



Kourosh Davoudi
kourosh@uoit.ca

Classification: Basic Concepts
and Techniques

CSCI 4150U: Data Mining

Outline and Learning Outcomes

- Classification (Part I)
 - Understanding basic classification **concepts** and definitions
 - Explain detailed case study: **Decision Tree**
 - How to perform **model evaluation**?
 - What is **model overfitting**?
 - How to do **model selection**?

Classification: Definition

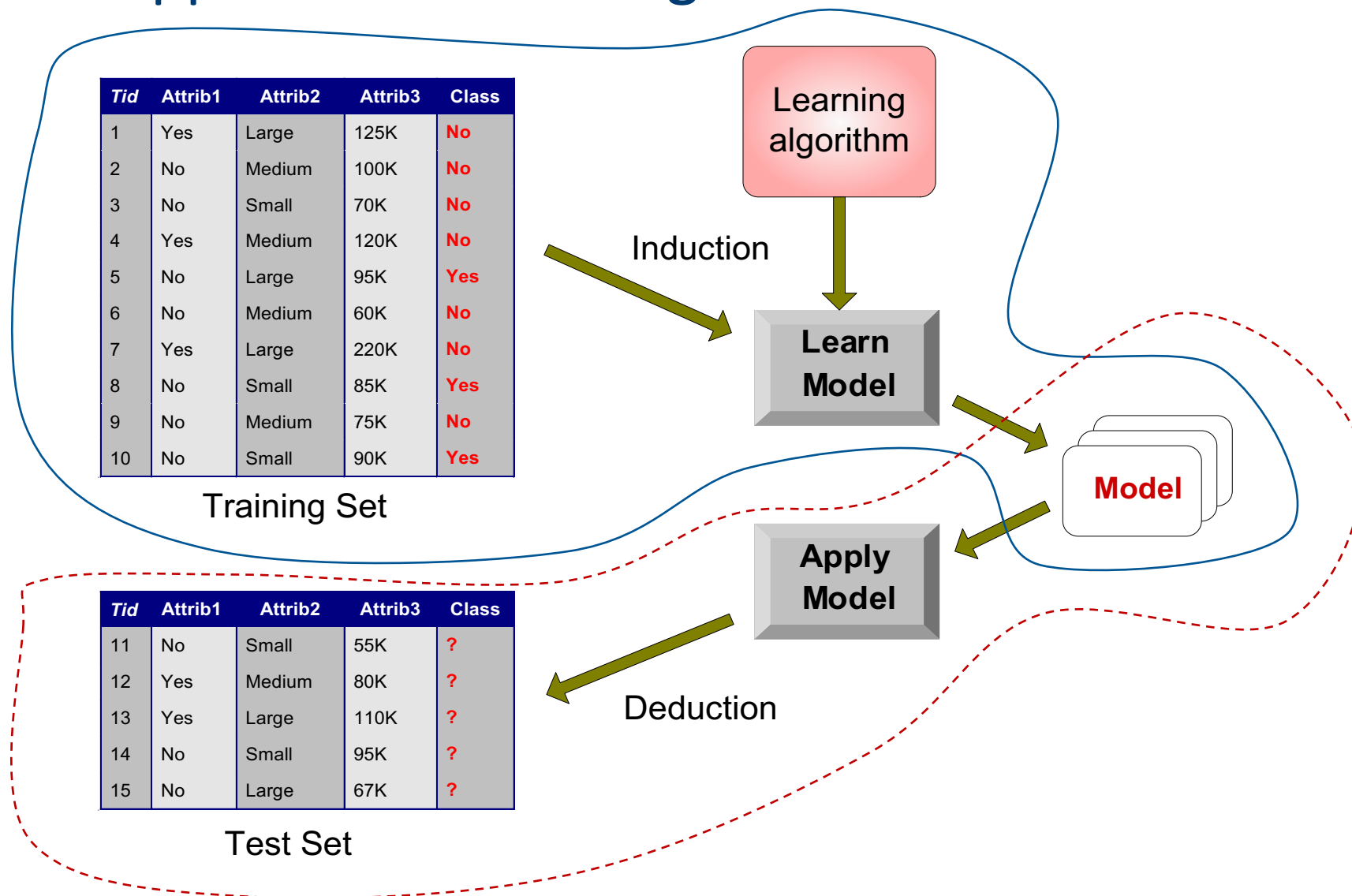
- Given a collection of records (training set)
 - Each record is characterized by a tuple (\mathbf{x}, y) , where \mathbf{x} is the attribute set and y is the class label
 - \mathbf{x} : attribute, predictor, independent variable, input
 - y : class, response, dependent variable, output
- Task:
 - Learn a model that maps each attribute set \mathbf{x} into one of the predefined class labels y



Examples of Classification Task

Task	Attribute set, \mathbf{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 x-rays or 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



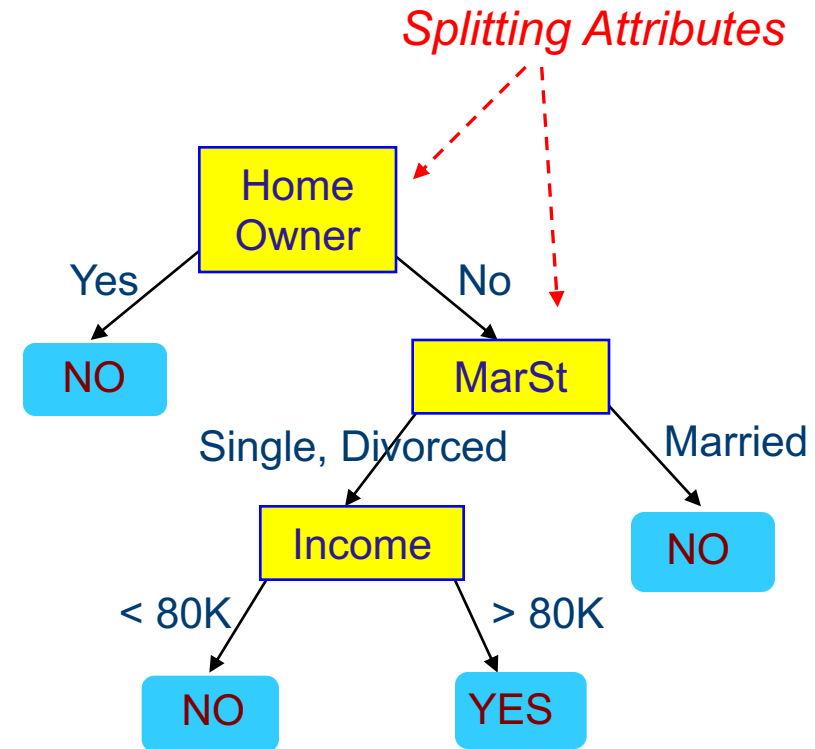
Classification Techniques

- Base Classifiers
 - Decision Tree based Methods ←
 - Rule-based Methods
 - Nearest-neighbor
 - Neural Networks
 - Deep Learning
 - Naïve Bayes and Bayesian Belief Networks
 - Support Vector Machines
- Ensemble Classifiers
 - Boosting, Bagging, Random Forests

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

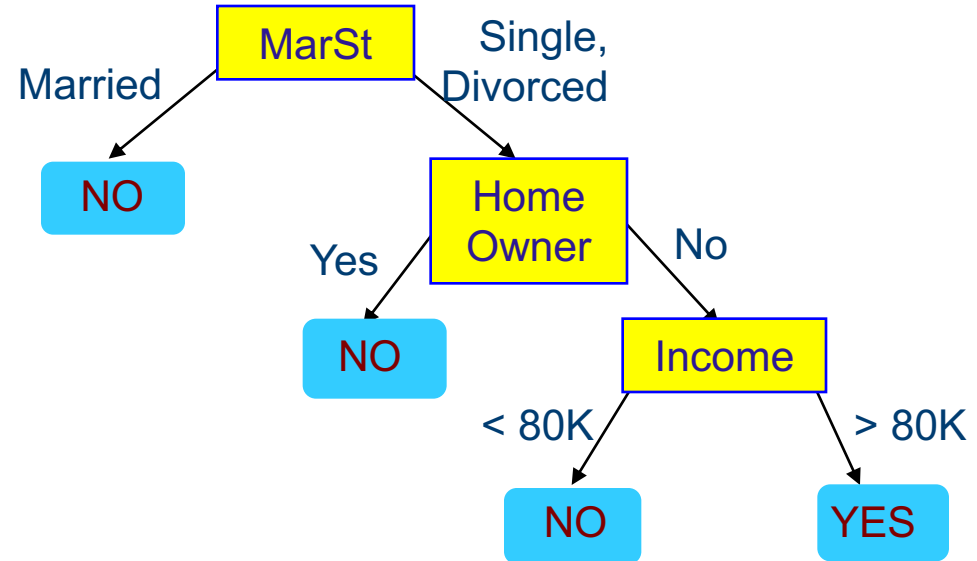


Model: Decision Tree

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

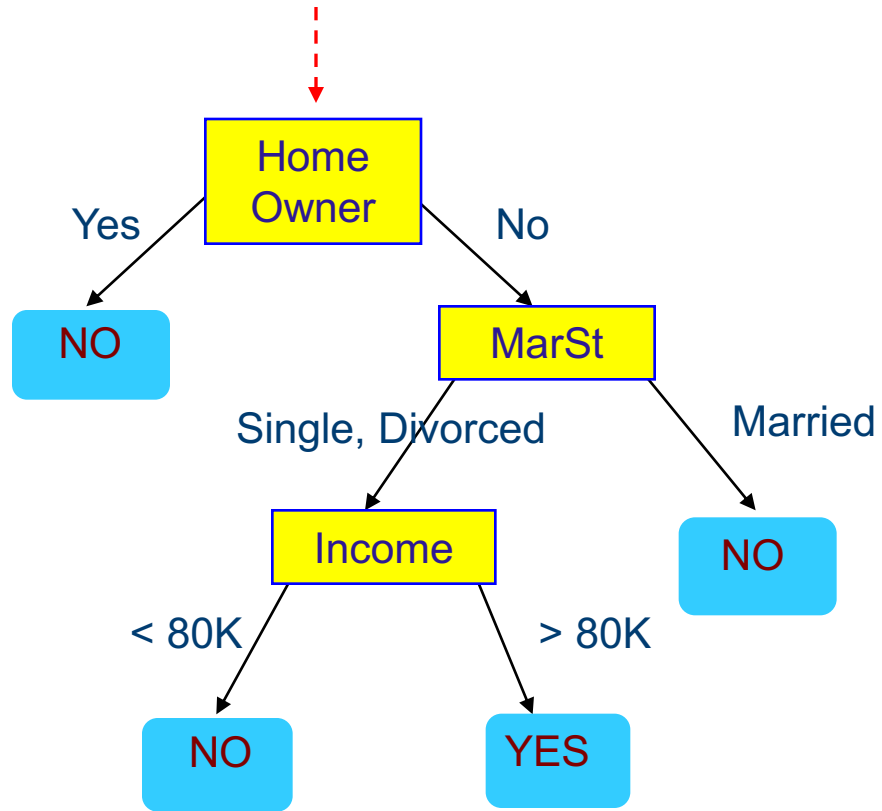
categorical
categorical
continuous
class



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

Apply Model to Test Data

Start from the root of tree.



