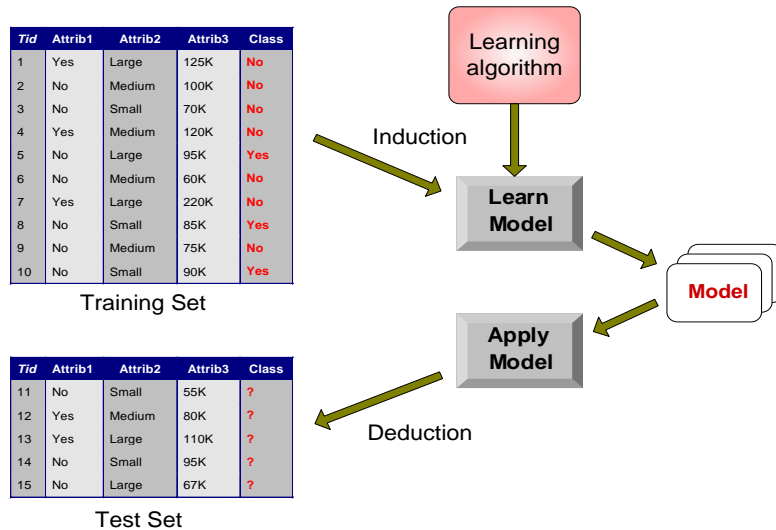

Kuliah 6

Classification: Basic Concepts and Decision Trees

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.
-

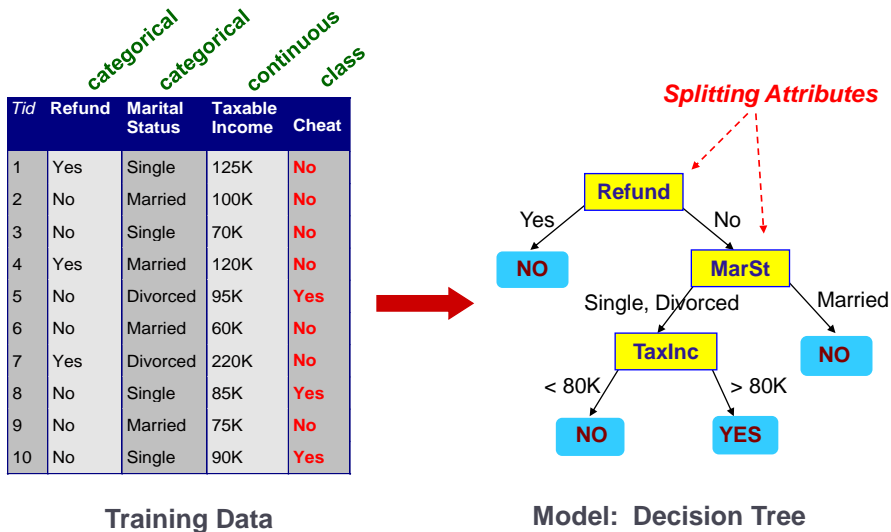
Illustrating Classification Task



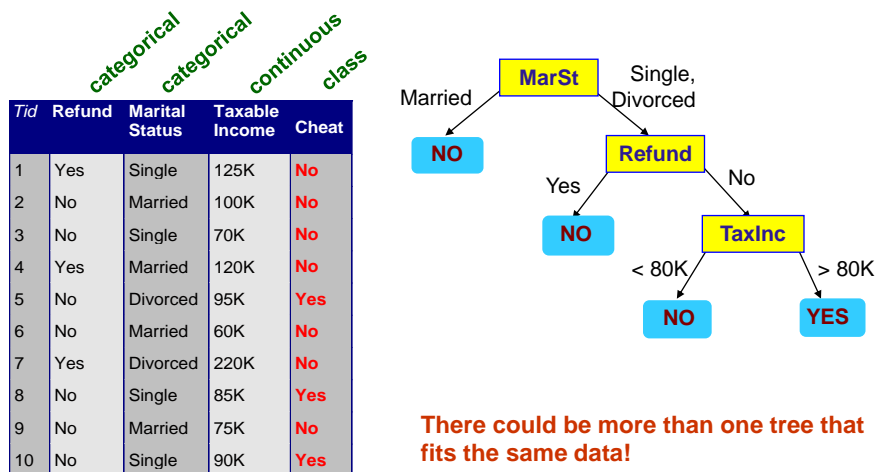
Classification Techniques

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

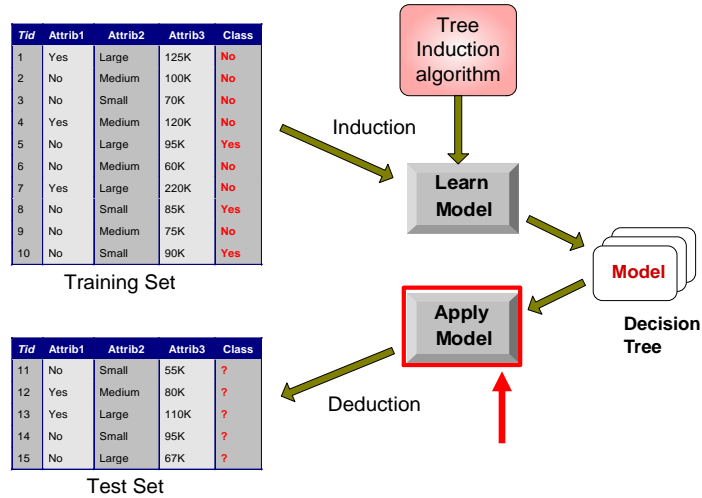
