

Yardgreen - Data Science Integrated Internship

Week 2:

project link: [Click here](#)

Data Cleaning and Preparation

Objective: Ensure the dataset is clean and ready for analysis

1. Missing Values Report:

Rows: 48,120

Columns: 4

Columns:

- DateTime → Timestamps (currently as object/string)
- Junction → ID of traffic junction
- Vehicles → Number of vehicles
- ID → Record ID

Cleaning Plan

- No missing values found in any column.

2. Standardized Date and Time Format:

Standardize Date/Time

Convert DateTime from string to proper datetime format YYYY-MM-DD HH:MM:SS

Add derived columns:

- Date
- Hour
- DayOfWeek
- Month

Clean and consistent DateTime

Added columns for flexible time-based analysis

Visual-ready structure for grouping, filtering, and trend analysis

3. Outlier Detection and Removal Summary:

Eliminate Outliers (Using IQR Method)

$$\text{IQR} = Q3 - Q1$$

$$\text{Lower Bound} = Q1 - 1.5 \times \text{IQR}$$

$$\text{Upper Bound} = Q3 + 1.5 \times \text{IQR}$$

Upper Bound ($Q3 + 1.5 \cdot IQR$): 59 vehicles

Remove the outliers in vehicles

4. Refined, Clean Dataset:

Accurate Time-Based Analysis

- With standardized DateTime, you can now group data by hour, day, week, or month to see traffic trends clearly.
- Enables robust visuals like line charts, area charts, and heatmaps.

Cleaner Statistical Summaries

- Outlier removal ensures means, medians, and totals reflect actual traffic behavior.
- Prevents skewed insights from rare traffic spikes or data errors.

Reliable Filtering and Forecasting

- No missing values = no NULL errors in DAX calculations or visuals.
- Clean structure allows smoother integration into forecasting models or Power BI Q&A.

Ready for Dashboards & KPIs

- Supports creation of accurate KPIs, cards, trend lines, and alerts

Visualize