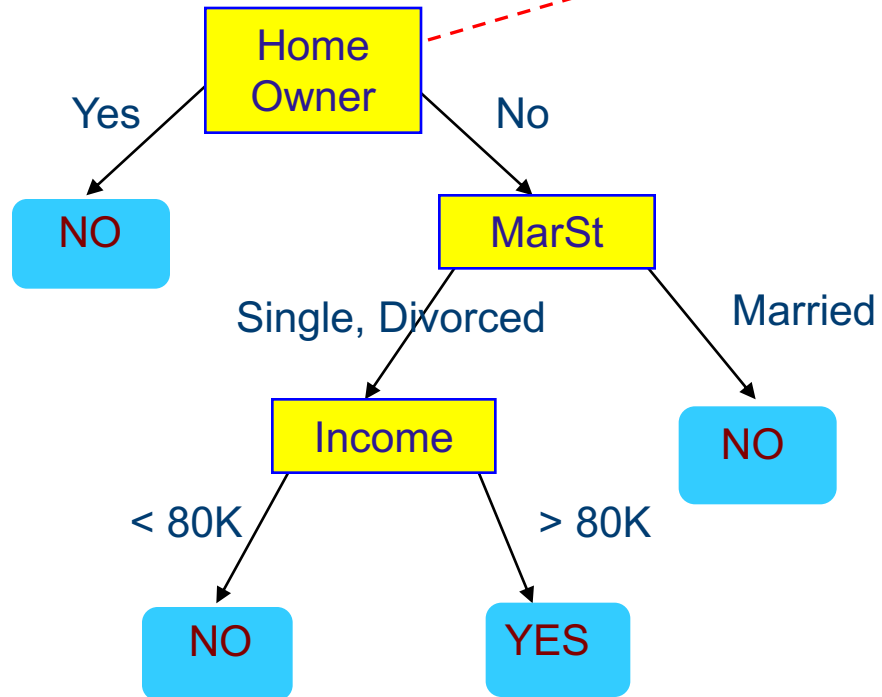
Test Data

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

Apply Model to Test Data

Test Data

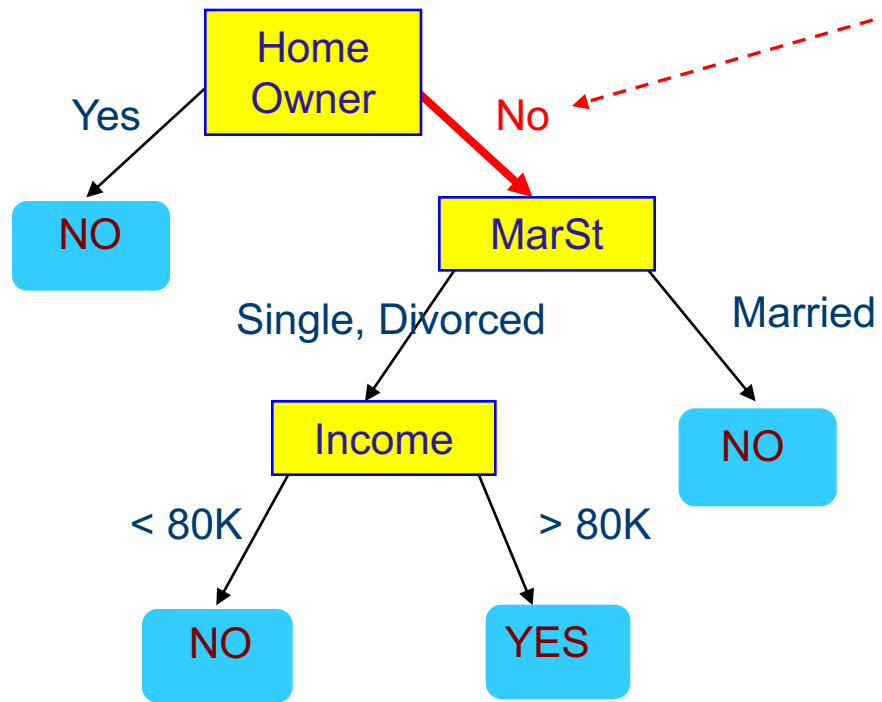
Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



Apply Model to Test Data

Test Data

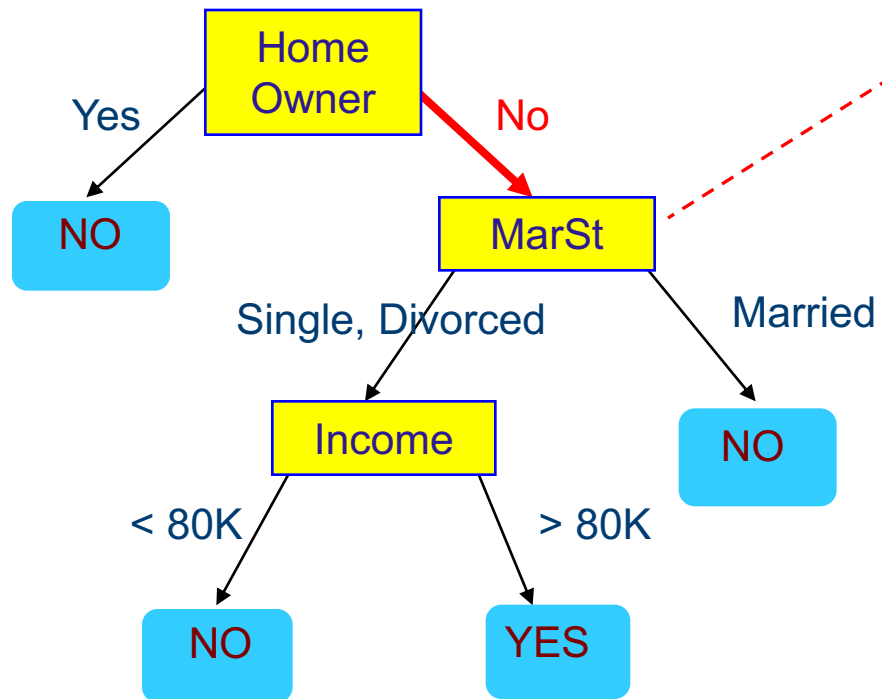
Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



Apply Model to Test Data

Test Data

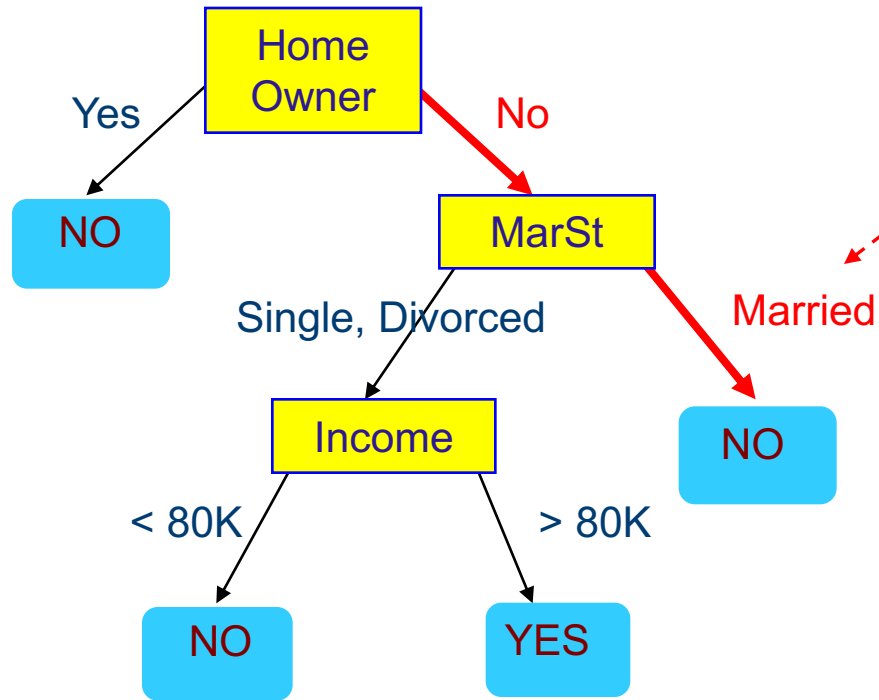
Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



Apply Model to Test Data

Test Data

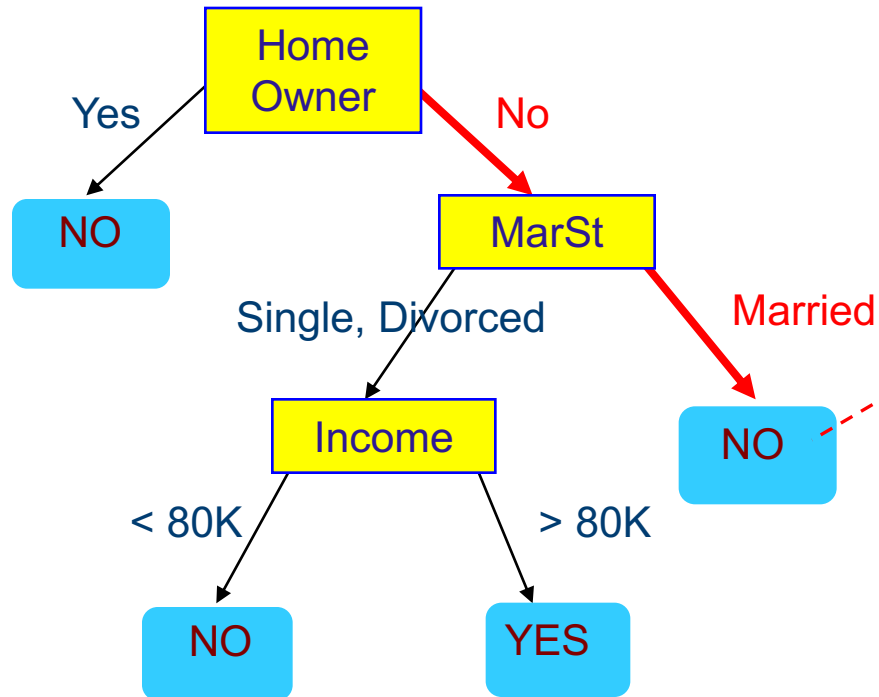
Home Owner	Marital Status	Annual Income	Defaulted Borrower
No	Married	80K	?



Apply Model to Test Data

Test Data

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



Assign Defaulted to "No"

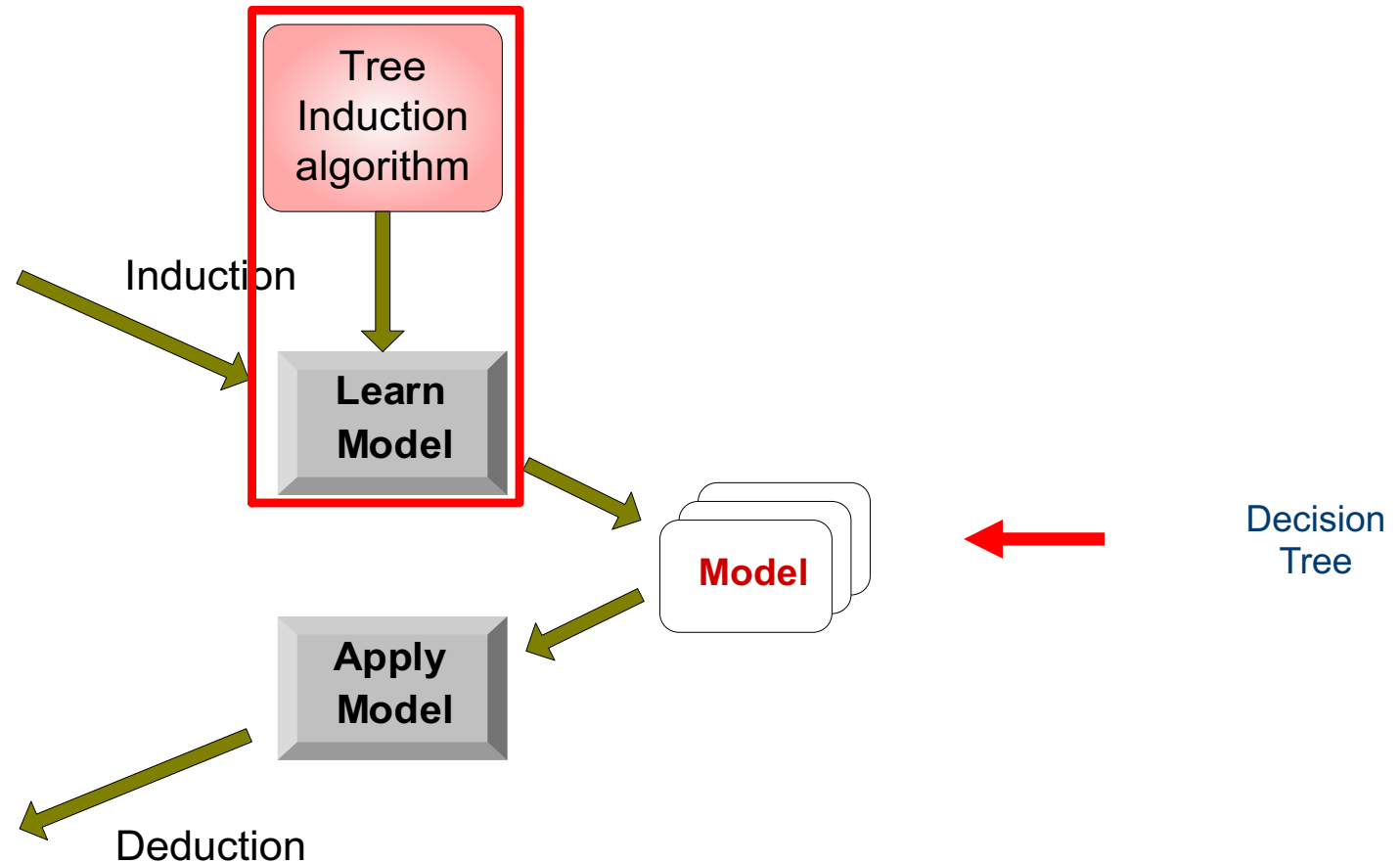
Decision Tree Classification Task

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

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



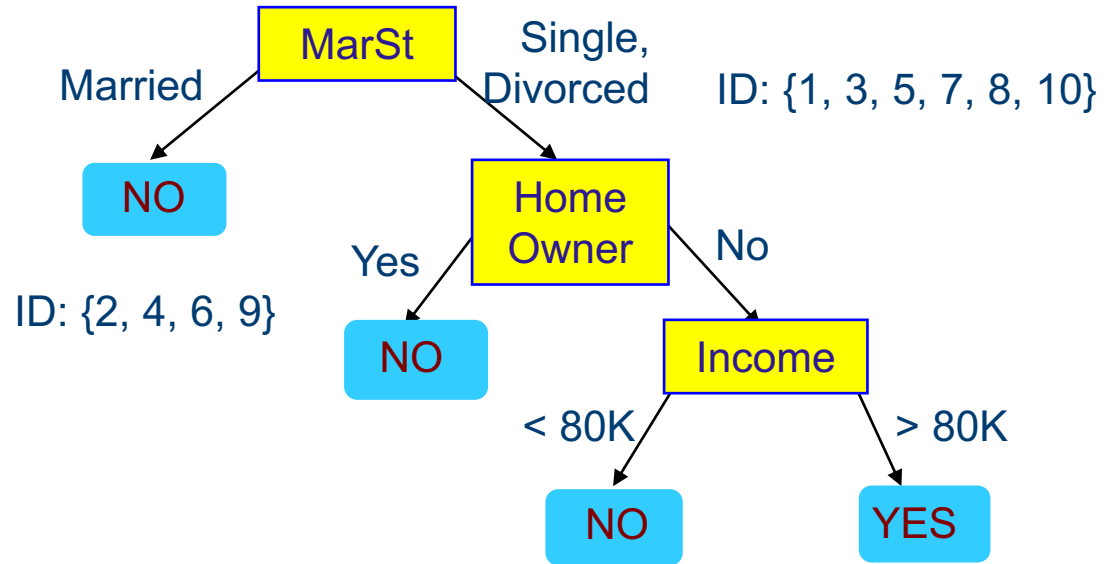
Decision Tree Induction

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

Main Questions

categorical
categorical
continuous
class

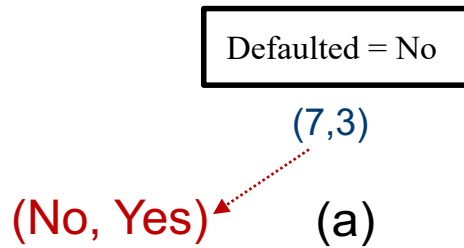
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Basic Questions

