

# Data Quality Report - Initial Findings

## 1. Overview

### Introduction to Data Quality Report for TLC Taxi Dataset.

This data quality plan for our dataset aims to ensure that the information we gather, process, and analyse meets the highest standards of quality which can be used for predictive modelling. This report outlines the strategies and methodologies that will be employed to assess, monitor, and enhance the quality of the data related to taxi ride trip data gathered from TLC.

Our dataset includes information on taxi trips around New York city from 01-2023 to 12-2024. It provides us informative features such as pickup times, drop off times, pickup locations, drop off locations, date and time of trips and most importantly the passenger count of each trip.

## 2. Summary

The parquet files containing monthly trip data were all joined together to create a dataset with 89 million rows of data. In order to be able to clean and view this dataset we must find ways of condensing the dataset while maintaining important information.

We may combine the pickup and drop off locations together in order to improve our predictive accuracy. We are not interested in one over the other as they both tell us important information about the areas. Analysing pickup locations helps understand commuting patterns, popular destinations, and areas of high activity, which can be correlated with events, attractions, or daily routines. Analysing drop-off locations can help identify popular hotspots, areas of high tourist activity, and places with significant economic and social activity. By combining pickup and dropoff locations, we can analyse the overall taxi activity in Manhattan, including both the origins and destinations of trips, equally.

Analysing the null/missing/0 values of our dataset, we can see that only passenger\_count contained any null values (around 4%). We decided that the best course of action would be to replace these rows with the median value of passenger\_count feature. This way we can still use these rows within our dataset and do not need to remove them.

## 3.2. Visualise Data

All histograms and boxplots can be found on the appendix. Outliers can evidently be seen in the passenger\_count box plot also. However, the points are only seen as 'outliers' because

of the huge percentage of data with value 1 and will not skew our data as we will be combining all passenger\_counts together by taxi\_zone.

## 4.1. Descriptive Statistics

There are 5 categorical and one continuous feature in the dataset. They can essentially be divided into a few groups.

### Time:

- Hour: The hour of the day in 24 hour standard (e.g. 13)
- Day: Day of the month (e.g. 2 - 2nd of the month)
- Week: Day of the week (e.g. Saturday)
- Year\_Month: Year and Month of the year ( e.g. 01-2022 - January 2022)

### Taxi Zone:

Shows the TaxiZoneID number for each zone in which the trips occurred.

### Passenger Count:

Shows the number of passengers in the taxi zone within that time frame.

## 5. Action to take

### Passenger Count:

Replace the null values with the median value.

### Combine Features:

Combine the existing features from the raw dataset to condense the data into useful features in which we can extract useful information which may help our predictions.

## 6. Appendix





