



INT3405 - Machine Learning

Lecture 5: Classification (P2)

Decision Tree

Duc-Trong Le & Viet-Cuong Ta

Hanoi, 09/2023

Outline

- Decision Tree Examples
- Decision Tree Algorithms
- Methods for Expressing Test Conditions
- Measures of Node Impurity
 - Gini index
 - Entropy
 - Gain Ratio
 - Classification Error

SVM: Optimization Formulation (2)

- SVM as a Quadratic Programming (QP) problem

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{subject to} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \\ & i = 1, \dots, l. \end{aligned}$$

- Convex problem, has unique minimum
- Quadratic objective function
- Linear equality and inequality constraints

Bài tập

Cho 6 điểm data dùng cho bài toán phân lớp $\{+1, -1\}$ như bảng. Vẽ các điểm, tìm các support vectors, độ rộng đường biên lớn nhất theo phương pháp SVM tuyến tính

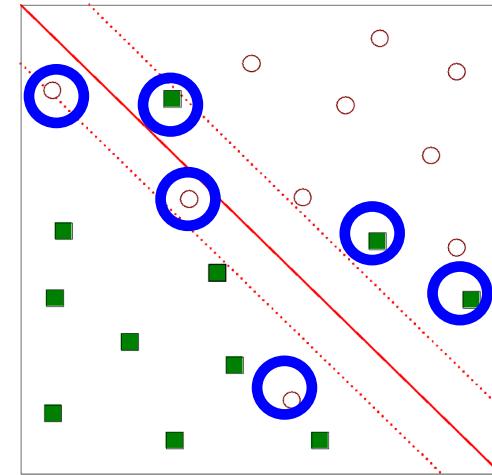
x1	x2	y
1	9	-1
5	5	-1
1	1	-1
8	5	+1
13	1	+1
13	9	+1

Soft Margin SVM

- Standard Linear SVM
 - Introduce slack variables
 - Relax the constraints
 - Penalize the relaxation

Primal Problem:
$$\min_{\mathbf{w}, b, \xi} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{i=1}^N \xi_i$$

subject to
$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, i = 1, \dots, N$$



C is a regularization parameter. Soft margin SVM trade off between maximizing the margin and minimizing the misclassification error rate

SVM: Nonlinear Case

- The dual problem

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C, i = 1, \dots, l \\ & \mathbf{y}^T \alpha = 0, \end{aligned}$$

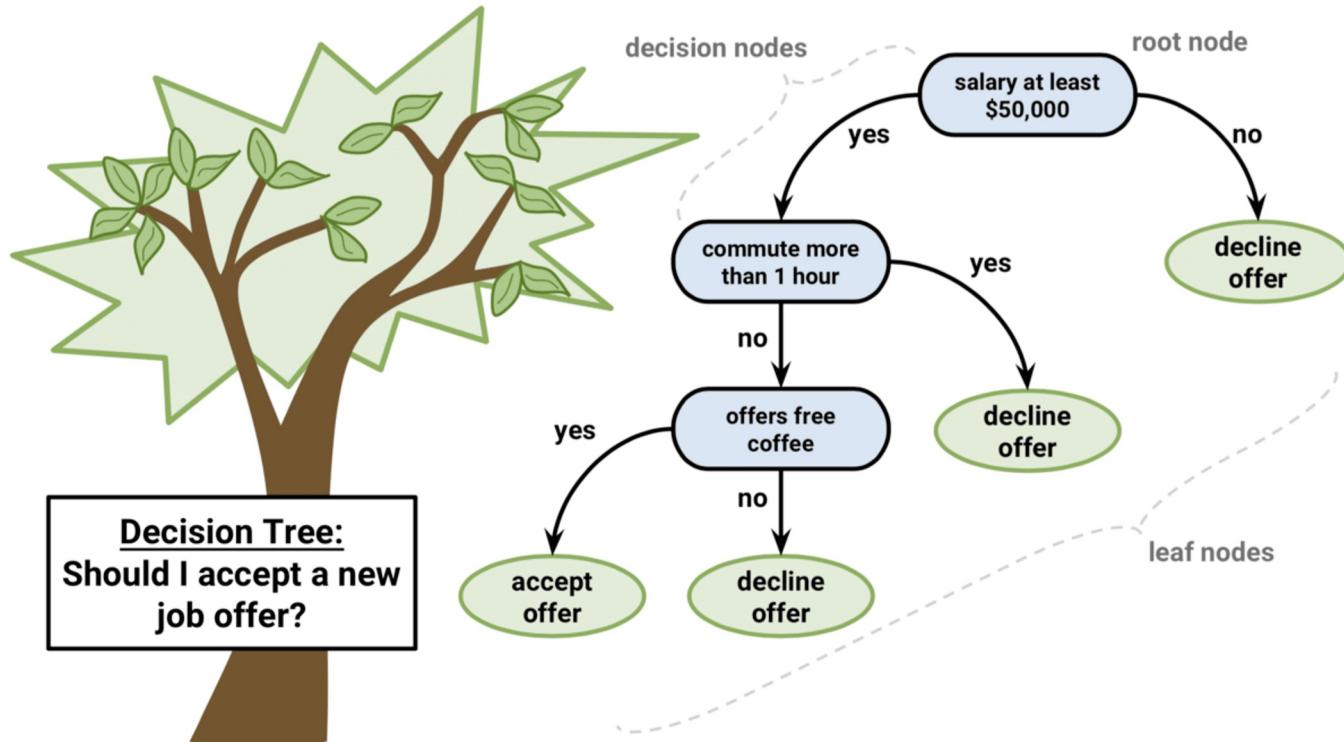
where $Q_{ij} = y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ and $\mathbf{e} = [1, \dots, 1]^T$

