# A systematic evaluation of regressions and loss functions for the prediction of monetary value in RFM analysis

MSc Research Project
Data Analytics

## Shiva Prasad Aruva
Student ID: x22115188

School of Computing
National College of Ireland

Supervisor: Dr Giovani Estrada

| Student Name: | Shiva Prasad Aruva |
|---|---|
| Student ID: | x22115188 |
| Programme: | Data Analytics |
| Year: | 2023 |
| Module: | MSc Research Project |
| Supervisor: | Dr Giovani Estrada |
| Submission Due Date: | 4/12/2023 |
| Project Title: | A systematic evaluation of regressions and loss functions for the prediction of monetary value in RFM analysis |
| Word Count: | 789 |
| Page Count: | 6 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

**ALL** internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| Signature: | Shiva Prasad Aruva |
|---|---|
| Date: | 4th December 2023 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

| Attach a completed copy of this sheet to each project (including multiple copies). | ☐ |
|---|---|
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| Office Use Only | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# A systematic evaluation of regressions and loss functions for the prediction of monetary value in RFM analysis

Shiva Prasad Aruva

x22115188

# 1  Introduction

This configuration handbook contains Technical specifications and a Description of the hardware and software utilized in the project for Customer Segmentation. Follow the guidelines in this manual to reproduce the results and valuation of regressions and loss functions for the prediction of monetary value in RFM analysis.

# 2  System Specification

## 2.1  Hardware Requirements:

Table 1:

| Processer | Intel CORE i5 x64 |
|---|---|
| RAM | 8GB |
| DISK Storage | 1GB Approx |

## 2.2  Software Requirements:

Table 2:

| Operating System | Windows 11 |
|---|---|
| Programming Language | Python version 3.10 |
| Web-Broser | Google Chrome |
| Other Softwares | Google Colaborator, Excel |

# 3  Environment Embedding:

## 3.1  Google Colaborator Setup

This is the initial stage where we Run the Google Colab. With Colab, anyone can write and run any Python code through a browser, making it particularly useful for data

analysis, machine learning. And it has default run-time to Python version 3.10.



Figure 1: Google Collaborator

## 3.2 Data Collection

The information can be found for download in CSV format and is taken from the source: `https://archive.ics.uci.edu/dataset/352/online+retail`

## 3.3 Imported Libraries

These are the list of libraries used for this entire research project 2 3



Figure 2:

## 3.4 Data Pre-processing

The pre-processed dataset is uploaded to the Google Collab environment.As seen in the illustration, A selection of the data cleaning procedures is displayed in the figures below 4 , for later usage.

Figure 3:

The RFM score is computed  5 for every client in the "retail data" data frame by this line of code. The reliability, frequency, and monetary worth of a client are measured by the RFM score.



Figure 4:

## 3.5    Plotting Elbow meathod

Here We Plot the inertia against the number of clusters is a technique known as the elbow method  6, which helps determine the ideal number of clusters. The within-cluster variation is measured by the inertia, which gets smaller as the number of clusters rises. The elbow's location on the plot indicates the ideal number of clusters.

Figure 5:

For k-values between 2 and 40, you are charting the inertia versus the number of clusters in your code. This is the ideal amount of clusters for your data, as the elbow seems to be around k=5.



Figure 6:

## 3.6 Evaluating the performance

The figure 7 below evaluates the performance of k-means clustering in predicting the Monetary Value variable.



Figure 7:

# 4 Performance of various machine learning algorithms using Loss Functions

Here The performance of various machine learning algorithms 8 is evaluated using the Loss functions for the negative mean absolute error, negative mean squared error, and negative median absolute error. In this, we Always consider the Lower values are better suggested fit for the model.

# 5 Defining Neural Network For Regression

In this below figure 9 we Define a neural network model for regression and evaluate its performance using cross-validation.The neural network's efficiency can be assessed using the obtained scores.

Figure 8:



Figure 9: