# Build data model, Data cleaning and preprocessing

# Comprehensive Report on the Work Submitted for Week One

1. Overview

The work submitted for the first week consists of :

- Python code in a Jupyter Notebook for processing and cleaning the "railway.csv" dataset related to UK train journeys.
- A set of SQL queries to analyze data stored in the "railway" database.
- Exploratory Data Analysis (EDA) performed in Excel, supplemented with Python analysis.
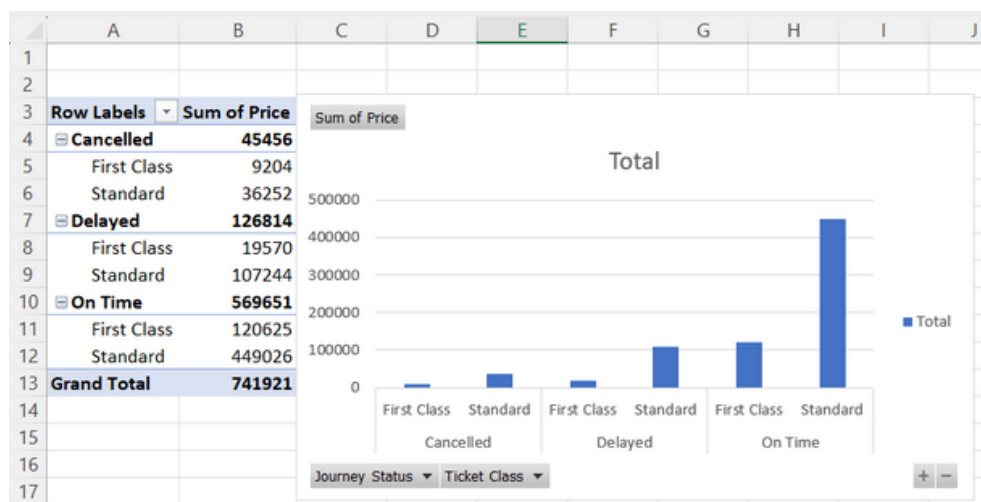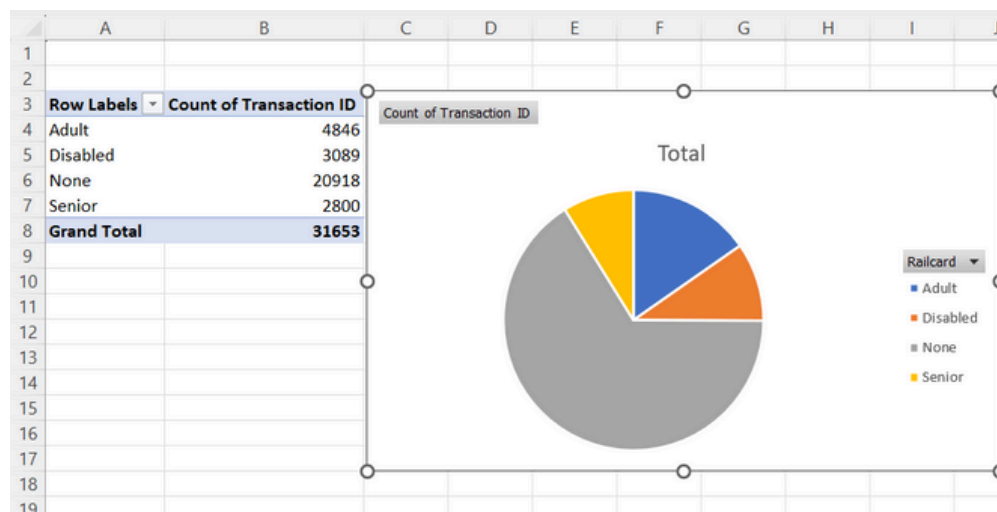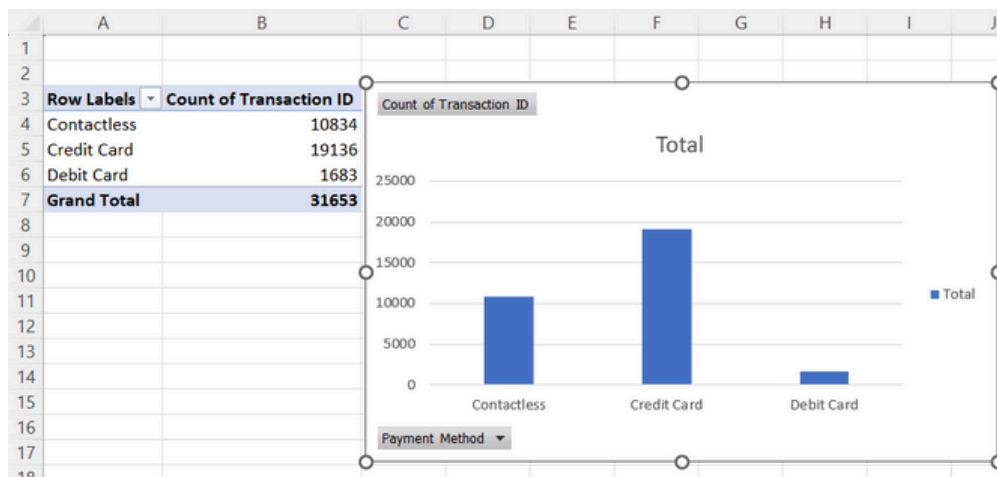
The objective is to clean the raw data, perform initial exploratory analysis, and extract analytical insights using Excel, Python, and SQL to understand patterns in train journeys, pricing, delays, and customer preferences.

