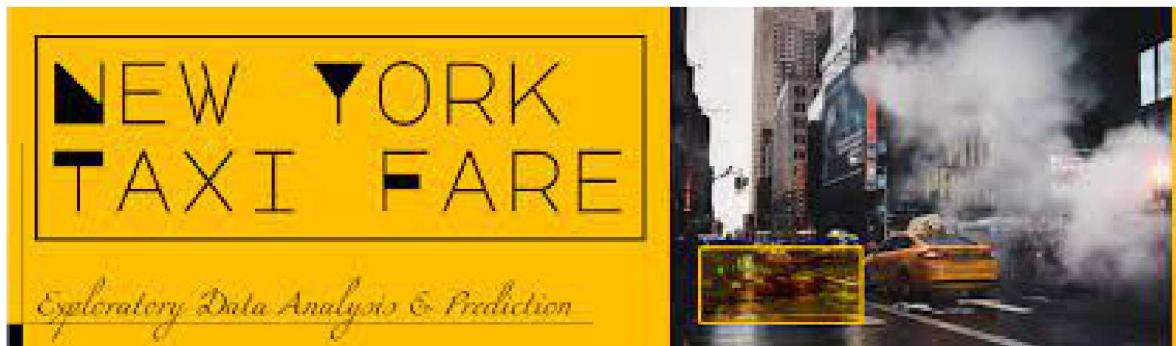


In [83]:

```

1 import matplotlib.pyplot as plt
2 import matplotlib.image as mpimg
3
4 # Load the image
5 img_path = r"C:\Users\jangi\Downloads\taxi.jpg"
6 img = mpimg.imread(img_path)
7
8 # Set the desired size (adjust width and height as needed)
9 fig, ax = plt.subplots(figsize=(20, 16))
10
11 # Display the image
12 ax.imshow(img)
13 ax.axis('off') # Turn off axis labels
14 plt.show()
15

```



The objective of this dataset could be to analyze and model factors influencing

- car cancellations in a transportation service. This analysis might involve understanding patterns in booking behavior,
- the impact of online and mobile site bookings, the influence of travel types and areas,
- and whether cancellations are correlated with certain time periods or locations.
- The ultimate goal could be to develop a predictive model that can identify factors contributing to cancellations
- and potentially reduce them, thereby improving the efficiency of the cab service.

Here's a brief interpretation of each column:

- 1, user_id: Identifies the unique user associated with the booking.
- 2, vehicle_model_id: Represents the model of the vehicle used for the transportation service.
- 3, travel_type_id: Indicates the type of travel (e.g., business, personal) associated with the booking.
- 4, from_area_id: Denotes the area or location from which the journey originates.
- 5, to_area_id: Represents the destination area or location for the journey.
- 6, from_date: Specifies the date and time when the journey starts.

- 7, online_booking: Binary indicator (0 or 1) of whether the booking was made online.
 - 8, mobile_site_booking: Binary indicator of whether the booking was made through a mobile site.
 - 9, booking_created: Indicates the date and time when the booking was created.
 - 10, from_lat: Latitude of the starting location.
 - 11, from_long: Longitude of the starting location.
 - 12, to_lat: Latitude of the destination location.
 - 13, to_long: Longitude of the destination location.
 - 14, Car_Cancellation: Binary indicator (0 or 1) of whether the car booking was canceled
- The dataset seems structured to capture various aspects of cab bookings, including user details, vehicle information,
- travel type, locations, timing, and whether a booking was canceled.
 - Analyzing this dataset could involve exploring patterns in booking behavior,
 - understanding factors influencing cancellations, and potentially developing models to predict
 - and mitigate cancellations for improved service efficiency.

