

Data Mining

Study Assignment Set #3

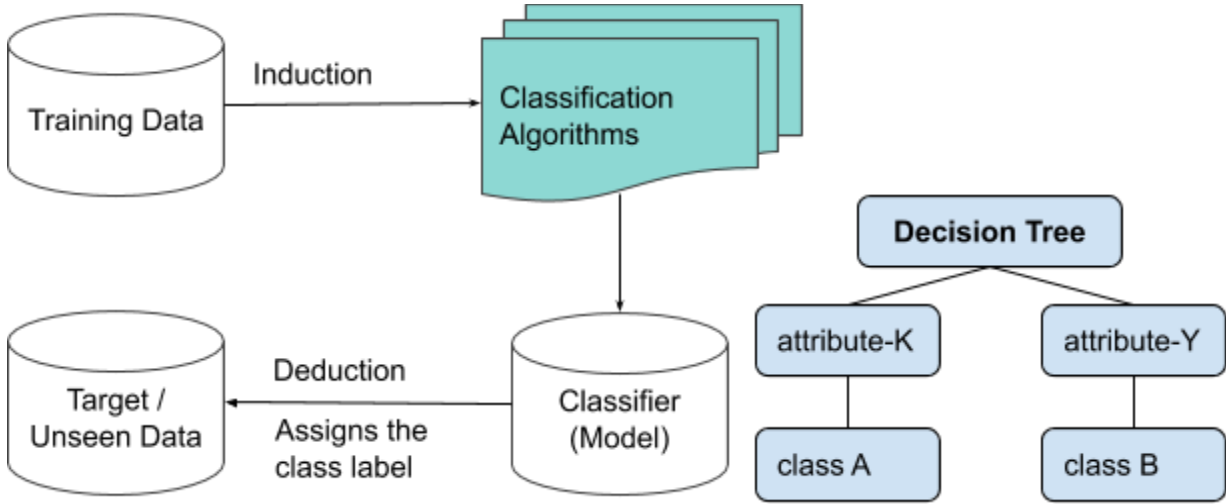
Prepared by: P V S Maruthi Rao | pvsmaruthirao@wilp.bits-pilani.ac.in

Reference Books:

Introduction to Data Mining by Tan P. N., Steinbach M and Kumar V.
Pearson Education, 2006

Data Mining: Concepts and Techniques, Second Edition by Jiawei Han and Micheline Kamber
Morgan Kaufmann Publishers, 2006

Topic: Classification of Data, Decision Trees, Gain Ratio

Classification of Data, Decision Trees	Question 1
<p>Learning objectives:</p> <ul style="list-style-type: none">- Basics of statistical learning with Decision Trees.- Decision Tree algorithm, and attribute selection methods.- Attribute selection by ‘Gain Ratio’- C4.5, a supervised learning algorithm proposed by J Ross Quinlan, uses attribute selection by Gain Ratio method. C4.5 is a successor of ID3, a supervised learning algorithm proposed by J Ross Quinlan, that uses ‘Information Gain’ as an attribute selection method. <p>Prerequisites:</p> <ul style="list-style-type: none">- Study Assignment Set #1 (Conditional probability)- Study Assignment Set #2 (Entropy, Information, Information Gain). <p>Basics of statistical learning learning with Decision Trees:</p>  <p>Some of the formulae are given as below.</p>	

$$Info(D) = - \sum_{i=1}^m p_i \log_2(p_i)$$

D: Training Data

m: Distinct values of the class label attribute.

p_i : non-zero probability that an **attribute tuple** in D belongs to a **class Y_i** and is estimated by $|Y_i, D| / |D|$

**** $P(Y_i | D) = P(Y_i, D) / P(D) = |Y_i, D| / |D|$ ****

[Some use C_i for class.]

How much more information would we still need (after partitioning) to arrive at an exact classification? Measure $Info_A(D)$ for attribute A as below.

$$Info_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j)$$

$Info_A(D)$ is the expected information required to classify a tuple from D based on the partition by the attribute A. The smaller the information (still) required, the greater the purity of the partition.

Gain (A) = Info (D) - $Info_A(D)$

Gain (A) is an indication of how much would be gained by branching on A (attribute A).

**** Branch on the attribute that gives highest gain ****

The C4.5 supervised learning algorithm applies a kind of **normalization to information gain** using a “**split information**” value defined as below.

$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \frac{|D_j|}{|D|}$$

v is a set of possible partitions on split attribute A.