# 1. Exploratory Data Analysis (EDA) Using Excel

**by Basem Fady**

A. Visualizations and Insights

The EDA in Excel utilized Pivot Tables and Charts to summarize key aspects of the dataset. The following visualizations were created:

# Run-Test Algorithm and Results

**by Basem Fady**

```
Run Test Results:
Number of Runs: 13966
Expected Runs: 15824.62
Standard Deviation of Runs: 0.94
Z-score: -1977.11
p-value: 0.00000
The data is not randomly distributed (Reject Null Hypothesis).
```

Conclusion: The price values in the dataset are not randomly distributed, suggesting the presence of a pattern in ticket prices.

# 2. Cleaning & Analysis Using SQL
**by Mariam Ahmed and Alaa Ebrahim.**

1. Checking for Null Values
- Checked for missing values in multiple columns such as Payment Method, Ticket Type, Price, Stations, Journey Date, Departure & Arrival Times, Journey Status, and Reason for Delay.
- Performed an additional check to identify null values in Actual_Arrival_Time when the journey status is not "Cancelled".



2. Handling Null Values
- Reason_for_Delay: Replaced null values with "No Delay".
- Railcard: Replaced "None" with "Without Rail Card".
- Actual_Arrival_Time: Set null values to "1900-01-01 00:00:00" for cancelled journeys only.

## 3. Checking for Duplicates

- Compared the total row count with the number of distinct Transaction_ID values to detect potential duplicates.

1. Average Price by Purchase Type
   - Calculates the average ticket price based on the Purchase_Type.
   - Orders the results in ascending order of price.
2. Ticket Price & Count by Journey Status
   - Computes the average ticket price and total number of tickets for journeys that were either Delayed or On Time.
   - Results are sorted in descending order by average ticket price.
3. Ticket Count by Payment Method & Ticket Type
   - Counts the number of tickets sold based on Payment_Method and Ticket_Type.
   - Orders the results by Payment_Method and ticket count (descending order).
4. Journey Frequency by Departure & Arrival Stations
   - Counts the total number of journeys between each departure and arrival station.
   - Orders the results by the most frequent journeys.
5. Departures & Arrivals Count
   - Counts the total number of departures per station.
   - Counts the total number of arrivals per destination.
   - Orders the results by the most frequent departure and arrival locations.

```sql
93 GROUP BY Payment_Method, Ticket_Type
94 ORDER BY Payment_Method, Avg_Ticket_Price DESC;
95
96 SELECT
97     Reason_for_Delay,
98     COUNT(*) AS Delay_Count,
99     ROUND((COUNT(*) * 100.0 / (SELECT COUNT(*) FROM railway WHERE Journey_Status = 'Delayed')), 2) AS Percentage
100 FROM railway
101 WHERE Journey_Status = 'Delayed'
102 GROUP BY Reason_for_Delay
103 ORDER BY Delay_Count DESC;
104
105
```