Example of a Decision Tree



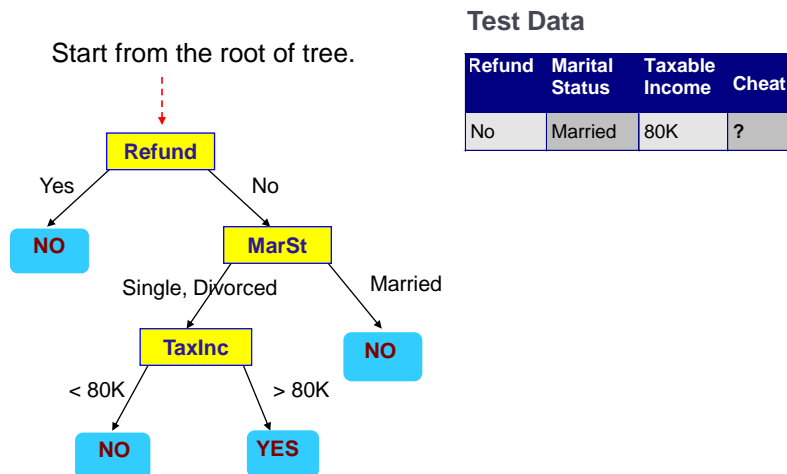
Another Example of Decision Tree



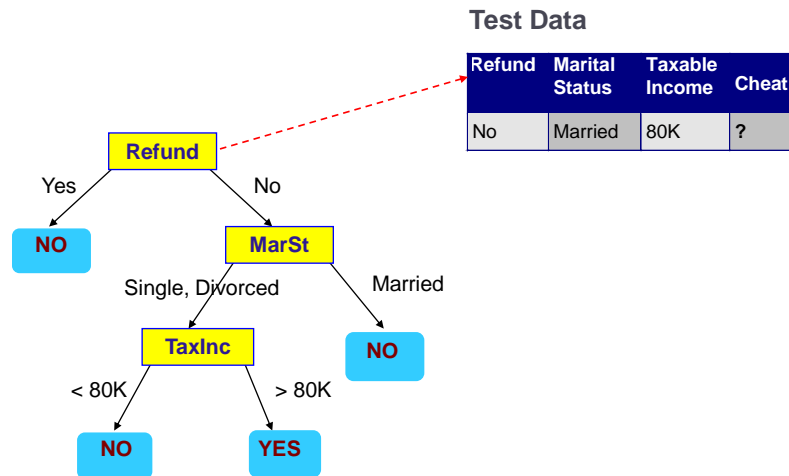
Decision Tree Classification Task



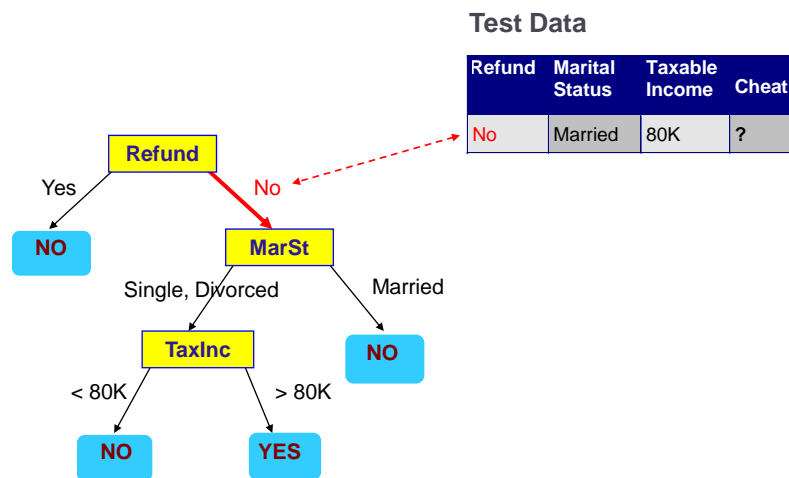
Apply Model to Test Data



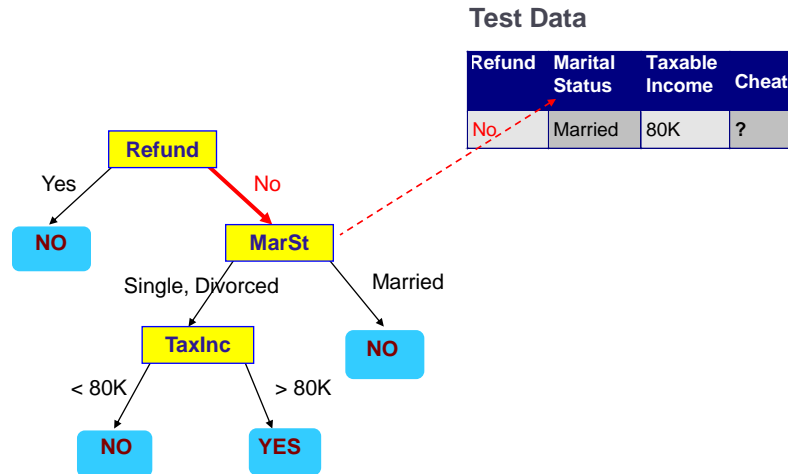
Apply Model to Test Data



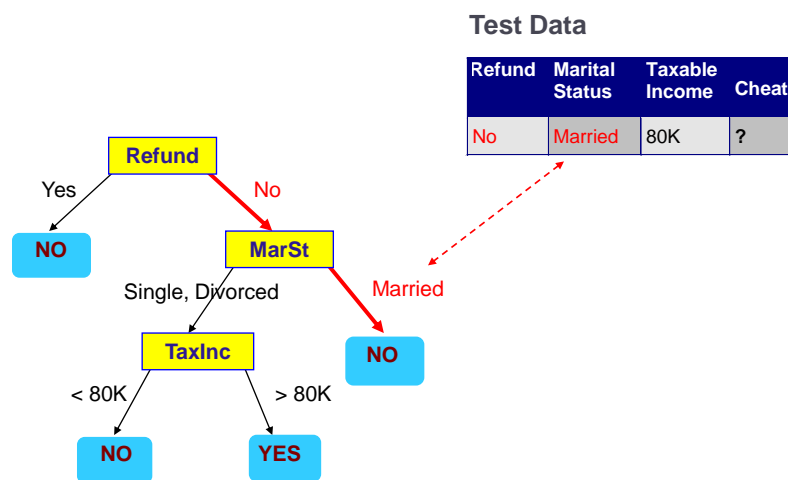
Apply Model to Test Data



