + 0.5)/30 = 10.5/30**
**Training Error (After splitting) = 9/30**
**Pessimistic error (After splitting)**
**= (9 + 4 × 0.5)/30 = 11/30**
**PRUNE!**

A?

A1    A4
A2    A3

| Class = Yes | 8 | Class = Yes | 3 | Class = Yes | 4 | Class = Yes | 5 |
| --- | --- | --- | --- | --- | --- | --- | --- |
| Class = No | 4 | Class = No | 4 | Class = No | 1 | Class = No | 1 |

---

## 4.2. Decision Tree based Methods

4.2.5. Other Issues

---

## Other Issues

- Data Fragmentation
- Search Strategy
- Expressiveness
- Tree Replication

## Data Fragmentation

- Number of instances gets smaller as you traverse down the tree

- Number of instances at the leaf nodes could be too small to make any statistically significant decision
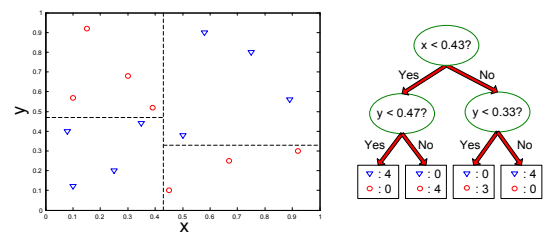
## Search Strategy

- Finding an optimal decision tree is NP-hard

- The algorithm presented so far uses a greedy, top-down, recursive partitioning strategy to induce a reasonable solution

- Other strategies?
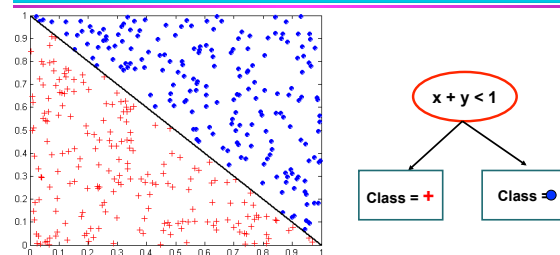  - Bottom-up
  - Bi-directional

## Expressiveness

- Decision tree provides expressive representation for learning discrete-valued function
  - But they do not generalize well to certain types of Boolean functions
    - Example: parity function:
      - Class = 1 if there is an even number of Boolean attributes with truth value = True
      - Class = 0 if there is an odd number of Boolean attributes with truth value = True
    - For accurate modeling, must have a complete tree

- Not expressive enough for modeling continuous variables
  - Particularly when test condition involves only a single attribute at-a-time

## 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**

## Oblique Decision Trees



- **Test condition may involve multiple attributes**

- **More expressive representation**

- **Finding optimal test condition is computationally expensive**