
Data Mining

Assignment Project Exam Help

Classification:

<https://powcoder.com>

Decision Tree

Add WeChat powcoder

Sample solutions to in-class exercises

- Design record data attributes for students at CityU. Think what attributes you want to choose? What are the possible attribute values for each attribute? Design the attributes to cover nominal, ordinal, interval, and ratio.

Assignment Project Exam Help

<https://powcoder.com>

[illegible]

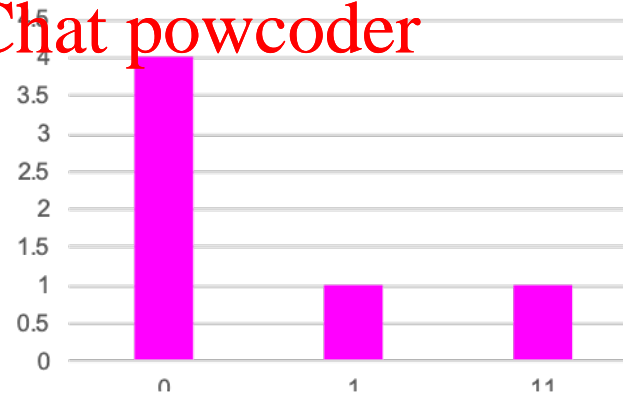
In-class exercise

□ Provide a set of integers so that:

- Its mode is 5
 - Its 10th percentile is 1
 - Its 50th percentile is 5
 - Its 90th percentile is 8
- Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder
- One possible answer: 1 5 5 5 8
 - Any other answer?

Sample solutions to exercises in note2

- Give me an input A such as A's mean = A's median = A's mode e.g. $\langle 1, 2, 3, 4, 5 \rangle$ $\langle 1, 1, 2, 3, 4, 4 \rangle$
- Give me an input A such as A's mean < A's median
 - e.g. $\langle 0, 0, 10, 10, 10 \rangle$
- Give me an input A such as A's mean > A's median
 - $0, 0, 0, 0, 1, 11$



Review exercises

1. Below is the rounded petal length values for 11 iris flowers

$L = \langle 3 \ 4 \ 2 \ 8 \ 4 \ 2 \ 3 \ 1 \ 3 \ 5 \ 2 \rangle$

Answer the following questions:

Q1: what is the mode of L?

Q2: plot the histogram of L (use bin width 1)

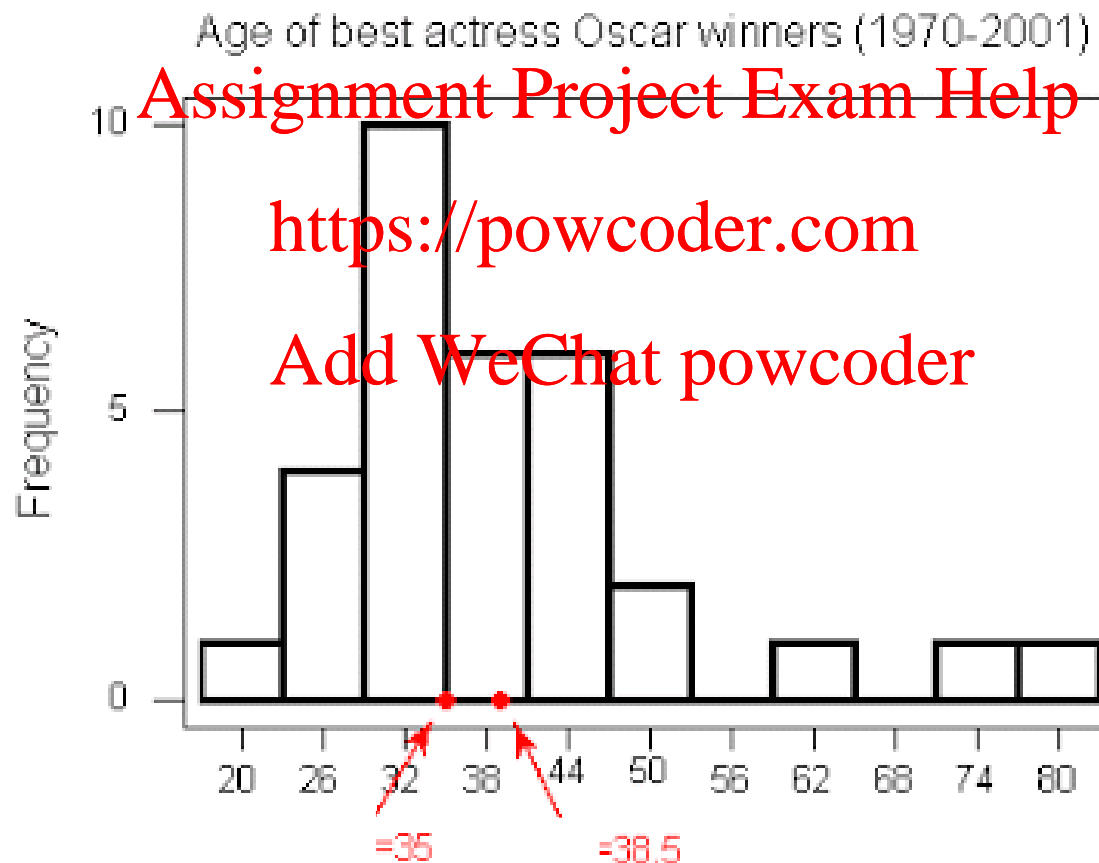
Q3: what is the 25th, 50th, and 75th percentile of L?

Q4: what is the mean and median of L?

Q5: what is the range of L?

Review exercises

- For the following histograms, what is the relationship between mean and median?



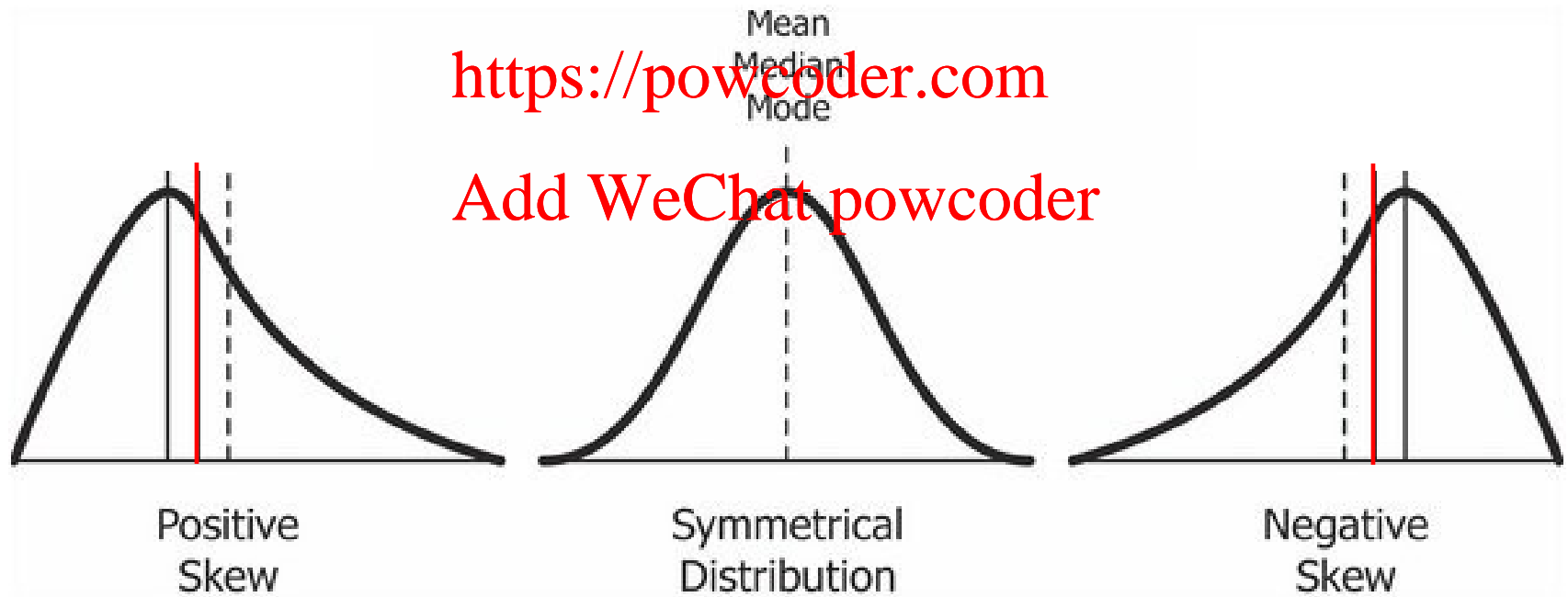
Review exercises

For the following distributions, what is the relationship between mean and median?

