

EDA & REGRESSION COURSEWORK 2024

The dataset for this coursework assignment, *MavenRail.csv*, is a synthetic dataset created by Maven Analytics for an online data analytics challenge. Participants were tasked with building an exploratory dashboard based on fictional traveler behavior and operational performance on the UK National Rail network.

According to Spiegelhalter’s ‘star rating’ system, the dataset would earn a “4*” rating—if it were factual—indicating high confidence as “numbers we can believe” (Muldoon, 2024).

The dataset contains 31,645 rows of individual traveler journeys, described across 13 columns:

1	Payment Method	Contactless, Credit Card, Debit Card
2	Railcard Type	Adult, Disabled, Senior, None
3	Ticket Class	First Class, Standard
4	Ticket Type	Advance, Anytime, Off-Peak
5	Price	GBP
6	Departure Station	Various UK National Rail Stations
7	Arrival Station	Various UK National Rail Stations
8	Departure Time	YYYY/mm/dd HH:MM
9	Scheduled Arrival	YYYY/mm/dd HH:MM
10	Actual Arrival	
11	Journey Status	On Time, Delayed, Cancelled
12	Reason for Delay	Signal Failure, Technical Issue, Weather, Staffing, Staff, Traffic
13	Refund Request	Yes/No

These details suggest the dataset was designed to study how journey delays and cancellations impact refund requests, a critical performance metric for service providers like National Rail.

Data Preparation

To begin, I examined the dataset structure. Most values were categorical, including Departure, Scheduled Arrival, and Actual Arrival, with some missing values in key columns: Railcard (67%), Actual Arrival (5.9%), and Reason for Delay (86.8%).

- **Railcard:** Null values were replaced with “None,” a valid category alongside “Adult,” “Disabled,” and “Senior.”
- **Reason for Delay:** Null values were similarly replaced with “None.”

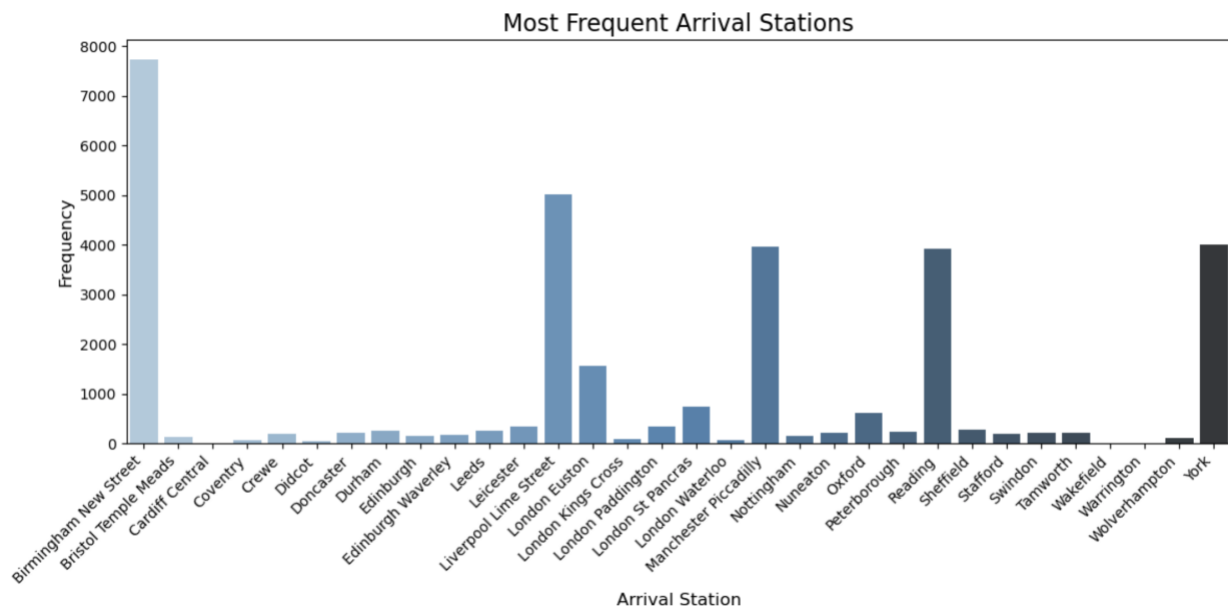
- **Scheduled Arrival and Departure:** Null values were removed due to their low frequency.
- **Actual Arrival:** Null values were replaced with “NA,” indicating no delay.

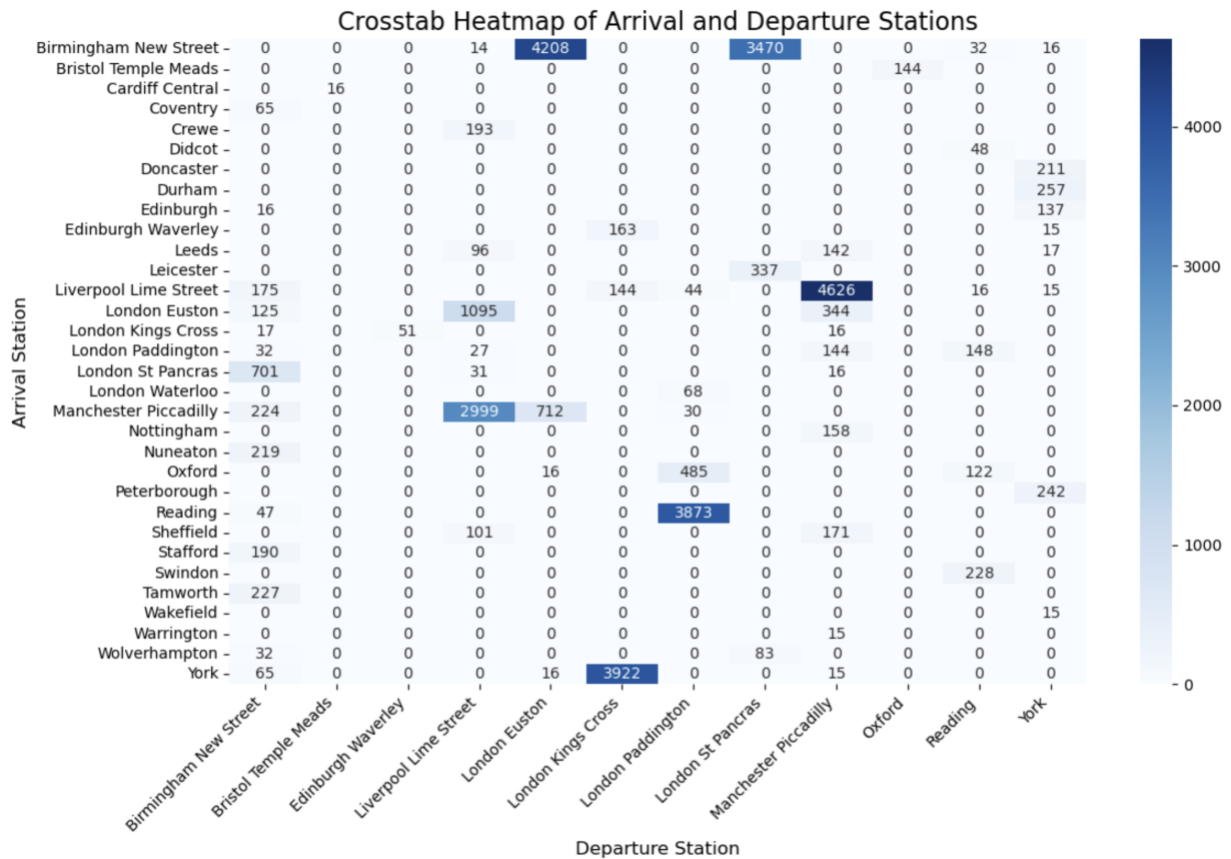
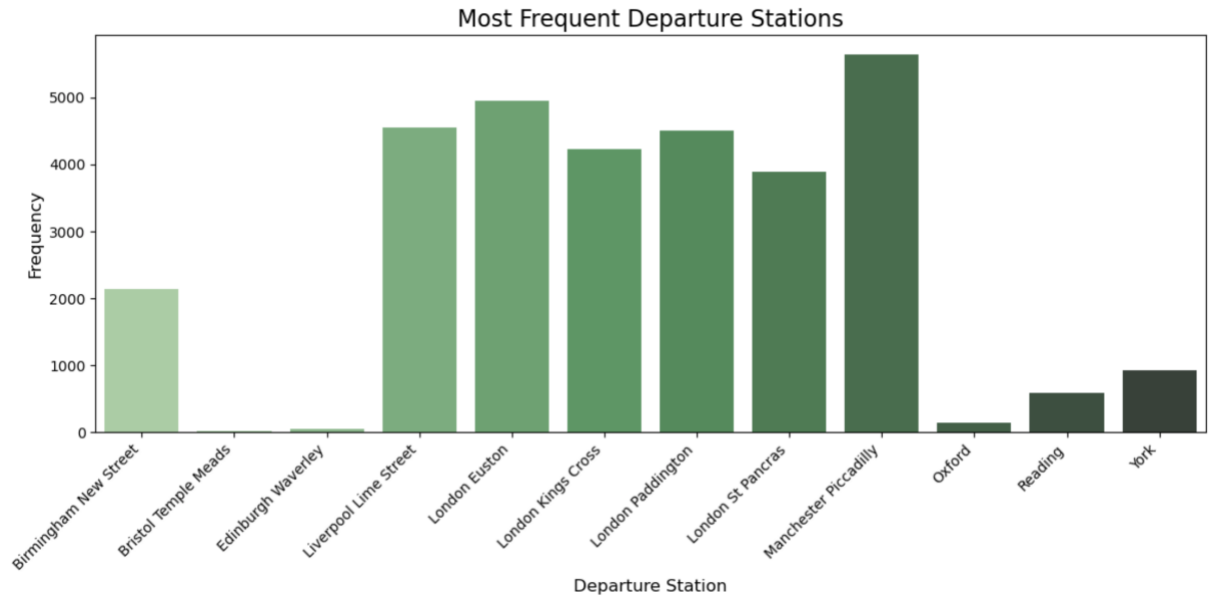
For analysis, departure and arrival times were converted to datetime64. Null values in Actual Arrival were marked as “Not a Time” to avoid exceptions during conversion.

