





### ← Go Back to Data Mining

## **:≡** Course Content

# **Project - Data Mining**

Submission type : File Upload

Due Date : Sep 01, 5:00 PM

Total Score : 60

Available from : Aug 13, 8:00 AM

Your Score : 49/60

## Description

^

Dear Participants,

Please find below the Data Mining Project instructions:

- You have to submit 2 files:
  - 1. Answer Report: In this, you need to submit all the answers to all the questions in a sequential manner. It should include the detailed explanation of the approach used, insights, inferences, all outputs of codes like graphs, tables etc. Your report should not be filled with codes. You will be evaluated based on the business report.
  - 2. Jupyter Notebook file: This is a must and will be used for reference while evaluating

# **Problem 1: Clustering**

A leading bank wants to develop a customer segmentation to give promotional offers to its customers. They collected a sample that summarizes the activities of users during the past few months. You are given the task to identify the segments based on credit card usage.

- **1.1** Read the data, do the necessary initial steps, and exploratory data analysis (Univariate, Bivariate, and multivariate analysis).
- 1.2 Do you think scaling is necessary for clustering in this case? Justify
- **1.3** Apply hierarchical clustering to scaled data. Identify the number of optimum clusters using Dendrogram and briefly describe them

- **1.4** Apply K-Means clustering on scaled data and determine optimum clusters. Apply elbow curve and silhouette score. Explain the results properly. Interpret and write inferences on the finalized clusters.
- **1.5** Describe cluster profiles for the clusters defined. Recommend different promotional strategies for different clusters.

Dataset for Problem 1: bank\_marketing\_part1\_Data.csv

### **Data Dictionary for Market Segmentation:**

- 1. spending: Amount spent by the customer per month (in 1000s)
- 2. advance\_payments: Amount paid by the customer in advance by cash (in 100s)
- 3. probability\_of\_full\_payment: Probability of payment done in full by the customer to the bank
- 4. current\_balance: Balance amount left in the account to make purchases (in 1000s)
- 5. credit\_limit: Limit of the amount in credit card (10000s)
- 6. min\_payment\_amt : minimum paid by the customer while making payments for purchases made monthly (in 100s)
- 7. max\_spent\_in\_single\_shopping: Maximum amount spent in one purchase (in 1000s)

### **Problem 2: CART-RF-ANN**

An Insurance firm providing tour insurance is facing higher claim frequency. The management decides to collect data from the past few years. You are assigned the task to make a model which predicts the claim status and provide recommendations to management. Use CART, RF & ANN and compare the models' performances in train and test sets.

- **2.1** Read the data, do the necessary initial steps, and exploratory data analysis (Univariate, Bivariate, and multivariate analysis).
- **2.2** Data Split: Split the data into test and train, build classification model CART, Random Forest, Artificial Neural Network
- **2.3** Performance Metrics: Comment and Check the performance of Predictions on Train and Test sets using Accuracy, Confusion Matrix, Plot ROC curve and get ROC\_AUC score, classification reports for each model.
- 2.4 Final Model: Compare all the models and write an inference which model is best/optimized.
- **2.5** Inference: Based on the whole Analysis, what are the business insights and recommendations

Dataset for Problem 2: insurance\_part2\_data-1.csv

#### Attribute Information:

- 1. Target: Claim Status (Claimed)
- 2. Code of tour firm (Agency\_Code)
- 3. Type of tour insurance firms (Type)
- 4. Distribution channel of tour insurance agencies (Channel)
- 5. Name of the tour insurance products (Product)
- 6. Duration of the tour (Duration)
- 7. Destination of the tour (Destination)
- 8. Amount of sales of tour insurance policies (Sales)
- 9. The commission received for tour insurance firm (Commission)
- 10. Age of insured (Age)

Important Note: Please reflect on all that you have learned while working on this project. This step is critical in cementing all your concepts and closing the loop. Please write down your thoughts here.

All the very best!