The Gain Ratio:

$$GainRatio(A) = \frac{Gain(A)}{SplitInfo_A(D)}$$

A is an Attribute, D is a training data set.

**** Select the attribute with highest ‘Gain Ratio’ ****

If Split Info is approaching zero, the gain ratio is unstable. So a constraint is added to avoid this, whereby the information gain of the test selected must be large - at least as great as the average gain over all tests examined.

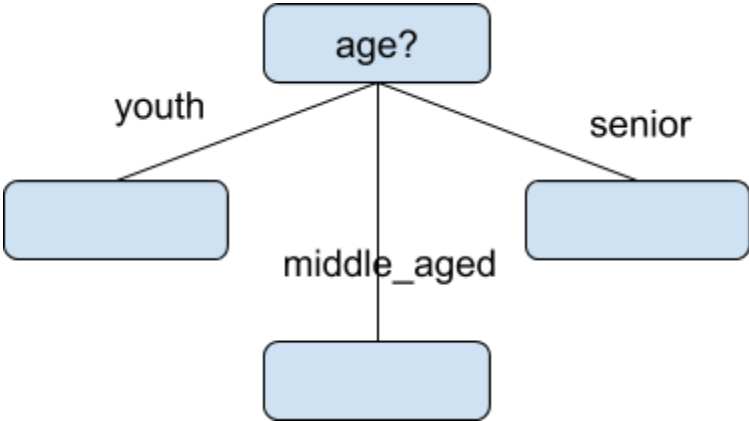
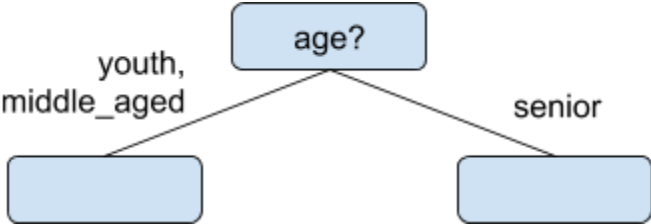
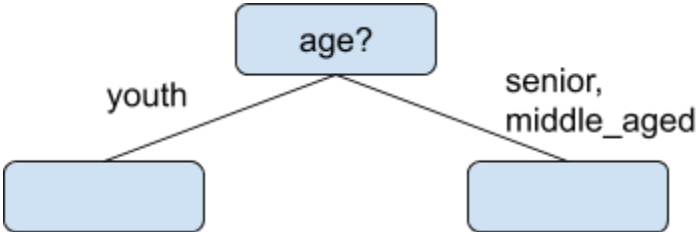
Note: When a calculation, system or subsystem behavior is tending towards unstable, then

design a constraint to avoid such instability.