- The optimal solution

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \Phi(\mathbf{x}_i)$$

$$f(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) + b = \sum_{i=1}^N \alpha_i y_i \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}) + b$$

Example of a Decision Tree

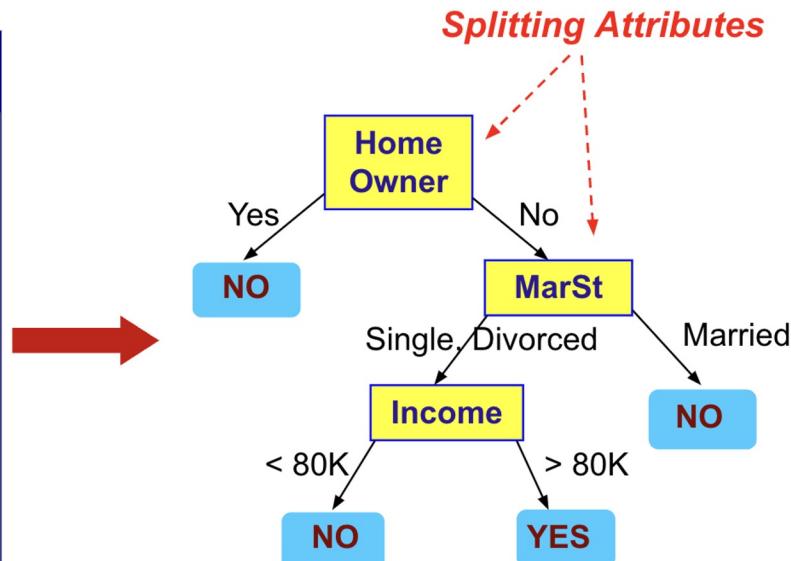


Source: <https://regenerativetoday.com/simple-explanation-on-how-decision-tree-algorithm-makes-decisions/>

Example of a Decision Tree

ID	categorical		continuous		class
	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	

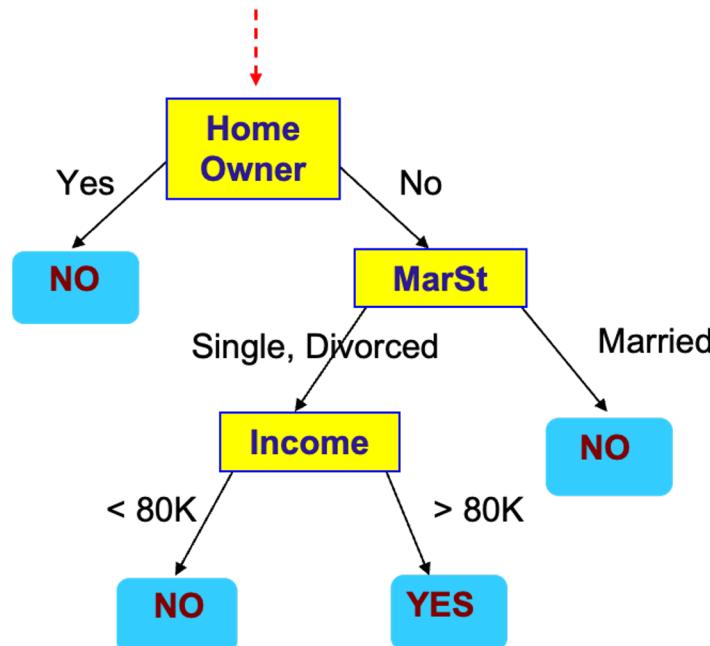
Training Data



Model: Decision Tree

Example of Model Prediction (1)

Start from the root tree.



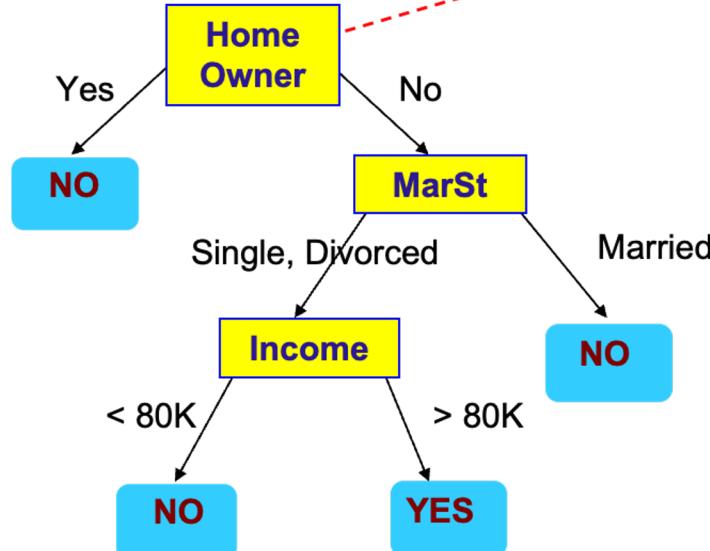
Test Data

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

Example of Model Prediction (2)

Test Data

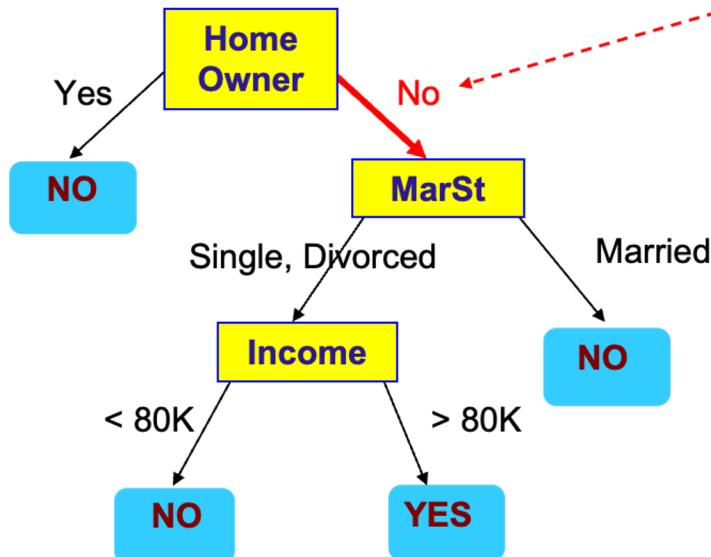
Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