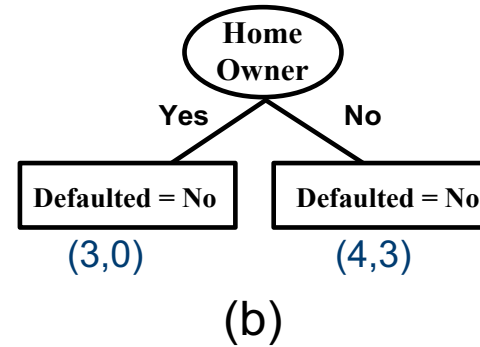
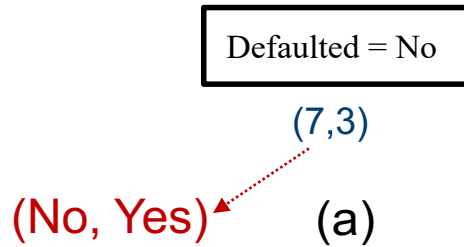
- How to select an attribute for a node?
- How many branches (multi-way or binary)?
- How to deal with continuous attributes?
- When we introduce a leaf?

Example: Hunt's Algorithm General Structure



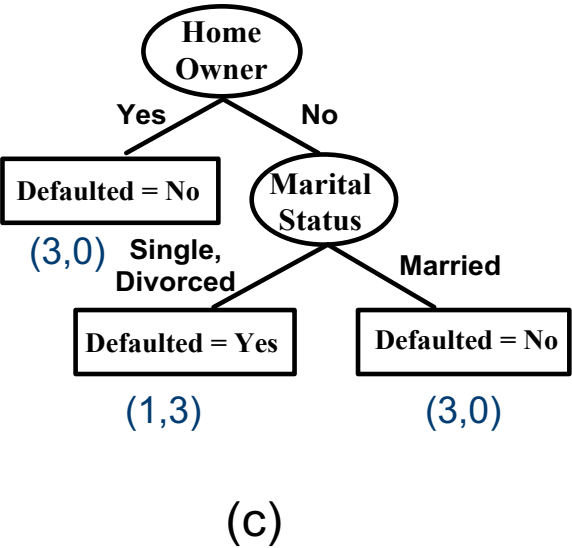
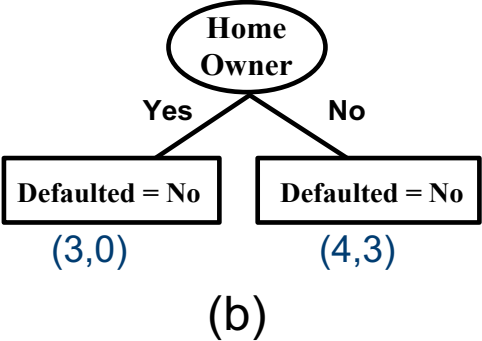
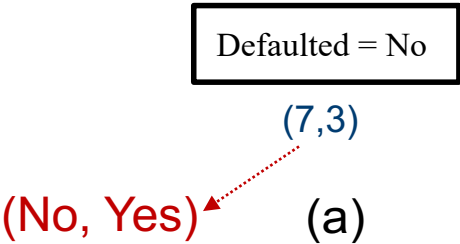
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Example: Hunt's Algorithm General Structure



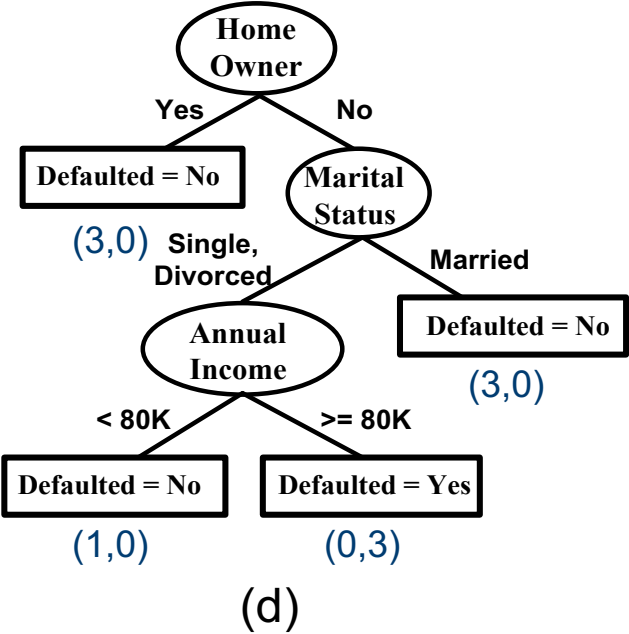
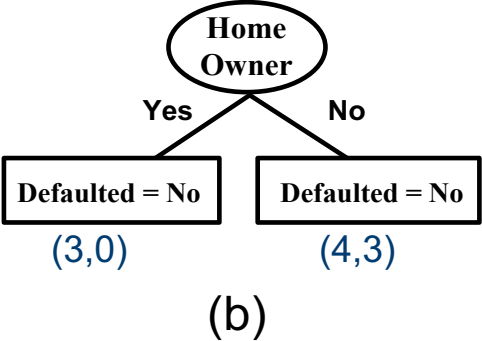
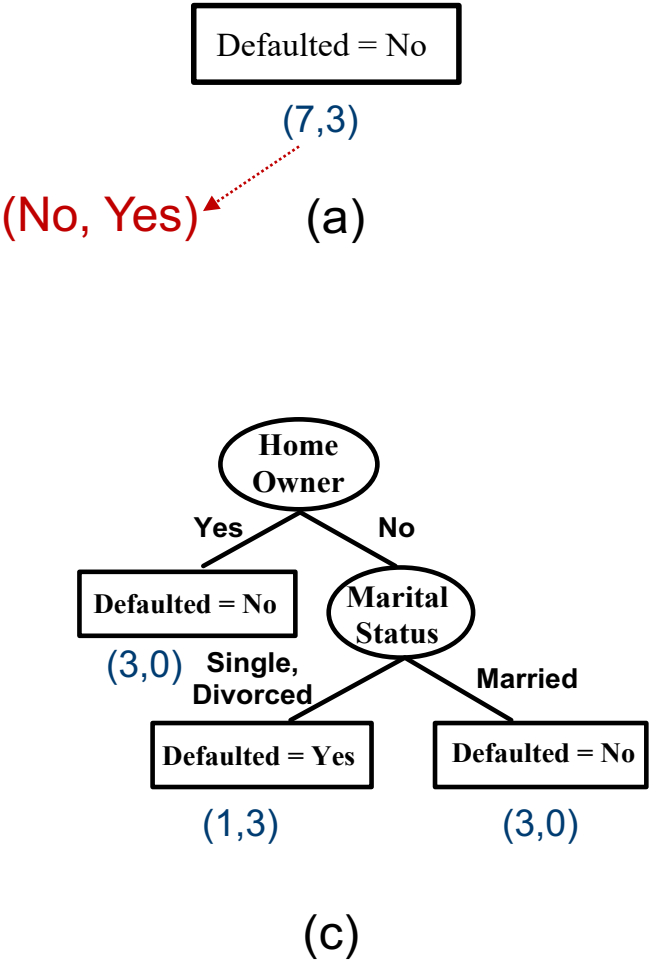
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Example: Hunt's Algorithm General Structure



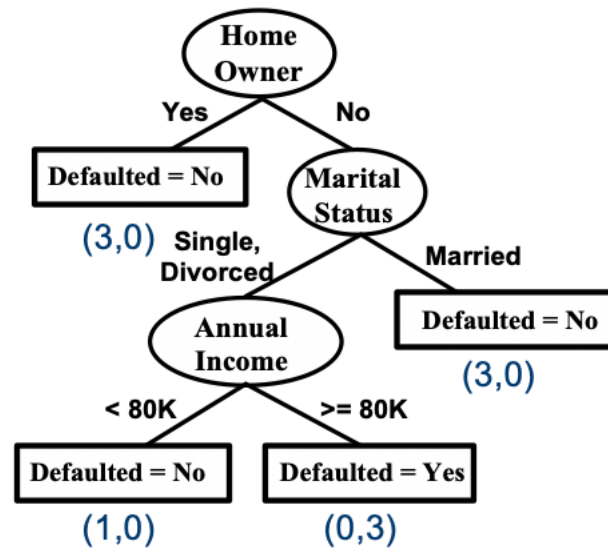
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Example: Hunt's Algorithm General Structure



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

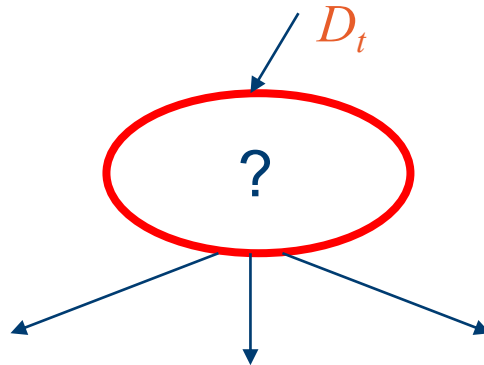
What is the error that model makes if I use **training** to evaluate the model?



- A. 50 %
- B. 100 %
- C. It is not known

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.



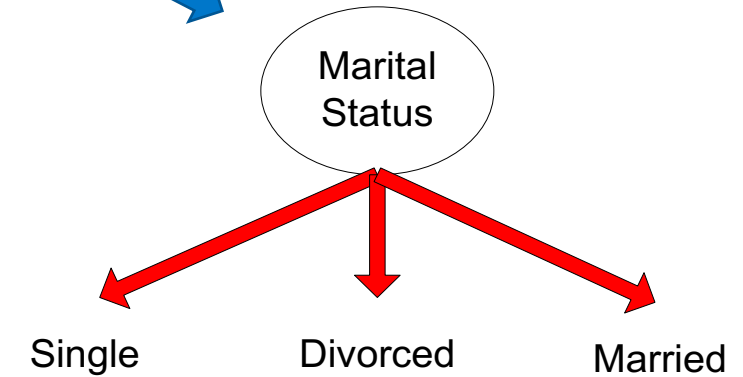
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



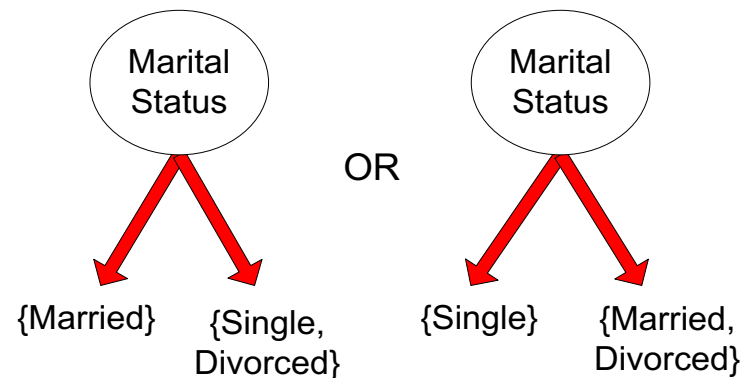
How to deal with different attribute types?

Test Condition for Nominal Attributes

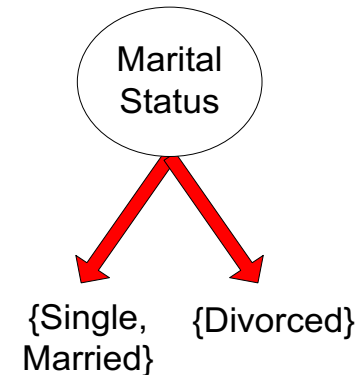
Our focus is on this method!