Classification of Data	Question 2																																																																																				
Learning objectives: <ul style="list-style-type: none">- Basics of statistical learning with Decision Trees.- Decision Tree algorithm, and attribute selection methods.- Attribute selection by ‘Gain Ratio’- C4.5, a supervised learning algorithm proposed by J Ross Quinlan, uses attribute selection by Gain Ratio method. C4.5 is a successor of ID3, a supervised learning algorithm proposed by J Ross Quinlan, that uses ‘Information Gain’ as an attribute selection method.																																																																																					
Prerequisites: <ul style="list-style-type: none">- Study Assignment Set #1 (Conditional probability)- Study Assignment Set #2 (Entropy, Information, Information Gain).- Study Assignment Set #3 (Question 1).																																																																																					
<hr/>																																																																																					
<p>An online computer store uses a Decision Tree classifier with ‘Gain Ratio’ as a method of attribute selection method. Please see the Question #1 above for Gain Ratio.</p>																																																																																					
<p>Let X is a set of attributes of the registered user.</p> <p>X = {id, age, income, student, credit_rating}</p>																																																																																					
<p>Let Y is the class variable</p> <p>Y = buys_computer = {yes, no}</p>																																																																																					
<p>The training dataset, D, is as below.</p>																																																																																					
<table><tr><th>id</th><th>age</th><th>income</th><th>student</th><th>credit_rating</th><th>buys_computer</th></tr><tr><td>1</td><td>youth</td><td>high</td><td>no</td><td>fair</td><td>no</td></tr><tr><td>2</td><td>youth</td><td>high</td><td>no</td><td>excellent</td><td>no</td></tr><tr><td>3</td><td>middle_aged</td><td>high</td><td>no</td><td>fair</td><td>yes</td></tr><tr><td>4</td><td>senior</td><td>medium</td><td>no</td><td>fair</td><td>yes</td></tr><tr><td>5</td><td>senior</td><td>low</td><td>yes</td><td>fair</td><td>yes</td></tr><tr><td>6</td><td>senior</td><td>low</td><td>yes</td><td>excellent</td><td>no</td></tr><tr><td>7</td><td>middle_aged</td><td>low</td><td>yes</td><td>excellent</td><td>yes</td></tr><tr><td>8</td><td>youth</td><td>medium</td><td>no</td><td>fair</td><td>no</td></tr><tr><td>9</td><td>youth</td><td>low</td><td>yes</td><td>fair</td><td>yes</td></tr><tr><td>10</td><td>senior</td><td>medium</td><td>yes</td><td>fair</td><td>yes</td></tr><tr><td>11</td><td>youth</td><td>medium</td><td>yes</td><td>excellent</td><td>yes</td></tr><tr><td>12</td><td>middle_aged</td><td>medium</td><td>no</td><td>excellent</td><td>yes</td></tr><tr><td>13</td><td>middle_aged</td><td>high</td><td>yes</td><td>fair</td><td>yes</td></tr></table>	id	age	income	student	credit_rating	buys_computer	1	youth	high	no	fair	no	2	youth	high	no	excellent	no	3	middle_aged	high	no	fair	yes	4	senior	medium	no	fair	yes	5	senior	low	yes	fair	yes	6	senior	low	yes	excellent	no	7	middle_aged	low	yes	excellent	yes	8	youth	medium	no	fair	no	9	youth	low	yes	fair	yes	10	senior	medium	yes	fair	yes	11	youth	medium	yes	excellent	yes	12	middle_aged	medium	no	excellent	yes	13	middle_aged	high	yes	fair	yes	
id	age	income	student	credit_rating	buys_computer																																																																																
1	youth	high	no	fair	no																																																																																
2	youth	high	no	excellent	no																																																																																
3	middle_aged	high	no	fair	yes																																																																																
4	senior	medium	no	fair	yes																																																																																
5	senior	low	yes	fair	yes																																																																																
6	senior	low	yes	excellent	no																																																																																
7	middle_aged	low	yes	excellent	yes																																																																																
8	youth	medium	no	fair	no																																																																																
9	youth	low	yes	fair	yes																																																																																
10	senior	medium	yes	fair	yes																																																																																
11	youth	medium	yes	excellent	yes																																																																																
12	middle_aged	medium	no	excellent	yes																																																																																
13	middle_aged	high	yes	fair	yes																																																																																

14	senior	medium	no	excellent	no
----	--------	--------	----	-----------	----

Questions:

A	<p>Find</p> <ul style="list-style-type: none"> - Info (D) (Entropy of the D) for the class label attribute, Y = buys_computer = {yes, no}.
B	<ul style="list-style-type: none"> - $Info_{age}(D)$ - Gain (age) # Let age is chosen as a splitting attribute. - Gain Ratio (age = {youth, middle_aged, senior}). <p>Apply v = {youth, middle_aged, senior} in the SplitInfo formula.</p> <div>  <pre> graph TD A[age?] -- youth --> B[] A -- middle_aged --> C[] A -- senior --> D[] </pre> </div>
C	<ul style="list-style-type: none"> - Gain (age) # Let age is chosen as a splitting attribute. <p>Please note, the split is {{youth, middle_aged}, senior}. It means, we have to calculate Gain (age = {{youth, middle_aged}, senior}).</p> <ul style="list-style-type: none"> - Gain Ratio (age = {{youth, middle_aged}, senior}). <p>Apply v = {youth, middle_aged}, senior} in the SplitInfo formula.</p> <p>Please note age will be a binary split.</p> <div>  <pre> graph TD A[age?] -- "youth, middle_aged" --> B[] A -- senior --> C[] </pre> </div>
D	<ul style="list-style-type: none"> - Gain (age) # Let age is chosen as a splitting attribute. <p>Please note, the split is {youth, {middle_aged, senior}}. It means, we have to calculate Gain (age = {youth, {middle_aged, senior}}).</p> <ul style="list-style-type: none"> - Gain Ratio (age = {youth, {middle_aged, senior}}). <p>Apply v = {youth, {middle_aged, senior}} in the SplitInfo formula.</p> <p>Please note age will be a binary split.</p> <div>  <pre> graph TD A[age?] -- youth --> B[] A -- "senior, middle_aged" --> C[] </pre> </div>

Answers:

A	<p>D has a total 14 tuples (training data).</p> <p>m: Distinct values of the class label attribute = 2. buys_computer has two distinct value {yes, no}.</p>
---	---