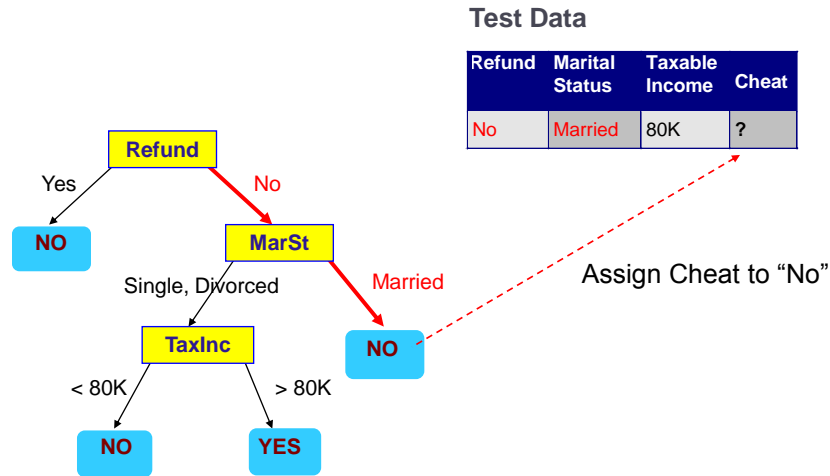
Apply Model to Test Data



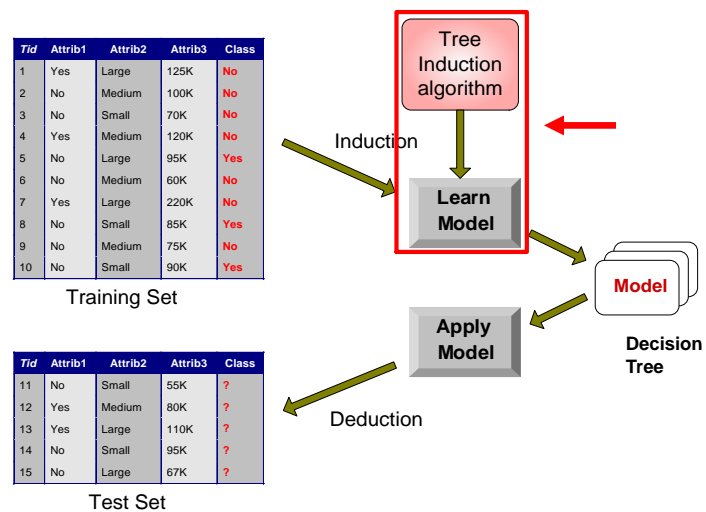
Apply Model to Test Data



Apply Model to Test Data



Decision Tree Classification Task



Decision Tree Induction

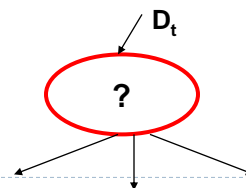
- ▶ **Many Algorithms:**
 - ▶ Hunt's Algorithm (one of the earliest)
 - ▶ CART
 - ▶ ID3, C4.5
 - ▶ SLIQ, SPRINT