| Reason_for_Delay | Delay_Count | Percentage |
|---|---|---|
| Weather | 758 | 33.07 |
| Technical Issue | 472 | 20.59 |
| Signal Failure | 242 | 10.56 |
| Signal failure | 209 | 9.12 |
| Staff Shortage | 183 | 7.98 |
| Staffing | 172 | 7.5 |
| Weather Conditions | 169 | 7.37 |
| Traffic | 87 | 3.8 |

# 3. Data Processing Using Python
**by Alaa Ebrahim, Saif AboElmagd and Basem Fady.**

**Library Imports:**
- Imported pandas, numpy, matplotlib.pyplot, and seaborn for data manipulation and visualization.

**Data Loading:**
- Loaded "railway.csv" into a DataFrame using pd.read_csv().



**Initial Data Exploration**
- Displayed the first few rows using df.head().
- Checked dataset structure with df.info() and missing values using df.isnull().sum().
- The dataset has 31,653 rows and 18 columns.

**Checking for Duplicates**
- Verified if duplicate rows exist using df.duplicated().sum(), and confirmed there were no duplicates.

**Handling Missing Values**
- The notebook appears to include steps to manage missing data, but further details would be in later cells.

**Normalization**
- Min-Max Scaling for normalization using MinMaxScaler from sklearn.preprocessing:

## Price Distribution



## Feature Correlation Matrix

# 4. Data Modeling by Power BI
**by Alaa Ebrahim.**



To analyze UK train data efficiently, we used Power BI's data modeling to structure and relate datasets. The model follows a star schema, linking key tables:

- dim_location (stations with Location_ID)
- dim_ticket_type (ticket categories)
- fact_journey (train journeys)
- dim_payment (payment methods)
- dim_date (time dimension)

This structure ensures data integrity, optimized queries, and efficient visualization. The example table shows station mappings, enabling seamless analysis of travel patterns and station usage.

| Departure Station | Arrival Destination | Location_ID |
|---|---|---|
| London Paddington | Liverpool Lime Street | 1 |
| London Kings Cross | York | 2 |
| Liverpool Lime Street | Manchester Piccadilly | 3 |
| London Paddington | Reading | 4 |
| Liverpool Lime Street | London Euston | 5 |
| London Euston | Oxford | 6 |
| London Euston | York | 7 |
| York | Durham | 8 |
| Manchester Piccadilly | Liverpool Lime Street | 9 |
| Birmingham New Street | London St Pancras | 10 |
| London St Pancras | Birmingham New Street | 11 |
| Birmingham New Street | Manchester Piccadilly | 12 |
| London Euston | Birmingham New Street | 13 |
| Manchester Piccadilly | London Paddington | 14 |
| Oxford | Bristol Temple Meads | 15 |
| Birmingham New Street | Tamworth | 16 |
| Manchester Piccadilly | London Euston | 17 |
| London Paddington | London Waterloo | 18 |
| Manchester Piccadilly | Sheffield | 19 |
| London St Pancras | Wolverhampton | 20 |
| Liverpool Lime Street | Leeds | 21 |
| Birmingham New Street | Stafford | 22 |
| Birmingham New Street | London Euston | 23 |
| York | Doncaster | 24 |
| London Euston | Manchester Piccadilly | 25 |
| Reading | Swindon | 26 |
| London Paddington | Oxford | 27 |
| Manchester Piccadilly | Nottingham | 28 |

dim_location (66 rows)