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Classification: Definition (review)


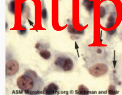
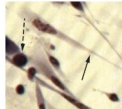
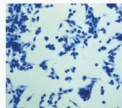
- Given a collection of records (training set)
 - Each record is characterized by a tuple (x,y) , where x is the attribute set and y is the class label
- ◆ x : attribute, predictor, independent variable, input
- ◆ y : class, response, dependent variable, output
- Task:
 - Learn a model that maps each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

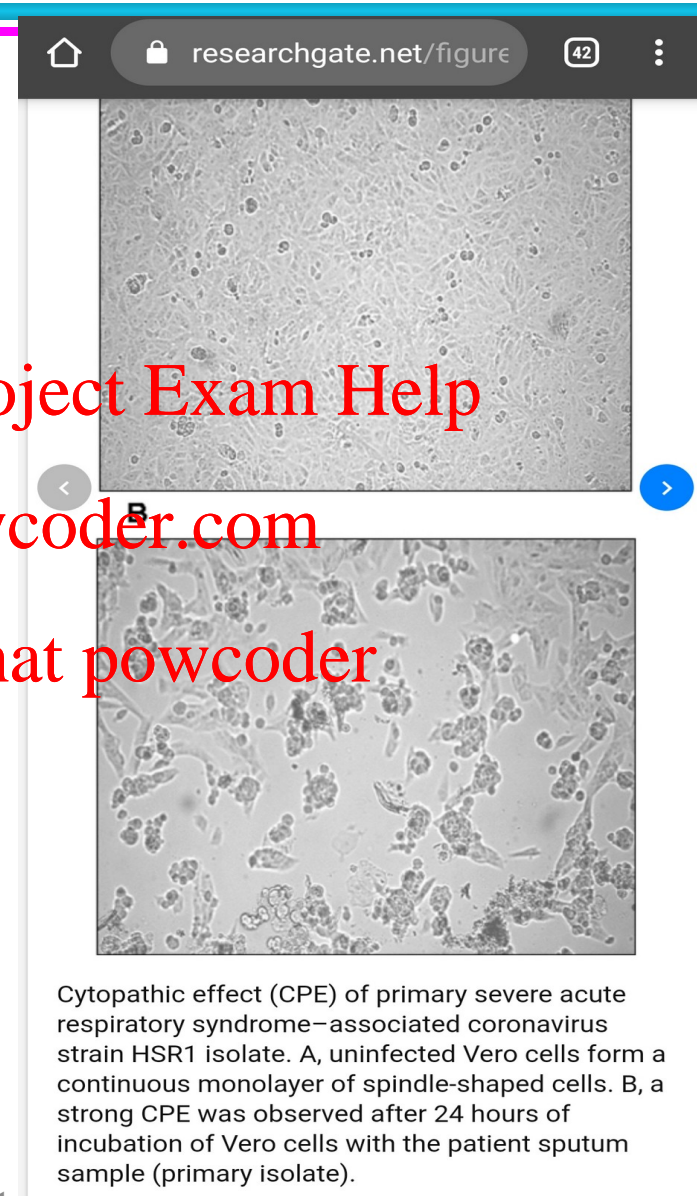
Examples of Classification Task (review)

Task	Attribute set, 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 MRI scans	malignant or benign cells
Cataloging galaxies	Features extracted from telescope images	Elliptical, spiral, or irregular-shaped galaxies

An example image classification problem

Table 1. Cytopathic Effects of Specific Viruses

Virus	Cytopathic Effect	Example
Paramyxovirus	Syncytium and faint basophilic cytoplasmic inclusion bodies	
Poxvirus	Pink eosinophilic cytoplasmic inclusion bodies (arrows) and cell swelling	
Herpesvirus	Cytoplasmic stranding (arrows) and nuclear inclusion bodies (dashed arrow)	
Adenovirus	Cell enlargement, rounding, and distinctive grape-like clusters	



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

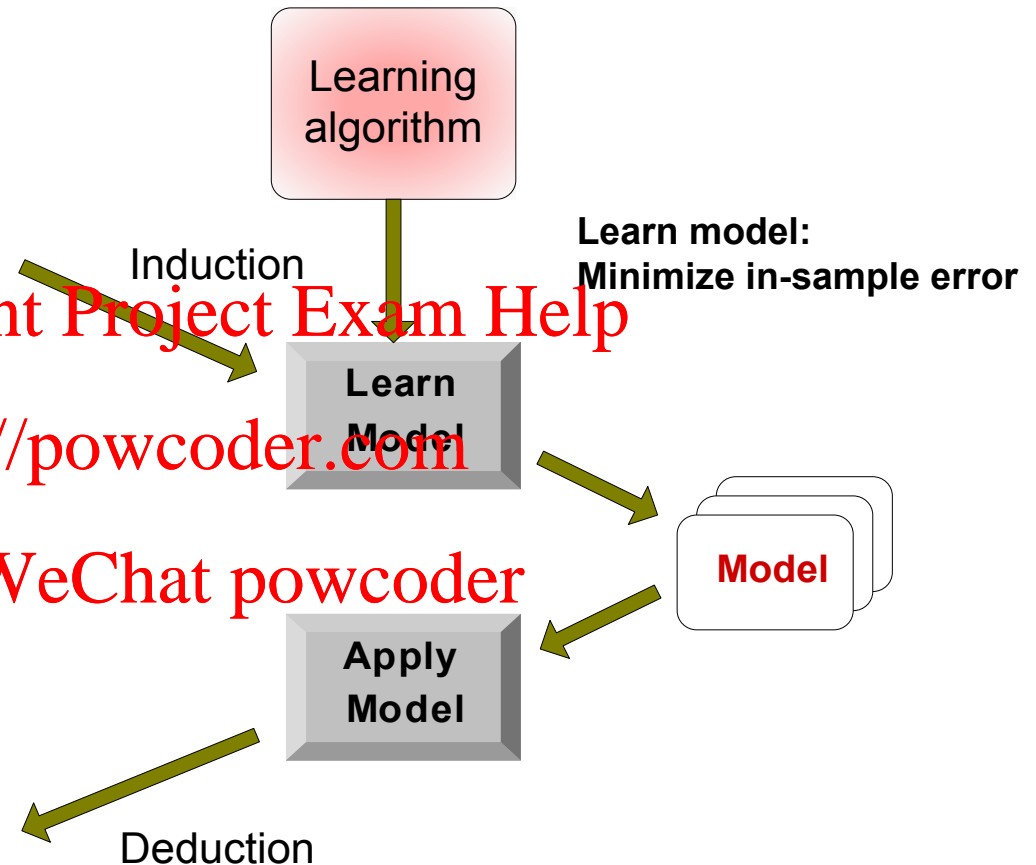
General Approach for Building Classification Model (review)

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



Classification Techniques

□ Base Classifiers

- **Decision Tree based Methods** (concept, prediction using a decision tree, and the training)
- Nearest-neighbor
- Neural Networks
- Deep Learning
- Naïve Bayes and Bayesian Belief Networks

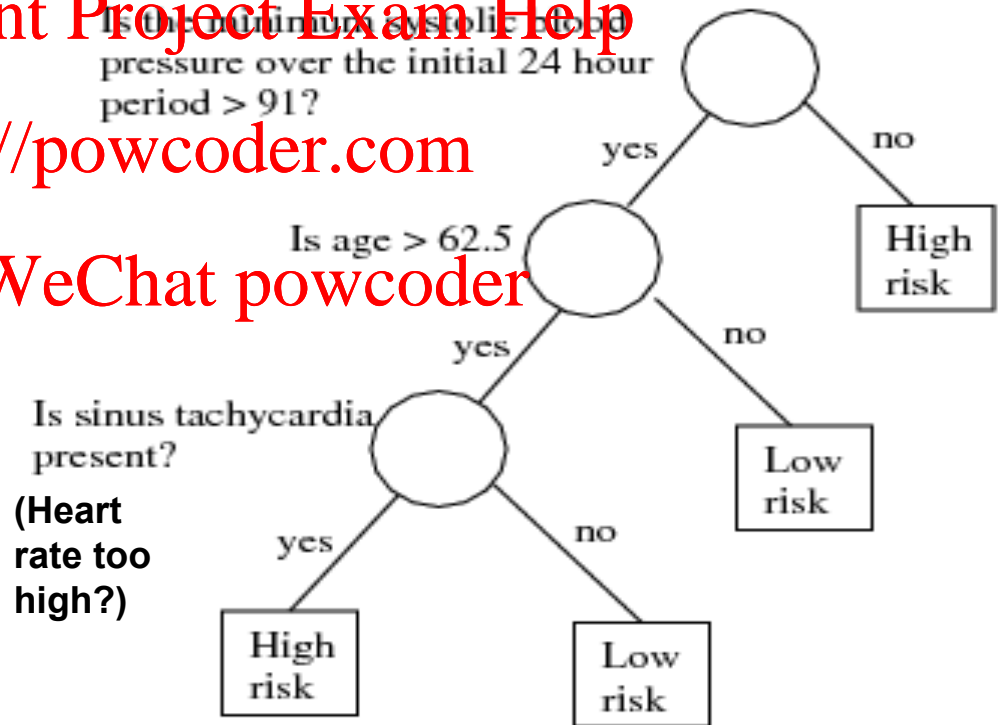
□ Ensemble Classifiers

- Boosting, Bagging, Random Forests

An example of classification using decision tree model

- Reference: *Classification and Regression Trees* by L. Breiman,
- J. H. Friedman, R. A. Olshen, and C. J. Stone, Chapman & Hall, 1984.
- A Medical Example: predictor, independent

- Predict high risk patients who will not survive at least 30 days on the basis of the initial 24-hour data
- 19 variables are measured during the first 24 hours, including blood pressure, age etc.