Exploratory Data Analysis

Departure and Arrival Patterns

Using cross-tabulations, I found that Manchester Piccadilly was the most frequent departure station, and Liverpool Lime Street the most common arrival station. The most frequent journey was from Manchester Piccadilly to Liverpool Lime Street.





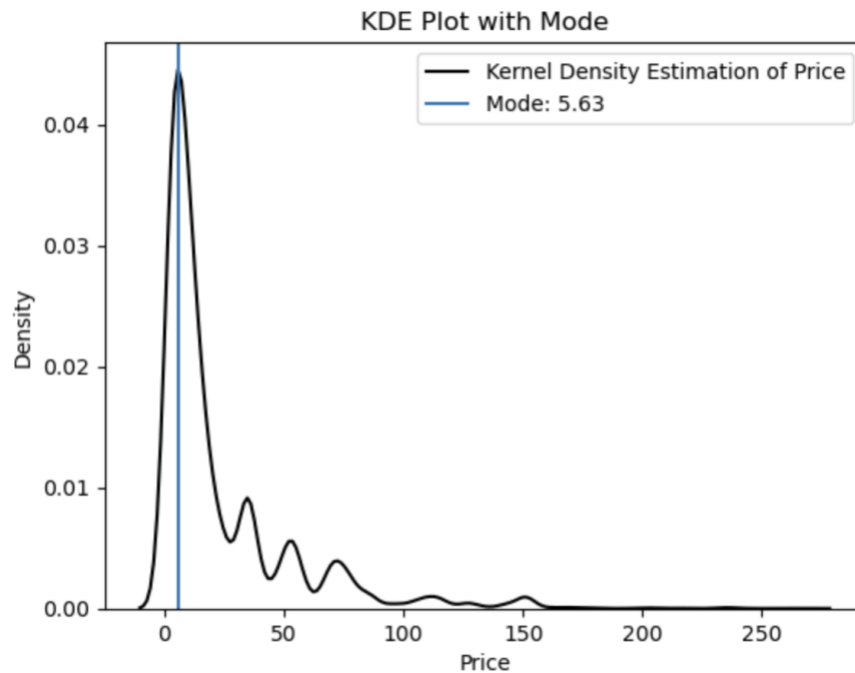
Ticket Prices

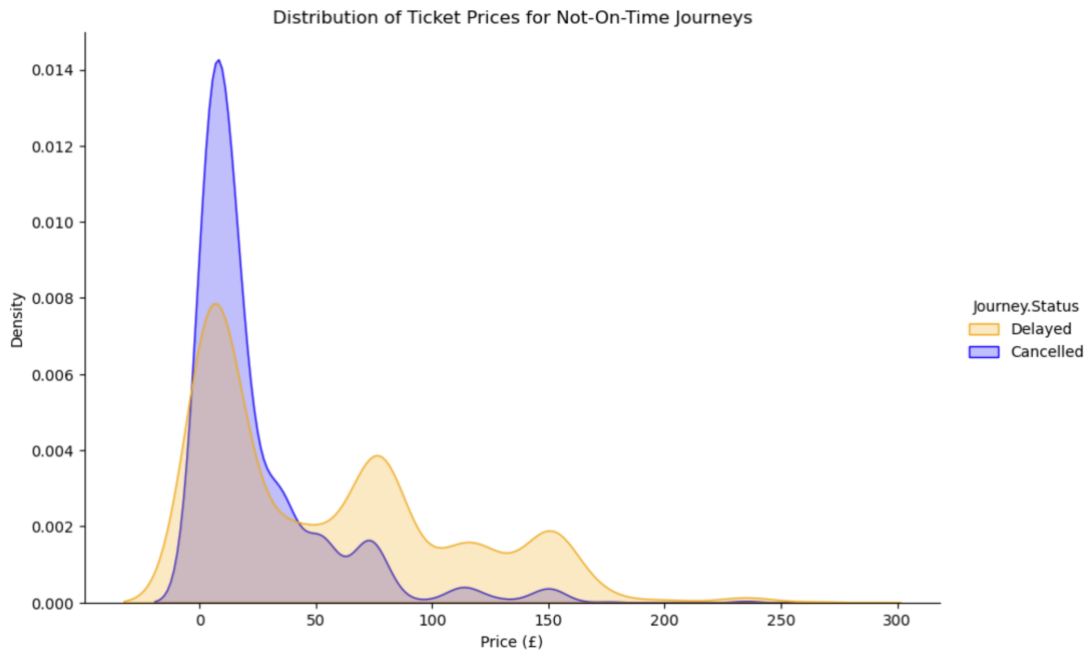
Univariate analysis of ticket prices revealed:

- Mean: £23.43

- Standard deviation: £29.98
- Minimum: £1
- Maximum: £267
- Median: £11
- Mode: £3

A Kernel Density Estimation (KDE) plot showed a KDE mode of £5.63, differing from the raw data mode likely due to the influence of outliers (e.g., tickets above £250). Price analysis by journey status showed that delayed trips were costlier, while cancelled trips had a higher density at lower prices. This might indicate that cheaper, more frequent routes are cancelled more often, while costlier routes are delayed to avoid cancellations.





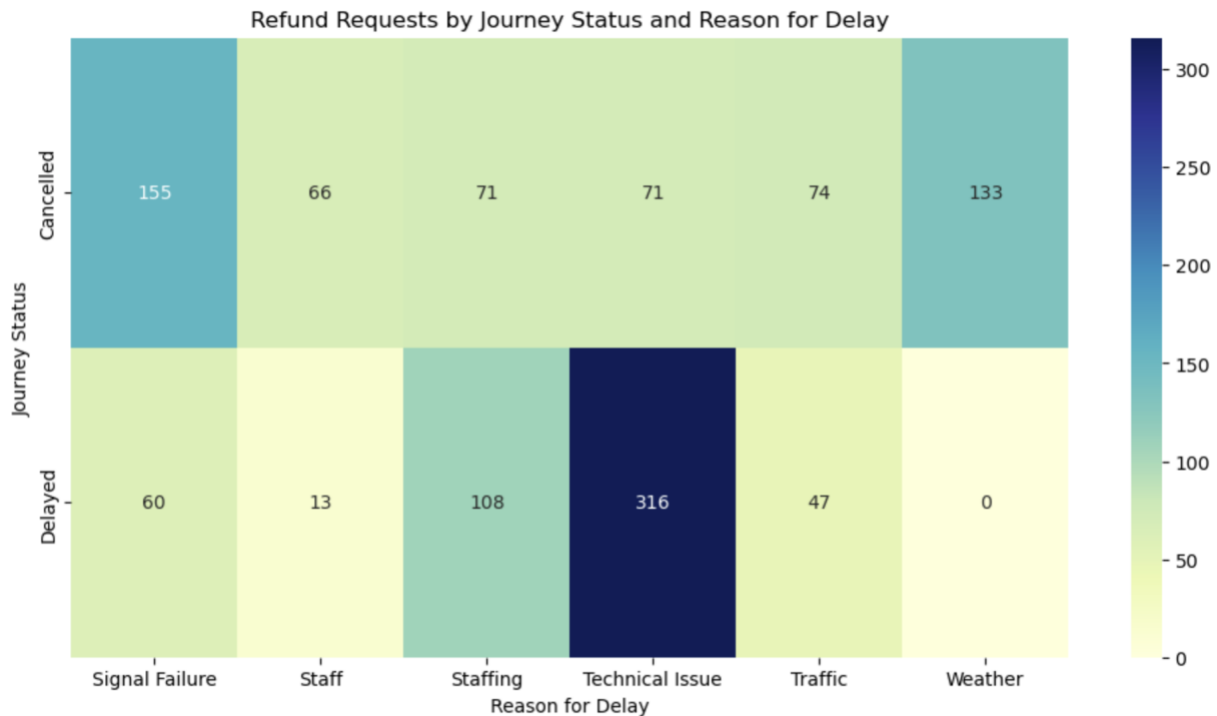
Analysis of Refund Requests

To examine refund requests, I filtered the dataset to include only delayed or cancelled journeys and cross-tabulated them with Reason for Delay. Results showed that:

- **Delayed Journeys:** Refund requests were most frequent when delays were due to technical issues.
- **Cancelled Journeys:** Signal failure was the leading cause, drawing the most refund requests.