- **Multi-way split:**
 - Use as many partitions as distinct values.
- **Binary split:**
 - Divides values into two subsets



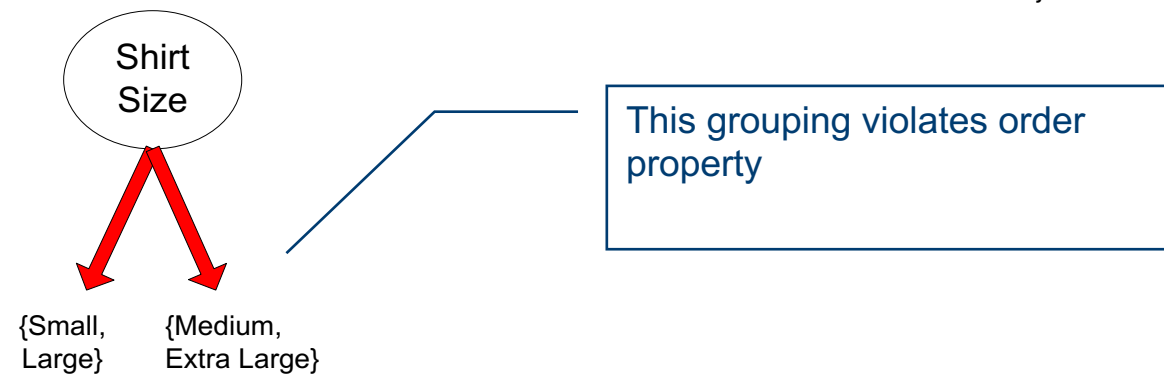
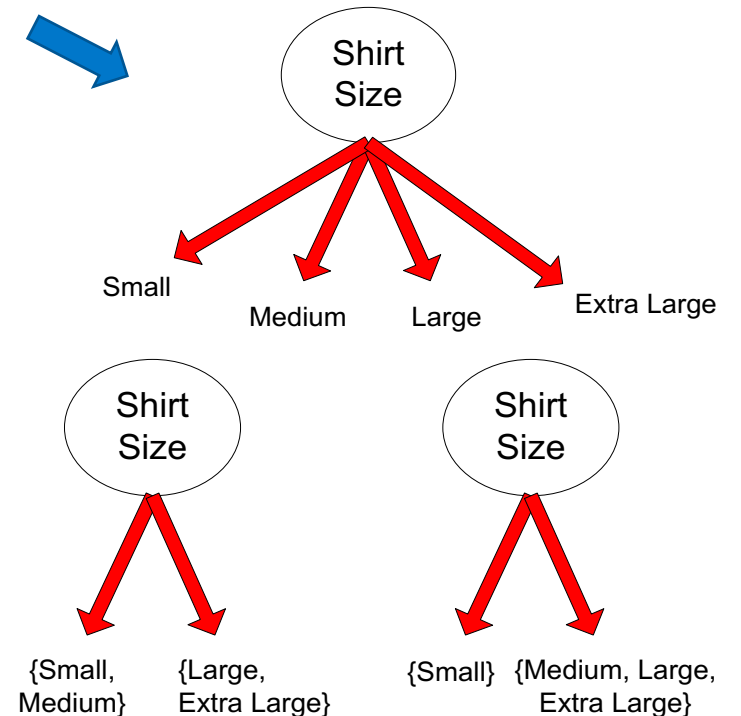
OR



Test Condition for Ordinal Attributes

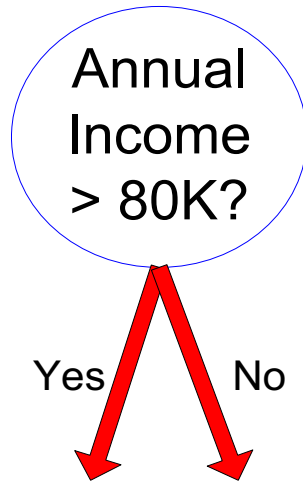
Our focus is on this method!

- **Multi-way split:**
 - Use as many partitions as distinct values
- **Binary split:**
 - Divides values into two subsets
 - Preserve order property among attribute values

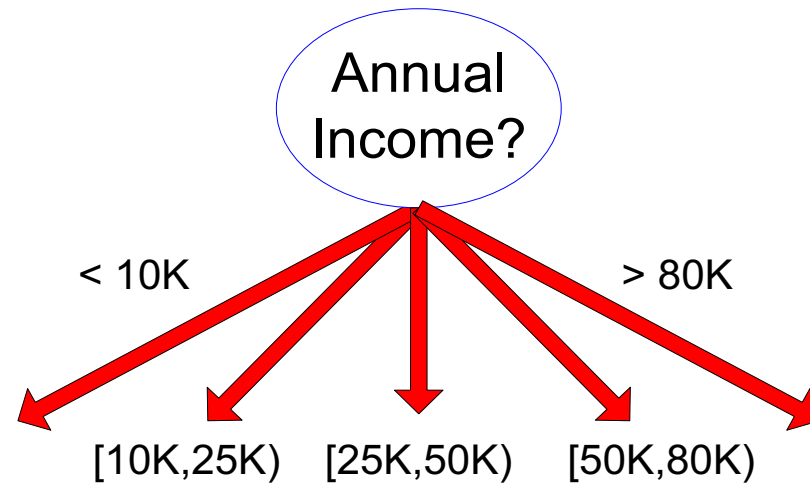


Test Condition for Continuous Attributes

Our focus is on this method!



(i) Binary split



(ii) Multi-way split



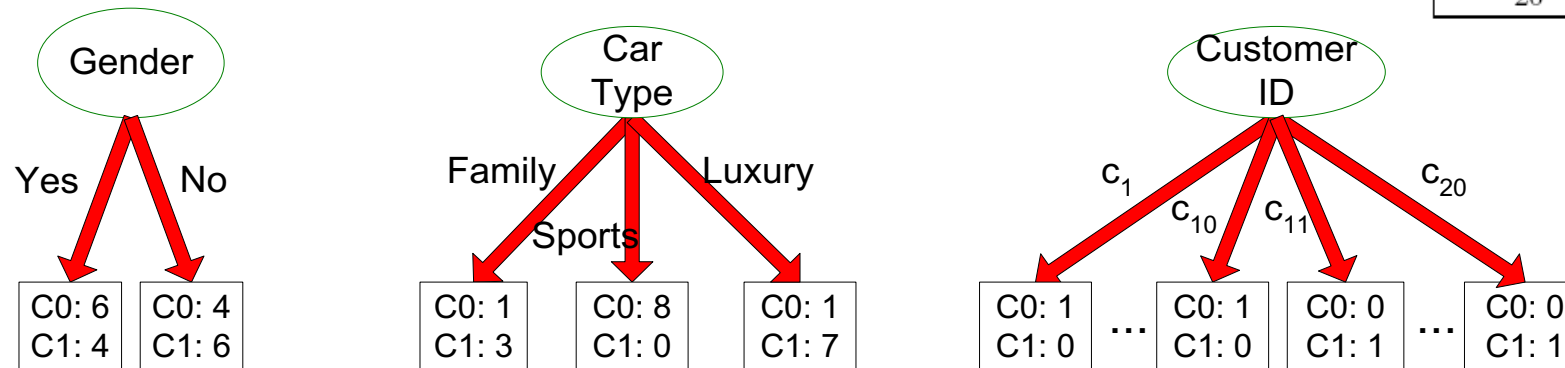
How to deal with different attribute types?

How to determine the best Split?

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


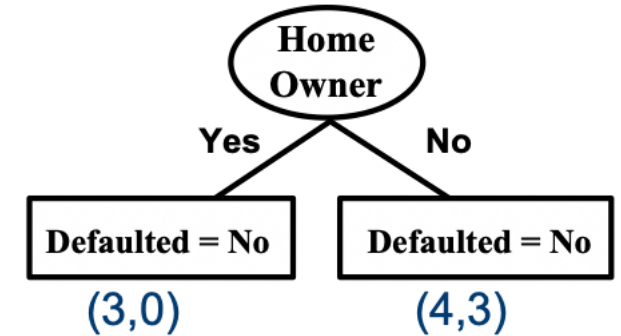
Customer Id	Gender	Car Type	Shirt Size	Class
1	M	Family	Small	C0
2	M	Sports	Medium	C0
3	M	Sports	Medium	C0
4	M	Sports	Large	C0
5	M	Sports	Extra Large	C0
6	M	Sports	Extra Large	C0
7	F	Sports	Small	C0
8	F	Sports	Small	C0
9	F	Sports	Medium	C0
10	F	Luxury	Large	C0
11	M	Family	Large	C1
12	M	Family	Extra Large	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

Which test condition is the best?




How to determine the best Split?

- Greedy approach:
 - Nodes with **pur**er 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 probability 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

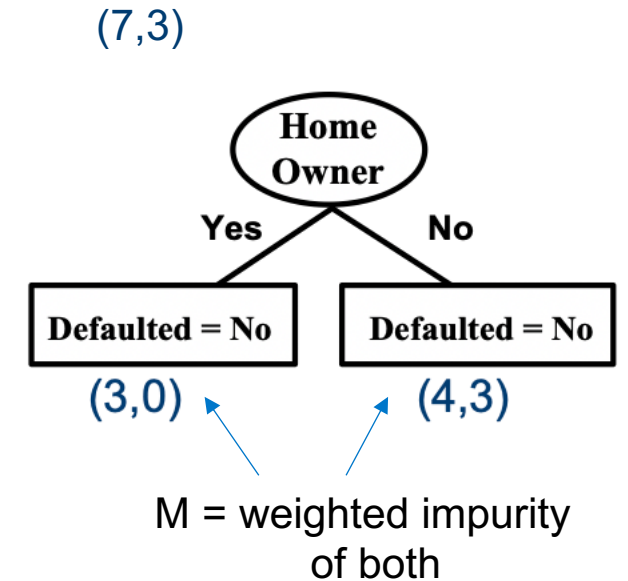
$$Classification\ error = 1 - \max[p_i(t)]$$

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**
 - M is the **weighted impurity** of **child nodes**
3. Choose the attribute test condition that produces the highest gain

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

or equivalently, lowest impurity measure after splitting (M)

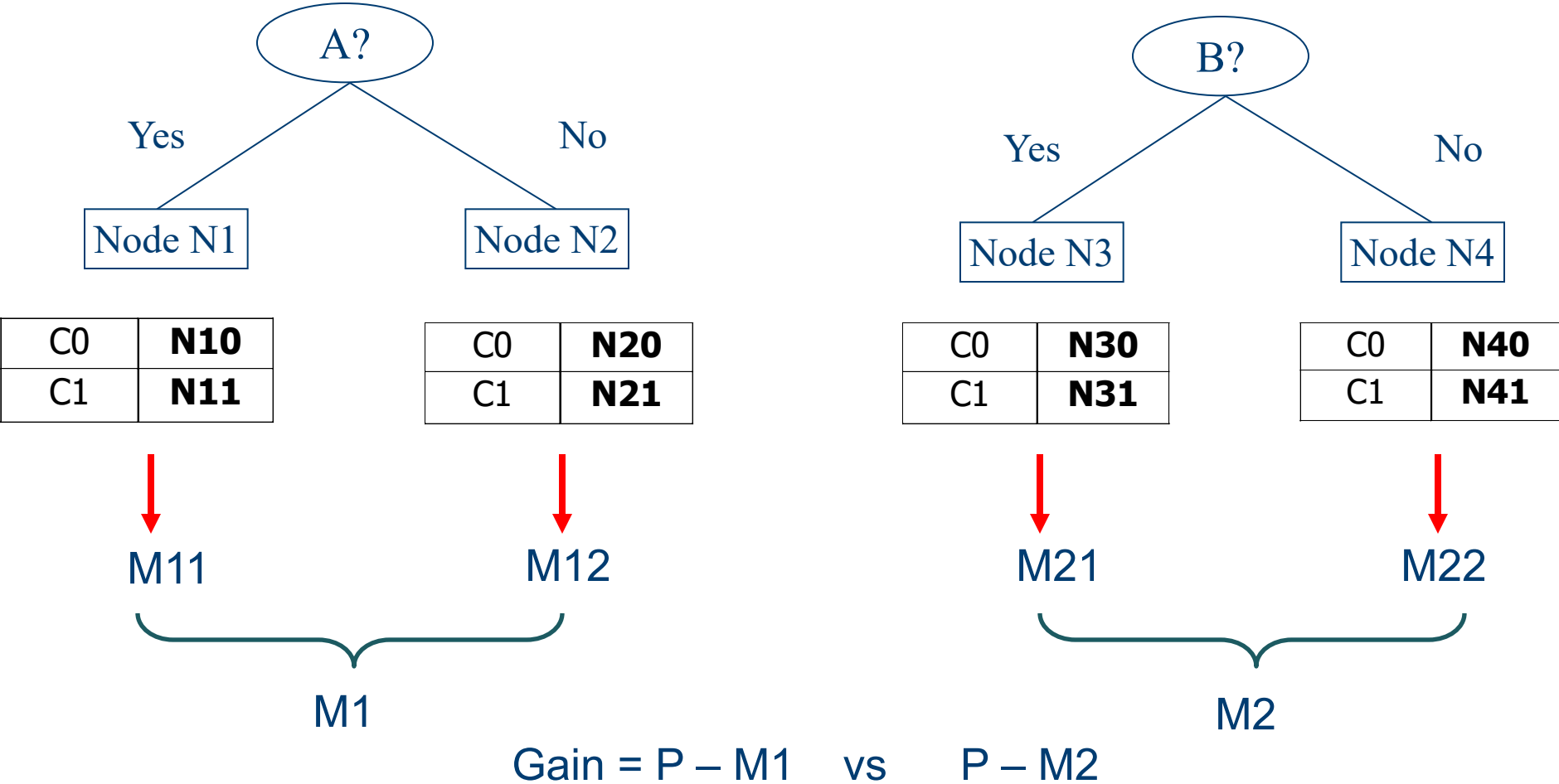


Finding the Best Split (Attribute A or B?)