Machine Learning Algorithms

Model Building :

- Linear Regression
- Decision Tree
- Random Forest
- SVM
- k-fold cross-validation

In [1]:

```

1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6 from sklearn import linear_model
7 import warnings
8 warnings.filterwarnings('ignore')

```

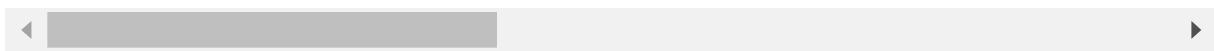
In [2]:

```
1 df = pd.read_csv(r"C:\Users\jangi\Downloads\YourCabs.csv")
```

In [3]: 1 df.head()

Out[3]:

	id	user_id	vehicle_model_id	package_id	travel_type_id	from_area_id	to_area_id	from
0	132512	22177		28	NaN	2	83.0	448.0
1	132513	21413		12	NaN	2	1010.0	540.0
2	132514	22178		12	NaN	2	1301.0	1034.0
3	132515	13034		12	NaN	2	768.0	398.0
4	132517	22180		12	NaN	2	1365.0	849.0



In [4]: 1 df.shape

Out[4]: (43431, 18)

In [5]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 43431 entries, 0 to 43430
Data columns (total 18 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   id               43431 non-null   int64  
 1   user_id          43431 non-null   int64  
 2   vehicle_model_id 43431 non-null   int64  
 3   package_id        7550 non-null    float64 
 4   travel_type_id   43431 non-null   int64  
 5   from_area_id      43343 non-null   float64 
 6   to_area_id        34293 non-null   float64 
 7   from_city_id     16345 non-null   float64 
 8   to_city_id       1588 non-null    float64  
 9   from_date         43431 non-null   object  
 10  online_booking   43431 non-null   int64  
 11  mobile_site_booking 43431 non-null   int64  
 12  booking_created  43431 non-null   object  
 13  from_lat          43338 non-null   float64 
 14  from_long         43338 non-null   float64 
 15  to_lat            34293 non-null   float64 
 16  to_long           34293 non-null   float64 
 17  Car_Cancellation 43431 non-null   int64  
dtypes: float64(9), int64(7), object(2)
memory usage: 6.0+ MB
```

In [6]: 1 df.describe().T.style.background_gradient(cmap = 'Reds')

Out[6]:

		count	mean	std	min	25%
	id	43431.000000	159206.473556	15442.386279	132512.000000	145778.000000
	user_id	43431.000000	30739.198153	10996.476709	16.000000	24614.000000
	vehicle_model_id	43431.000000	25.717230	26.798250	1.000000	12.000000
	package_id	7550.000000	2.030066	1.461756	1.000000	1.000000
	travel_type_id	43431.000000	2.137252	0.437712	1.000000	2.000000
	from_area_id	43343.000000	714.544494	419.883553	2.000000	393.000000
	to_area_id	34293.000000	669.490917	400.638225	2.000000	393.000000
	from_city_id	16345.000000	14.915081	1.165306	1.000000	15.000000
	to_city_id	1588.000000	68.537783	49.880732	4.000000	32.000000
	online_booking	43431.000000	0.351592	0.477473	0.000000	0.000000
	mobile_site_booking	43431.000000	0.043241	0.203402	0.000000	0.000000
	from_lat	43338.000000	12.982461	0.085933	12.776630	12.926450
	from_long	43338.000000	77.636255	0.059391	77.386930	77.593661
	to_lat	34293.000000	13.026648	0.113487	12.776630	12.951850
	to_long	34293.000000	77.640595	0.064045	77.386930	77.582030
	Car_Cancellation	43431.000000	0.072114	0.258680	0.000000	0.000000

◀ ▶

In [7]: 1 df.isna().sum().sort_values(ascending=True)

Out[7]:

id	0
booking_created	0
mobile_site_booking	0
online_booking	0
from_date	0
Car_Cancellation	0
travel_type_id	0
vehicle_model_id	0
user_id	0
from_area_id	88
from_lat	93
from_long	93
to_area_id	9138
to_long	9138
to_lat	9138
from_city_id	27086
package_id	35881
to_city_id	41843
dtype: int64	

```
In [8]: 1 df.isna().sum()/len(df)*100
```

```
Out[8]: id           0.000000
         user_id       0.000000
         vehicle_model_id 0.000000
         package_id      82.616104
         travel_type_id   0.000000
         from_area_id     0.202620
         to_area_id       21.040271
         from_city_id     62.365591
         to_city_id       96.343626
         from_date        0.000000
         online_booking    0.000000
         mobile_site_booking 0.000000
         booking_created   0.000000
         from_lat          0.214133
         from_long          0.214133
         to_lat             21.040271
         to_long            21.040271
         Car_Cancellation  0.000000
         dtype: float64
```