When analyzing the relationship between Reason for Delay and Minutes Delayed:

- Travelers on delayed trains waited most often due to weather.
- For cancelled trains, the longest delays were caused by signal failure



Data Transformation

To prepare for modeling, I calculated the delay in minutes by subtracting Scheduled Arrival from Actual Arrival, with no-delay values marked as 0.0 and then replaced with “NA.” A new column, *MediumPrice*, was created to indicate whether ticket prices fell between £10 and £30, mapped as binary values (“Yes” or “No”).

Outliers in ticket prices were removed using a Z-score threshold, reducing the maximum price to £113 and the mean to £20.43 across 30,881 values. I converted categorical columns to numerical or categorical types as needed, with *DelayInMinutes* set to numeric, replacing “NA” values with 0.

Price	Departure.Station	Arrival.Station	Departure	Scheduled.Arrival	Actual.Arrival	Journey.Status	Reason.for.Delay	Refund.Request	DelayInMinutes
43	London Paddington	Liverpool Lime Street	2024-01-01 11:00:00	2024-01-01 13:30:00	2024-01-01 13:30:00	On Time	NA	No	NA
23	London Kings Cross	York	2024-01-01 09:45:00	2024-01-01 11:35:00	2024-01-01 11:40:00	Delayed	Signal Failure	No	5.0
3	Liverpool Lime Street	Manchester Piccadilly	2024-01-02 18:15:00	2024-01-02 18:45:00	2024-01-02 18:45:00	On Time	NA	No	NA
13	London Paddington	Reading	2024-01-01 21:30:00	2024-01-01 22:30:00	2024-01-01 22:30:00	On Time	NA	No	NA
76	Liverpool Lime Street	London Euston	2024-01-01 16:45:00	2024-01-01 19:00:00	2024-01-01 19:00:00	On Time	NA	No	NA

Logistic Regression Models

Single-Predictor Model

```
#fit small model with one predictor
X_1 = journey_final['MediumPrice'].values.reshape(-1, 1)
X_1 = sm.add_constant(X_1)
#set response variable
y = journey_final.loc[:, journey_final.columns == 'Refund.Request_Yes']

logit_model=sm.Logit(y,X_1)
result1=logit_model.fit()
print(result1.summary2())
```

Optimization terminated successfully.
Current function value: 0.578797
Iterations 5

Results: Logit

Model:	Logit	Method:	MLE
Dependent Variable:	Refund.Request_Yes	Pseudo R-squared:	0.003
Date:	2024-11-21 00:14	AIC:	4825.3797
No. Observations:	4165	BIC:	4838.0487
Df Model:	1	Log-Likelihood:	-2410.7
Df Residuals:	4163	LL-Null:	-2418.7
Converged:	1.0000	LLR p-value:	6.1887e-05
No. Iterations:	5.0000	Scale:	1.0000

	Coef.	Std.Err.	z	P> z	[0.025	0.975]
const	-1.0753	0.0393	-27.3720	0.0000	-1.1523	-0.9983
x1	0.3540	0.0872	4.0583	0.0000	0.1830	0.5250

Using *MediumPrice* as the sole predictor, I developed a logistic regression model to predict refund requests. In this model:

- Travelers paying less than £10 or more than £30 had lower refund probabilities.
- For a £5 ticket, the refund probability was 25%.

$$p = \frac{1}{1 + e^{-\theta}} = \frac{1}{1 + e^{-(-1.0753 + 0.3540 \times 0)}} = \frac{1}{1 + e^{-(-1.0753)}} = \frac{1}{3.93} = 0.25 \times 100 = 25\%$$

- For a £25 ticket, the refund probability was 33%.

$$p = \frac{1}{1 + e^{-\theta}} = \frac{1}{1 + e^{-(-1.0753 + 0.3540 \times 0)}} = \frac{1}{1 + e^{-(-0.7213)}} = \frac{1}{3.057} = 0.33 \times 100 = 33\%$$

However, this model performed poorly, with a Pseudo R-squared value of 0.003, indicating minimal explanatory power.

```
# Predict refund request using test data from to_predict

X_1 = predict_journey_final['MediumPrice'].values.reshape(-1, 1)
X_1 = sm.add_constant(X_1)

yhat = result1.predict(X_1)
prediction = list(map(round, yhat))

print("Predicted probabilities:", yhat)
print("Binary predictions (0 or 1):", prediction)
```

Predicted probabilities: [0.2543911 0.2543911 0.2543911 0.2543911 0.2543911 0.3271028] Binary predictions (0 or 1): [0, 0, 0, 0, 0, 0]

Multiple Logistic Regression

To improve predictions, I used the following predictors: *MediumPrice*, *Price*, *DelayInMinutes*, *Journey Status Delayed*, *Reason for Delay (Staffing)*, and *Reason for Delay (Technical Issue)*. This model achieved a better Pseudo R-squared value of 0.179. Key results:

```
Optimization terminated successfully.
Current function value: 0.476949
Iterations 7

Results: Logit
=====
Model:                Logit                Method:                MLE
Dependent Variable:    Refund.Request_Yes    Pseudo R-squared:       0.179
Date:                  2024-11-22 00:08      AIC:                    3986.9828
No. Observations:      4165                  BIC:                    4031.3241
Df Model:              6                      Log-Likelihood:         -1986.5
Df Residuals:          4158                  LL-Null:                -2418.7
Converged:             1.0000                LLR p-value:            1.8257e-183
No. Iterations:        7.0000                Scale:                  1.0000
=====
```

	Coef.	Std.Err.	z	P> z	[0.025	0.975]
const	-1.0594	0.0696	-15.2188	0.0000	-1.1958	-0.9230
MediumPrice	-0.0408	0.0985	-0.4141	0.6788	-0.2338	0.1523
Price	-0.0048	0.0010	-4.5999	0.0000	-0.0068	-0.0028
DelayInMinutes	-0.0593	0.0038	-15.5963	0.0000	-0.0667	-0.0518
Journey.Status_Delayed	1.4501	0.1293	11.2189	0.0000	1.1968	1.7035
Reason.for.Delay_Staffing	0.9982	0.1172	8.5144	0.0000	0.7684	1.2280
Reason.for.Delay_Technical_Issue	1.4509	0.0965	15.0341	0.0000	1.2617	1.6400

- **Journey Status Delayed (coefficient: 1.45):** Strongly associated with refund requests.
- **Reason for Delay - Technical Issue (coefficient: 1.45):** Significant contributor.
- **MediumPrice (coefficient: -0.0408):** Not a meaningful predictor compared to other variables.

Predictions on Test Data

Using the multiple logistic model, I predicted refund probabilities for the *To_Predict* dataset:

- Journey 2: 1%
- Journey 3: 14.8%
- Journey 4: 55.7%
- Journey 5: 45.4%
- Journey 6: 0.1%
- Journey 7: 12.3%

These probabilities reflect that journey delays, especially those caused by technical issues, significantly increase refund requests, far more than ticket price or whether the price falls within a specific range.

Conclusion

Delays and cancellations are critical drivers of refund requests, with technical issues and staffing delays being the most problematic causes. While pricing influences refund probabilities to a small extent, operational performance, particularly avoiding delays, has a more substantial impact. National Rail should focus on addressing the root causes of delays, especially technical issues, to improve customer satisfaction and reduce operational costs.

REFERENCES

Agresti, A. (2018). 'Multiple Logistic Regression', *Statistical Methods for the Social Sciences*. Fifth edn: Pearson, pp. 477-479.

Maven Rail Challenge (2024). Maven Analytics. Available at:
<https://bit.ly/MavenRailChallenge> (Accessed: 25/11/2024).

Muldoon, M. (2024). 'Lecture 2A: Univariate Exploratory Data Analysis', *Statistics and Machine Learning, Week 2, 2024*(15/11/2024) Available at:
https://online.manchester.ac.uk/ultra/courses/_83866_1/cl/outline.