Before Splitting:

C0	N00
C1	N01

→ P



What is the Gini index for the following class distribution?

C1	0
C2	6
Gini=?	

- A. 0
- B. 0.25
- C. 0.5
- D. 1

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

Measure of Impurity: GINI

- Gini Index for a given node t :

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


- Examples:

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	


$$1 - (1/6)^2 - (5/6)^2$$

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 .

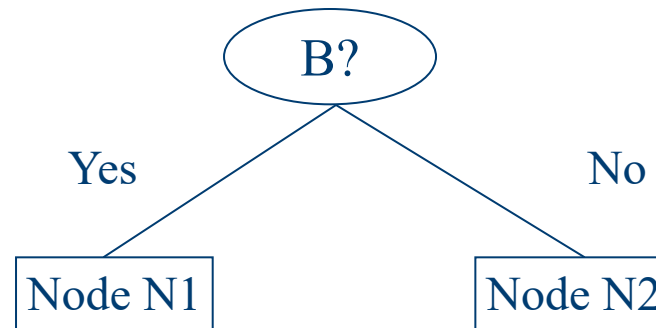
- Choose the attribute that minimizes weighted average Gini index of the children
- Gini index is used in decision tree algorithms such as CART, SLIQ, SPRINT

Computing GINI Index

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

$$\begin{aligned} \text{Gini}(N1) &= 1 - (5/6)^2 - (1/6)^2 \\ &= 0.278 \end{aligned}$$

$$\begin{aligned} \text{Gini}(N2) &= 1 - (2/6)^2 - (4/6)^2 \\ &= 0.444 \end{aligned}$$



	N1	N2
C1	5	2
C2	1	4
Gini=0.361		

	Parent
C1	7
C2	5
Gini = 0.486	

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

$$\begin{aligned} \text{Weighted Gini of N1 N2} &= 6/12 * 0.278 + \\ &\quad 6/12 * 0.444 \\ &= 0.361 \end{aligned}$$

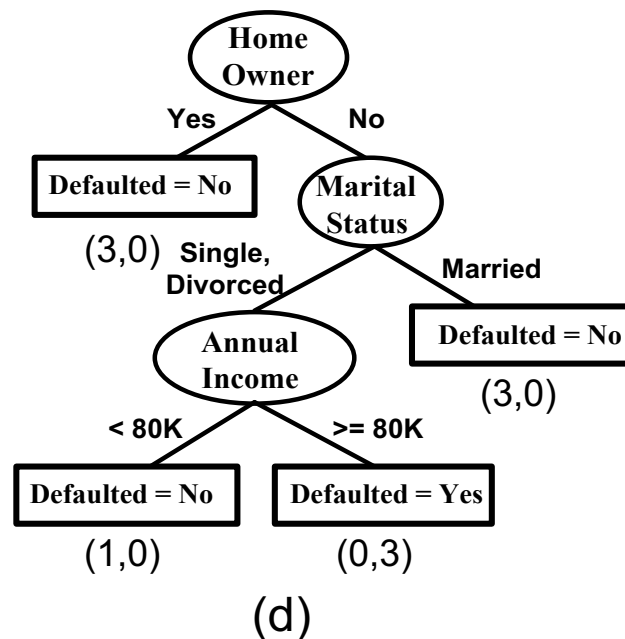
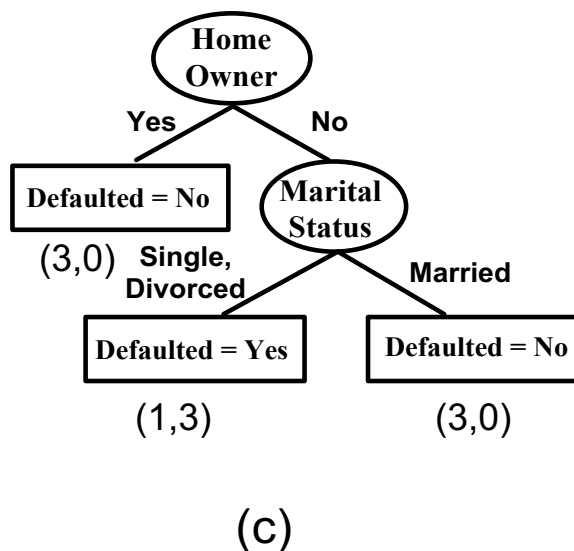
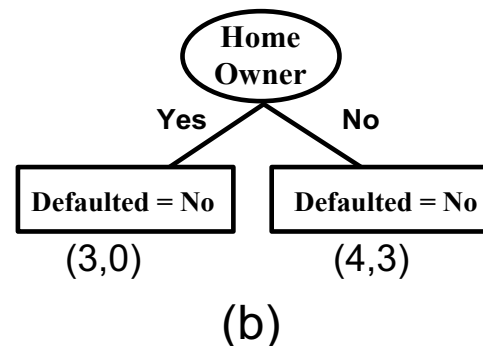
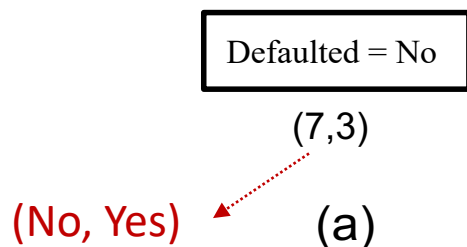
$$\text{Gain} = 0.486 - 0.361 = 0.125$$

What is the impurity measure (M) for “CarType”?

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

- A. 0.163
- B. 0.343
- C. 0.485
- D. 0.677

Example: Hunt's Algorithm General Structure Recap)



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

Now we now how to select the attributes!

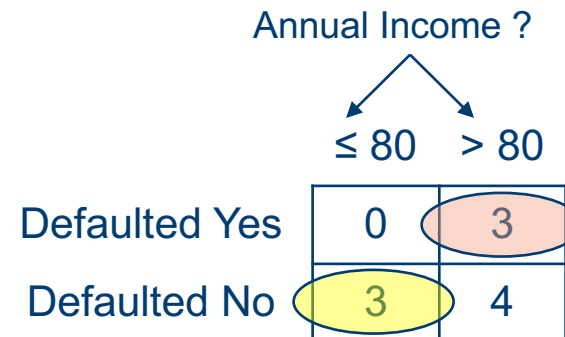
One with maximum information gain

Continuous Attributes: Computing Gini Index

- Use Binary Decisions based on one value
- Each splitting value has a count matrix associated with it
 - Class counts in each of the partitions, $A \leq v$ and $A > v$

	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

How to choose v ?



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

Sorted Values →

Cheat	No	No	No	Yes	Yes	Yes	No	No	No	No
Annual Income										
	60	70	75	85	90	95	100	120	125	220

Continuous Attributes

Computing Gini Index...

Sorted Values Split Positions	→	Cheat	No	No	No	Yes	Yes	Yes	No	No	No	No
		Annual Income										
	→	60	70	75	85	90	95	100	120	125	220	
	→	55	65	72	80	87	92	97	110	122	172	230
		<=	>	<=	>	<=	>	<=	>	<=	>	<=

Continuous Attributes:

Computing Gini Index...

Sorted Values

Split Positions

Cheat

No

No

No

Yes

Yes

Yes

No

No

No

No

Annual Income

60

70

75

85

90

95

100

120

125

220

55

65

72

80

87

92

97

110

122

172

230

<=

>

<=

>

<=

>

<=

>

<=

>

<=

>

<=

>

<=

>

<=

>

Yes

No

Gini

0

3

3

4

0.343

Continuous Attributes:

Computing 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	1	2					
	No				3	4	3	4					
	Gini				0.343		0.417						

Continuous Attributes:

Sorted Values Split Positions		Cheat	No		No		No		Yes		Yes		Yes		No		No		No		No		
		Annual Income																					
		60		70		75		85		90		95		100		120		125		220			
		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		<u>0.300</u>		0.343		0.375		0.400		0.420	

Measure of Impurity: Entropy

- Entropy at a given node t

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

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

- Maximum of $\log_2 c$ when records are equally distributed among all classes, implying the least beneficial situation for classification
- Minimum of 0 when all records belong to one class, implying most beneficial situation for classification
- Entropy based computations are quite similar to the GINI index computations

Computing Entropy of a Single Node

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

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:

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

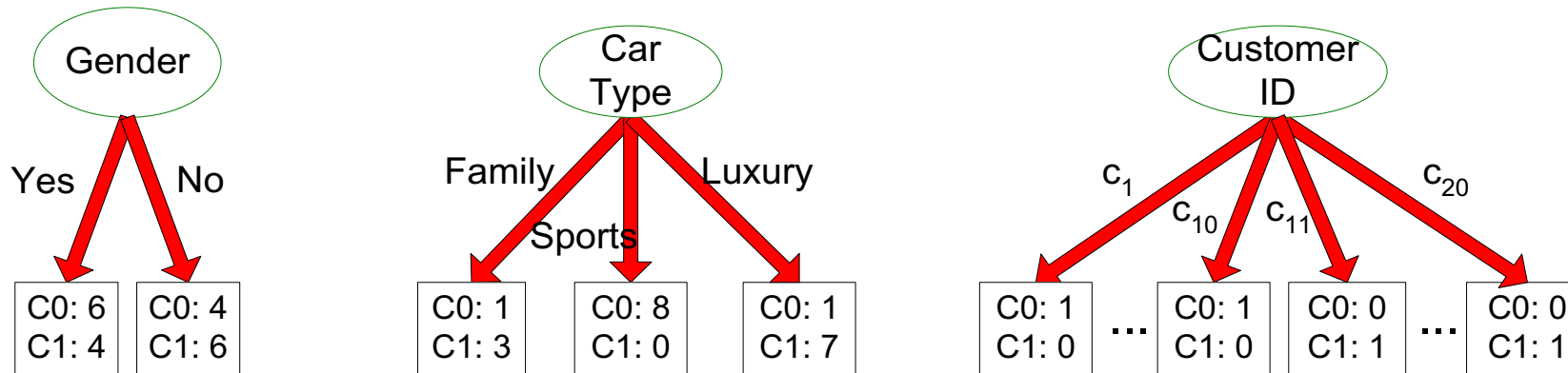
Parent Node, p is split into k partitions (children)

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:

$$\text{Gain Ratio} = \frac{\text{Gain}_{\text{split}}}{\text{Split Info}} \qquad \text{Split Info} = - \sum_{i=1}^k \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

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:

$$\text{Gain Ratio} = \frac{\text{Gain}_{\text{split}}}{\text{Split Info}}$$

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

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

n_i is number of records in child node i

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

SplitINFO = 1.52

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

SplitINFO = 0.72

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

SplitINFO = 0.97

$$-(16/20) \log_2 (16/20) - (4/20) \log_2 (4/20)$$

Decision Tree Based Classification

- **Advantages:**
 - Inexpensive to construct
 - Extremely fast at classifying unknown records
 - Easy to interpret for small-sized trees
 - Robust to noise (especially when methods to avoid overfitting are employed)
- **Disadvantages:**
 - Space of possible decision trees is exponentially large. Greedy approaches are often unable to find the best tree.
 - Does not take into account interactions between attributes
 - Each decision boundary involves only a single attribute



How to evaluate a model?

Model Evaluation

- Purpose:
 - To estimate performance of classifier on previously unseen data (test set)
- **Metrics** for Performance Evaluation
 - How to evaluate the performance of a model?
- **Methods** for Performance Evaluation
 - How to obtain reliable estimates?

Metrics for Performance Evaluation

- Focus on the **predictive capability** of a model
 - Rather than how fast it takes to classify or build models, scalability, etc.
- **Confusion Matrix:**

	PREDICTED CLASS		
		Class=Yes	Class=No
	Class=Yes	a	b
	Class=No	c	d

a: TP (true positive)

b: FN (false negative)

c: FP (false positive)

d: TN (true negative)

Metrics for Performance Evaluation

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

- Most widely-used metric:

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

Other Measures

$$\text{Precision (p)} = \frac{a}{a + c}$$

$$\text{Recall (r)} = \frac{a}{a + b}$$

$$\text{F - measure (F)} = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

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

Consider a 2-class problem, and what is accuracy if model always predicts YES?

