

Name: Arvind Murali

College: ~~Amrita~~ Amrita University Coimbatore (Amrita School of Engineering)

A Report on - VARIOUS DATA CLEANING AND ANALYSIS TECHNIQUES ON A MULTI- SERVICE BUSINESS DATASET

Aim: To generate a report on the various data cleaning and analysis processes involved while using a multi-service business dataset, with documentation on discrepancies or inaccuracies in the dataset (if any).

Procedure undertaken: The multi-service business dataset is provided in the form of an excel sheet (extension: .xlsx), with column names being 'Booking Type', 'Booking Date', 'Status', 'Class Type', 'Instructor', 'Time Slot', 'Duration (mins)', 'Price', 'Facility', 'Theme', 'Service Name', 'Subscription Type', 'Service Name', 'Service Type', 'Customer Email' and 'Customer Phone'.

The dataset is ~~first~~ ^{first} loaded into ^{front-end of} Microsoft Power BI after the process of transformation of data in the backend interface called as Power Query Editor, where all the cleaning and ~~pre~~ preprocessing of data takes place. ~~before~~

Datatype name: 'Data Analyst - Assessment - Dataset'

The various data cleaning and transformation are ^{done} as follows -

- Since Booking ID, Customer ID and Customer Name should not contain duplicate values, 'remove duplicates' feature is used to remove any redundant values.
- The first row containing column names are promoted as headers using 'use first row as headers' feature.
- Since no column should contain null values, the columns 'Class Type', 'Instructor', 'Facility' & 'Theme' ~~and~~ have their null values replaced with 'Not Applicable / Not Specified / (N/A)' using 'replace values' feature.
- Since 'subscription type' has column has empty values throughout, the column is removed completely, as it does not have ~~contribute~~ information for any part of analysis.
- The columns 'Customer Email' and 'Customer Phone' have their null values replaced with (N/A).
- ~~Empty~~ Empty spaces are removed, if any, using 'remove empty' feature of all columns.
- 'Time Slot' column datatype is changed from date/time to time.
- ~~Column~~ Column Quality, Column Profile and Column Distribution features are used to evaluate columns and the information they contain; no empty values were found and all the information was given as needed.
- The 'customer phone' column is split into 'customer phone number' and 'country code' to separate country code from phone number. This is done using 'split columns using delimiter' '-' (leftmost).

Conclusion: The cleaned dataset is evaluated once again and loaded to port end using 'Close and Apply' feature.