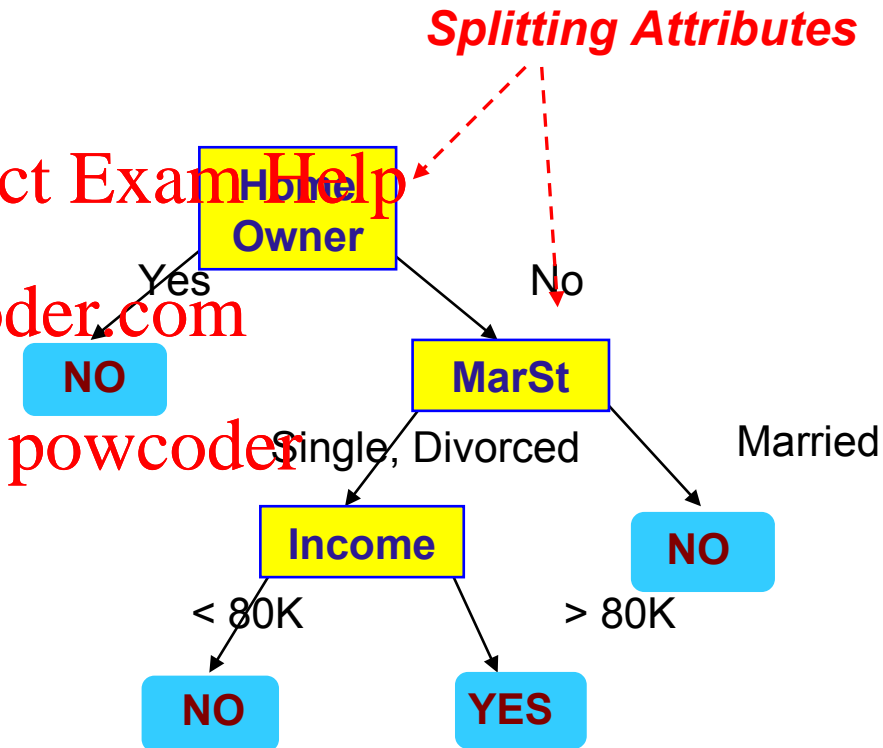


Toy Example of a Decision Tree

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	No
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

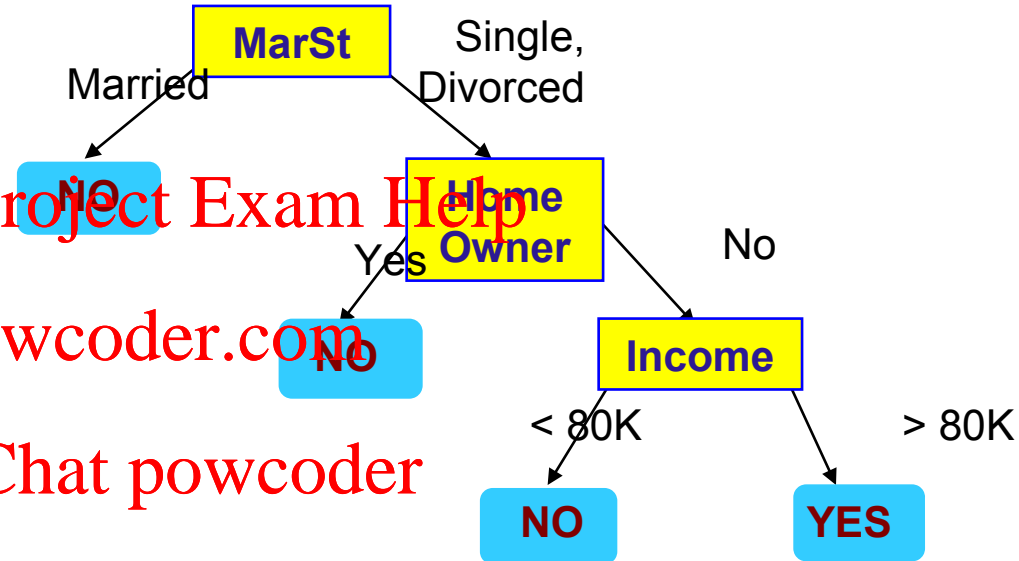


Model: Decision Tree

Another Example of Decision Tree

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



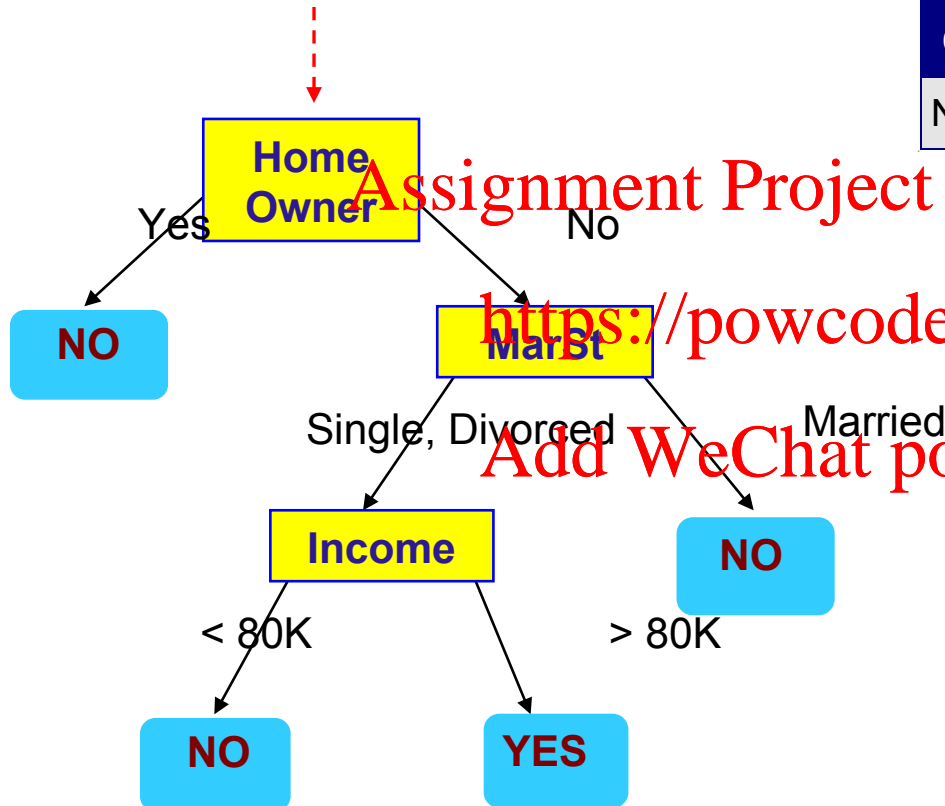
There could be more than one tree that fits the same data!

Apply Model to Test Data

Test Data

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

Start from the root of tree.



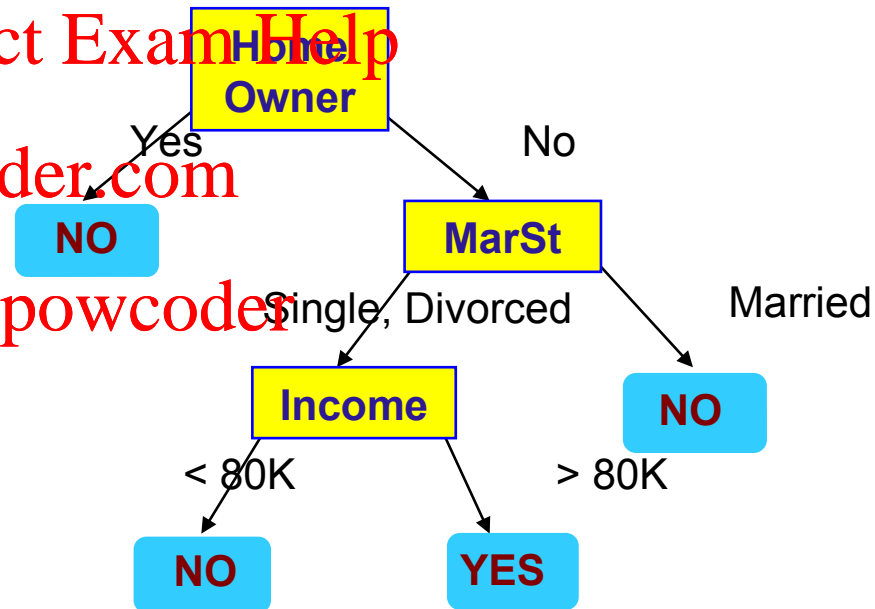
In-class exercise (to submit)

Please scan your solutions to the in-class exercise problem and upload it to Canvas. Please do so right after class. **Don't forget to write your name and ID.**

Problem: Derive a different decision tree using Income as the root.