Number of Class YES examples = 99

Number of Class NO examples = 1

- A. 0 %
- B. 25%
- C. 50%
- D. 99%

Limitation of Accuracy

- Consider a 2-class problem
 - Number of Class 0 examples = 9990
 - Number of Class 1 examples = 10
- If model predicts everything to be class 0, accuracy is $9990/10000 = 99.9 \%$

Accuracy is misleading because model
does not detect any class 1 example !!

Cost Matrix

	PREDICTED CLASS		
	$C(i j)$	Class=Yes	Class=No
	Class=Yes	$C(\text{Yes} \text{Yes})$	$C(\text{No} \text{Yes})$
	Class=No	$C(\text{Yes} \text{No})$	$C(\text{No} \text{No})$

$C(i|j)$: Cost of misclassifying class j example as class i

Computing Cost of Classification

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

Cost Matrix	PREDICTED CLASS		
ACTUAL CLASS	C(i j)	+	-
	+	-1	100
	-	1	0

Will be define based on different problems

Model M_1	PREDICTED CLASS		
ACTUAL CLASS		+	-
	+	150	40
	-	60	250

Accuracy = 80%

Cost = 3910

Model M_2	PREDICTED CLASS		
ACTUAL CLASS		+	-
	+	250	45
	-	5	200

Accuracy = 90%

Cost = 4255

What if I don't have test data (almost always!)

How to estimate performance of classifier on previously unseen data (test set)?

- Holdout
 - Reserve $k\%$ for training and $(100-k)\%$ for testing
 - Random subsampling: repeated holdout

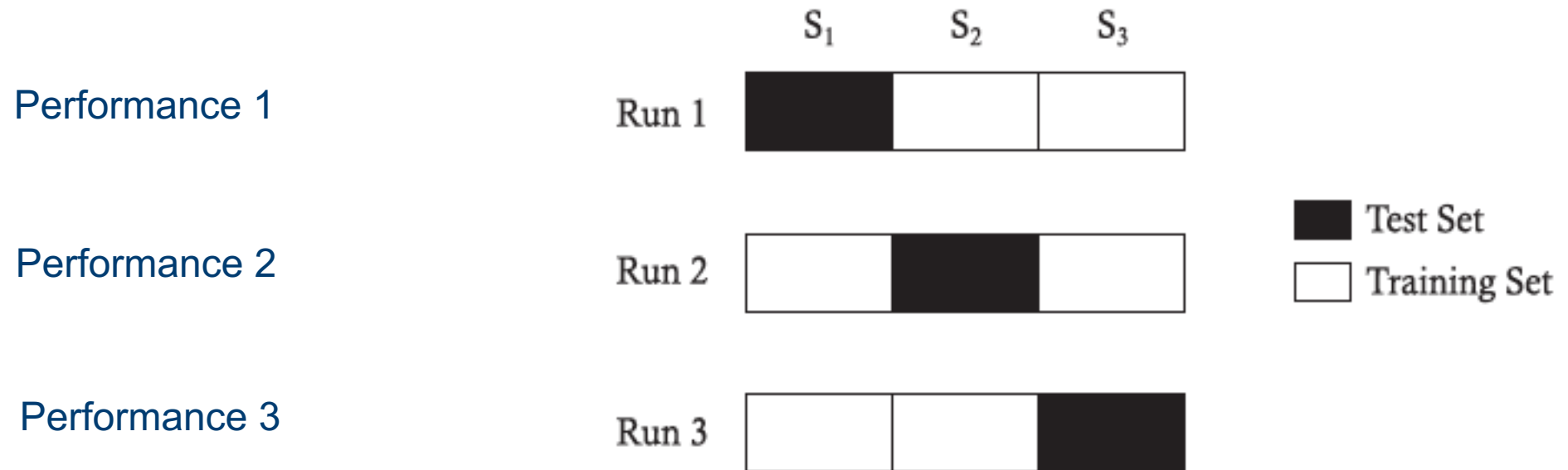
What if I don't have test data (almost always!)

How to estimate performance of classifier on previously unseen data (test set)?

- Cross validation
 - Partition data into k disjoint subsets
 - k-fold: **train** on $k-1$ partitions, **test** on the remaining one
 - Leave-one-out: $k = n$ (number of instances)

Cross-validation Example

- 3-fold cross-validation



Final Performance = **Average** of All Performance

What is over the fitting problem?

Example Data Set

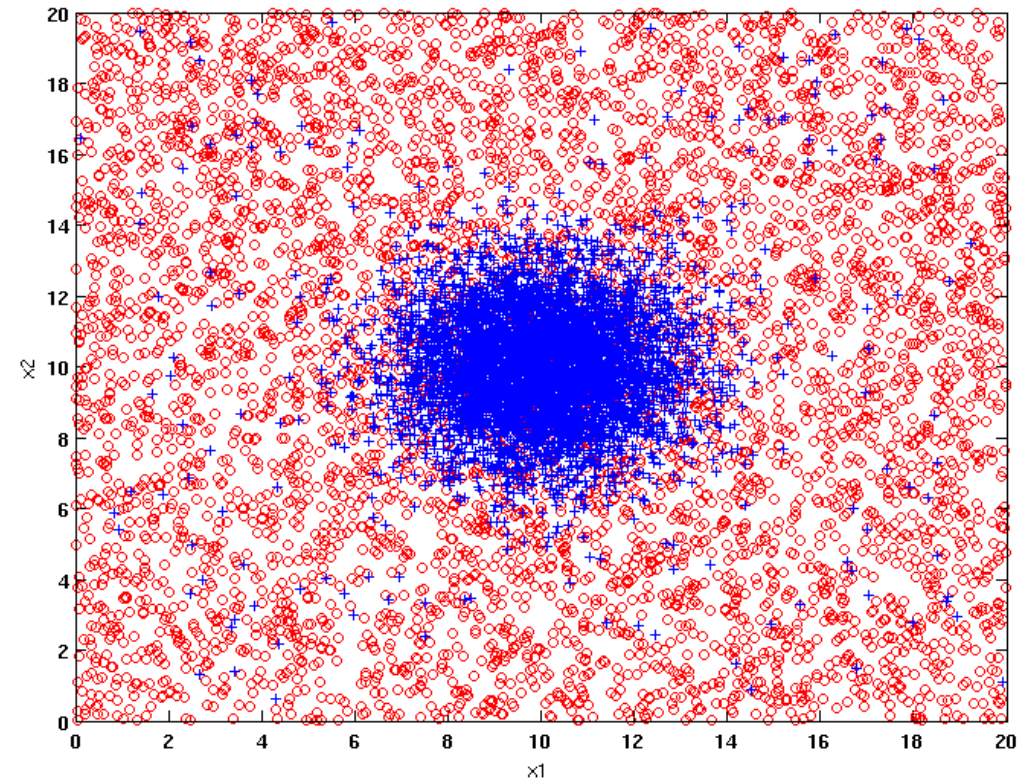
Two class problem:

+ : 5400 instances

- 5000 instances generated from a Gaussian centered at (10,10)
- 400 noisy instances added