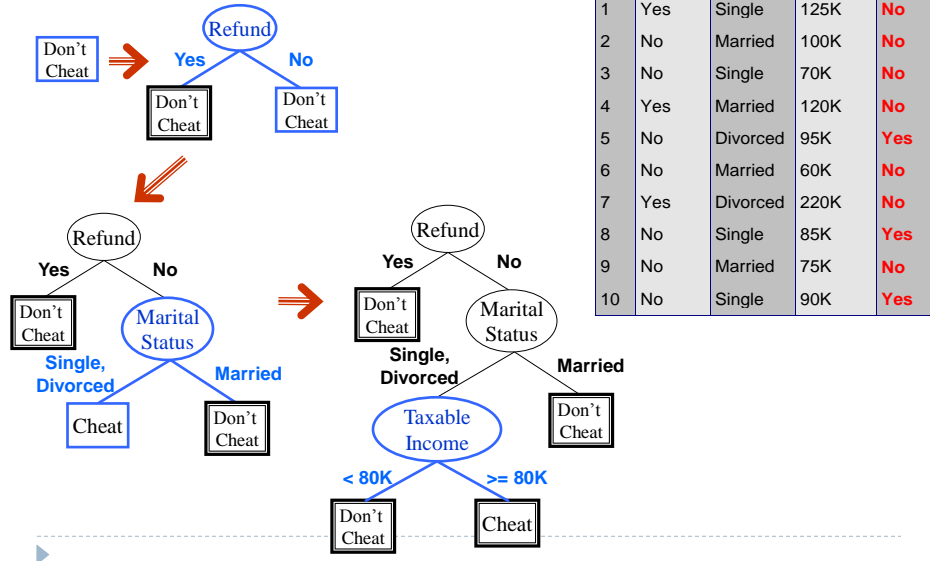
General Structure of Hunt's Algorithm

- ▶ Let D_t be the set of training records that reach a node t
- ▶ **General Procedure:**
 - ▶ If D_t contains records that belong to the same class y_t , then t is a leaf node labeled as y_t
 - ▶ If D_t is an empty set, then t is a leaf node labeled by the default class, y_d
 - ▶ If D_t contains records that belong to more than one class, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.

<i>Tid</i>	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



Hunt's Algorithm



Tree Induction

- ▶ **Greedy strategy.**
 - ▶ Split the records based on an attribute test that optimizes certain criterion.
- ▶ **Issues**
 - ▶ Determine how to split the records
 - ▶ How to specify the attribute test condition?
 - ▶ How to determine the best split?
 - ▶ Determine when to stop splitting

Tree Induction

- ▶ **Greedy strategy.**
 - ▶ Split the records based on an attribute test that optimizes certain criterion.
- ▶ **Issues**
 - ▶ Determine how to split the records
 - ▶ **How to specify the attribute test condition?**
 - ▶ How to determine the best split?
 - ▶ Determine when to stop splitting



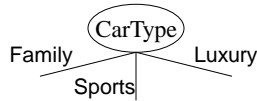
How to Specify Test Condition?

- ▶ **Depends on attribute types**
 - ▶ Nominal
 - ▶ Ordinal
 - ▶ Continuous
- ▶ **Depends on number of ways to split**
 - ▶ 2-way split
 - ▶ Multi-way split

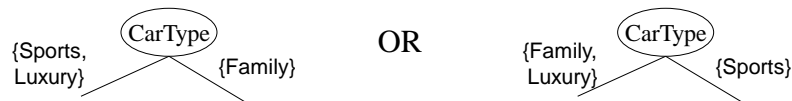


Splitting Based on Nominal Attributes

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

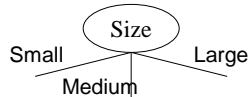


- ▶ **Binary split:** Divides values into two subsets.
Need to find optimal partitioning.

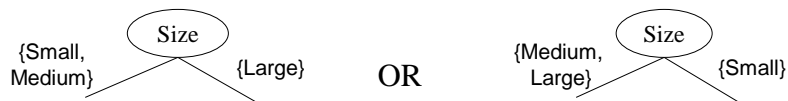


Splitting Based on Ordinal Attributes

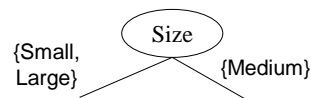
- ▶ **Multi-way split:** Use as many partitions as distinct values.



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



- ▶ What about this split?

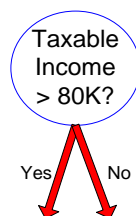


Splitting Based on Continuous Attributes

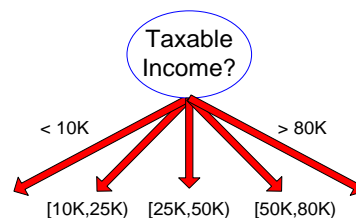
- ▶ Different ways of handling
 - ▶ **Discretization** to form an ordinal categorical attribute
 - ▶ Static – discretize once at the beginning
 - ▶ Dynamic – ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.
 - ▶ **Binary Decision**: ($A < v$) or ($A \geq v$)
 - ▶ consider all possible splits and finds the best cut
 - ▶ can be more compute intensive



Splitting Based on Continuous Attributes



(i) Binary split



(ii) Multi-way split



Tree Induction

- ▶ Greedy strategy.