Example of Model Prediction (3)

Test Data

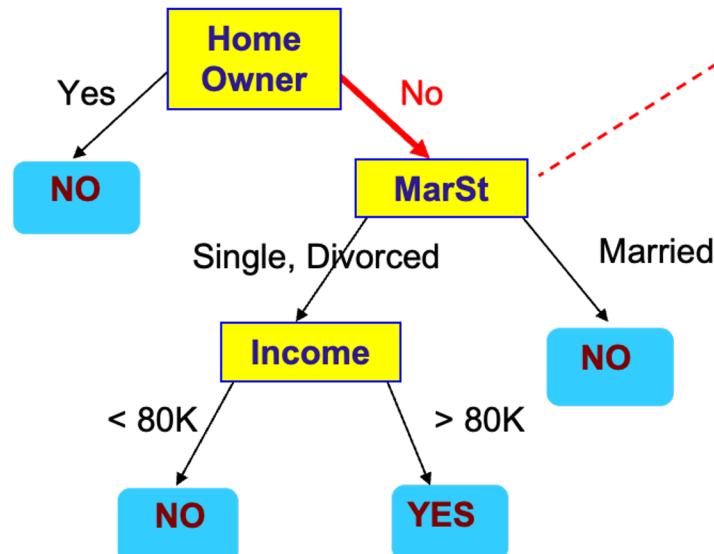
Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



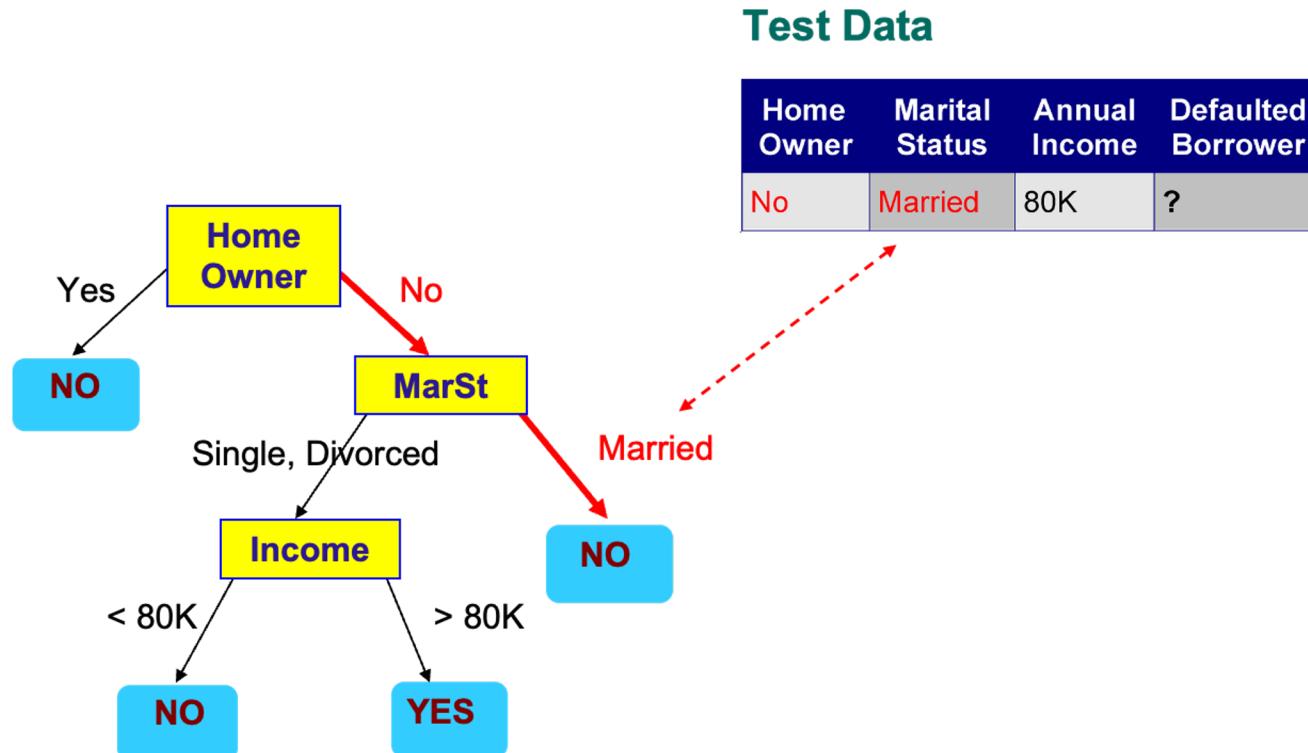
Example of Model Prediction (4)