1. MovenRail Data Description

```
In [2]: import pandas as pd
import numpy as np
df = pd.read_csv('MovenRail.csv')
df.head()
```

```
Out[2]: Payment.Method  Railcard  Ticket.Class  Ticket.Type  Price  Departure.Station  Arrival.Station  Departure  \
0      Contactless      Adult      Standard      Advance    43      London Paddington  Liverpool Lime  2024-01-01 10:00
1      Credit Card      Adult      Standard      Advance    23      London Kings Cross      York      2024-01-02 18:15
2      Credit Card      NaN      Standard      Advance     3      Liverpool Line Street  Manchester Piccadilly  2024-01-02 18:15
3      Credit Card      NaN      Standard      Advance    13      London Paddington      Reading      2024-01-01 21:30
4      Contactless      NaN      Standard      Advance    76      Liverpool Lime Street  London Euston  2024-01-01 16:45

In [3]: df.describe()
Out[3]: Payment.Method  Railcard  Ticket.Class  Ticket.Type  Price  \
0      Contactless      Adult      Standard      Advance    43
1      Credit Card      Adult      Standard      Advance    23
2      Credit Card      NaN      Standard      Advance     3
3      Credit Card      NaN      Standard      Advance    13
4      Contactless      NaN      Standard      Advance    76
...      ...      ...      ...      ...      ...
31640      Credit Card      NaN      Standard      Off-Peak     4
31641      Contactless      NaN      Standard      Off-Peak    10
31642      Credit Card      NaN      Standard      Off-Peak     3
31643      Credit Card      NaN      Standard      Off-Peak    10
31644      Credit Card      Adult      Standard      Off-Peak     3

Departure.Station  Arrival.Station  Departure  \
0      London Paddington  Liverpool Line Street  2024-01-01 11:00
1      London Kings Cross      York      2024-01-01 09:45
2      Liverpool Line Street  Manchester Piccadilly  2024-01-02 18:15
3      London Paddington      Reading      2024-01-01 21:30
4      Liverpool Line Street  London Euston  2024-01-01 16:45
...      ...      ...      ...
31640      Manchester Piccadilly  Liverpool Line Street  2024-04-30 20:00
31641      London Euston      Birmingham New Street  2024-04-30 20:15
31642      Manchester Piccadilly  Liverpool Line Street  2024-04-30 20:15
31643      London Euston      Birmingham New Street  2024-04-30 21:15
31644      Liverpool Line Street  Manchester Piccadilly  2024-04-30 21:30

Scheduled.Arrival  Actual.Arrival  Journey.Status  Reason.for.Delay  \
0      2024-01-01 13:30  2024-01-01 13:30      On Time      Signal Failure
1      2024-01-02 18:45  2024-01-02 18:45      On Time      NaN
2      2024-01-02 18:45  2024-01-02 18:45      On Time      NaN
3      2024-01-01 19:00  2024-01-01 19:00      On Time      NaN
4      2024-01-01 19:00  2024-01-01 19:00      On Time      NaN
...      ...      ...      ...
31640      2024-04-30 20:30  2024-04-30 20:30      On Time      NaN
31641      2024-04-30 21:35  2024-04-30 21:35      On Time      NaN
31642      2024-04-30 20:45  2024-04-30 20:45      On Time      NaN
31643      2024-04-30 22:35  2024-04-30 22:35      On Time      NaN
31644      2024-04-30 22:00  2024-04-30 22:00      On Time      NaN

Refund.Request      \
0      No
1      No
2      No
3      No
4      No
...      ...
31640      No
31641      No
31642      No
31643      No
31644      No

[31645 rows x 13 columns]
```

```
In [4]: df.shape
Out[4]: (31645, 13)
```

```
In [5]: df.columns
```

```
Out[5]: Index(['Payment.Method', 'Railcard', 'Ticket.Class', 'Ticket.Type', 'Price',
      'Departure.Station', 'Arrival.Station', 'Departure',
      'Scheduled.Arrival', 'Actual.Arrival', 'Journey.Status',
      'Reason.for.Delay', 'Refund.Request'],
      dtype='object')
```

```
In [6]: for column in ['Payment.Method', 'Railcard', 'Ticket.Class', 'Ticket.Type', 'Price',
      'Departure.Station', 'Arrival.Station', 'Departure',
      'Scheduled.Arrival', 'Actual.Arrival', 'Journey.Status',
      'Reason.for.Delay', 'Refund.Request']:
    print(f'{column}: {df[column].dtype}')
```

```
Payment.Method: object
Railcard: object
Ticket.Class: object
Ticket.Type: object
Price: int64
Departure.Station: object
Arrival.Station: object
Departure: object
Scheduled.Arrival: object
Actual.Arrival: object
Journey.Status: object
Reason.for.Delay: object
Refund.Request: object
```

```
In [7]: df.isnull().sum()
```

```
Out[7]: Payment.Method      0
Railcard      20911
Ticket.Class      0
Ticket.Type      0
Price      0
Departure.Station      0
Arrival.Station      0
Departure      3
Scheduled.Arrival      4
Actual.Arrival      1880
Journey.Status      0
Reason.for.Delay      27479
Refund.Request      0
dtype: int64
```

```
In [8]: (df.isnull().sum()/(len(df)))>100
```

```
Out[8]: Payment.Method      0.000000
Railcard      66.079949
Ticket.Class      0.000000
Ticket.Type      0.000000
Price      0.000000
Departure.Station      0.000000
Arrival.Station      0.000000
Departure      0.000000
Scheduled.Arrival      0.012640
Actual.Arrival      0.049007
Journey.Status      0.000000
Reason.for.Delay      86.835203
Refund.Request      0.000000
dtype: float64
```

```
In [9]: payment_methods = []
for value in df['Payment.Method']:
    if value in payment_methods:
        pass
    else:
        payment_methods.append(value)
print(payment_methods)
```

```
Out[9]: ['Contactless', 'Credit Card', 'Debit Card']
```

```
In [10]: journey_status = []
for value in df['Journey.Status']:
    if value in journey_status:
        pass
    else:
        journey_status.append(value)
print(journey_status)
```

```
Out[10]: ['On Time', 'Delayed', 'Cancelled']
```

```
In [11]: reasons_for_delay = []
for value in df['Reason.for.Delay']:
    if value in reasons_for_delay:
        pass
    else:
        reasons_for_delay.append(value)
print(reasons_for_delay)
```

```
Out[11]: ['Signal Failure', 'Technical Issue', 'Weather', 'Staffing', 'Staff', 'Traffic']
```

```
In [12]: railcard_type = []
for value in df['Railcard']:
    if value in railcard_type:
        pass
    else:
        railcard_type.append(value)
print(railcard_type)
```

```
Out[12]: ['Adult', 'nan', 'Disabled', 'Senior']
```

```
In [13]: df_clean = df.copy()
```

```
In [14]: df_clean.isnull().sum()
(df_clean.isnull().sum()/(len(df)))>100
```

```
Out[14]: Payment.Method      0.000000
Railcard      66.079949
Ticket.Class      0.000000
Ticket.Type      0.000000
Price      0.000000
Departure.Station      0.000000
Arrival.Station      0.000000
Departure      0.009480
Scheduled.Arrival      0.012640
Actual.Arrival      0.049007
Journey.Status      0.000000
Reason.for.Delay      86.835203
Refund.Request      0.000000
dtype: float64
```

```
In [15]: # Remove null values. Convert 'nan' for 'No reason' entry in 'Reason for delay' to 'none'.
df_clean['Reason.for.Delay'] = df_clean['Reason.for.Delay'].fillna('NA')
```

```
In [16]: # Remove null values. Convert 'nan' for 'no railcard' entry in 'Railcards' to 'none'.
df_clean['Railcard'] = df_clean['Railcard'].fillna('None')
```

```
In [17]: df_clean.isnull().sum()
```

```
Out[17]: Payment.Method      0
Railcard      0
Ticket.Class      0
Ticket.Type      0
Price      0
Departure.Station      0
Arrival.Station      0
Departure      3
Scheduled.Arrival      4
Actual.Arrival      1880
Journey.Status      0
Reason.for.Delay      0
Refund.Request      0
dtype: int64
```

```
In [18]: # Remove missing values in 'Scheduled Arrival' and 'Departure' due to low number and lack of alter
df_clean = df_clean[df_clean['Scheduled.Arrival'].isnull()<100]
df_clean = df_clean[df_clean['Departure'].isnull()<100]
```

```
In [19]: # Remove null values. Convert null values in 'Actual Arrival' to duplicate 'Scheduled Arrival' value
df_clean['Actual.Arrival'] = df_clean['Actual.Arrival'].fillna(df_clean['Scheduled.Arrival'])
```

```
Out[19]: df_clean['Actual.Arrival']
```

```
Out[20]: 0      2024-01-01 13:30
1      2024-01-01 11:40
2      2024-01-02 18:45
3      2024-01-01 22:30
4      2024-01-01 19:00
...      ...
31640      2024-04-30 20:30
31641      2024-04-30 21:35
31642      2024-04-30 20:45
31643      2024-04-30 22:35
31644      2024-04-30 22:00
Name: Actual.Arrival, Length: 31639, dtype: object
```

```
In [21]: df_clean.isnull().sum()
```

```
Out[21]: Payment.Method      0
Railcard      0
Ticket.Class      0
Ticket.Type      0
Price      0
Departure.Station      0
Arrival.Station      0
Departure      0
Scheduled.Arrival      0
Actual.Arrival      0
Journey.Status      0
Reason.for.Delay      0
Refund.Request      0
dtype: int64
```

```
In [22]: crosstab_result = pd.crosstab(df_clean['Ticket.Class'], df_clean['Ticket.Type'])
print(crosstab_result)
```

```
Out[22]: Ticket.Class  Advance  Anytime  Off-Peak
First class      1762      406      884
Standard      1793      4850      7940
```

```
In [23]: # Convert string values for 'Actual Arrival', 'Departure', and 'Scheduled Arrival' to datetime64
df_clean['Departure'] = pd.to_datetime(df_clean['Departure'], format='%Y-%m-%d %H:%M')
df_clean['Scheduled.Arrival'] = pd.to_datetime(df_clean['Scheduled.Arrival'], format='%Y-%m-%d %H:%M')
```

```
In [24]: df_clean.info()
```

```
Out[24]: <class 'pandas.core.frame.DataFrame'>
Int64Index: 31639 entries, 0 to 31644
Data columns (total 13 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   Payment.Method      31639 non-null  object
 1   Railcard            31639 non-null  object
 2   Ticket.Class        31639 non-null  object
 3   Ticket.Type         31639 non-null  object
 4   Price               31639 non-null  int64
 5   Departure.Station   31639 non-null  object
 6   Arrival.Station     31639 non-null  object
 7   Departure           31639 non-null  datetime64[ns]
 8   Scheduled.Arrival   31639 non-null  datetime64[ns]
 9   Actual.Arrival      31639 non-null  datetime64[ns]
10   Journey.Status      31639 non-null  object
11   Reason.for.Delay    31639 non-null  object
12   Refund.Request      31639 non-null  object
dtypes: datetime64[ns](3), int64(1), object(9)
memory usage: 3.4+ MB
```

```
In [25]: crosstab_result = pd.crosstab(df_clean['Arrival.Station'], df_clean['Departure.Station'], margins=True)
print(crosstab_result)
```

```
Out[25]: Departure.Station  Birmingham New Street  Bristol Temple Meads  \
Arrival.Station
Birmingham New Street      0      0
Bristol Temple Meads      0      0
Cardiff Central      0      16
Coventry      65      0
Crewe      0      0
Didcot      0      0
Doncaster      0      0
Durham      0      0
Edinburgh      16      0
Edinburgh Waverley      0      0
Leeds      0      0
Leicester      0      0
Liverpool Line Street      175      0
London Euston      125      0
London Kings Cross      172      0
London Paddington      32      0
London St Pancras      781      0
London Waterloo      224      0
Manchester Piccadilly      0      0
Nottingham      219      0
Oxford      0      0
Peterborough      0      0
Reading      47      0
Sheffield      0      0
Stafford      190      0
Swindon      0      0
Tamworth      227      0
Warrington      0      0
Wolverhampton      0      0
York      32      0
All      2135      16

Departure.Station  Edinburgh Waverley  Liverpool Line Street  \
Arrival.Station
Birmingham New Street      0      14
Bristol Temple Meads      0      0
Cardiff Central      0      0
Coventry      0      0
Crewe      0      0
Didcot      0      0
Doncaster      0      0
Durham      0      0
Edinburgh      0      0
Edinburgh Waverley      0      0
Leeds      0      0
Leicester      0      0
Liverpool Line Street      0      142
London Euston      0      30
London Kings Cross      51      0
London Paddington      0      27
London St Pancras      0      344
London Waterloo      0      0
Manchester Piccadilly      0      0
Nottingham      0      0
Oxford      0      0
Peterborough      0      0
Reading      0      0
Sheffield      0      0
Stafford      0      0
Swindon      0      0
Tamworth      0      0
Warrington      0      0
Wolverhampton      0      0
York      0      0
All      51      4556

Departure.Station  London Euston  London Kings Cross  London Paddington  \
Arrival.Station
Birmingham New Street      4208      0      0
Bristol Temple Meads      0      0      0
Cardiff Central      0      0      0
Coventry      0      0      0
Crewe      0      0      0
Didcot      0      0      0
Doncaster      0      0      0
Durham      0      0      0
Edinburgh      0      0      0
Edinburgh Waverley      0      163      0
Leeds      0      0      0
Leicester      0      0      0
Liverpool Line Street      0      144      44
London Euston      0      0      0
London Kings Cross      0      0      0
London Paddington      0      0      0
London St Pancras      0      0      68
London Waterloo      712      0      0
Manchester Piccadilly      0      0      0
Nottingham      219      0      0
Oxford      16      0      485
Peterborough      0      0      0
Reading      0      0      3873
Sheffield      0      0      0
Stafford      0      0      0
Swindon      0      0      0
Tamworth      0      0      0
Warrington      0      0      0
Wolverhampton      0      0      0
York      16      3922      4500
All      4952      4229      4500

Departure.Station  London St Pancras  Manchester Piccadilly  Oxford  \
Arrival.Station
Birmingham New Street      3470      0      0
Bristol Temple Meads      0      0      144
Cardiff Central      0      0      0
Coventry      0      0      0
Crewe      0      0      0
Didcot      0      0      0
Doncaster      0      0      0
Durham      0      0      0
Edinburgh      0      0      0
Edinburgh Waverley      0      0      0
Leeds      0      144      36
Leicester      0      142      30
Liverpool Line Street      0      0      0
London Euston      0      1905      0
London Kings Cross      51      0      0
London Paddington      0      27      0
London St Pancras      0      344      0
London Waterloo      0      0      0
Manchester Piccadilly      0      2999      0
Nottingham      0      0      0
Nuneaton      0      0      0
Oxford      0      0      0
Peterborough      0      0      0
Reading      0      0      0
Sheffield      0      101      0
Stafford      0      0      0
Swindon      0      0      0
Tamworth      0      0      0
Warrington      0      0      0
Wolverhampton      0      15      0
York      0      0      0
All      51      4556      144

Departure.Station  Reading  York  All
Arrival.Station
Birmingham New Street      32      16      7740
Bristol Temple Meads      0      0      144
Cardiff Central      0      0      65
Coventry      0      0      16
Crewe      0      0      193
Didcot      48      0      48
Doncaster      0      211      211
Durham      0      257      257
Edinburgh      0      137      153
Edinburgh Waverley      0      15      178
Leeds      0      17      255
Leicester      0      0      337
Liverpool Line Street      16      15      200
London Euston      0      0      1504
London Kings Cross      0      84      84
London Paddington      32      0      351
London St Pancras      0      0      748
London Waterloo      0      0      68
Manchester Piccadilly      0      0      3965
Nottingham      0      0      219
Nuneaton      0      0      158
Oxford      122      0      623
Peterborough      0      242      242
Reading      0      0      3920
Sheffield      0      0      272
Stafford      0      0      190
Swindon      228      0      228
Tamworth      0      0      227
Warrington      0      15      15
Wolverhampton      0      0      83
York      0      0      115
All      594      925      31639
```