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

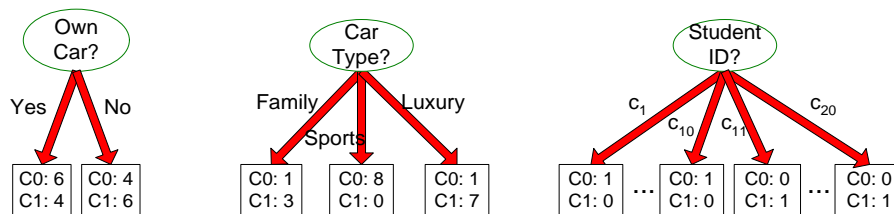
- ▶ Issues

- ▶ 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**



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**

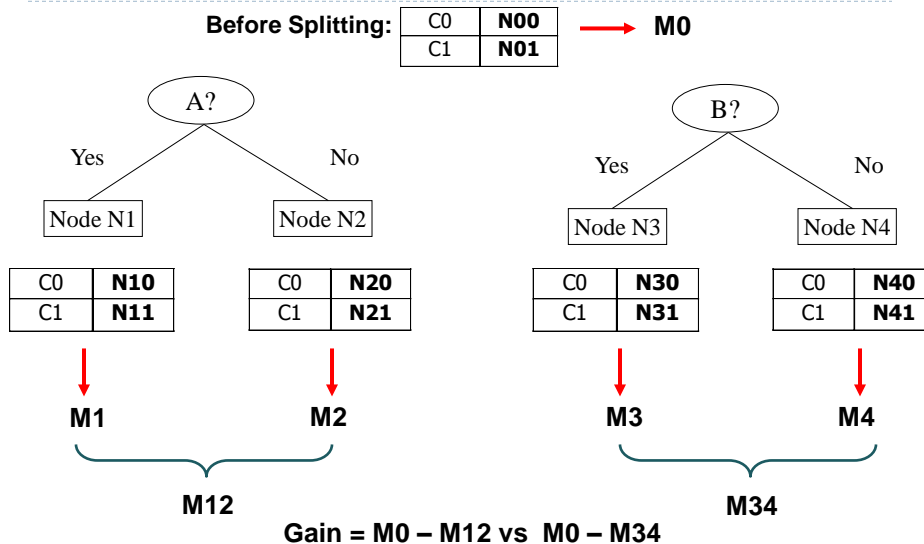


Measures of Node Impurity

- ▶ Gini Index
- ▶ Entropy
- ▶ Misclassification error



How to Find the Best Split



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, 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

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

C1	0
C2	6

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Gini = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$$

C1	1
C2	5

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

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

C1	2
C2	4

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

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



Splitting Based on GINI

- ▶ Used in CART, SLIQ, SPRINT.
- ▶ When a node p is split into k partitions (children), the quality of split is computed as,

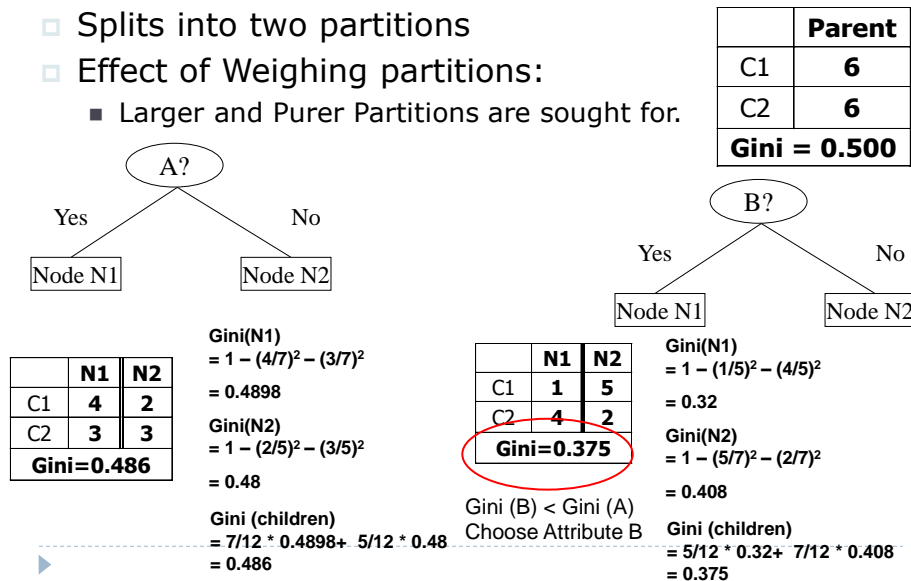
$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

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



Binary Attributes: Computing GINI Index

- Splits into two partitions
- Effect of Weighing partitions:
 - Larger and Purer Partitions are sought for.