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

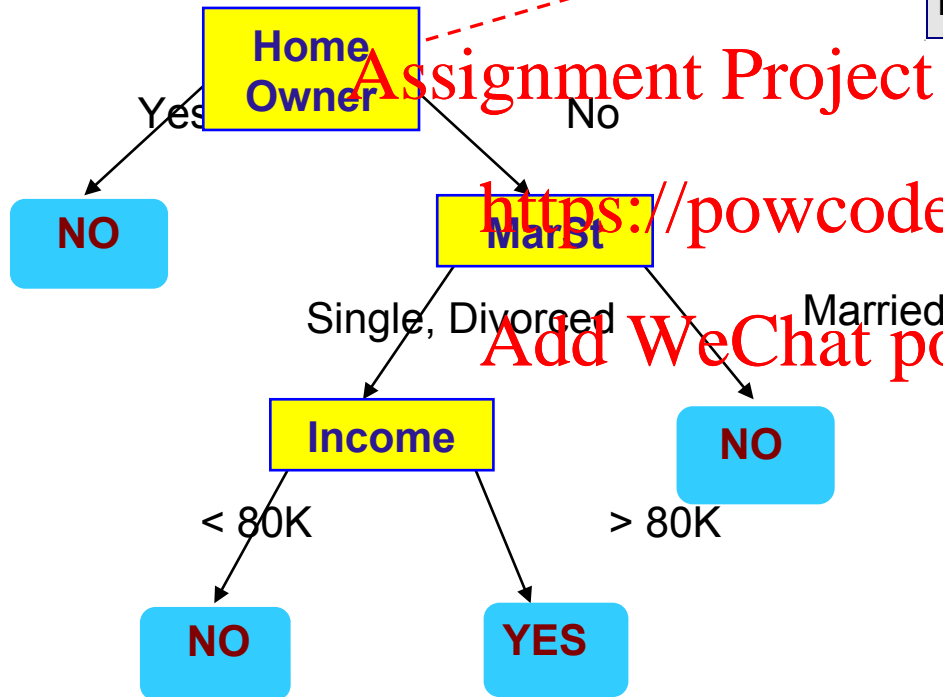


Model: Decision Tree

Apply Model to Test Data

Test Data

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



Assignment Project Exam Help

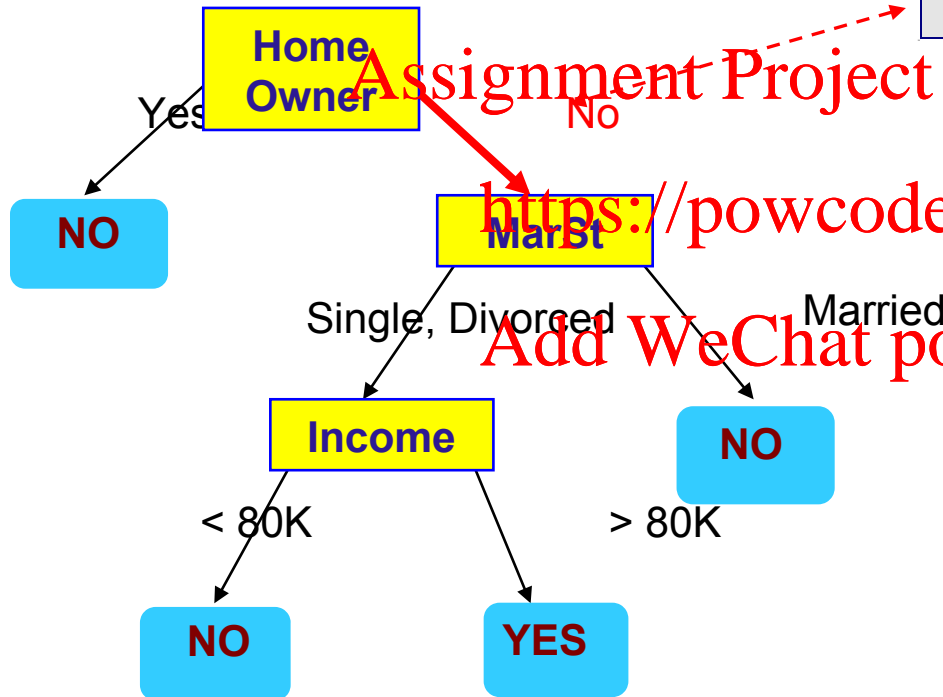
<https://powcoder.com>

Add WeChat powcoder

Apply Model to Test Data

Test Data

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



Assignment Project Exam Help

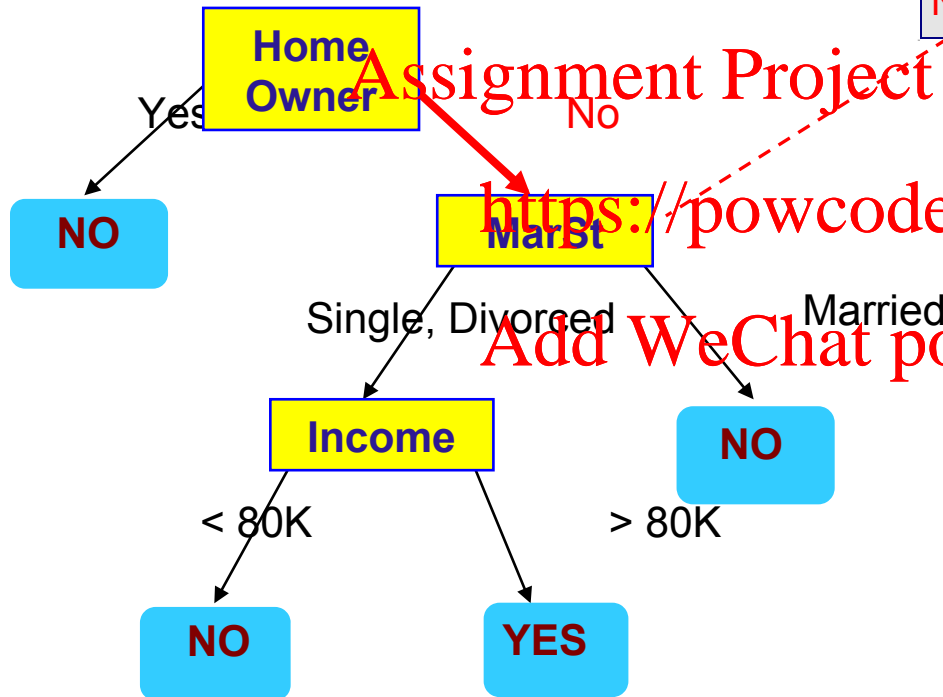
<https://powcoder.com>

Add WeChat powcoder

Apply Model to Test Data

Test Data

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



Assignment Project Exam Help

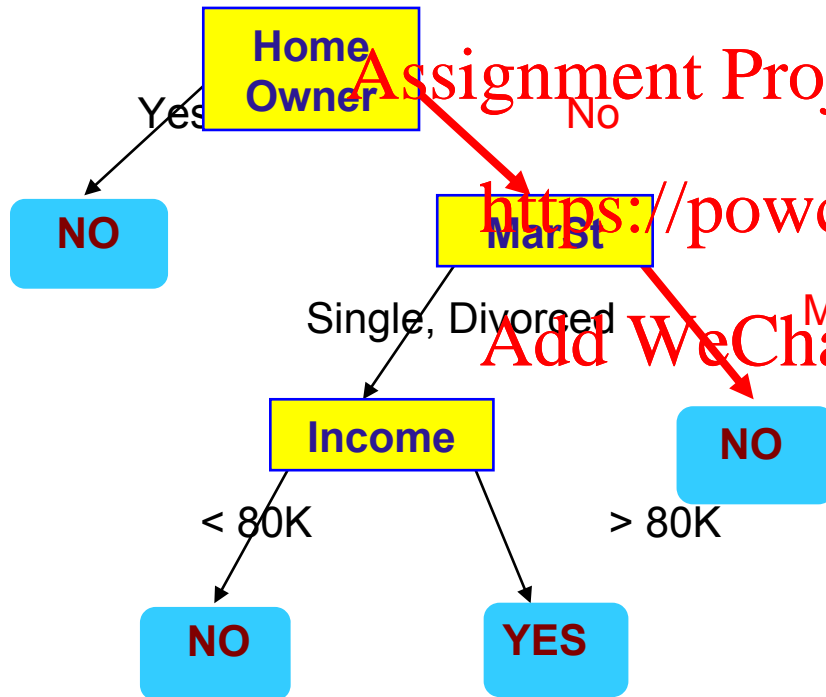
<https://powcoder.com>

Add WeChat powcoder

Apply Model to Test Data

Test Data

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



Assignment Project Exam Help

<https://powcoder.com>

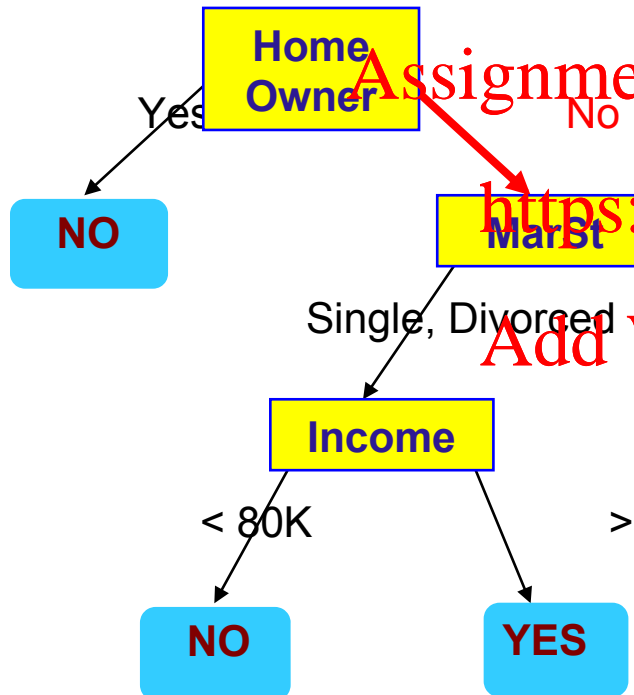
Add WeChat powcoder

How many tests do you need to reach a decision?