Regards,

**Program Office** 

Scoring guide (Rubric) - Project - Data Mi	ning Evaluated		^
Criteria	Ratings	Points	

Criteria Ratings **Points** 1.1 Read the data and do exploratory data analysis (3 pts). Describe the data briefly. Interpret the inferences for each (3 pts). Initial steps like head() .info(), Data Types, etc. Null value check. Distribution plots(histogram) or similar plots for the continuous columns. Box plots, Correlation plots. Appropriate plots for 6/6 categorical variables. Inferences on each plot. Summary stats, Skewness, Outliers proportion should be discussed, and inferences from above used plots should be there. There is no restriction on how the learner wishes to implement this but the code should be able to represent the correct output and inferences should be logical and correct. 1.2 Do you think scaling is necessary for clustering in this case? Justify The learner is expected to check and comment about the difference in scale of different features on the bases of appropriate measure for example 2/2 std dev, variance, etc. Should justify whether there is a necessity for

works.

scaling and which method is he/she

using to do the scaling. Can also comment on how that method

1.3 Apply hierarchical clustering to scaled data (3 pts). Identify the number of optimum clusters using Dendrogram and briefly describe them (4). Students are expected to apply hierarchical clustering. It can be obtained via Fclusters or Agglomerative Clustering. Report should talk about the used criterion, affinity and linkage. Report must contain a Dendrogram and a logical reason behind choosing the optimum number of clusters and Inferences on the dendrogram. Customer segmentation can be visualized using limited features or whole data but it should be clear, correct and logical. Use appropriate plots to visualize the clusters.

good effort but how optimum no. of clusters are chosen from dendogram to be explained on details. Customer Segments can be visualized using appropriate graphs( scatter plot)

1.4 Apply K-Means clustering on scaled data and determine optimum clusters (2 pts). Apply elbow curve and silhouette score (3 pts). Interpret the inferences from the model (2.5 pts). K-means clustering code application with different number of clusters. Calculation of WSS(inertia for each value of k) Elbow Method must be applied and visualized with different values of K. Reasoning behind the selection of the optimal value of K must be explained properly. Silhouette Score must be calculated for the same values of K taken above and commented on. Report must contain logical and correct explanations for choosing the optimum clusters using both elbow method and silhouette scores. Append cluster labels obtained from K-means clustering into the original data frame. Customer Segmentation can be visualized using appropriate graphs.

good effort but cluster profiles can be visualized using appropriate graphs( scatter plot)

1.5 Describe cluster profiles for the clusters defined (2.5 pts).

Recommend different promotional strategies for different clusters in context to the business problem inhand (2.5 pts). After adding the final clusters to the original dataframe, do the cluster profiling. Divide the data in the finalyzed groups and check their means. Explain each of the group briefly. There should be at least 3-4 Recommendations.

understandable and business specific, students should not give any technical suggestions. Full marks will only be allotted if the recommendations are correct and business specific. variable means. Students to explain the profiles and

suggest a mechanism to approach

each cluster. Any logical explanation is acceptable.

Recommendations should be easily

Criteria **Ratings Points** 2.1 Read the data and do exploratory data analysis (4 pts). Describe the data briefly. Interpret the inferences for each (2 pts). Initial steps like head() .info(), Data Types, etc. Null value check. Distribution plots(histogram) or similar plots for the continuous columns. Box plots, Correlation plots. Appropriate plots for categorical variables. Inferences on 6/6 each plot. Summary stats, Skewness, Outliers proportion should be discussed, and inferences from above used plots should be there. There is no restriction on how the learner wishes to implement this but the code should be able to represent the correct output and inferences should be logical and correct.