Categorical Attributes: Computing Gini Index

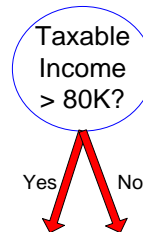
- ▶ For each distinct value, gather counts for each class in the dataset
- ▶ Use the count matrix to make decisions

Multi-way split				Two-way split (find best partition of values)					
				CarType		CarType			
				{Sports, Luxury}		{Sports}			
				{Family}		{Family, Luxury}			
C1	1	2	1	C1	3	2	2		
C2	4	1	1	C2	2	1	5		
Gini	0.393			Gini	0.400		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 $A \geq v$
- ▶ 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.

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



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

		Taxable Income																				
Cheat		No		No		No		Yes		Yes		No		No		No		No				
Sorted Values	→	Taxable Income																				
		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	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	4	3	5	2	6	1	7	0	
Gini		0.420		0.400		0.375		0.343		0.417		0.400		0.300		0.343		0.375		0.400		0.420

Alternative Splitting Criteria based on INFO

- ▶ Entropy at a given node t :

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

(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

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Entropy = -0 \log 0 - 1 \log 1 = -0 - 0 = 0$$

C1	1
C2	5

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$Entropy = -(1/6) \log_2 (1/6) - (5/6) \log_2 (5/6) = 0.65$$

C1	2
C2	4

$$P(C1) = 2/6 \quad P(C2) = 4/6$$

$$Entropy = -(2/6) \log_2 (2/6) - (4/6) \log_2 (4/6) = 0.92$$



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)
- 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 INFO...

► 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



Splitting Criteria based on 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

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

C1	0
C2	6

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$
$$Error = 1 - \max(0, 1) = 1 - 1 = 0$$

C1	1
C2	5

$$P(C1) = 1/6 \quad P(C2) = 5/6$$
$$Error = 1 - \max(1/6, 5/6) = 1 - 5/6 = 1/6$$

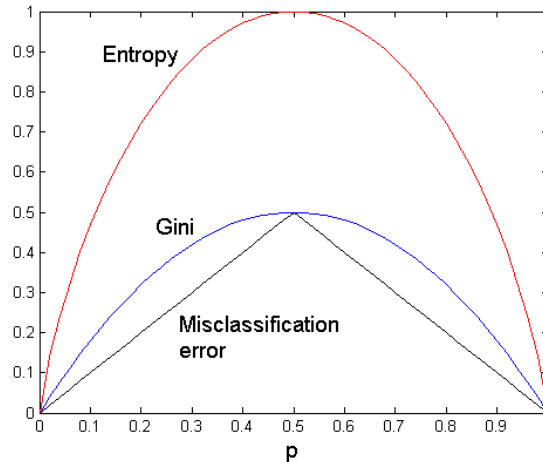
C1	2
C2	4

$$P(C1) = 2/6 \quad P(C2) = 4/6$$
$$Error = 1 - \max(2/6, 4/6) = 1 - 4/6 = 1/3$$



Comparison among Splitting Criteria

For a 2-class problem:



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 similar 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



Example: C4.5

- ▶ Simple depth-first construction.
- ▶ Uses Information Gain
- ▶ Sorts Continuous Attributes at each node.
- ▶ Needs entire data to fit in memory.
- ▶ Unsuitable for Large Datasets.
 - ▶ Needs out-of-core sorting to handle massive amounts of data.

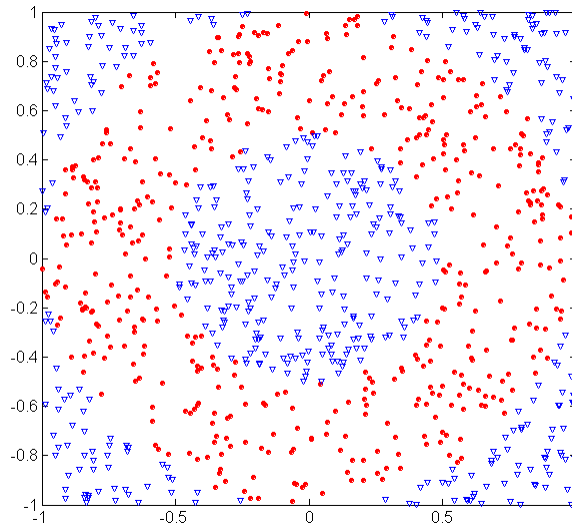


Practical Issues of Classification

- ▶ Underfitting and Overfitting
- ▶ Missing Values
- ▶ Costs of Classification



Underfitting and Overfitting (Example)



500 circular and 500 triangular data points.