2. Exploratory Data Analysis - Univariate

```
In [53]: df_clean['Price'] = df_clean['Price'].astype(float)
```

```
In [55]: import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [57]: # Univariate analysis of Price
df_clean['Price'].describe()
df_clean['Price'].mode()
```

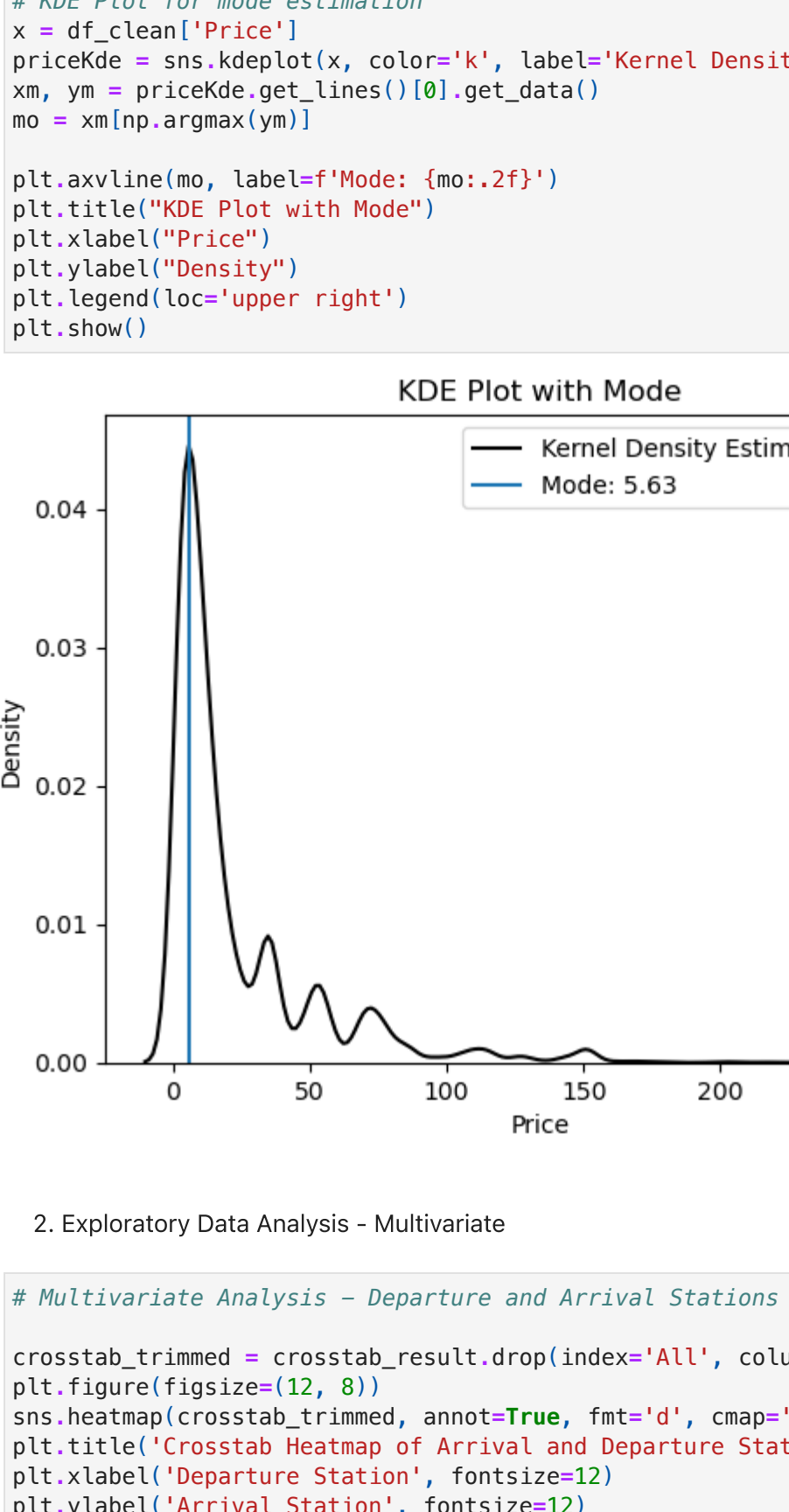
```
Out[57]: 0      3.0
Name: Price, dtype: float64
```

```
In [59]: df_clean['Price'].median()
```

```
Out[59]: 11.0
```

```
In [61]: # Histogram of railcard prices
sns.displot(df_clean['Price'], label='Price Distribution')
plt.title('Distribution of Ticket Prices')
plt.xlabel('Price')
plt.ylabel('Count')
plt.legend()
plt.show()
```