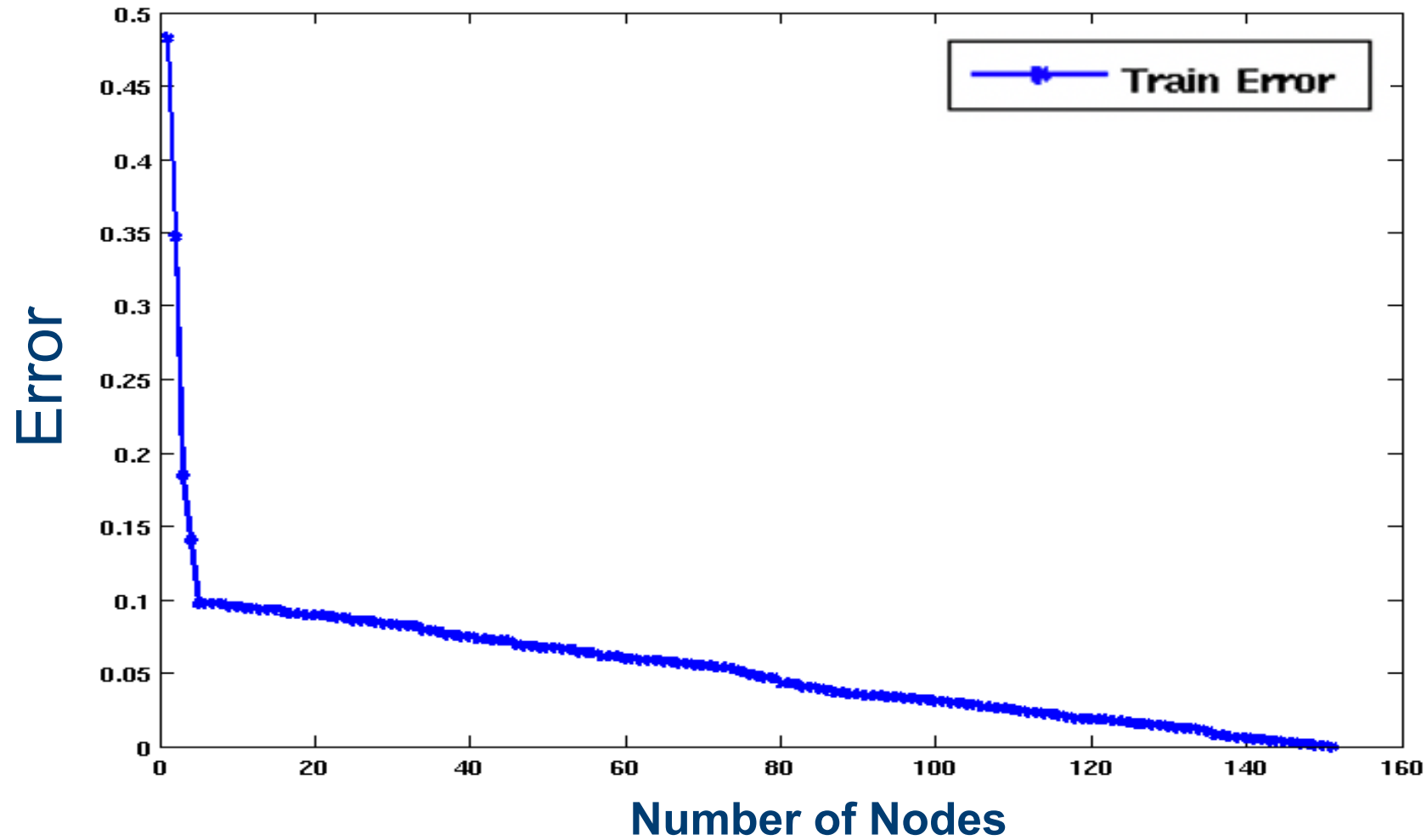
o : 5400 instances

- Generated from a uniform distribution

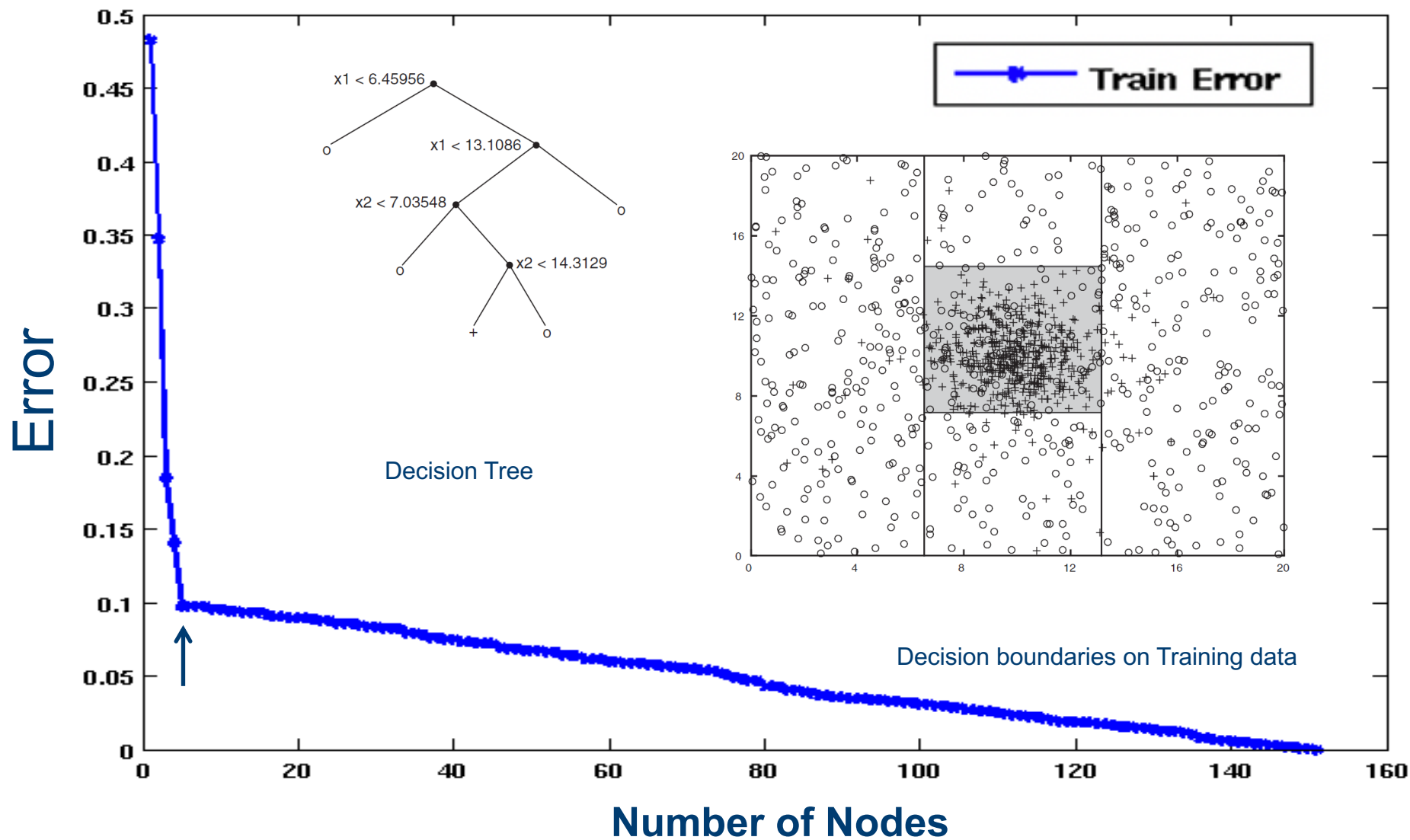


10 % of the data used for training and 90%
of the data used for testing

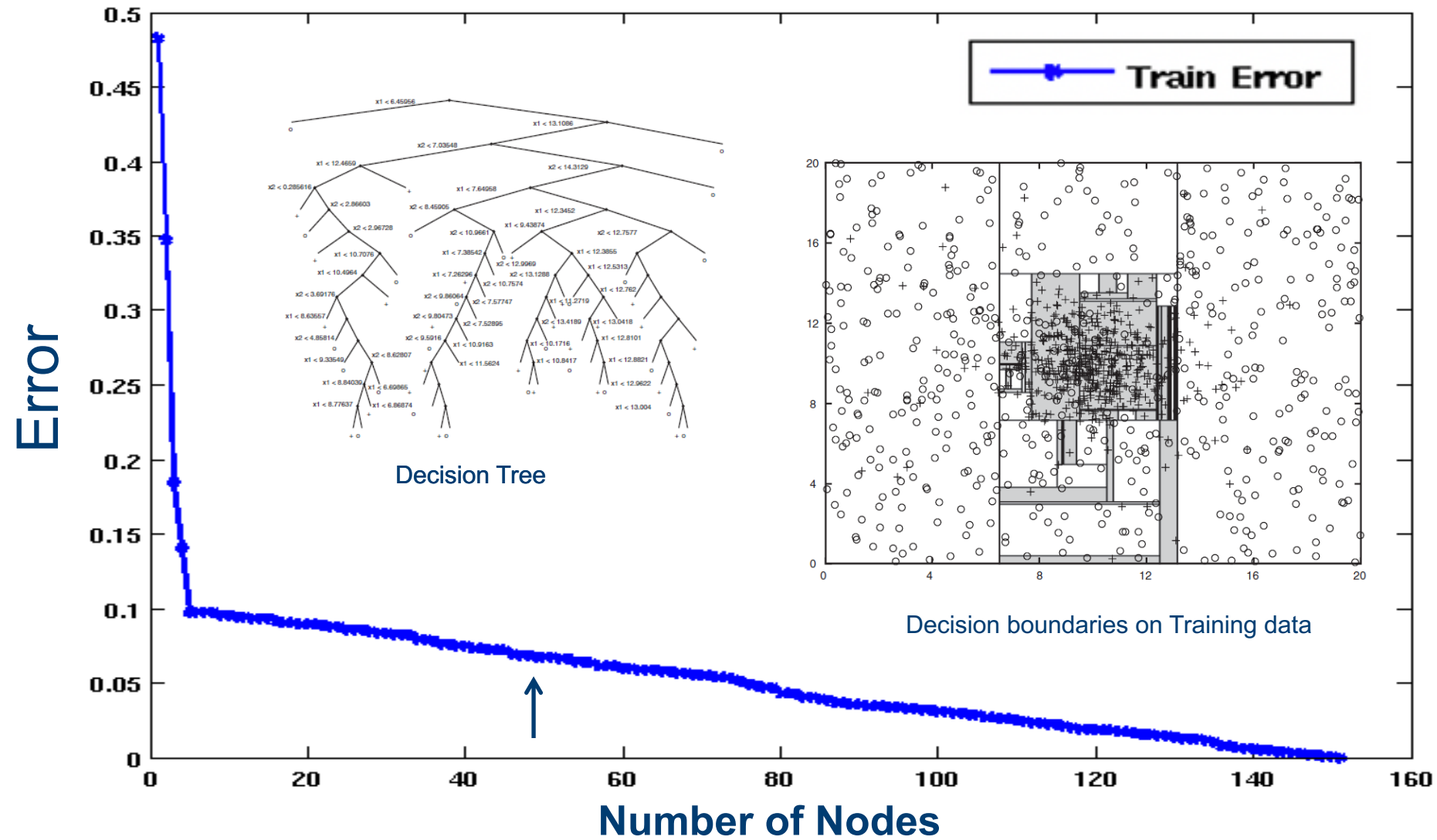
Increasing number of nodes in Decision Trees



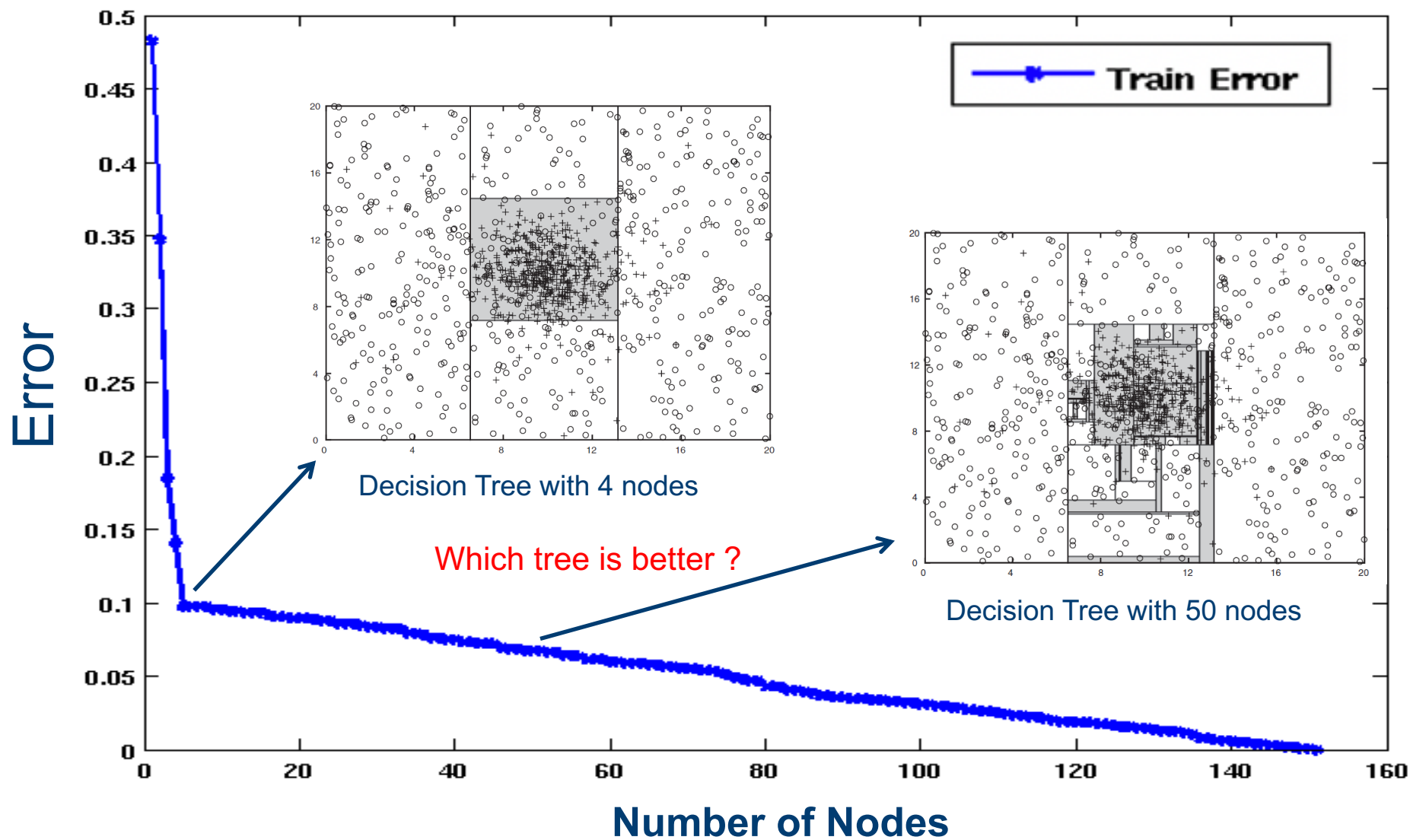
Decision Tree with 4 nodes



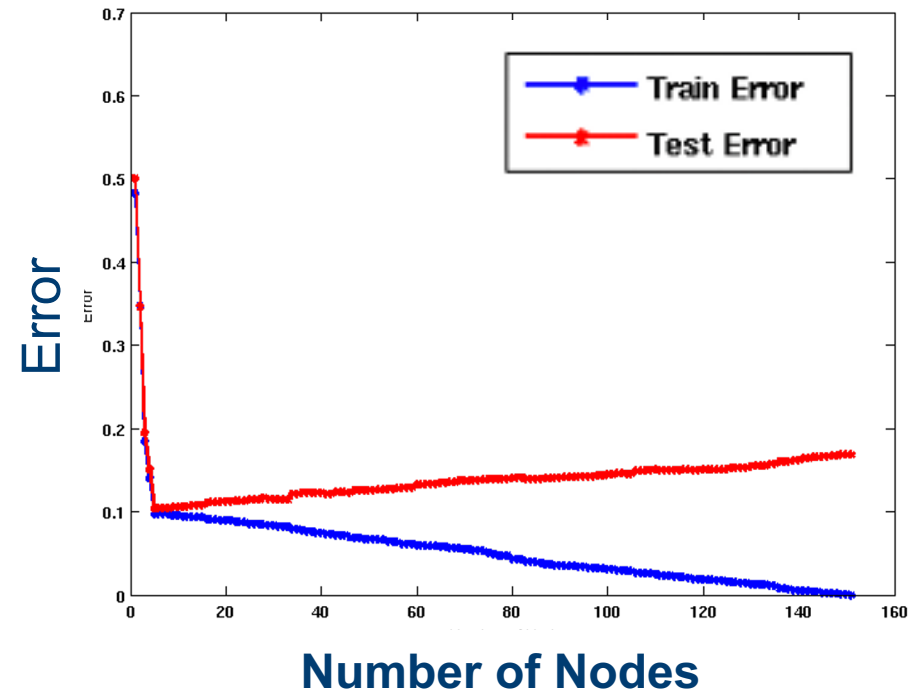
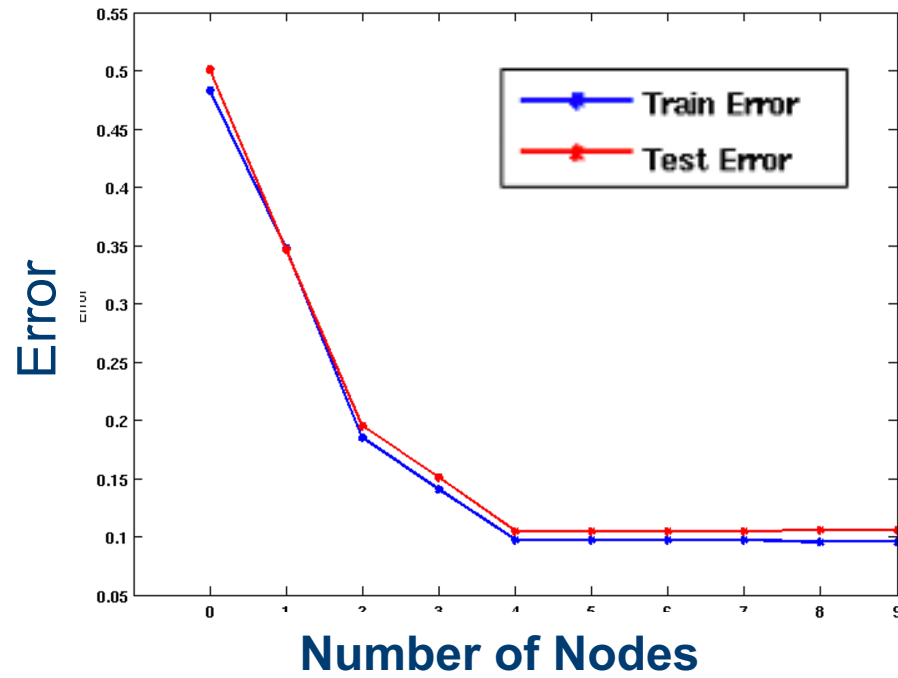
Decision Tree with 50 nodes



Which tree is better?