	<p>p(buys_computer = yes D) = 9/14 p(buys_computer = no D) = 5/14</p> $Info(D) = - \sum_{i=1}^m p_i \log_2(p_i)$ <p>Info (D) = -9/14 log₂ (9/14) - 5/14 log₂ (5/14) = 0.940 bits. It is also known as Entropy of D.</p>
B	<p>Let's select age as a splitting attribute. D: Training data set. Class label Y: buys_computer = {yes, no}. A: age v_{age}: {youth, middle_aged, senior}</p> <p>Info (D) = -9/14 log₂ (9/14) - 5/14 log₂ (5/14) = 0.940 bits. It is also known as Entropy of D.</p> <p>Info (D_{age}) = (5/14 * Info(D_{youth})) + (4/14 * Info(D_{middle_aged})) + (5/14 * Info(D_{senior}))</p> <p>Gain (age) = Info (D) - Info (D_{age}) = 0.246 bits.</p> <p>(From the Assignment Set #2 OR you may calculate here)</p> $SplitInfo_A(D) = - \sum_{j=1}^v \frac{ D_j }{ D } \times \log_2 \frac{ D_j }{ D }$ $SplitInfo_A(D) = - \frac{ D_{youth} }{ D } \times \log_2 \frac{ D_{youth} }{ D } - \frac{ D_{middle_aged} }{ D } \times \log_2 \frac{ D_{middle_aged} }{ D } - \frac{ D_{senior} }{ D } \times \log_2 \frac{ D_{senior} }{ D }$ <p>SplitInfo_{age} (D) = -5/14log₂(5/14) - 4/14log₂(4/14) - 5/14log₂(5/14) = 1.5447 bits</p> $GainRatio(A) = \frac{Gain(A)}{SplitInfo_A(D)}$ <p>GainRaio (age) = Gain (age) / SplitInfo_{age} (D)</p> <p>GainRaio (age) = 0.246 bits / 1.5447 bits = 0.1559</p>
C	<p>Let's select age as a splitting attribute. D: Training data set. Class label Y: buys_computer = {yes, no}. A: age v_{age}: {{youth, middle_aged}, senior}</p> <p>Info (D) = -9/14 log₂ (9/14) - 5/14 log₂ (5/14) = 0.940 bits.</p>

It is also known as Entropy of D.				
From the data set D,				
	age = youth	age = middle_aged	age = senior	
buys_computer = yes	2	4	3	SUM = 9
buys_computer = no	3	0	2	SUM = 5
	SUM = 5	SUM = 4	SUM = 5	

$$\text{Info}(D_{\text{age}}) = (9/14 * \text{Info}(D_{\text{youth, middle_aged}})) + (5/14 * \text{Info}(D_{\text{senior}}))$$
$$= 9/14 (-6/9 \log_2 (6/9) - 3/9 \log_2 (3/9)) + 5/14 (-3/5 \log_2 (3/5) - 2/5 \log_2 (2/5))$$
$$= \mathbf{0.9371}$$

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}(D_{\text{age}})$$
$$= 0.940 - 0.9371$$
$$= \mathbf{0.0029}$$

$$\text{SplitInfo}_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \frac{|D_j|}{|D|}$$

Let's select age as a splitting attribute.

D: Training data set.

A: age

v_{age} : {{youth, middle_aged}, senior}

$$\text{SplitInfo}_A(D) = - \frac{|D_{\text{youth, middle_aged}}|}{|D|} \times \log_2 \frac{|D_{\text{youth, middle_aged}}|}{|D|} - \frac{|D_{\text{senior}}|}{|D|} \times \log_2 \frac{|D_{\text{senior}}|}{|D|}$$

$$\text{SplitInfo}_{\text{age}}(D) = -9/14 \log_2 (9/14) - 5/14 \log_2 (5/14)$$
$$= 0.9403 \text{ bits}$$

$$\text{GainRatio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}_A(D)}$$

$$\text{GainRaio}(\text{age}) = \text{Gain}(\text{age}) / \text{SplitInfo}_{\text{age}}(D)$$

$$\text{GainRaio}(\text{age}) = 0.0029 \text{ bits} / 0.9403 \text{ bits}$$
$$= \mathbf{0.0031}$$

| D | Follow the answer B) above and calculate the Gain Ratio for the spilt age = {youth, {middle_aged, senior}}. | | | |

Remarks	<p>In the above example, the root node is split on age.</p> <p>The tree gets added with new nodes or split partitions recursively.</p>
---------	--