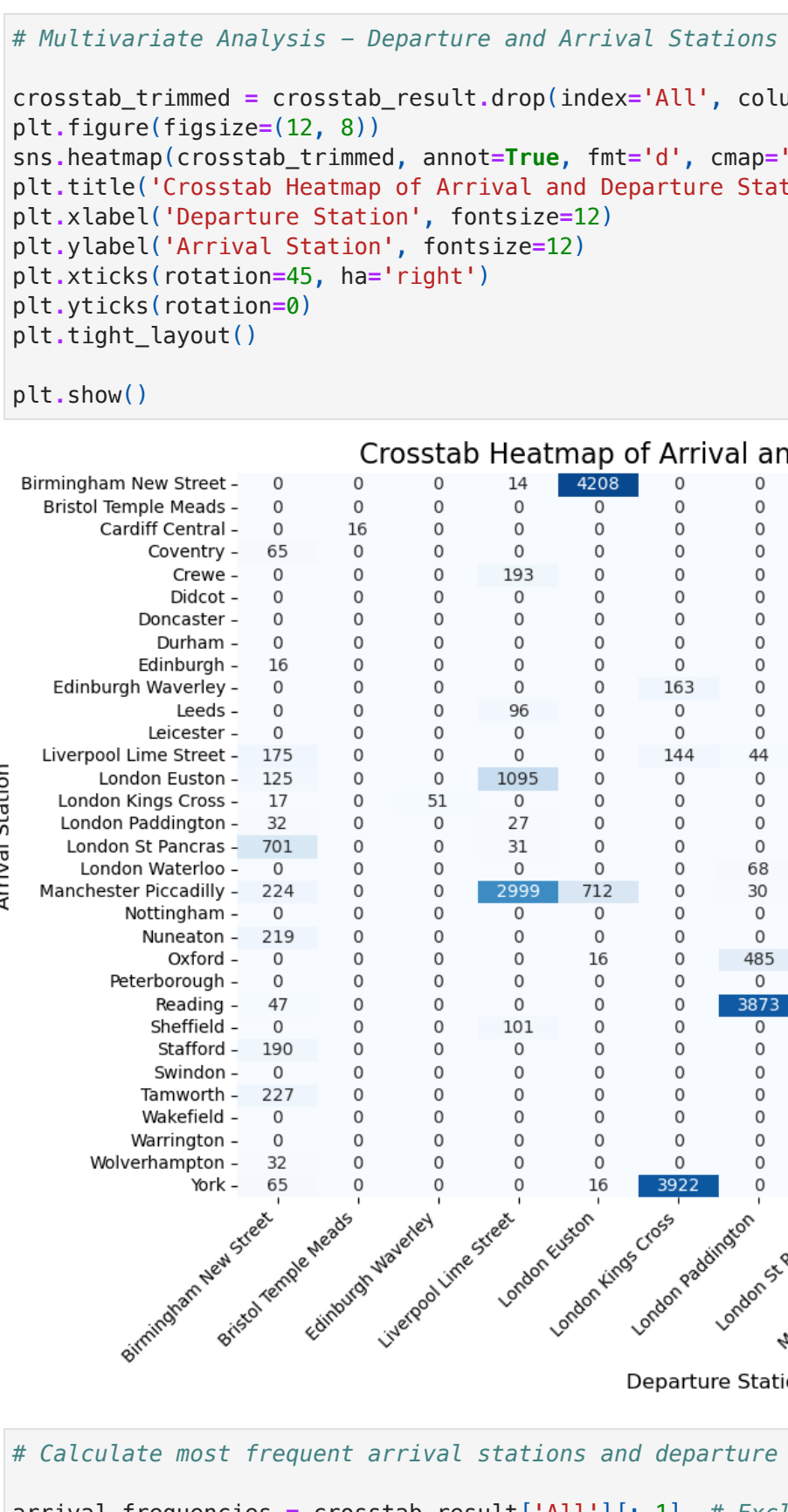
Distribution of Ticket Prices



```
In [63]: # Histogram, kernel density plot and rug plot of Railcard prices
sns.displot(df_clean['Price'], kde=True, rug=True)
```

```
Out[63]: <seaborn.axisgrid.FacetGrid at 0x170dbbce0>
```

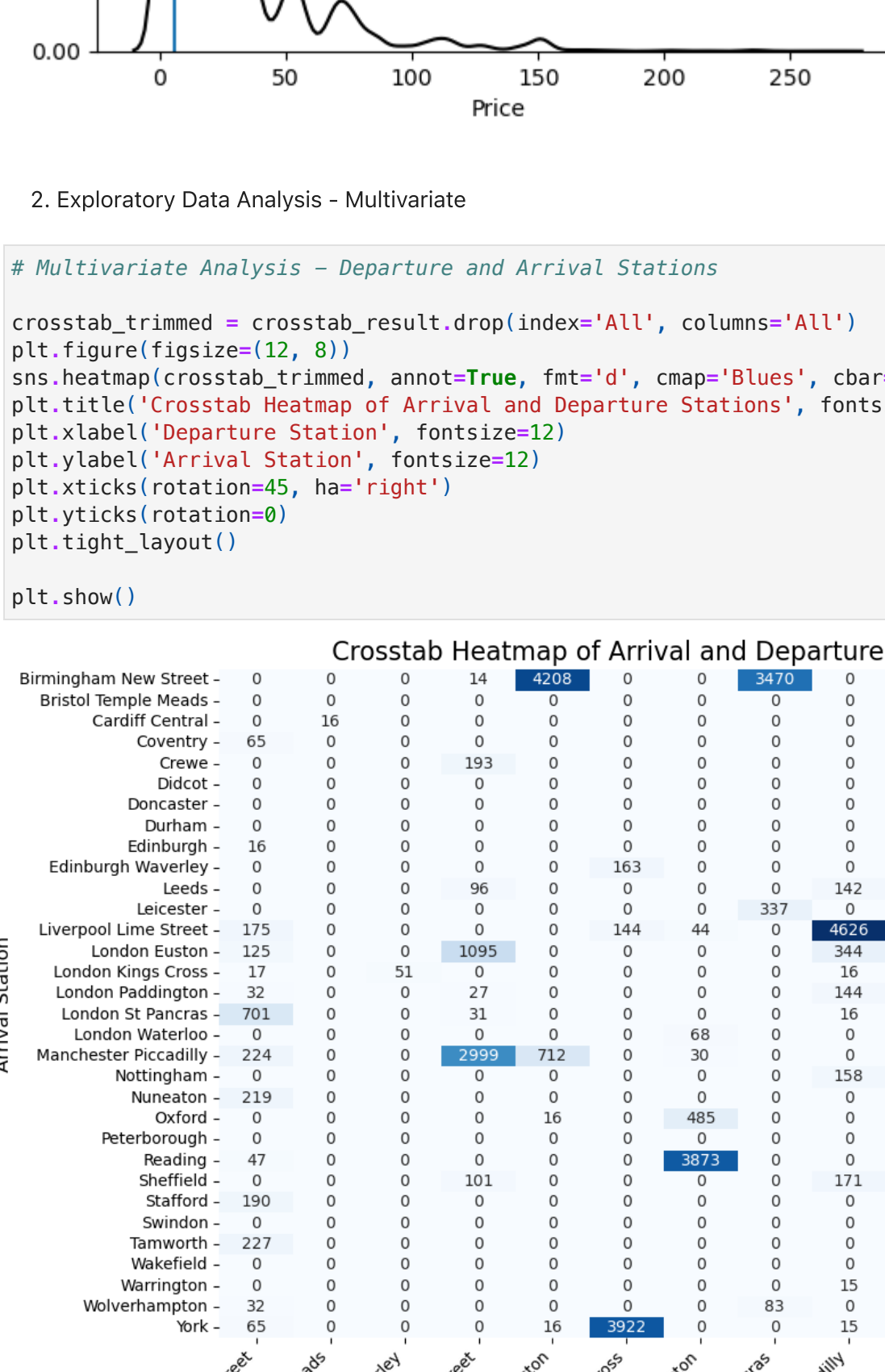
Distribution of Ticket Prices



```
In [65]: # KDE Plot for mode estimation
x = df_clean['Price']
price_kde = sns.kdeplot(x, color='k', label='Kernel Density Estimation of Price')
xg, yg = price_kde.get_lines()[0].get_data()
xg = xg[np.argmax(yg)]

plt.axvline(xg, label=f'Mode: {xg:.2f}')
plt.title('KDE plot with Mode')
plt.xlabel('Price')
plt.ylabel('Density')
plt.legend(loc='upper right')
plt.show()
```