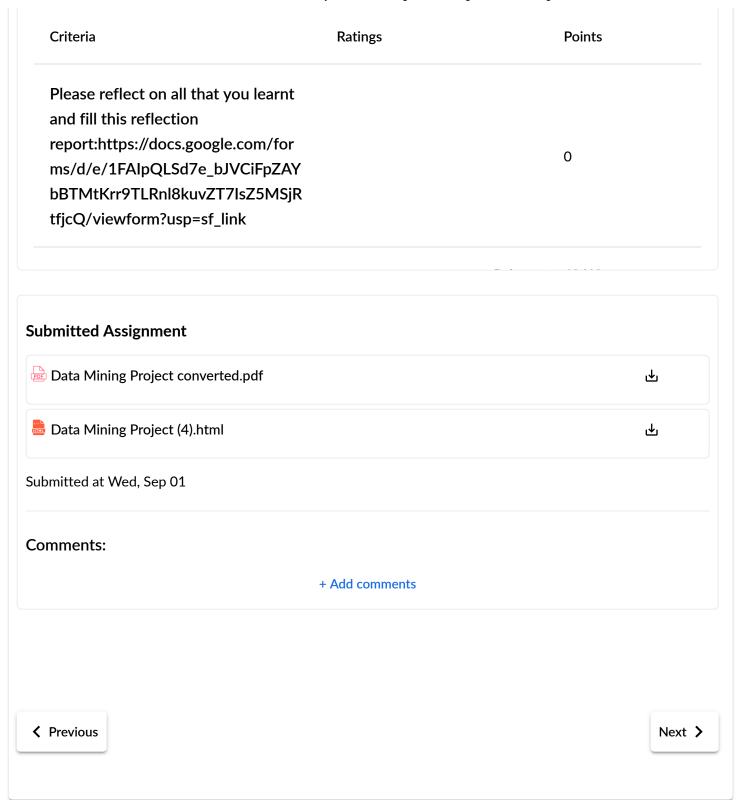
2.2 Data Split: Split the data into test and train(1 pts), build classification model CART (1.5 pts), Random Forest (1.5 pts), Artificial Neural Network(1.5 pts). Object data should be converted into categorical/numerical data to fit in the models. (pd.categorical().codes(), pd.get\_dummies(drop\_first=True)) Data split, ratio defined for the split, train-test split should be discussed. Any reasonable split is acceptable. Use of random state is mandatory. Successful implementation of each model. Logical reason behind the selection of different values for the parameters involved in each model. Apply grid search for each model and make models on best\_params. Feature importance for each model.

5.5/5.5

2.3 Performance Metrics: Check the performance of Predictions on Train and Test sets using Accuracy (1 pts), Confusion Matrix (2 pts), Plot ROC curve and get ROC\_AUC score for each model (2 pts), Make classification reports for each model. Write inferences on each model (2 pts). Calculate Train and Test Accuracies for each model. Comment on the validness of models (overfitting or underfitting) Build confusion matrix for each model. Comment on the positive class in hand. Must clearly show obs/pred in row/col Plot roc\_curve for each model. Calculate roc\_auc\_score for each model. Comment on the above calculated scores and plots. Build classification reports for each model. Comment on f1 score, precision and recall, which one is important here.

confusion matrix to be explained on details. Need to comment on Comment on f1 score, precision and recall, which one is important here and why

Criteria	Ratings	Points
2.4 Final Model - Compare all		
models on the basis of the		
performance metrics in a structured		
tabular manner (2.5 pts). Describe		
on which model is best/optimized		
(1.5 pts ). A table containing all the		
values of accuracies, precision,		
recall, auc_roc_score, f1 score.		4/4
Comparison between the different		
models(final) on the basis of above		
table values. After comparison		
which model suits the best for the		
problem in hand on the basis of		
different measures. Comment on		
the final model.		
2.5 Based on your analysis and		
working on the business problem,		
detail out appropriate insights and		
recommendations to help the		
management solve the business		
objective. There should be at least		
3-4 Recommendations and insights		45/45
in total. Recommendations should		4.5/4.5
be easily understandable and		
business specific, students should		
not give any technical suggestions.		
Full marks should only be allotted if		
the recommendations are correct		
and business specific.		
Quality of Business Report	good, but detailed description of answers to be written	5/6



Proprietary content. ©Great Learning. All Rights Reserved. Unauthorized use or distribution prohibited.

© 2021 All rights reserved

Privacy Terms of service Help