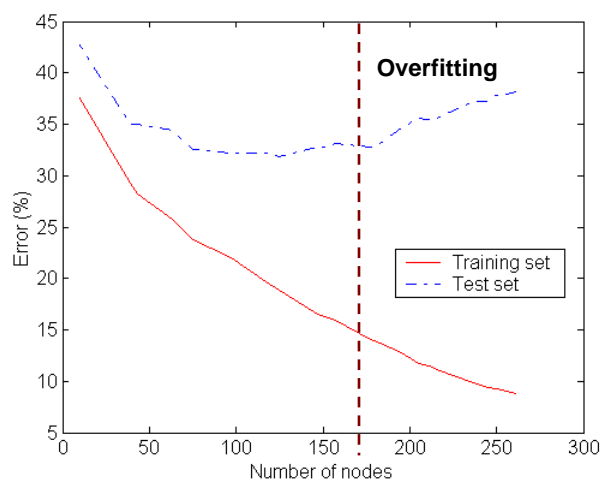
Circular points:

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

Triangular points:

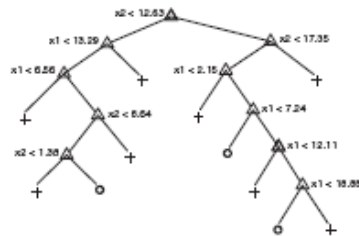
$$\sqrt{x_1^2 + x_2^2} > 0.5 \text{ or } \sqrt{x_1^2 + x_2^2} < 1$$

Underfitting and Overfitting

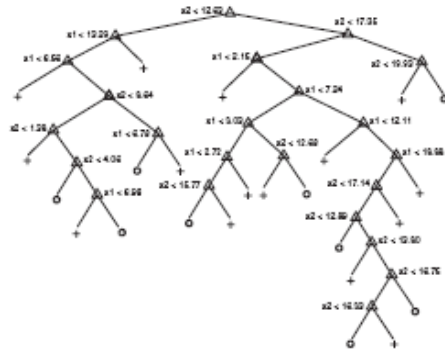


Underfitting: when model is too simple, both training and test errors are large

Example

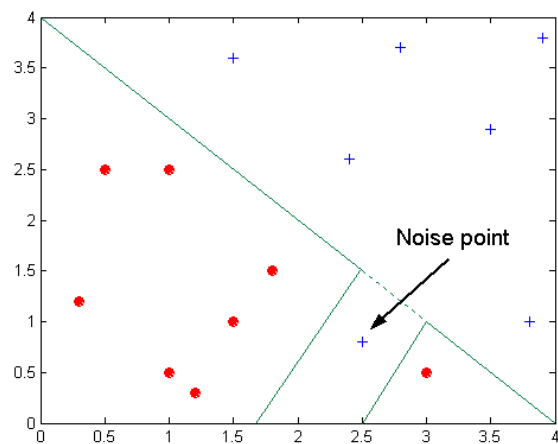


(a) Decision tree with 11 leaf nodes.



(b) Decision tree with 24 leaf nodes.

Overfitting due to Noise



Decision boundary is distorted by noise point

Example

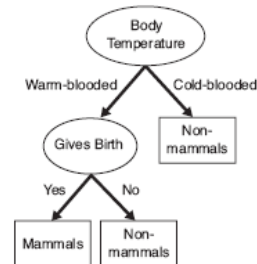
Training set

Name	Body Temperature	Gives Birth	Four-legged	Hibernates	Class Label
porcupine	warm-blooded	yes	yes	yes	yes
cat	warm-blooded	yes	yes	no	yes
bat	warm-blooded	yes	no	yes	no*
whale	warm-blooded	yes	no	no	no*
salamander	cold-blooded	no	yes	yes	no
komodo dragon	cold-blooded	no	yes	no	no
python	cold-blooded	no	no	yes	no
salmon	cold-blooded	no	no	no	no
eagle	warm-blooded	no	no	no	no
guppy	cold-blooded	yes	no	no	no

* mislabeled

Test set

Name	Body Temperature	Gives Birth	Four-legged	Hibernates	Class Label
human	warm-blooded	yes	no	no	yes
pigeon	warm-blooded	no	no	no	no
elephant	warm-blooded	yes	yes	no	yes
leopard shark	cold-blooded	yes	no	no	no
turtle	cold-blooded	no	yes	no	no
penguin	cold-blooded	no	no	no	no
eel	cold-blooded	no	no	no	no
dolphin	warm-blooded	yes	no	no	yes
spiny anteater	warm-blooded	no	yes	yes	yes
gila monster	cold-blooded	no	yes	yes	no



(b) Model M2

Example

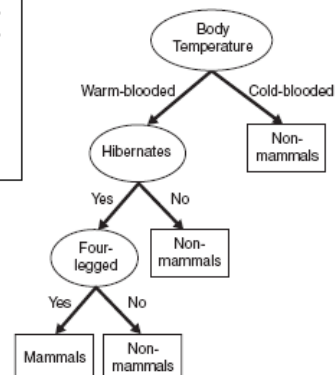
Training set

Name	Body Temperature	Gives Birth	Four-legged	Hibernates	Class Label
porcupine	warm-blooded	yes	yes	yes	yes
cat	warm-blooded	yes	yes	no	yes
bat	warm-blooded	yes	no	yes	no*
whale	warm-blooded	yes	no	no	no*
salamander	cold-blooded	no	yes	yes	no
komodo dragon	cold-blooded	no	yes	no	no
python	cold-blooded	no	no	yes	no
salmon	cold-blooded	no	no	no	no
eagle	warm-blooded	no	no	no	no
guppy	cold-blooded	yes	no	no	no