Test Data

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



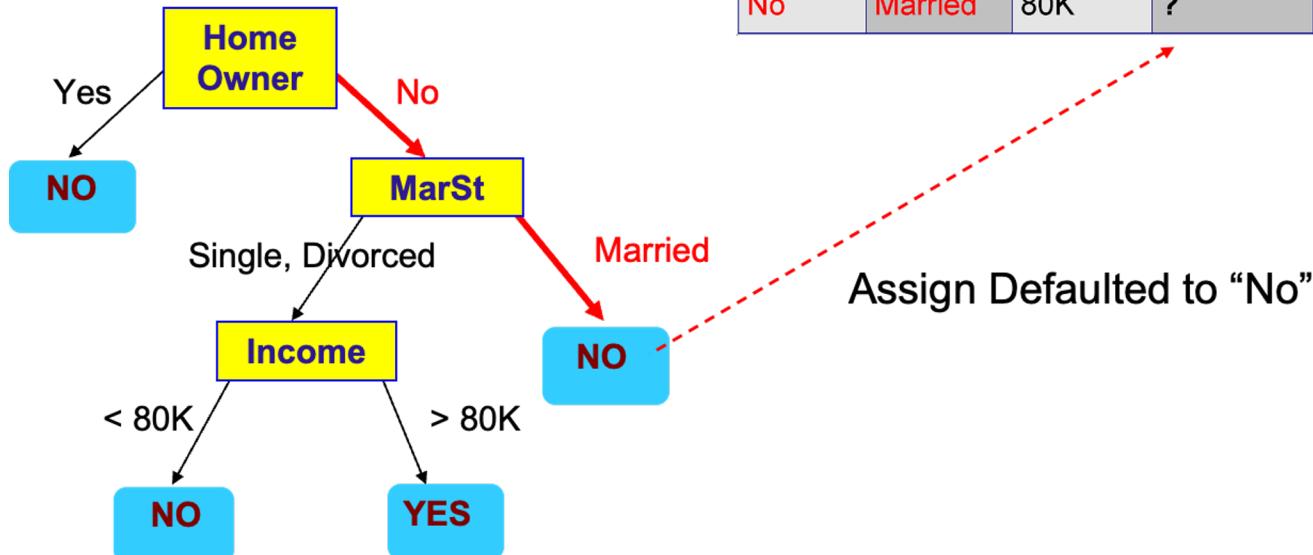
Example of Model Prediction (5)



Example of Model Prediction (6)

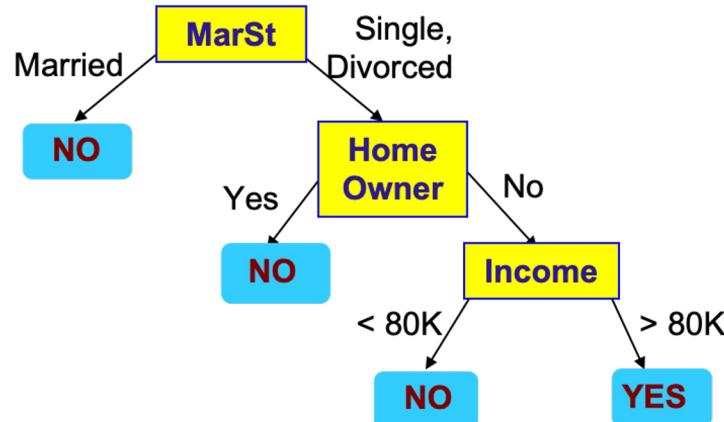
Test Data

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



Decision Tree - Another Solution

ID	categorical		Annual Income	Defaulted Borrower	class
	Home Owner	Marital Status			
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!

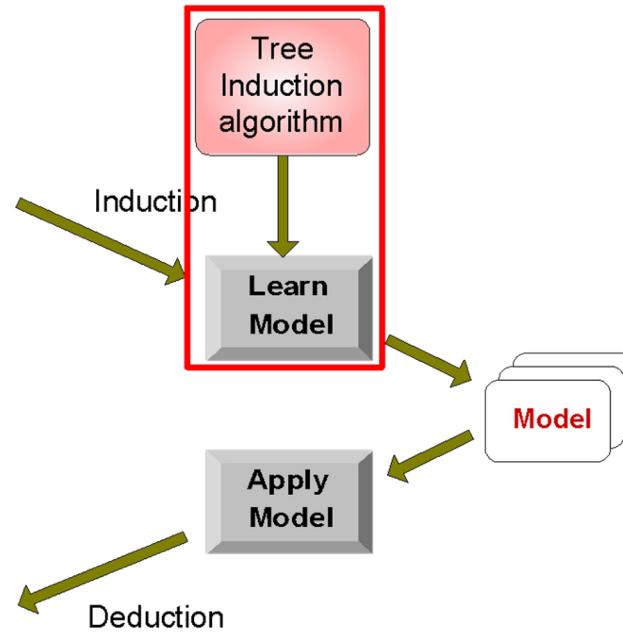
Decision Tree - Classification Task

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

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



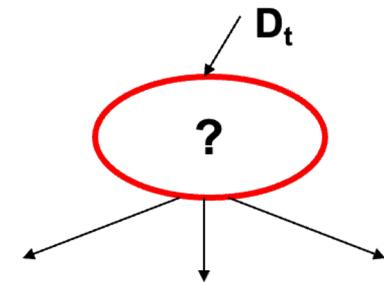
Decision Tree Induction

- Various 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 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



Hunt's Algorithm

Defaulted = No

(7,3)

(a)

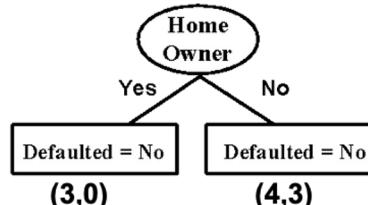
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Hunt's Algorithm

Defaulted = No

(7,3)

(a)



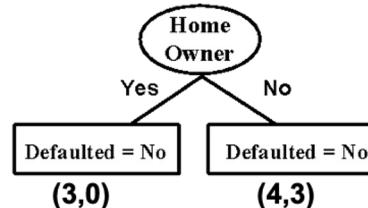
(b)

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

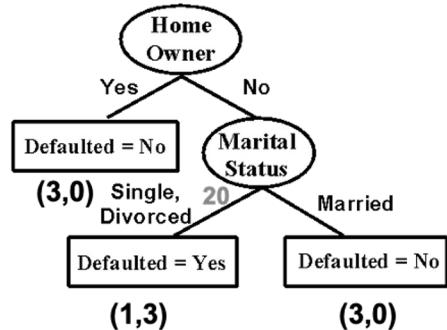
Hunt's Algorithm

Defaulted = No
(7,3)