KDE Plot with Mode



2. Exploratory Data Analysis - Multivariate

```
In [68]: # Multivariate Analysis - Departure and Arrival Stations
crosstab_trimmed = crosstab_result.drop('All', columns='All')
plt.figure(figsize=(12, 8))
sns.heatmap(crosstab_trimmed, annot=True, fmt='d', cmap='Blues', cbar=True)
plt.title('Crosstab Heatmap of Arrival and Departure Stations', fontsize=16)
plt.xlabel('Departure Station', fontsize=12)
plt.ylabel('Arrival Station', fontsize=12)
plt.xticks(rotation=45, ha='right')
plt.yticks(rotation=0)
plt.tight_layout()
plt.show()
```

Crosstab Heatmap of Arrival and Departure Stations



```
In [70]: # Calculate most frequent arrival stations and departure stations
arrival_frequencies = crosstab_result['All'][-1:] # Exclude the 'All' row
departure_frequencies = crosstab_result.loc['All'][-1:] # Exclude the 'All' column

plt.figure(figsize=(12, 6))
sns.barplot(
    x=arrival_frequencies.index,
    y=arrival_frequencies.values,
    palette='Blues_r',
    hue=arrival_frequencies.index,
    legend=False
)
plt.title('Most Frequent Arrival Stations', fontsize=16)
plt.xlabel('Frequency', fontsize=12)
plt.ylabel('Frequency', fontsize=12)
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
plt.show()
```

Most Frequent Arrival Stations



```
In [72]: # Plot most frequent arrival and departure stations
plt.figure(figsize=(12, 6))
sns.barplot(
    x=departure_frequencies.index,
    y=departure_frequencies.values,
    palette='Blues_r',
    hue=departure_frequencies.index
)
plt.title('Most Frequent Departure Stations', fontsize=16)
plt.xlabel('Frequency', fontsize=12)
plt.ylabel('Frequency', fontsize=12)
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
plt.show()
```

Most Frequent Departure Stations



```
In [74]: # Identify most popular journeys by grouping 'Departure Station' and 'Arrival Station'
popular_journeys = (
    df_clean.groupby(['Departure.Station', 'Arrival.Station'])
    .size()
    .reset_index(name='Frequency')
    .sort_values(by='Frequency', ascending=False)
)
print(popular_journeys)
```

```
Out[74]: Departure.Station  Arrival.Station  Frequency
40  Manchester Piccadilly  Liverpool Line Street  4208
24  London Euston      Birmingham New Street  4208
36  London Kings Cross      York      3922
35  London Paddington      Reading      3873
36  London St Pancras      Birmingham New Street  3470
...      ...      ...
40  Manchester Piccadilly      York      15
60  London Euston      Edinburgh Waverley      15
47  Manchester Piccadilly      Warrington      15
64  Liverpool Line Street      Wakefield      15
16  Liverpool Line Street      Birmingham New Street  14

[65 rows x 3 columns]
```

```
In [76]: # Plot most popular journeys
top_journeys = popular_journeys.head(1)
plt.figure(figsize=(12, 6))
sns.barplot(
    x=top_journeys.index,
    y=top_journeys.values,
    palette='viridis',
    hue=top_journeys.index
)
plt.title('Most Popular Journeys', fontsize=16)
plt.xlabel('Journey (Departure - Arrival)', fontsize=12)
plt.ylabel('Frequency', fontsize=12)
plt.tight_layout()
plt.show()
```