Apply Model to Test Data

Test Data

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



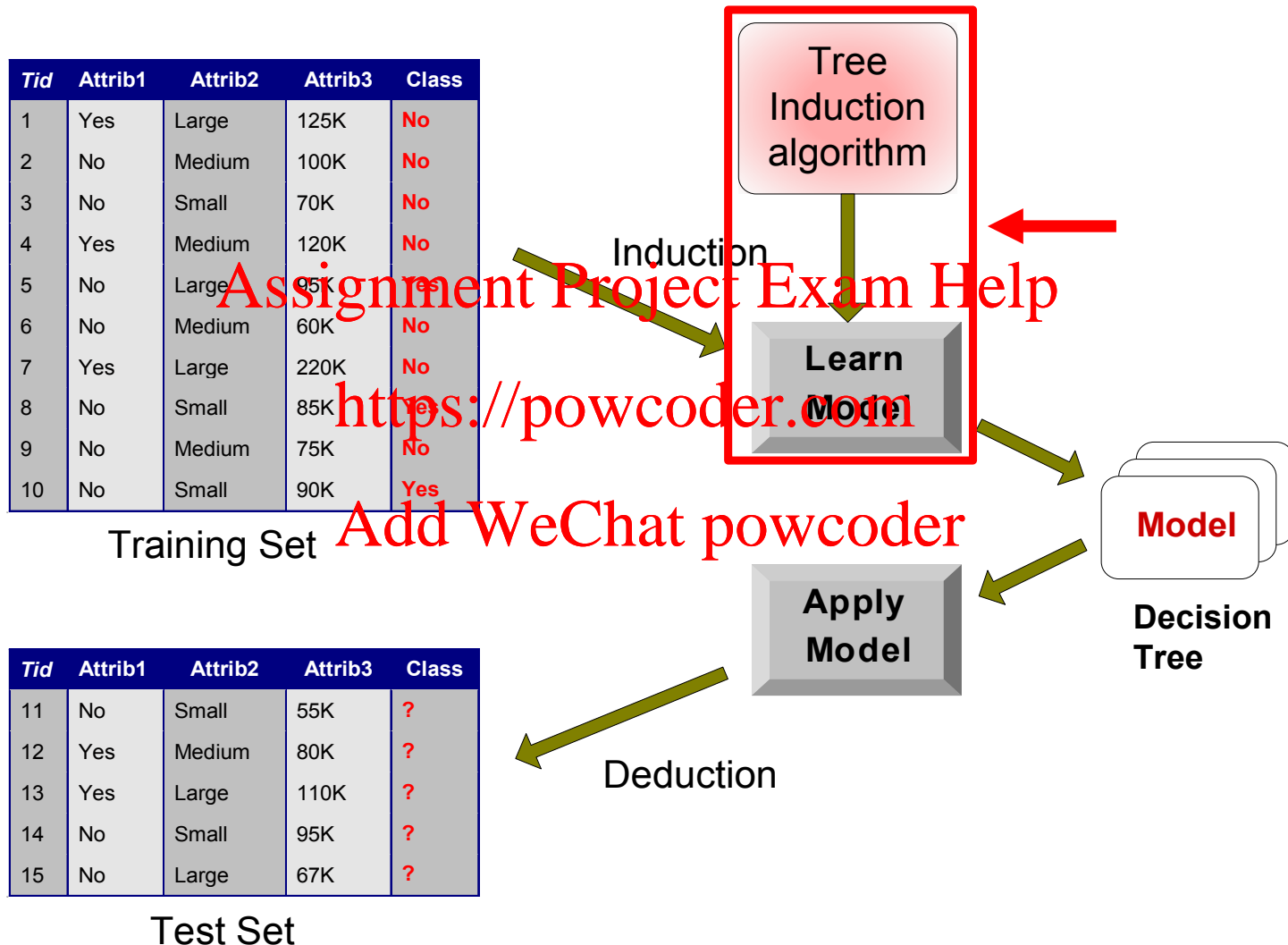
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Assign Defaulted to "No"

Decision Tree Classification Task



How do you define “optimal” decision tree

CONSTRUCTING OPTIMAL BINARY DECISION TREES IS NP-COMPLETE*

Laurent HYAFIL

IRIA – Laboria, 78150 Rocquencourt, France

and

Ronald L. RIVEST

Dept. of Electrical Engineering and Computer Science, M.I.T., Cambridge, Massachusetts 02139, USA

Received 7 November 1975, revised version received 26 January 1976

Add WeChat powcoder

We demonstrate that constructing optimal binary decision trees is an NP-complete problem, where an optimal tree is one which minimizes the expected number of tests required to identify the unknown object.

Decision Tree Induction

□ Many Algorithms:

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

Assignment Project Exam Help

<https://powcoder.com>

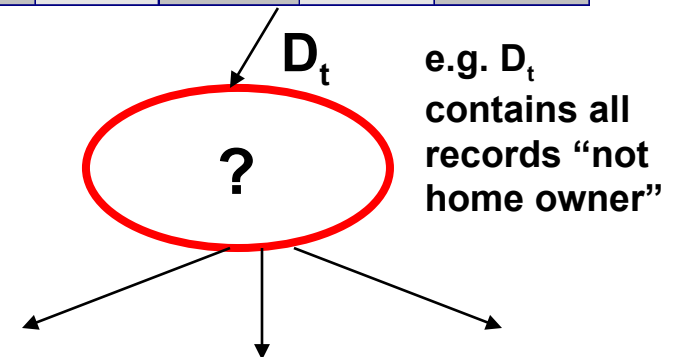
Add WeChat powcoder

General Structure of Hunt's Algorithm: how to create a node

- 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.

Defaulted borrower: fail to pay back

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)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

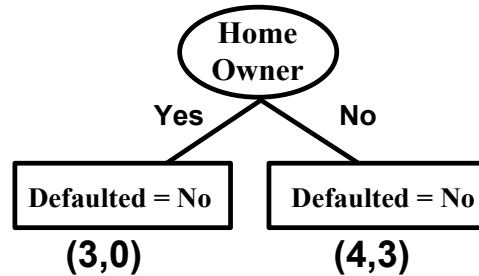
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Hunt's Algorithm

Defaulted = No

(7,3)

(a)



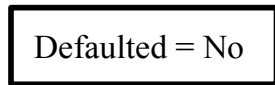
Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

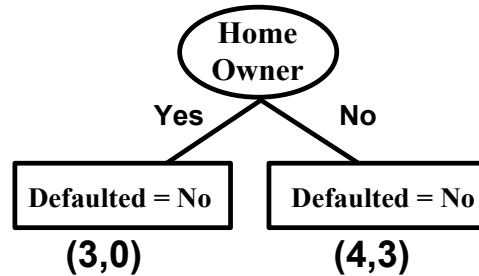
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Hunt's Algorithm



(7,3)

(a)

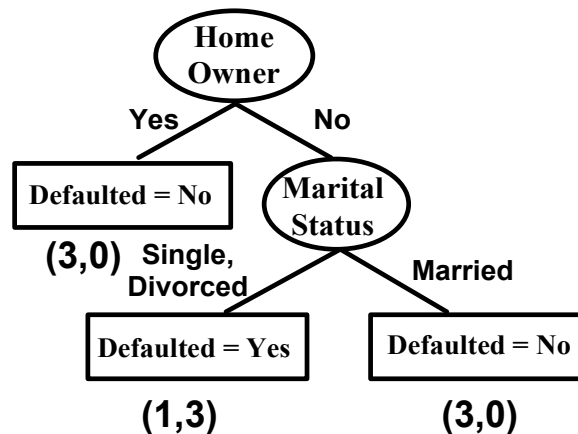


Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

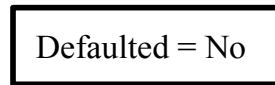
ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



(c)

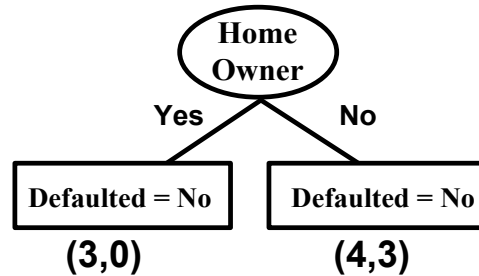
Hunt's Algorithm

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



(7,3)

(a)



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



(3,0)

Single, Divorced

Annual Income

< 80K

>= 80K

Defaulted = No

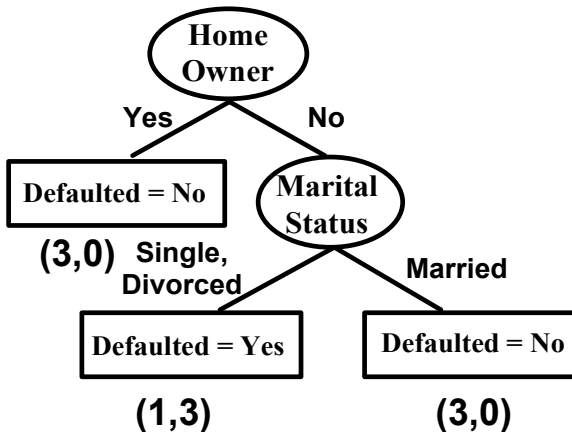
Defaulted = Yes

(1,0)

(0,3)

(3,0)

(d)



(3,0)

Single, Divorced

Married

Defaulted = Yes

Defaulted = No

(1,3)

(3,0)

(c)

Design Issues of Decision Tree Induction

- How should training records be split?
 - Method for specifying test condition
 - ◆ depending on attribute types
 - Measure for evaluating the goodness of a test condition <https://powcoder.com>

[Add WeChat powcoder](https://powcoder.com)

- 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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

□ Depends on number of ways to split

- 2-way split

- Multi-way split

Test Condition for Nominal Attributes

Multi-way split:

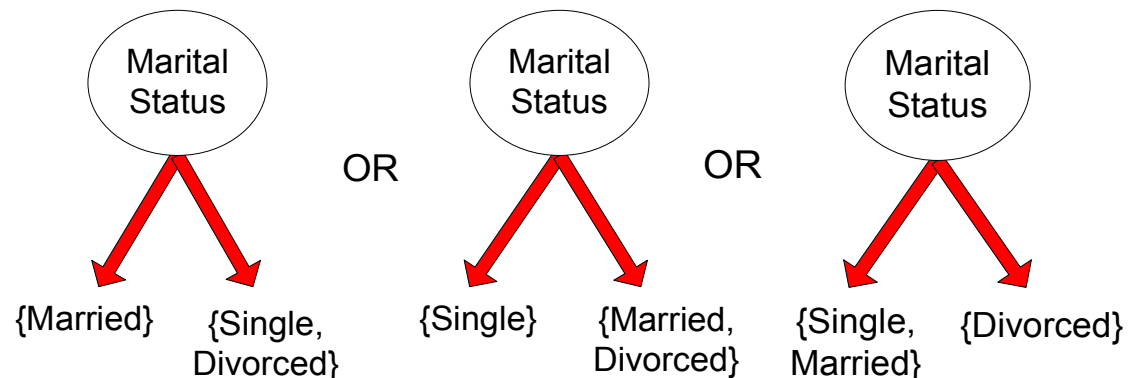
- Use as many partitions as distinct values.

