

PES University, Bangalore

Established under the Karnataka Act No. 16 of 2013

UE20CS312 - Data Analytics

Worksheet 1a - Part 2: EDA with R | ANOVA

Harshith Mohan Kumar - harshithmohankumar@pesu.pes.edu

Yashas Kadambi - yashasks@pesu.pes.edu

Nishanth M S - nishanthmsathish.23@gmail.com

Anushka Hebbar - anushkahebbar@pesu.pes.edu

Prerequisites

To download the data required for this worksheet, visit [this Github link](#). This worksheet has two parts, the first focuses on the basics of dealing with data and exploratory data analysis using R. The second deals with ANOVA. To help guide you through the worksheet, here are a few resources:

- Revise how to deal with DataFrames in R [here](#).
- [This online book](#) has everything you need to get started with visualizations in R.
- Check out [this](#) resource for an excellent deep-dive of visualizations using the `ggplot2` library (**optional**).
- The following are resources to learn about ANOVA:
 - [Anova in Python](#)
 - [Anova in R](#)

Part I. Exploratory Data Analysis with R

Book Club Marketing Dataset

Charles Book Club (CBC) is a book club that has an active database of 500,000 subscribers. The organization sends out monthly mailings to its database of members with the latest promotional offerings. Its marketing team would like to see if customer data can be used to reduce the cost of marketing activities to improve the profitability of their marketing operations. For an initial pilot of a predictive analytics solution, CBC decided to focus on its strongest customers and run a marketing test for a new book release of *'The Art History of Florence'*.

The dataset provided consists of information about customer purchases CBC has as its disposal after conducting the marketing test. Use the `CharlesBookClubDataset.csv` for Part I of the worksheet. This data was adapted from a famous business database called the 'Charles Book Club', dealt with in more detail in a case study from the 'Data Mining for Business Analytics' book.

Data Dictionary

ID#: Customer Identification number

Gender: Male, Female

M: Monetary - Total money spent on books

R: Recency - Months since last purchase

F: Frequency - Total number of purchases
FirstPurch: Months since first purchase
ChildBks: Number of purchases from category of child books |
YouthBks: Number of purchases from category of youth books
CookBks: Number of purchases from category of cook books
DoItYBks: Number of purchases from category of DIY books
RefBks: Number of purchases from category of reference books
ArtBks: Number of purchases from category of art books
GeoBks: Number of purchases from category of geography books
ItalCook: Number of purchases of book title 'Secrets of Italian Cooking'
ItalAtlas: Number of purchases of book title 'Historical Atlas of Italy'
ItalArt: Number of purchases of book title 'Italian Art'
Related Purchase: Number of related books purchased
Florence: = 1 if 'Art History of Florence' was purchased; = 0 if not

Loading the Dataset

Use the following commands to load the dataset from CSV format and get a high-level overview of its fields:

```
library(tidyverse)
cbc_df <- read_csv(path_to_csv)
head(cbc_df)
```

Points

The problems for this part of the worksheet are for a total of 8 points, with a non-uniform weightage.

- *Problem 1* : 1 point
- *Problem 2* : 2 points
- *Problem 3* : 2 points
- *Problem 4.1* : 1 point
- *Problem 4.2* : 1 point
- *Problem 4.3* : 1 point

Problems

Problem 1 (1 point)

Generate an understanding of the dataset via a summary of its features. Find the count, missing count, minimum, 1st quartile, median, mean, 3rd quartile, max and standard deviation of all relevant columns. Separately, print the total number of missing values in each column.

Problem 2 (2 points)

Replace missing values within the Recency, Frequency, and Monetary features with suitable values. Explain your reasoning behind the method of substitution used. *Hint*: Try plotting the distribution of the values in each feature using the `hist` function. Think about how to best deal with data imputation. Also, plot the distribution of feature values after imputation.

Problem 3 (2 points)

Discretize the continuous values of Monetary, Recency, and Frequency into appropriate bins, and create three new columns `Mcode`, `Rcode` and `Fcode` respectively, for the discretized values. Explicitly mention the number of bins used and explain the choice for the bin size. Print out the summary of the newly created columns. *Hint*: Use the `cut` function to break on preset breakpoints. What are the most optimum breakpoints you can choose? Try to think of a statistical function that provides these breakpoints for optimum binning.

Problem 4

The marketing team heavily relies on the RFM variables of the recency of last purchase, total number of purchases, and total money spent on purchases to gauge the health of the members of the book club. Increases in either the frequency of purchases or monetary spend and decreases in time since last purchase across the customer base, will intuitively lead to more sales for the business.