(a)



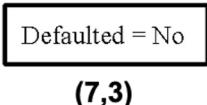
(b)



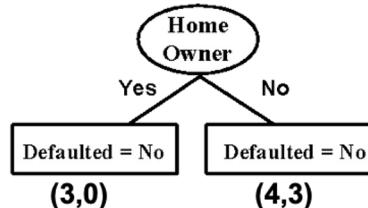
(c)

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

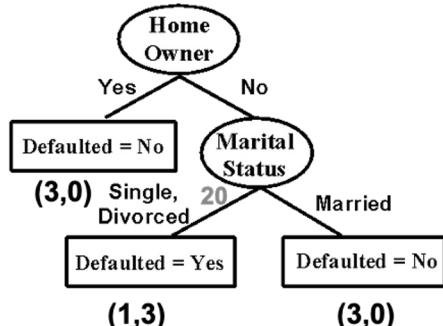
Hunt's Algorithm



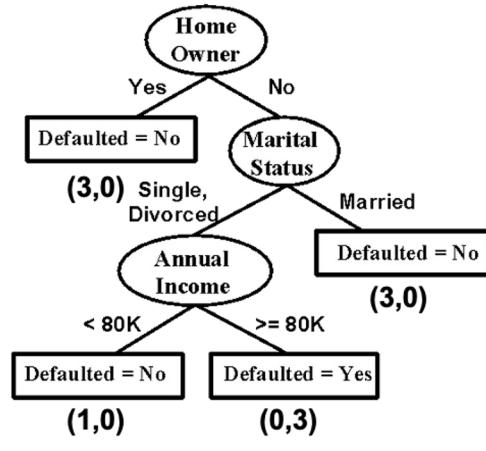
(a)



(b)



(c)



(d)

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

Design Issues of Decision Tree Induction

- **How should training records be split?**
 - Method for expressing 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

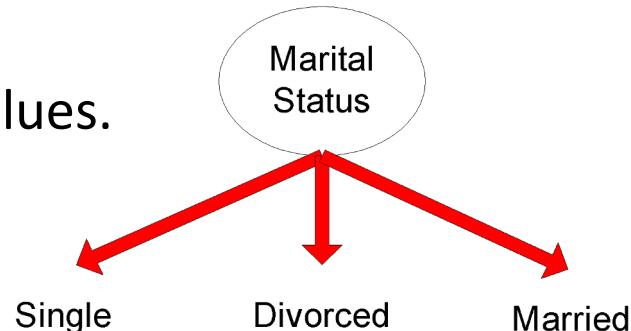
Methods for Expressing Test Conditions

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

Test Conditions for Nominal Attributes

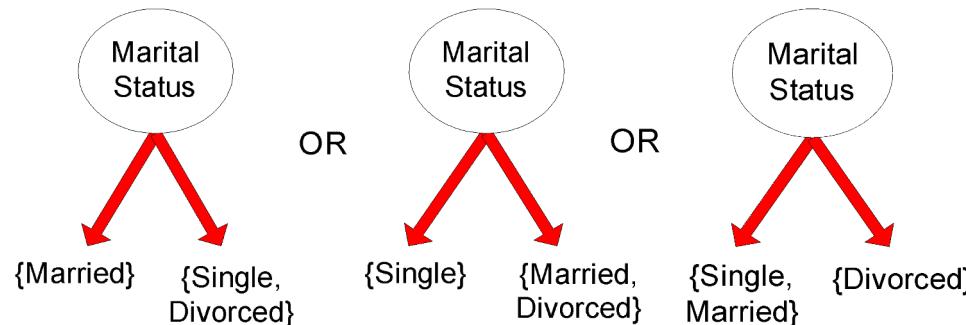
- **Multi-way split:**

- Use as many partitions as distinct values.



- **Binary split:**

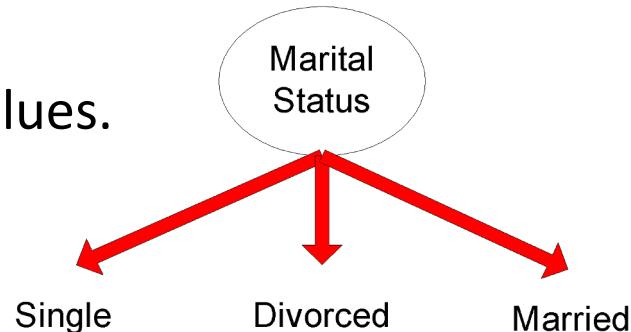
- Divides values into two subsets



Test Conditions for Nominal Attributes

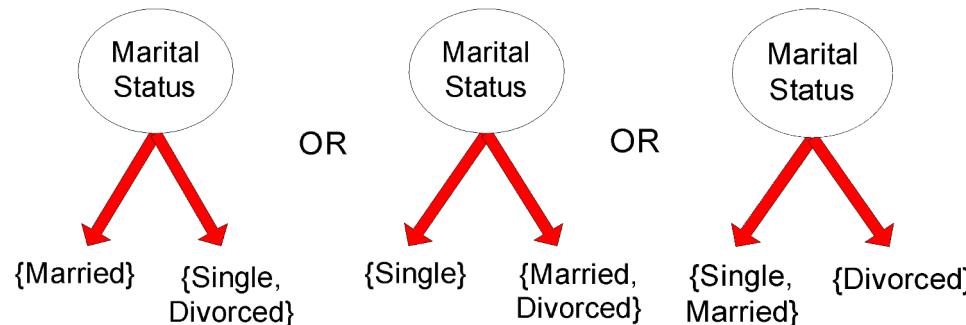
- **Multi-way split:**

- Use as many partitions as distinct values.



- **Binary split:**

- Divides values into two subsets



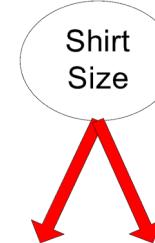
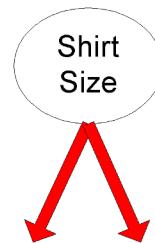
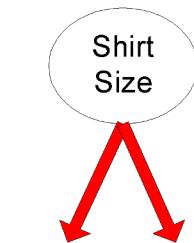
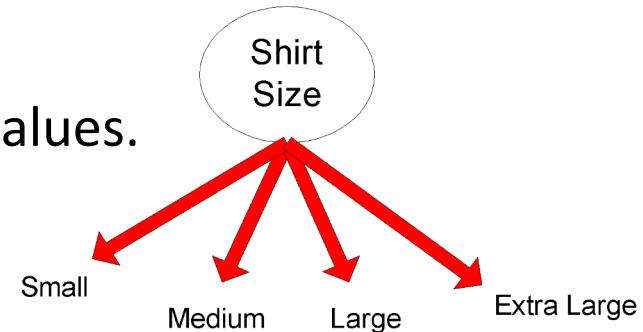
Test Conditions for Ordinal Attributes

- **Multi-way split:**

- Use as many partitions as distinct values.

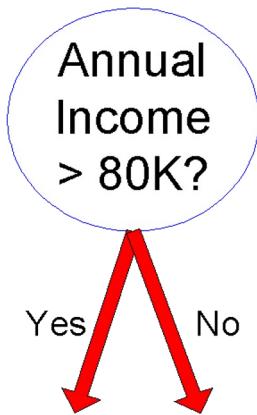
- **Binary split:**

- Divides values into two subsets
- Preserve order property among attribute values

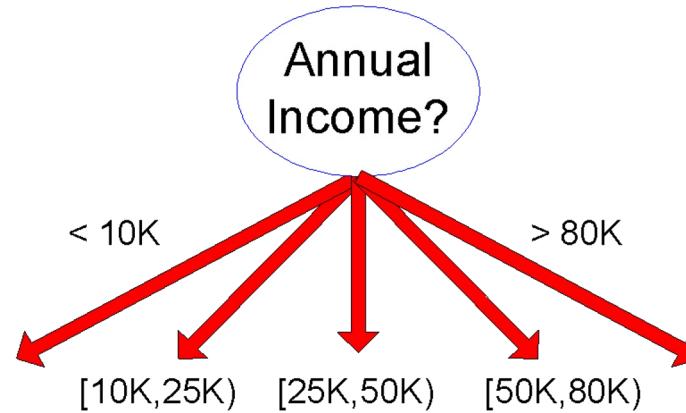


This grouping violates order property

Test Conditions for Continuous Attributes



(i) Binary split



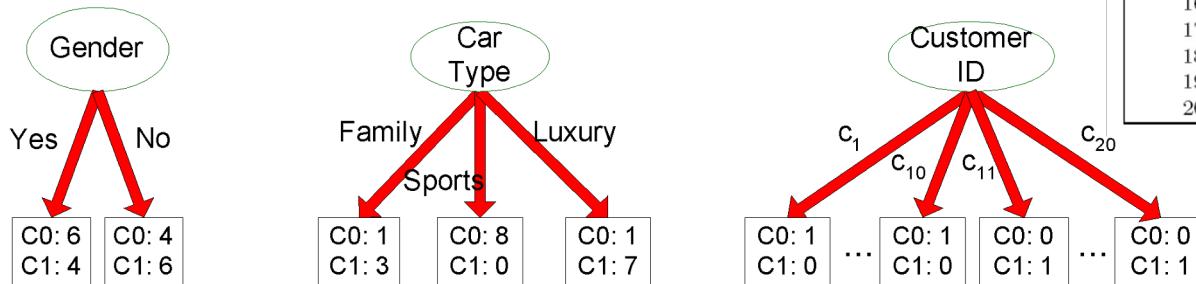
(ii) Multi-way split

Splitting based on Continuous Attributes

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

How to determine the Best Split

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



Which test condition is the best?

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

How to determine the Best Split

- Greedy approach:
 - Nodes with **purer** 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

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

Where $p_i(t)$ is the frequency of class i at node t , and c is the total number of classes