Most Popular Journeys



```
In [78]: # Multivariate analysis: what is the most common cause for a refund request?
# What are the reasons for delay when a train is cancelled?
reasons_for_cancel = []
cancelled_journeys = df_clean[df_clean['Journey.Status'] == 'Cancelled']
for reason in cancelled_journeys['Reason.for.Delay']:
    if reason not in reasons_for_cancel:
        reasons_for_cancel.append(reason)
print(reasons_for_cancel)
```

```
Out[78]: ['Technical Issue', 'Staffing', 'Staff', 'Signal Failure', 'Weather', 'Traffic']
```



```

In [142.], Payment.Method      0
Railcard                 0
Ticket.Class             0
Ticket.Type              0
Price                   0
Departure.Station        0
Arrival.Station          0
Departure                0
Scheduled.Arrival        0
Actual.Arrival           0
Journey.Status           0
Reason.for.Delay         0
Refund.Request           0
DelayInMinutes           0
MediumPrice              0
dtype: int64

In [142.], journey_not_ontime.head()

Out[142.], Payment.Method      Railcard  Ticket.Class  Ticket.Type  Price  Departure.Station  Arrival.Station  Departure
1      Credit Card      Adult      Standard      Advance  23.0      London Kings Cross      York      2024-01-01 09:45:00
8      Credit Card      None      Standard      Advance  37.0      London Euston          York      2024-01-01 00:00:00
20     Debit Card      Adult      Standard      Advance  7.0      Birmingham New Street  Manchester Piccadilly  2024-01-01 11:15:00
26     Credit Card      Senior    First Class      Advance  34.0      Oxford                 Bristol Temple Meads   2024-01-01 14:15:00
39     Credit Card      None      Standard      Advance  7.0      London Euston          Birmingham New Street  2024-01-02 02:15:00

In [144.], ## Logistic regression model using sklearn

# from sklearn.model_selection import train_test_split
# from sklearn.linear_model import LogisticRegression
# from sklearn.metrics import classification_report, accuracy_score

# X = journey_not_ontime[['MediumPrice']]
# y = journey_not_ontime['Refund.Request']

# X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)

# Fit the logistic regression model
# log_model = LogisticRegression()
# log_model.fit(X_train, y_train)

# Make predictions
# y_pred = log_model.predict(X_test)

# Evaluate the model
# print("Accuracy:", accuracy_score(y_test, y_pred))
# print(classification_report(y_test, y_pred, zero_division=0))

4. Logistic Regression Model with Single Predictor 'MediumPrice'

In [147.], ## Logistic regression model using statsmodels

import statsmodels.api as sm

In [149.], # Convert object types to categorical

journey_not_ontime = journey_not_ontime.apply(Lambda col: col.astype('category') if col.dtypes ==

In [151.], print(journey_not_ontime.dtypes)

Payment.Method      category
Railcard            category
Ticket.Class        category
Ticket.Type         category
Price              float64
Departure.Station   category
Arrival.Station     category
Departure           datetime64[ns]
Scheduled.Arrival   datetime64[ns]
Actual.Arrival      datetime64[ns]
Journey.Status      category
Reason.for.Delay    category
Refund.Request      category
DelayInMinutes      float64
MediumPrice         int64
dtype: object

In [153.], # Concatenate category values to dataset
cat_vars = journey_not_ontime.select_dtypes(include=['category']).columns

for var in cat_vars:
    cat_list = pd.get_dummies(journey_not_ontime[var], prefix=var, drop_first=True)
    journey_not_ontime = pd.concat([journey_not_ontime, cat_list], axis=1)

In [155.], data_vars=journey_not_ontime.columns.values.tolist()
to_keep=[i for i in data_vars if i not in cat_vars]

In [157.], # Final dataset for logistic regression model

journey_final=journey_not_ontime[to_keep]
journey_final.columns.values

Out[157.], array(['Price', 'Departure', 'Scheduled.Arrival', 'Actual.Arrival',
'DelayInMinutes', 'MediumPrice', 'Payment.Method.Credit Card',
'Payment.Method.Debit Card', 'Railcard.Disabled', 'Railcard.None',
'Railcard.Senior', 'Ticket.Class.Standard', 'Ticket.Type.Anytime',
'Ticket.Type.Off-Peak', 'Departure.Station.Edinburgh Waverley',
'Departure.Station.Liverpool Lime Street',
'Departure.Station.London Euston',
'Departure.Station.London Kings Cross',
'Departure.Station.London Paddington',
'Departure.Station.London St Pancras',
'Departure.Station.Manchester Piccadilly',
'Departure.Station.Oxford', 'Departure.Station.Reading',
'Departure.Station.York', 'Arrival.Station.Bristol Temple Meads',
'Arrival.Station.Covey', 'Arrival.Station.Crews',
'Arrival.Station.Didcot', 'Arrival.Station.Doncaster',
'Arrival.Station.Durham', 'Arrival.Station.Edinburgh',
'Arrival.Station.Edinburgh Waverley', 'Arrival.Station.Leeds',
'Arrival.Station.Leicester',
'Arrival.Station.Liverpool Lime Street',
'Arrival.Station.London Euston',
'Arrival.Station.London Kings Cross',
'Arrival.Station.London Paddington',
'Arrival.Station.London St Pancras',
'Arrival.Station.London Waterloo',
'Arrival.Station.Manchester Piccadilly',
'Arrival.Station.Notttingham', 'Arrival.Station.Nuneaton',
'Arrival.Station.Oxford', 'Arrival.Station.Peterborough',
'Arrival.Station.Reading', 'Arrival.Station.Sheffield',
'Arrival.Station.Stafford', 'Arrival.Station.Swindon',
'Arrival.Station.Tamworth', 'Arrival.Station.Wakefield',
'Arrival.Station.Wolverhampton', 'Arrival.Station.York',
'Journey.Status.Delayed', 'Reason.for.Delay.Staff',
'Reason.for.Delay.Staffing', 'Reason.for.Delay.Technical Issue',
'Reason.for.Delay.Traffic', 'Reason.for.Delay.Weather',
'Refund.Request.Yes'], dtype=object)

In [159.], journey_final.head()

Out[159.], Price  Departure  Scheduled.Arrival  Actual.Arrival  DelayInMinutes  MediumPrice  Payment.Method.Credit Card  Refund
1      23.0      2024-01-01 09:45:00      2024-01-01 11:35:00      2024-01-01 11:40:00      5.0      1      True
8      37.0      2024-01-01 00:00:00      2024-01-01 01:50:00      2024-01-01 02:07:00      17.0      0      True
20     7.0      2024-01-01 11:15:00      2024-01-01 12:35:00      2024-01-01 13:06:00      31.0      0      False
26     34.0      2024-01-01 14:15:00      2024-01-01 15:30:00      2024-01-01 15:54:00      24.0      0      True
39     7.0      2024-01-02 02:15:00      2024-01-02 03:35:00      2024-01-02 03:35:00      0.0      0      True
5 rows x 60 columns

In [161.], # Fit small model with one predictor
X_1 = journey_final['MediumPrice'].values.reshape(-1, 1)
X_1 = sm.add_constant(X_1)

# Set response variable
y = journey_final.loc[:, journey_final.columns == 'Refund.Request.Yes']

In [163.], logit_model=sm.Logit(y,X_1)
result=logit_model.fit()
print(result.summary2())

Optimization terminated successfully.
Current function value: 0.578797
Iterations: 5

Results: Logit
=====
Model:                Logit                Method:                MLE
Dependent Variable:    Refund.Request.Yes      Pseudo R-squared:  0.003
Date:                 2024-11-22 00:08          AIC:               4825.3197
No. Observations:     4165                    BIC:               4838.0487
Df Model:             1                      Log-Likelihood:    -2410.7
Df Residuals:         4163                    LL-Null:           -2418.7
Converged:             1.0000                 LLR p-value:       6.1887e-05
No. Iterations:       5.0000                 Scale:            1.0000

=====
Coef.      Std.Err.      z      P>|z|      [0.025      0.975]
-----
const      -1.0753      0.0393      -27.3720      0.0000      -1.1523      -0.9983
x1          0.3548      0.0072      4.8583      0.0000      0.1830      0.5258
=====

In [165.], print(result.params)

const      -1.075328
x1          0.354818
dtype: float64

In [167.], # Plot model results

# Import numpy as np
# Import matplotlib.pyplot as plt

# Model coefficients
Intercept = -1.0753
coef_x1 = 0.3548

# Define the logistic function
def logistic_function(x):
    return 1 / (1 + np.exp(-(Intercept + coef_x1 * x)))

# Generate a range of values for the predictor (e.g., ticket prices)
x1_values = np.linspace(0, 50, 100) # You can change the range based on your data

# Calculate the predicted probabilities for these x1 values
predicted_probabilities = logistic_function(x1_values)

# Create the plot
plt.figure(figsize=(8, 6))
plt.plot(x1_values, predicted_probabilities, label='Probability of Refund Request', color='blue')
plt.title('Logistic Regression: Probability of Refund Request vs. Predictor (MediumPrice)')
plt.xlabel('Predictor (e.g., Ticket Price)')
plt.ylabel('Probability of Refund Request')
plt.grid(True)
plt.legend()
plt.show()

Logistic Regression: Probability of Refund Request vs. Predictor (MediumPrice)

Probability of Refund Request

1.0
0.9
0.8
0.7
0.6
0.5
0.4
0.3
0.2
0.1
0.0

0      10      20      30      40      50

Predictor (e.g., Ticket Price)

Probability of Refund Request

5. Fit Regression Models to To_Predict using MavennRail Predictors

In [170.], # Clean and format To_Predict data to conform with MavennRail data

In [172.], to_predict = pd.read_csv('ToPredict.csv')
to_predict.head()

Out[172.], Payment.Method      Railcard  Ticket.Class  Ticket.Type  Price  Departure.Station  Arrival.Station  Departure
0      Debit Card      NaN      First Class      Advance  54      London St Pancras  Birmingham New Street  2024-01-01 17:45:00
1      Credit Card      NaN      Standard      Advance  7      London Euston      Birmingham New Street  2024-01-09 05:15:00
2      Debit Card      NaN      Standard      Off-Peak  113     Liverpool Lime Street  London Euston          2024-01-09 15:30:00
3      Contactless     Adult      Standard      Off-Peak  3      Liverpool Lime Street  Manchester Piccadilly  2024-01-31 05:45:00
4      Credit Card      NaN      Standard      Off-Peak  4      Manchester Piccadilly  Liverpool Lime Street  2024-02-10 16:00:00

In [174.], # Remove null values. Convert 'nan' for 'No reason' entry in 'Reason for delay' to 'none'.
to_predict['Reason.for.Delay'] = to_predict['Reason.for.Delay'].fillna('NA')

In [176.], # Remove null values. Convert 'nan' for 'no railcard' entry in 'Railcards' to 'none'.
to_predict['Railcard'] = to_predict['Railcard'].fillna('None')

In [178.], to_predict['Actual.Arrival'] = to_predict['Actual.Arrival'].fillna(to_predict['Scheduled.Arrival'])

In [180.], to_predict.isnull().sum()

Out[180.], Payment.Method      0
Railcard                 0
Ticket.Class             0
Ticket.Type              0
Price                   0
Departure.Station        0
Arrival.Station          0
Departure                0
Scheduled.Arrival        0
Actual.Arrival           0
Journey.Status           0
Reason.for.Delay         0
dtype: int64

In [182.], # Convert string values for 'Actual.Arrival', 'Departure', and 'Scheduled.Arrival' to datetime64[ns]
to_predict['Actual.Arrival'] = pd.to_datetime(to_predict['Actual.Arrival'], format='%Y-%m-%d %H:%M')
to_predict['Departure'] = pd.to_datetime(to_predict['Departure'], format='%Y-%m-%d %H:%M')
to_predict['Scheduled.Arrival'] = pd.to_datetime(to_predict['Scheduled.Arrival'], format='%Y-%m-%d %H:%M')

In [184.], to_predict.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8 entries, 0 to 7
Data columns (total 12 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Payment.Method      8 non-null      object
1   Railcard            8 non-null      object
2   Ticket.Class        8 non-null      object
3   Ticket.Type         8 non-null      object
4   Price              8 non-null      int64
5   Departure.Station   8 non-null      object
6   Arrival.Station     8 non-null      object
7   Departure           8 non-null      datetime64[ns]
8   Scheduled.Arrival   8 non-null      datetime64[ns]
9   Actual.Arrival      8 non-null      datetime64[ns]
10  Journey.Status       8 non-null      object
11  Reason.for.Delay     8 non-null      object
12  DelayInMinutes       8 non-null      float64
dtypes: datetime64[ns](3), float64(1), int64(1), object(8)
memory usage: 960.0+ bytes

In [186.], Payment.Method      Railcard  Ticket.Class  Ticket.Type  Price  Departure.Station  Arrival.Station  Departure
0      Debit Card      None      First Class      Advance  54      London St Pancras  Birmingham New Street  2024-01-01 17:45:00
1      Credit Card      None      Standard      Advance  7      London Euston      Birmingham New Street  2024-01-09 05:15:00
2      Debit Card      None      Standard      Off-Peak  113     Liverpool Lime Street  London Euston          2024-01-09 15:30:00
3      Contactless     Adult      Standard      Off-Peak  3      Liverpool Lime Street  Manchester Piccadilly  2024-01-31 05:45:00
4      Credit Card      None      Standard      Off-P
```