4.1 Bar Graphs (1 point) Create and visualize histograms for the discretized Recency, Frequency, Monetary features. Also create one for the `FirstPurch` feature.

4.2 Box Plot (1 point) Transform the `Florence` variable into a categorical feature that can take up the values `True` or `False`. Create and visualize horizontal box plots for the original Recency, Frequency, Monetary and `FirstPurch` features against the `Florence` variable. *Hint:* To transform `Florence`, use the concept of factors in R and set the labels `True` and `False`.

4.3 Density Plot (1 point) Create and visualize a density plot for Recency, Frequency, Monetary and `FirstPurch` features.

Part II. ANOVA

An Analysis of Variance Test, or ANOVA, can be thought of as a generalization of the t-tests for more than 2 groups. The independent t-test is used to compare the means of a condition between two groups. ANOVA is used when we want to compare the means of a condition between more than two groups. ANOVA tests if there is a difference in the mean somewhere in the model (testing if there was an overall effect), but it does not tell us where the difference is (if there is one). To find where the difference is between the groups, we have to conduct post-hoc tests.

To perform any tests, we first need to define the null and alternate hypothesis:

- **Null Hypothesis:** There is *no significant difference* among the groups.
- **Alternate Hypothesis:** There is a *significant difference* among the groups.

Points

The problems for this part of the worksheet are for a total of 6 points, with a non-uniform weightage.

- *Problem 1* : 2 points
- *Problem 2* : 3 points
- *Problem 3* : 1 point

Scenario 1

It's a brand new day in the 99th precinct of the New York Police Department. Lieutenant Terrance has had enough of Hitchcock and Scully being useless paper pushers and wanted to assign them work to help the investigations; they were assigned the duty of gaining insights from the different types of objects in the evidence log of an ongoing investigation focused on the New York Mafia.

Problems

Problem 1 (2 points)

Captain Holt provided a file containing the names of a few `People of Interest` and the number of items logged at various evidence lockers of various precincts pertaining to them. He also instructs Peralta and Diaz to look into the file as he was told it should contain more information.

Scully decided to use ANOVA.

For this problem, use the data file named **Scenario 1.csv** in the data repository. Load the following libraries before moving on and read the dataset,

```
library(ggpubr)
library(dplyr)
library(ggplot2)
library(ggpubr)
library(broom)
library(car)

data <- read.csv('Scenario 1.csv')
```

1. Consider the dataset. Which type of ANOVA can Scully use? (Justify why the particular test)
2. What function(s) could have been used by Scully for ANOVA if he uses the R programming language?
3. What does the output of this/these functions tell Scully? (Specify hypotheses and what each column in the summary of the output means considering 5% significance)

Problem 2 (3 points)

Peralta and Diaz find a member of the family, a certain Frank Pentangeli, through Doug Judy. They discovered that the *famiglia* had altered this file resulting in invalid results. The original file was then recovered by the squad and was sent to Scully and Hitchcock for analysis. To their surprise they discovered that the file also had additional column of which gives the priority.

The dataset has three columns:

- First column has the **Person of Interest(POI)** in the Mafia
- Second column has the number of evidence items collected in particular evidence locker (evidence lockers are present across the city and many precincts have multiple squads working on the mafia, so one POI has multiple entries).
- Third column gives the **Priority** given to collect the evidence by a particular squad with respect to a POI.

Read the dataset before moving on. For this problem, use the data file named **Scenario 2.csv** in the data repository.

```
data <- read.csv('Scenario 2.csv')
```

1. Consider the data. Which type of ANOVA can Scully use? (Justify why the particular test)
2. What function(s) could have been used by Scully for the ANOVA if he uses the R programming language?
3. What does the output of this/these functions tell Scully? (Specify hypotheses and what each column in the summary of the output means considering 5% significance)
4. Hitchcock thinks that Scully has missed a task which completes the ANOVA test. What should Scully have thought of? *Hint*: Philosophically, a hypothesis is a proposition made as a basis for reasoning, without any assumption of its truth.

Problem 3 (1 point)

Hitchcock also wanted to compare the number of items collected for each pair of Person of Interest and priority. He decided to follow the common practice of doing a **Tukey's HSD**. The [Tukey's Honestly-Significant-Difference](#)[TukeyHSD] test lets us see which groups are different from one another.

What insights did Hitchcock gain after doing the Tukey's HSD? (The **TukeyHSD** function can be used to do this test and the output of this function can be represented graphically using the **plot** function.)