Assignment Project Exam Help

<https://powcoder.com>

Binary split:

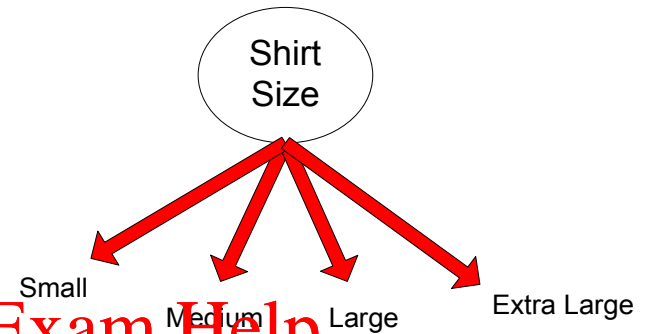
- Divides values into two subsets



Test Condition for Ordinal Attributes

Multi-way split:

- Use as many partitions as distinct values

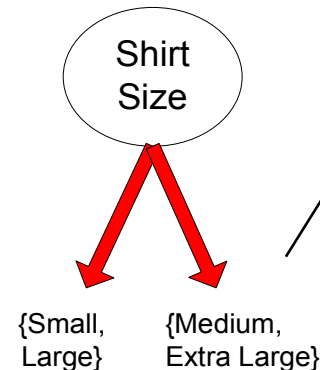
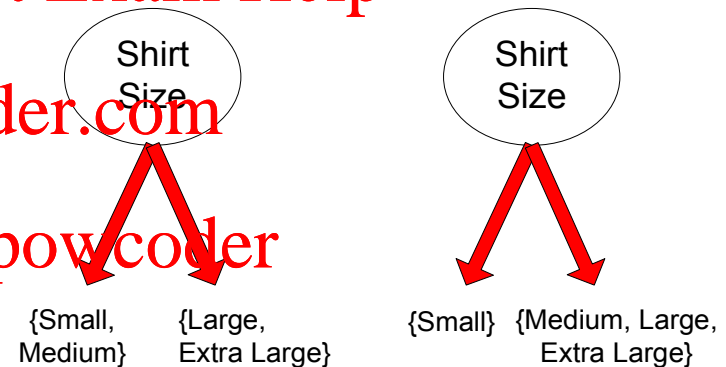


Assignment Project Exam Help

Binary split:

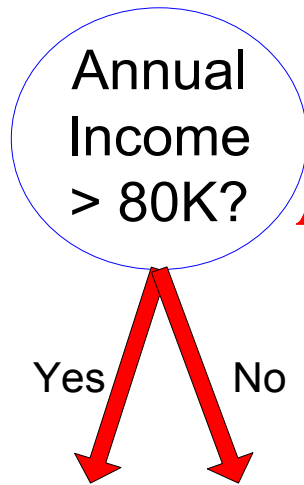
- Divides values into two subsets
- Preserve order property among attribute values

<https://powcoder.com>
Add WeChat powcoder



This grouping violates order property

Test Condition for Continuous Attributes

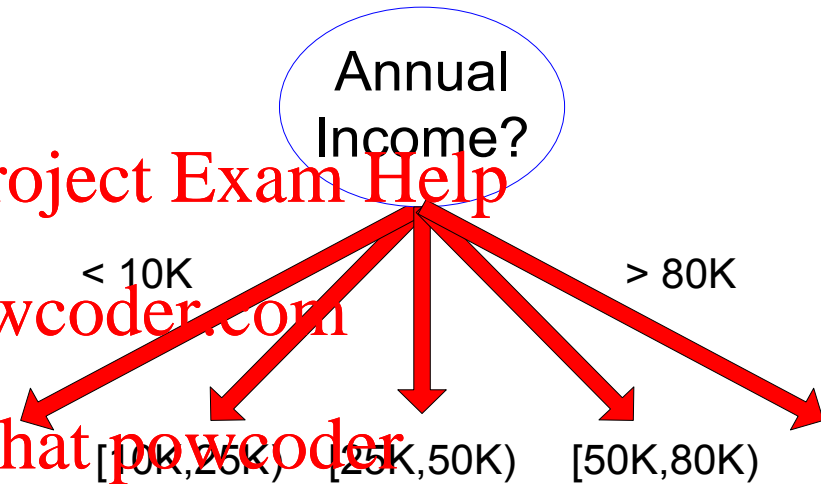


(i) Binary split

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



(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

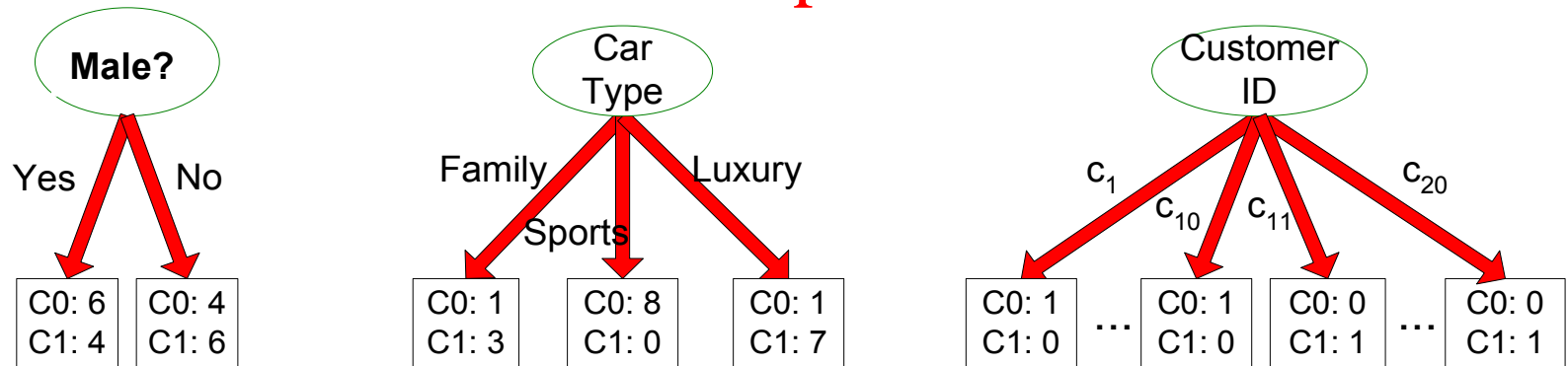
Before Splitting: 10 records of class 0,

10 records of class 1

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	F	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

<https://powcoder.com>

Add WeChat powcoder



Which test condition is the best?

How to determine the Best Split

- Greedy approach:
 - Nodes with **purser** class distribution are preferred

Assignment Project Exam Help

- Need a measure of node impurity:

Add WeChat powcoder

C0: 5
C1: 5

High degree of impurity

High uncertainty

C0: 9
C1: 1

Low degree of impurity

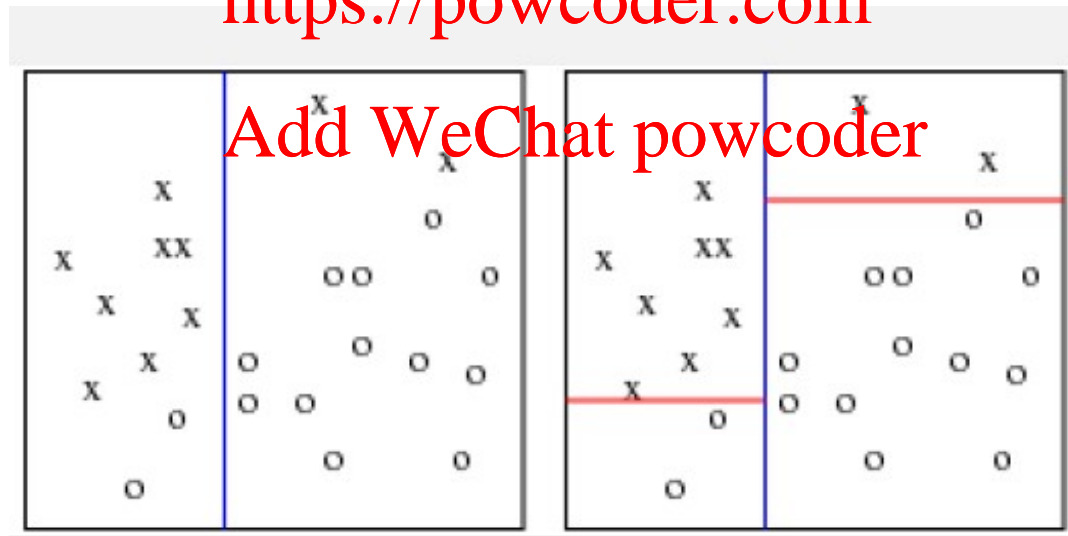
Low uncertainty

Goodness of Split

- The goodness of split is measured by an impurity function defined for each node.
- Intuitively, we want each leaf node to be “pure”, that is, one class dominates.

Assignment Project Exam Help

<https://powcoder.com>



The Impurity Function

- *Definition: An impurity function is a function φ defined on the set of all K -tuples of numbers (p_1, \dots, p_K) satisfying $p_j \geq 0, j = 1, \dots, K, \sum_{j=1}^K p_j = 1$ with the properties:*
Assignment Project Exam Help
- *1. φ is a maximum only at the point $(1/K, 1/K, \dots, 1/K)$*
https://powcoder.com
- *φ achieves its minimum only at the points $(1, 0, \dots, 0), (0, 1, 0, \dots, 0), \dots, (0, 0, \dots, 0, 1)$.*
Add WeChat powcoder
- *φ is a symmetric function of p_1, \dots, p_K , i.e., if you permute p_j , φ remains constant.*

Measures of Node Impurity

□ Gini Index

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

t is a node

Assignment Project Exam Help

□ Entropy

<https://powcoder.com>

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

Add WeChat powcoder

□ Misclassification error

**Entropy quantifies
uncertainty**

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

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 children
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

Before Splitting:

C0	N00
C1	N01

→ P

N00: number of records with label 0 at node N0

A?

Yes

No

Node N1

Node N2

B?

Yes

No

Node N3

Node N4

C0	N10
C1	N11

C0	N20
C1	N21

C0	N30
C1	N31

C0	N40
C1	N41

↓
M11

↓
M12

↓
M21

↓
M22

M1

M2

Gain = P – M1 vs P – M2

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Measure of Impurity: GINI

□ Gini Index for a given node t :

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

Assignment Project Exam Help

(NOTE: $p(j | t)$ is the relative frequency of class j at node t). <https://powcoder.com>

- 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

Measure of Impurity: GINI

□ Gini Index for a given node t :

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

Assignment Project Exam Help

(NOTE: $p(j | t)$ is the relative frequency of class j at node t).

<https://powcoder.com>

— 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

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

C1	0
C2	6

Assignment Project Exam Help

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

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

<https://powcoder.com>

C1	1
C2	5

Add WeChat powcoder

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

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

C1	2
C2	4

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

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

In-class exercise (to submit)

Consider the split of 3 classes in the following nodes. Compute their Gini index values:

C0: 5, C1: 5, C2: 5 [Assignment Project Exam Help](https://powcoder.com)

C0: 0, C1: 10, C2: 5 <https://powcoder.com>

C0: 15, C1: 0, C2: 0 [Add WeChat powcoder](https://powcoder.com)

If we have 20 data items with 4 labels, give a split with minimum and maximum Gini index values

Computing Gini Index for a Collection of Nodes

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

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

Assignment Project Exam Help

where, n_i = number of records at child i .

n = number of records at parent node p .

<https://powcoder.com>
Add WeChat powcoder

- 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

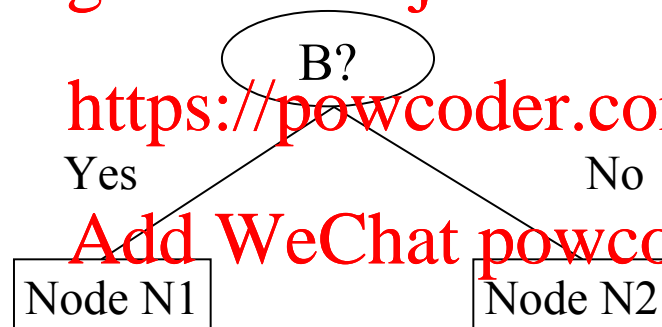
Binary Attributes: Computing GINI Index

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

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



	Parent
C1	7
C2	5
Gini = 0.486	

Gini(N1)

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

$$= 0.278$$

Gini(N2)

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

$$= 0.444$$

	N1	N2
C1	5	2
C2	1	4
Gini=0.361		

Weighted Gini of N1 N2

$$= 6/12 * 0.278 +$$

$$6/12 * 0.444$$

$$= 0.361$$

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

Categorical Attributes: Computing Gini Index

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

Assignment Project Exam Help

Multi-way split

<https://powcoder.com>

Two-way split

(find best partition of values)

	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	

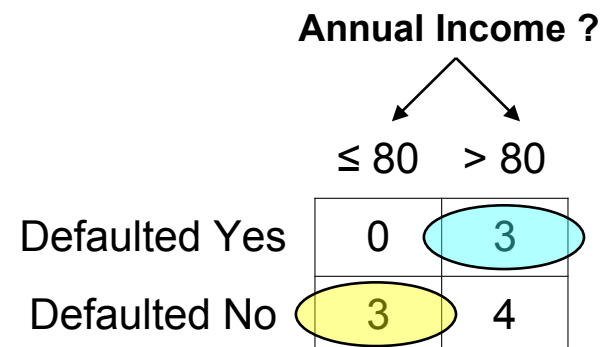
	CarType	
	{Sports}	{Family, Luxury}
C1	8	2
C2	0	10
Gini	0.167	

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, $A < v$ 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.

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



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

Assignment Project Exam Help

<https://powcoder.com>

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...

- 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

Assignment Project Exam Help

<https://powcoder.com>

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...

- 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

Assignment Project Exam Help

<https://powcoder.com>

	Cheat	No	No	No	Yes	Yes	Yes	No	No	No	No	
		Add WeChat powercoder										
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						
No					3	4						
Gini					0.343							

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

Assignment Project Exam Help

<https://powcoder.com>

	Cheat	No	No	No	Yes	Yes	Yes	No	No	No	No		
		Add WeChat powecoder											
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: 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

Assignment Project Exam Help

<https://powcoder.com>

	Cheat	No	No	No	Yes	Yes	Yes	No	No	No	No												
Sorted Values Split Positions		Add WeChat 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	1	2	2	1	3	0	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:

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

Assignment Project Exam Help

(NOTE: $p(j | t)$ is the relative frequency of class j at node t).

<https://powcoder.com>

- ◆ 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 quite similar to the GINI index computations

Supplementary reading

- Entropy
- There are many materials about entropy if you google it. I found the following is a good one. The author, Dr. Li, was a professor at Michigan State University. <https://powcoder.com>
- http://episte.math.ntu.edu.tw/articles/mm/mm_13_3_01/

Computing Entropy of a Single Node

$$Entropy(t) = -\sum_j p(j | t) \log_2 p(j | t)$$

C1	0
C2	6

Assignment Project Exam Help
<https://powcoder.com>

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

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

<https://powcoder.com>

C1	1
C2	5

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

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

C1	2
C2	4

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

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

Computing Information Gain After Splitting

Information Gain:

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

Assignment Project Exam Help

Parent Node, p is split into k partitions;

<https://powcoder.com>

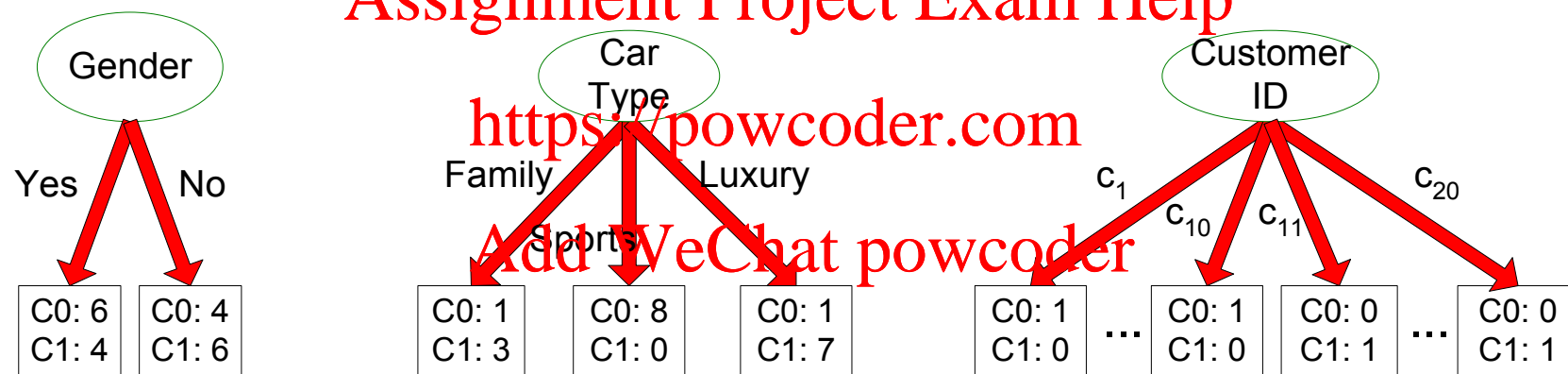
n_i is number of records in partition i

Add WeChat powcoder

- Choose the split that achieves most reduction (maximizes GAIN)
- Used in the ID3 and C4.5 decision tree algorithms

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:

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

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 algorithm
- Designed to overcome the disadvantage of Information Gain

Gain Ratio

□ Gain Ratio:

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

Assignment Project Exam Help

Parent Node, p is split into k partitions

n_i is the number of records in partition i

Add WeChat powcoder

	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

Measure of Impurity: Classification Error

- Classification error at a node t :

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

Assignment Project Exam Help

- Maximum ($1 - 1/n_c$) when records are equally distributed among all classes, implying least interesting information
- Minimum (0) when all records belong to one class, implying most interesting information

