

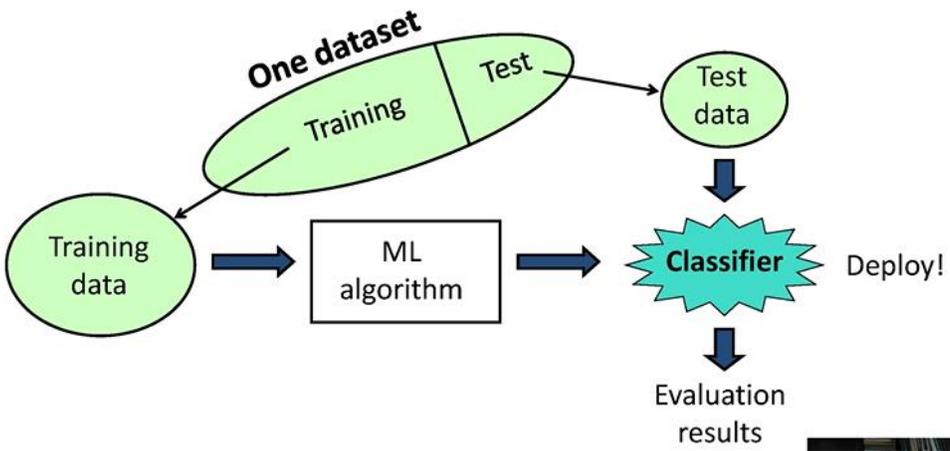
| | |
|---------------|--------------------------------------------------------------------------------------------------------------------------------|
| Name: |  UPES UNIVERSITY WITH A PURPOSE |
| Enrolment No: | |

UNIVERSITY OF PETROLEUM AND ENERGY STUDIES
End Semester Examination(Online) - July 2020

| | |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------|
| Course: Data Mining Program: MBA(BA) Course code: DSBA 7008 Instructions: Use Weka to solve questions wherever required. | Semester: II Time: 03 Hours Max. Marks: 100 |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------|

IMPORTANT INSTRUCTIONS

1. The student must write his/her name and enrolment no. in the space designated above.
2. The questions have to be answered in this MS Word document.
3. After attempting the questions in this document, the student has to upload this MS Word document on Blackboard.
4. Use Weka software wherever required and paste screen shot in document file for support of your answer

| | | Marks | CO |
|-----|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------|-----|
| Q1 | <p>A) Download and open the anneal dataset.</p> <ol style="list-style-type: none"> i. How many <i>attributes</i> does it have? ii. Apply the unsupervised attribute filter <i>RemoveUseless</i> and find how many <i>attributes</i> does the dataset have now? iii. Identify one of the attributes that was removed by clicking <i>Undo</i> and then <i>Apply</i>. Now figure out why it was removed. <p>B) Explain the below terms along with the process of Machine learning as shown in diagram:</p> <div style="text-align: center;">  </div> | 10+10 | CO2 |
| Q2. | <p>Classify the attribute ‘Type’ of glass.arff dataset using J48.</p> <p>a) No. of instances and Attributes</p> | (2+2+5+5+6=20) | CO2 |

| | | | |
|------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------|-------------------|
| | <ul style="list-style-type: none"> b) No of leaves and trees c) Write confusion matrix and explain it d) Write impact when Change minNumObj=15 e) Display Decision tree and interpret the model | | |
| <p>Q3.</p> | <p>Open the <i>diabetes.arff</i> dataset and answer the following:</p> <ul style="list-style-type: none"> a) Select <i>Percentage split</i> as test option and set <i>percentage for training</i> to 80%. How many instances will be used for training, and how many for testing? b) Select the J48 classifier (default options) and evaluate it with the following seed values (<i>More options</i>): 1, 2, 3, 4, 5 c) What are the minimum and maximum values for the number of incorrectly classified instances? What is the mean of the accuracy for these five seed values? <ul style="list-style-type: none"> i. 63.3% ii. 75.3% iii. 76.0% iv. 95.0% d) What is the standard deviation of the accuracy for these five seed values? <ul style="list-style-type: none"> i. 2.1 ii. 2.8 iii. 3.1 iv. 3.3 e) If you did the experiment of (b) with 10 different random seeds rather than 5, how would you expect this to affect the mean and standard deviation? <ul style="list-style-type: none"> i. They would both stay about the same. ii. The mean would be a bit bigger but the standard deviation would be about the same. iii. The mean would be about the same and the standard deviation would be a little smaller. iv. Both the mean and standard deviation would be a bit smaller. | <p>5X4=20</p> | <p>CO2</p> |
| <p>Q4.</p> | <ul style="list-style-type: none"> a) Open weather.normal.arff file. <ul style="list-style-type: none"> i. Write step to remove 3rd attribute ii. Write step to remove High values of Humidity attribute b) Open cpu.arff file. <ul style="list-style-type: none"> i. Run Non- linear regression. | <p>2.5X2+3 X5 =20)</p> | <p>CO3</p> |

| | | | |
|-----|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------|-----|
| | ii. Write all the linear equations with interpretation iii. Justify which one is better and why. | | |
| Q5. | <p>Answer the following based on below given file:</p> <pre>@relation weather.symbolic @attribute outlook {sunny, overcast, rainy} @attribute temperature {hot, mild, cool} @attribute humidity {high, normal} @attribute windy {TRUE, FALSE} @attribute play {yes, no} @data sunny,hot,high,FALSE,no sunny,hot,high,TRUE,no overcast,hot,high,FALSE,yes rainy,mild,high,FALSE,yes rainy,cool,normal,FALSE,yes rainy,cool,normal,TRUE,no overcast,cool,normal,TRUE,yes sunny,mild,high,FALSE,no sunny,cool,normal,FALSE,yes rainy,mild,normal,FALSE,yes sunny,mild,normal,TRUE,yes overcast,mild,high,TRUE,yes overcast,hot,normal,FALSE,yes rainy,mild,high,TRUE,no</pre> <p>a) Write the name of relation, attributes and number of instances. b) Name the class attribute and its labels. c) Write any two decisions list based on above data. d) Calculate number of possible instances. e) Draw decision tree based on above data.</p> | <p>(3+ 2+ 5+ 3+ 7=20)</p> | CO2 |

Answers