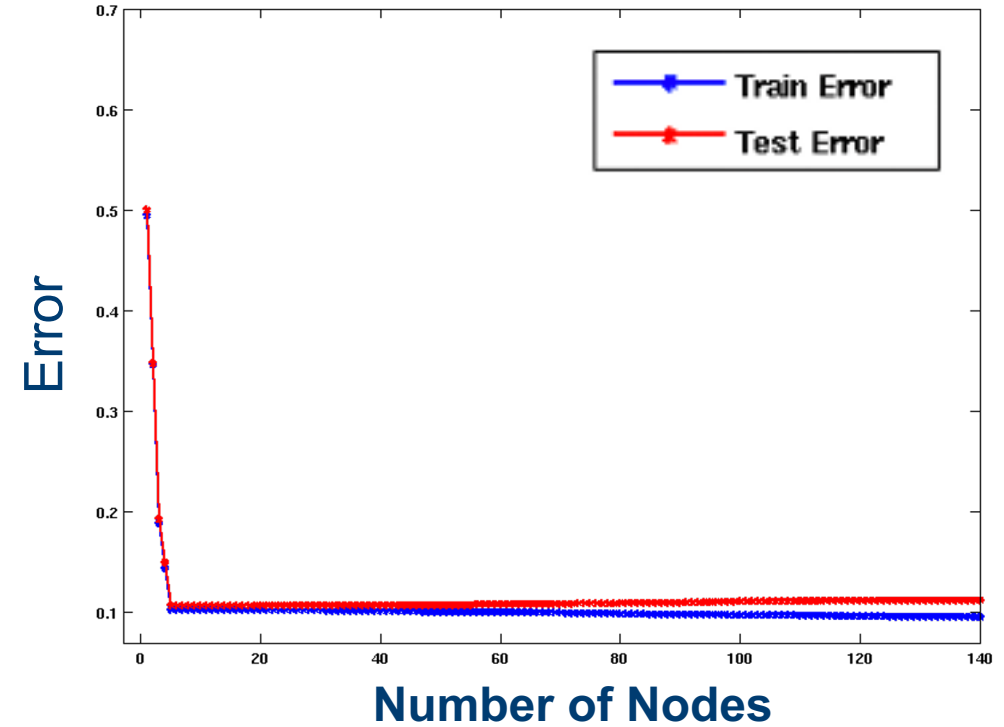
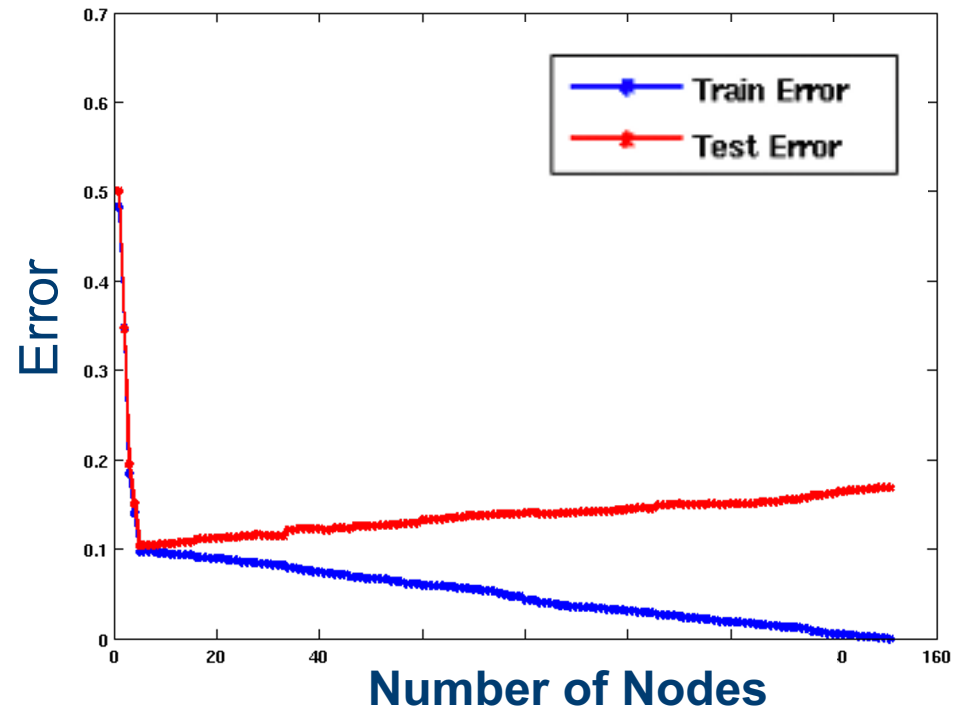


Model Overfitting



- As the model becomes more and more complex, test errors can start increasing even though training error may be decreasing
- **Underfitting**: when model is too simple, both training and test errors are large
- **Overfitting**: when model is too complex, training error is small but test error is large

Model Overfitting



(Using twice the number of data instances)

- Increasing the size of training data reduces the difference between training and testing errors at a given size of model

What are two important ways of dealing with over fitting?

- A. Increasing size of training set
- B. Reduce the complexity of model
- C. A and B



How to do model selection?

- Which one of my two/more model is better?
- I don't know which one is not overfitted

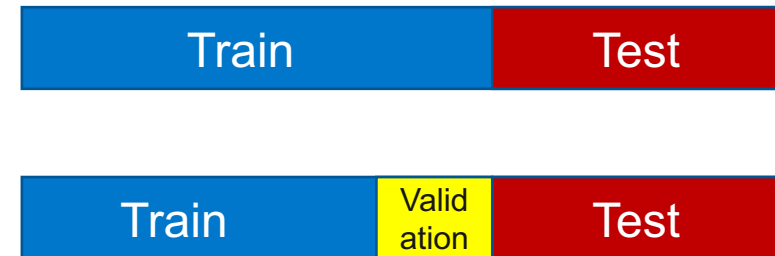


Model Selection

- Performed during model building
- Purpose is to ensure that model is not overly complex (to avoid overfitting)
- How we do model selection?
 - Using validation set method
 - Incorporating model complexity in learning/evaluation process

Model Selection: Using Validation Set

- Divide training data into two parts:
 - Training set:
 - use for model building
 - Validation set:
 - use for estimating generalization error
 - Note: validation set is not the same as test set
- Drawback:
 - Less data available for training



Model Selection: Incorporating Model Complexity

- Rationale: Occam's Razor
 - Given two models of similar **generalization errors**, one should prefer the **simpler** model over the more **complex** model
 - A complex model has a greater chance of being fitted accidentally
 - Therefore, one should include model complexity when evaluating a model

$$\text{Gen. Error}(\text{Model}) = \text{Train. Error}(\text{Model}, \text{Train. Data}) + \alpha \times \text{Complexity}(\text{Model})$$

Estimating the Complexity of Decision Trees

- Pessimistic Error Estimate of decision tree T with k leaf nodes:

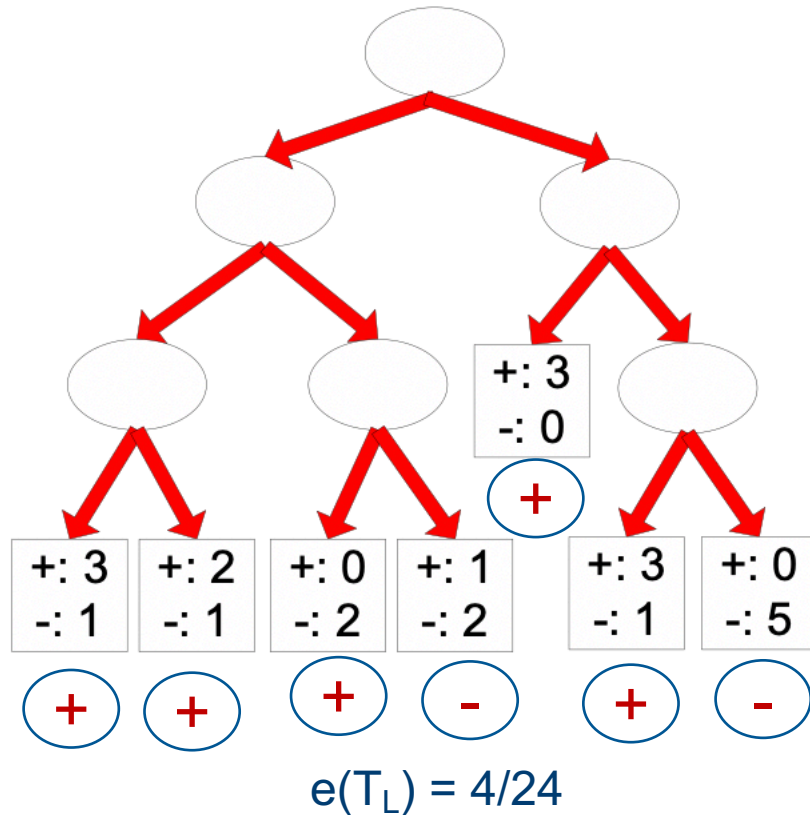
$$err_{gen}(T) = err(T) + \Omega \times \frac{k}{N_{train}}$$

- $err(T)$: error rate on all training records
- Ω : trade-off hyper-parameter (similar to α)
 - Relative cost of adding a leaf node
- k : number of leaf nodes
- N_{train} : total number of training records

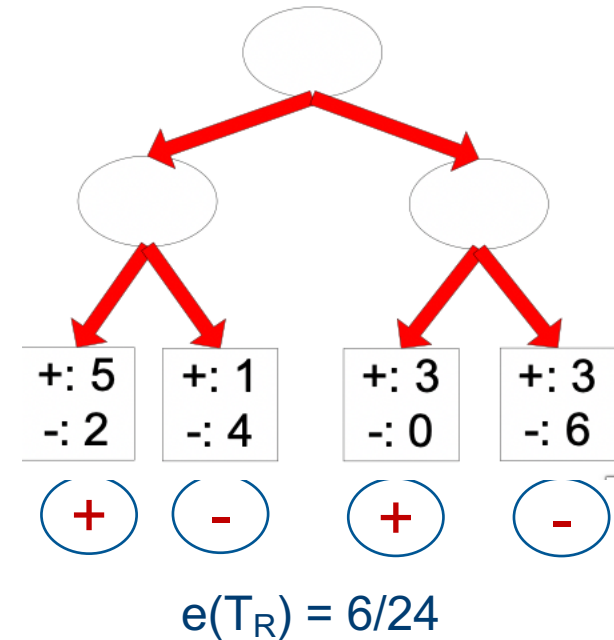
Estimating the Complexity of Decision Trees: Example

- Counts shows the class distribution of training instances and $\Omega = 1$

$$err_{gen}(T) = err(T) + \Omega \times \frac{k}{N_{train}}$$



$$err_{gen}(T_L) = 4/24 + 1 \cdot 7/24 = 11/24 = 0.458$$



$$err_{gen}(T_R) = 6/24 + 1 \cdot 4/24 = 10/24 = 0.417$$

Model Selection for Decision Trees

- 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
 - More restrictive conditions:
 - Stop if number of instances is less than some user-specified threshold
 - Stop if expanding the current node does not improve impurity measures (e.g., Gini or information gain).
 - Stop if estimated generalization error falls below certain threshold

Model Selection for Decision Trees

- Post-pruning
 - Grow decision tree to its entirety
 - Subtree replacement
 - 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
 - Subtree raising
 - Replace subtree with most frequently used branch

Examples of Post-pruning

Decision Tree:

```
depth = 1 :
| breadth > 7 : class 1
| breadth <= 7 :
| | breadth <= 3 :
| | | ImagePages > 0.375 : class 0
| | | ImagePages <= 0.375 :
| | | | totalPages <= 6 : class 1
| | | | totalPages > 6 :
| | | | | breadth <= 1 : class 1
| | | | | breadth > 1 : class 0
| | breadth > 3 :
| | | MultiIP = 0:
| | | | ImagePages <= 0.1333 : class 1
| | | | ImagePages > 0.1333 :
| | | | | breadth <= 6 : class 0
| | | | | breadth > 6 : class 1
| | | MultiIP = 1:
| | | | TotalTime <= 361 : class 0
| | | | TotalTime > 361 : class 1
depth > 1 :
| MultiAgent = 0:
| | depth > 2 : class 0
| | depth <= 2 :
| | | MultiIP = 1: class 0
| | | MultiIP = 0:
| | | | breadth <= 6 : class 0
| | | | breadth > 6 :
| | | | | RepeatedAccess <= 0.0322 : class 0
| | | | | RepeatedAccess > 0.0322 : class 1
| MultiAgent = 1:
| | totalPages <= 81 : class 0
| | totalPages > 81 : class 1
```

Simplified Decision Tree:

```
depth = 1 :
| ImagePages <= 0.1333 : class 1
| ImagePages > 0.1333 :
| | breadth <= 6 : class 0
| | breadth > 6 : class 1
depth > 1 :
| MultiAgent = 0: class 0
| MultiAgent = 1:
| | totalPages <= 81 : class 0
| | totalPages > 81 : class 1
```

Subtree
Raising

Subtree
Replacement