- Entropy

$$Entropy = - \sum_{i=0}^{c-1} p_i(t) \log_2 p_i(t)$$

- Misclassification error

$$Classification\ error = 1 - \max[p_i(t)]$$

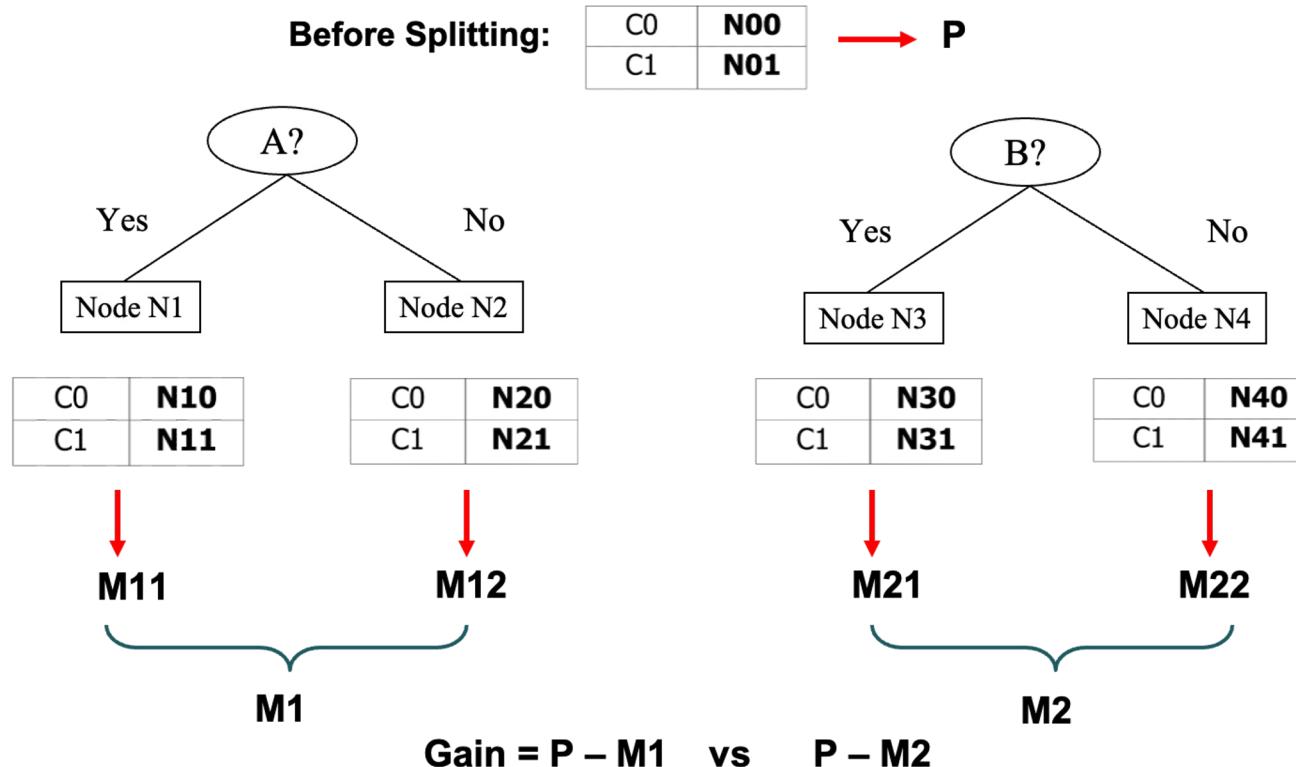
Find the Best Split

- Compute impurity measure (**P**) before splitting
- Compute impurity measure (**M**) after splitting
 - Compute impurity measure of each child node
 - M is the weighted impurity of child nodes
- Choose the attribute test condition that produces the highest gain

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

- or equivalently, lowest impurity measure after splitting (**M**)

Find the Best Split



Measure of Impurity: Gini (index)

- Gini Index for a given node t :

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

- For 2-class problem (p, 1 – p):

◆ $GINI = 1 - p^2 - (1 - p)^2 = 2p(1-p)$

C1	0
C2	6
Gini=0.000	

C1	1
C2	5
Gini=0.278	

C1	2
C2	4
Gini=0.444	

C1	3
C2	3
Gini=0.500	

Computing Gini Index of a Single Node

$$Gini\ Index = 1 - \sum_{i=0}^{c-1} p_i(t)^2$$

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

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

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

Computing Gini Index for a Collection of Nodes

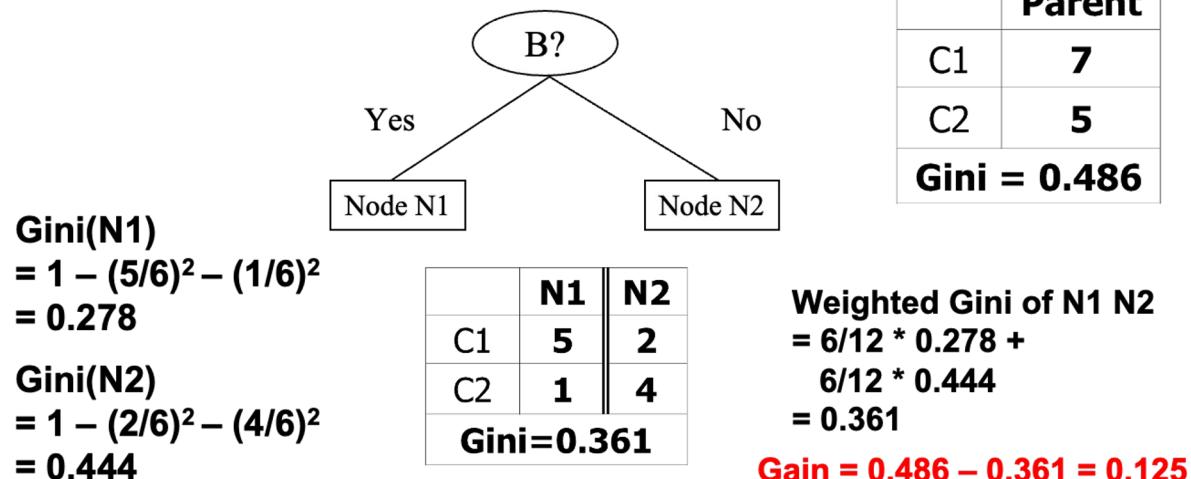
- When a node p is split into k partitions (children)

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

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

Binary Attributes: Computing GINI Index

- Splits into two partitions (child nodes)
- Effect of Weighing partitions:
 - Larger and purer partitions are sought



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	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		

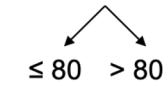
Which of these is the best?

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, $A \leq v$ 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.

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 ?