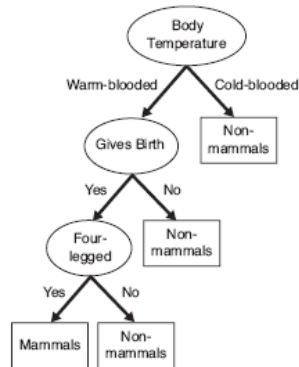
Test set

Name	Body Temperature	Gives Birth	Four-legged	Hibernates	Class Label
human	warm-blooded	yes	no	no	yes
pigeon	warm-blooded	no	no	no	no
elephant	warm-blooded	yes	yes	no	yes
leopard shark	cold-blooded	yes	no	no	no
turtle	cold-blooded	no	yes	no	no
penguin	cold-blooded	no	no	no	no
eel	cold-blooded	no	no	no	no
dolphin	warm-blooded	yes	no	no	yes
spiny anteater	warm-blooded	no	yes	yes	yes
gila monster	cold-blooded	no	yes	yes	no

Model M1

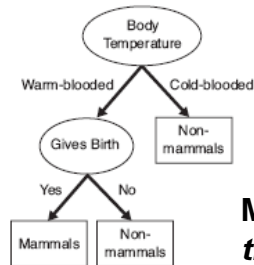


Tree



(a) Model M1

M1 overfitting training set

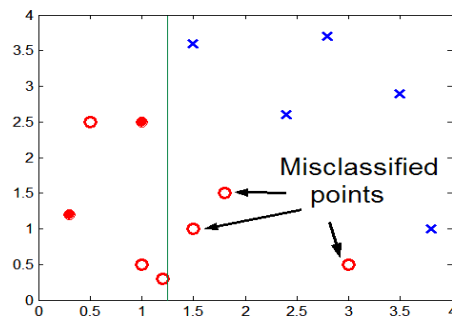


(b) Model M2

M1:
*training error : 0%,
test error : 30%*

M2:
*training error : 20%,
test error : 10%*

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



Estimating Generalization Errors

- ▶ **Re-substitution errors:** error on training ($\sum e(t)$)
- ▶ **Generalization errors:** error on testing ($\sum e'(t)$)
- ▶ 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: $e'(T) = e(T) + N \times 0.5$ (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 = $(10 + 30 \times 0.5)/1000 = 2.5\%$
 - ▶ **Reduced error pruning (REP):**
 - ▶ uses validation data set to estimate generalization error



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 χ^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

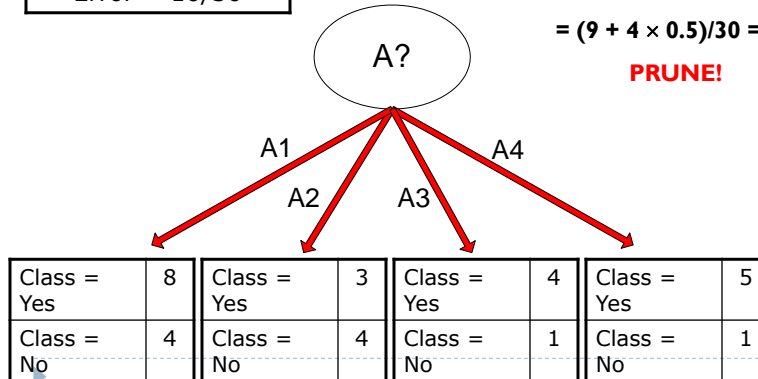


Example of Post-Pruning

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)
 $= (9 + 4 \times 0.5)/30 = 11/30$

PRUNE!



Scalable Decision Tree Induction Methods

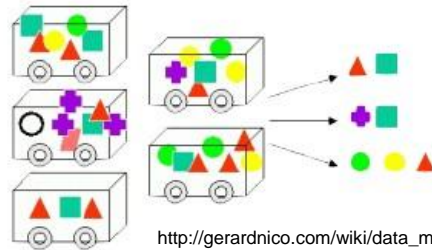
- ▶ **SLIQ** (EDBT'96 — Mehta et al.)
 - ▶ Builds an index for each attribute and only class list and the current attribute list reside in memory
- ▶ **SPRINT** (VLDB'96 — J. Shafer et al.)
 - ▶ Constructs an attribute list data structure
- ▶ **PUBLIC** (VLDB'98 — Rastogi & Shim)
 - ▶ Integrates tree splitting and tree pruning: stop growing the tree earlier
- ▶ **RainForest** (VLDB'98 — Gehrke, Ramakrishnan & Ganti)
 - ▶ Builds an AVC-list (attribute, value, class label)
- ▶ **BOAT** (PODS'99 — Gehrke, Ganti, Ramakrishnan & Loh)
 - ▶ Uses bootstrapping to create several small samples

Catatan

- ▶ Semua slide diambil dari

- ▶ Tan P, Michael S., & Vipin K. 2006. *Introduction to Data mining*. Pearson Education, Inc.
- ▶ Han J & Kamber M. 2006. *Data mining – Concept and Techniques*. Morgan-Kaufman, San Diego
- ▶ Dan sumber lainnya yang relevan

- ▶ Topik selanjutnya: *Association Rule Mining*



http://gerardnico.com/wiki/data_mining/association