Computing Error of a Single Node

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

C1	0
C2	6

Assignment Project Exam Help

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

$$Error = 1 - \max(0, 1) = 1 - 1 = 0$$

<https://powcoder.com>

C1	1
C2	5

Add WeChat powder

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

$$Error = 1 - \max(1/6, 5/6) = 1 - 5/6 = 1/6$$

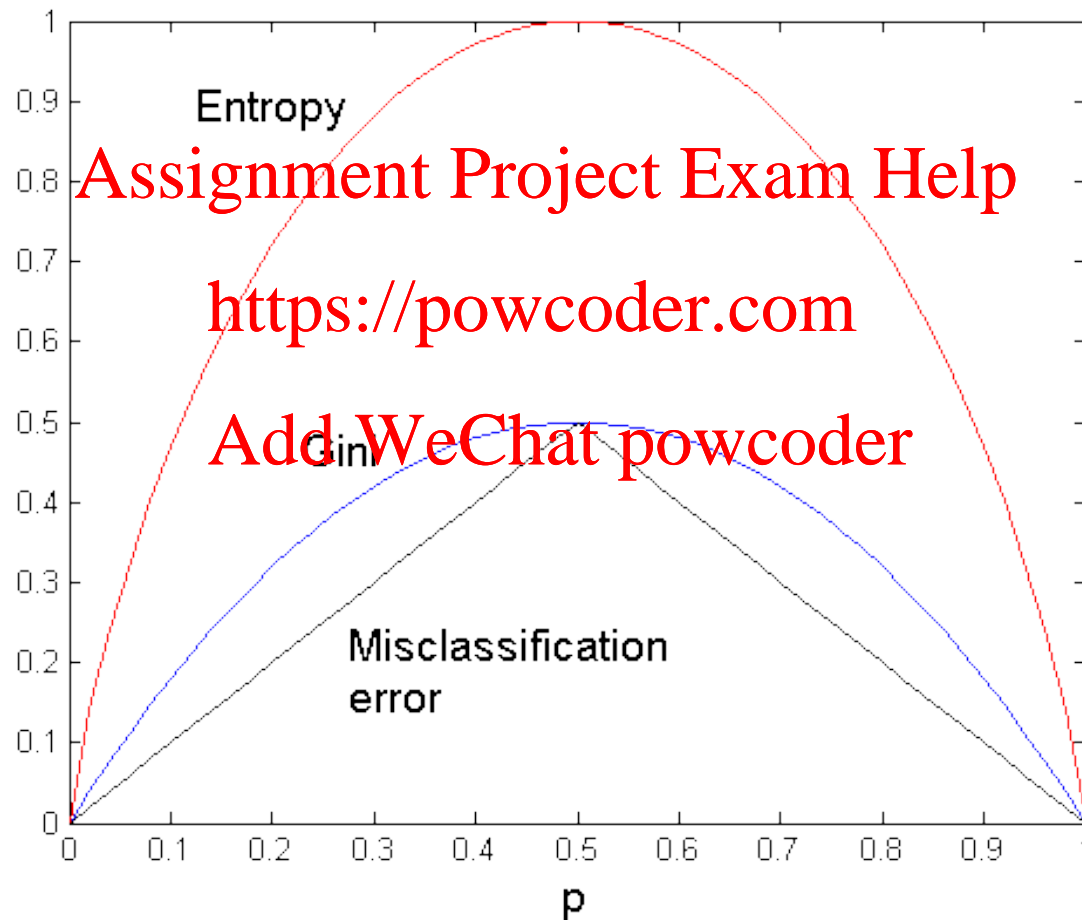
C1	2
C2	4

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

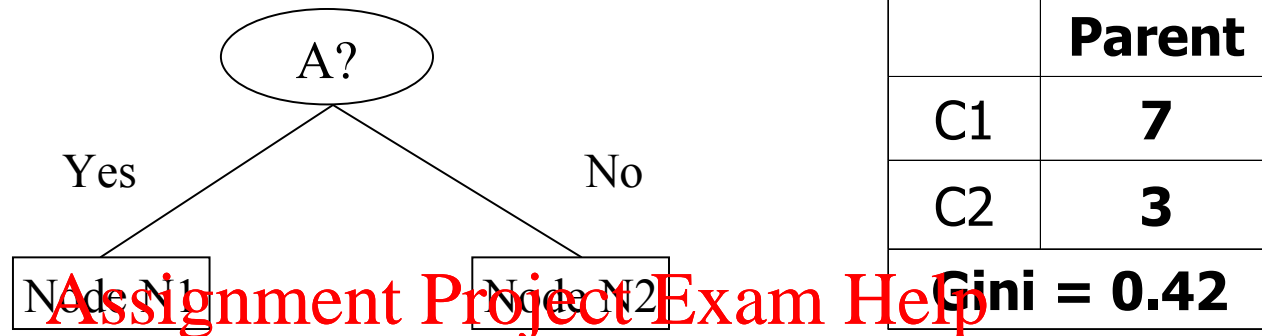
$$Error = 1 - \max(2/6, 4/6) = 1 - 4/6 = 1/3$$

Comparison among Impurity Measures

For a 2-class problem:



Misclassification Error vs Gini Index



<https://powcoder.com>

$$\begin{aligned} \text{Gini}(N1) &= 1 - (3/3)^2 - (0/3)^2 \\ &= 0 \end{aligned}$$

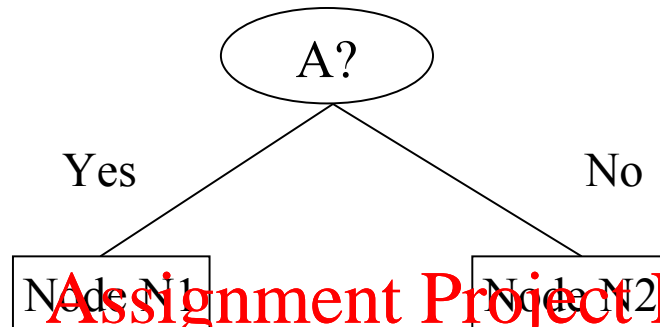
$$\begin{aligned} \text{Gini}(N2) &= 1 - (4/7)^2 - (3/7)^2 \\ &= 0.489 \end{aligned}$$

	N1	N2
C1	3	4
C2	0	3
Gini=0.342		

$$\begin{aligned} \text{Gini(Children)} &= 3/10 * 0 \\ &+ 7/10 * 0.489 \\ &= 0.342 \end{aligned}$$

Gini improves but error remains the same!!

Misclassification Error vs Gini Index



	Parent
C1	7
C2	3
Gini = 0.42	

Assignment Project Exam Help

<https://powcoder.com>

	N1	N2
C1	3	4
C2	0	3
Gini=0.342		

N1: $1 - \max(3/3, 0/3) = 0$

N2: $1 - \max(4/7, 3/7) = 3/7$

Weighted sum:

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

Add WeChat powcoder

	N1	N2
C1	3	4
C2	1	2
Gini=0.416		

Decision Tree Based Classification

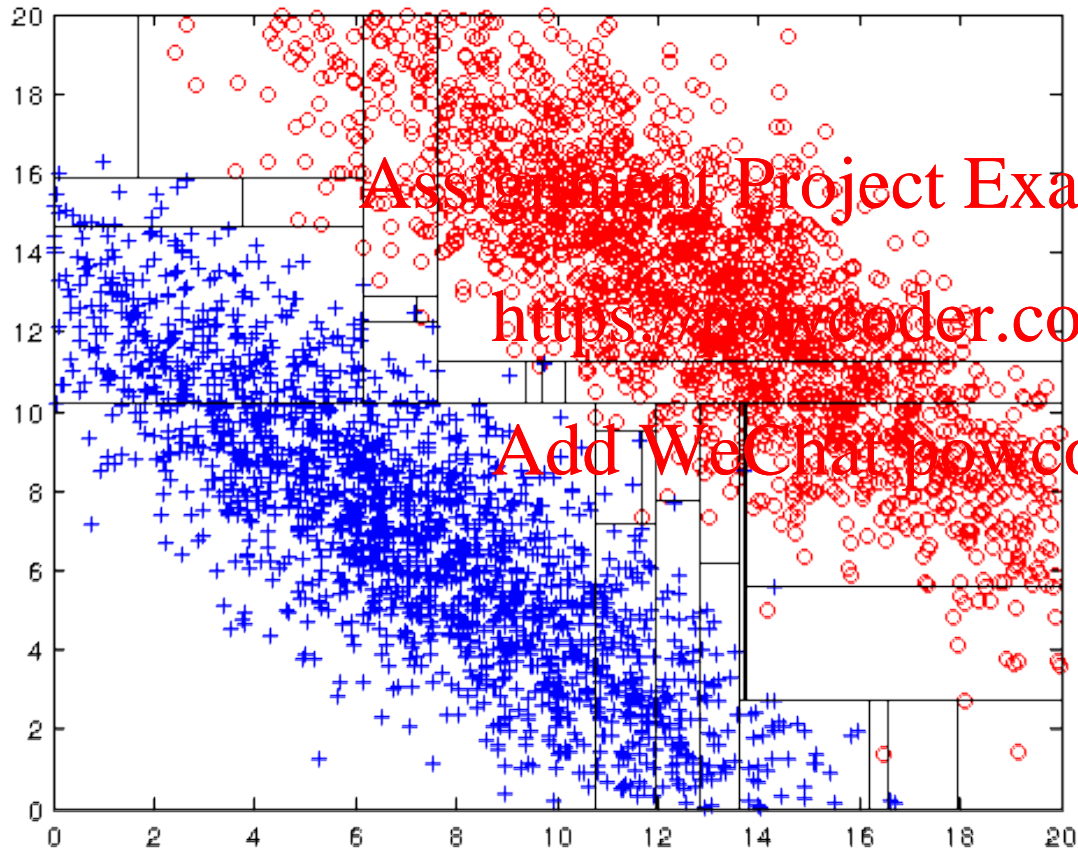
Advantages:

- 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)
- Can easily handle redundant or irrelevant attributes (unless the attributes are interacting)

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

Limitations of single attribute-based decision boundaries



Both **positive (+)** and **negative (o)** classes generated from skewed Gaussians with centers at (8,8) and (12,12) respectively.