Defaulted Yes	0	3
Defaulted No	3	4

Measure of Impurity: Entropy

- Entropy for a given node t :

$$Entropy = - \sum_{i=0}^{c-1} p_i(t) \log_2 p_i(t)$$

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

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

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

C1	2
C2	4

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

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

Computing Information Gain After Splitting

- Information Gain:

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

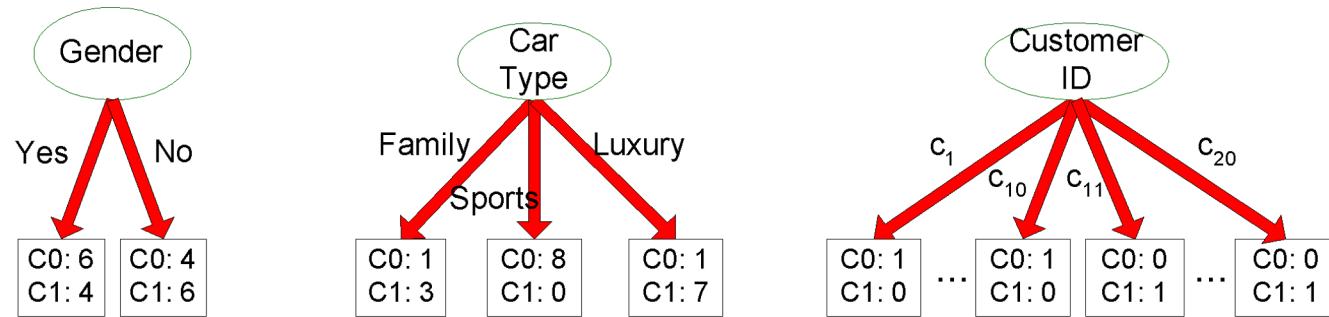
Parent Node, p is split into k partitions (children)

n_i is number of records in child node i

- Choose the split that achieves most reduction (maximizes GAIN)
- Used in the ID3 and C4.5 decision tree algorithms
- Information gain is the mutual information between the class variable and the splitting variable

Problem with large number of partitions

- Node impurity measures tend 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:

$$Gain\ Ratio = \frac{Gain_{split}}{Split\ Info}$$

$$Split\ Info = - \sum_{i=1}^k \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

Parent Node, p is split into k partitions (children)

n_i is number of records in child node i

- Adjusts Information Gain by the entropy of the partitioning (*Split Info*).
 - ◆ Higher entropy partitioning (large number of small partitions) is penalized!
- Used in C4.5 algorithm
- Designed to overcome the disadvantage of Information Gain

Gain Ratio

$$Gain\ Ratio = \frac{Gain_{split}}{Split\ Info}$$

$$Split\ Info = \sum_{i=1}^k \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

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

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		

SplitINFO = 1.52

SplitINFO = 0.72

SplitINFO = 0.97

Measure of Impurity: Classification Error

- Classification error at a node t

$$Error(t) = 1 - \max_i[p_i(t)]$$

- Maximum of $1 - 1/c$ when records are equally distributed among all classes, implying the least interesting situation
- Minimum of 0 when all records belong to one class, implying the most interesting situation

Computing Error of a Single Node

$$\text{Error}(t) = 1 - \max_i[p_i(t)]$$

C1	0
C2	6

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

C1	1
C2	5

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

C1	2
C2	4

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

Computing Error of a Single Node

$$\text{Error}(t) = 1 - \max_i[p_i(t)]$$

C1	0
C2	6

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

C1	1
C2	5

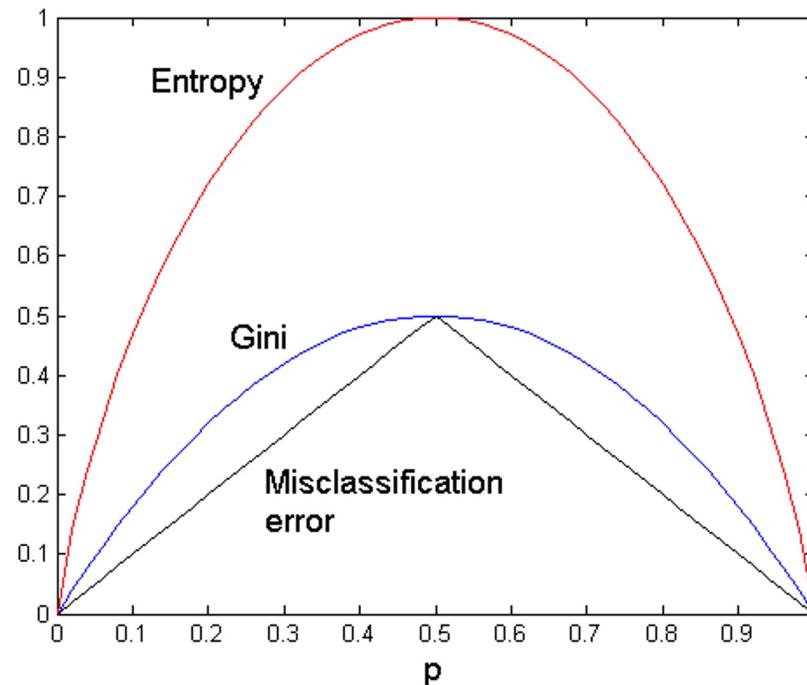
$$\mathbf{P(C1) = 1/6} \quad \mathbf{P(C2) = 5/6}$$
$$\mathbf{Error = 1 - max (1/6, 5/6) = 1 - 5/6 = 1/6}$$

C1	2
C2	4

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

Comparison among Impurity Measures

For a 2-class problem:



Decision Tree Classification

- **Advantages:**
 - Relatively inexpensive to construct
 - Extremely fast at classifying unknown records
 - Easy to interpret for small-sized trees
 - Robust to noise
 - Can easily handle redundant attributes, irrelevant attributes
- **Disadvantages:** .
 - Due to the greedy nature of splitting criterion, interacting attributes (that can distinguish between classes together but not individually) may be passed over in favor of other attributes that are less discriminating.
 - Each decision boundary involves only a single attribute

Summary

- Decision Tree Examples
- Decision Tree Algorithms
- Methods for Expressing Test Conditions
- Measures of Node Impurity
 - Gini index
 - Entropy
 - Gain Ratio
 - Classification Error