```
In [9]: 1 df.duplicated()
```

```
Out[9]: 0      False
        1      False
        2      False
        3      False
        4      False
        ...
       43426  False
       43427  False
       43428  False
       43429  False
       43430  False
Length: 43431, dtype: bool
```

```
In [10]: 1 df.nunique().sort_values(ascending=True)
```

```
Out[10]: Car_Cancellation      2
          mobile_site_booking    2
          online_booking         2
          travel_type_id         3
          from_city_id           3
          package_id              7
          vehicle_model_id        27
          to_city_id               116
          to_long                  447
          to_lat                   450
          from_long                462
          from_lat                 466
          to_area_id                568
          from_area_id              598
          from_date                 20377
          user_id                  22267
          booking_created            39349
          id                        43431
          dtype: int64
```

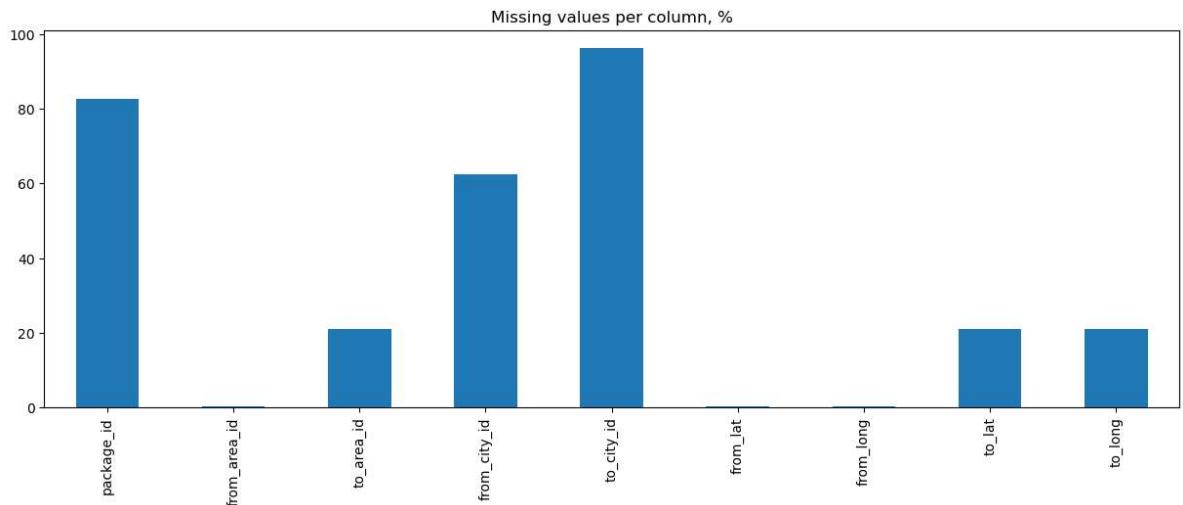
```
In [11]: 1 df.isna().sum()/len(df)*100
```

```
Out[11]: id                    0.000000
          user_id                0.000000
          vehicle_model_id        0.000000
          package_id              82.616104
          travel_type_id          0.000000
          from_area_id             0.202620
          to_area_id                21.040271
          from_city_id              62.365591
          to_city_id                 96.343626
          from_date                 0.000000
          online_booking            0.000000
          mobile_site_booking       0.000000
          booking_created            0.000000
          from_lat                  0.214133
          from_long                  0.214133
          to_lat                     21.040271
          to_long                     21.040271
          Car_Cancellation           0.000000
          dtype: float64
```

Data cleaning and preprocessing

In [12]:

```
1 #find percentage of missing values for each column
2 listings_missing_df = df.isnull().mean()*100
3
4 #filter out only columns, which have missing values
5 listings_columns_with_nan = listings_missing_df[listings_missing_df > 0]
6
7 #plot the results
8 listings_columns_with_nan.plot.bar(title='Missing values per column, %',x=
```



Removing Unnecessary Columns:

In [13]:

```
1 df.drop(['package_id', 'from_city_id', 'to_city_id'], axis="columns", inplace=True)
2
```

In [14]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 43431 entries, 0 to 43430
Data columns (total 15 columns):
 #   Column           Non-Null Count Dtype  
 --- 
 0   id               43431 non-null  int64   
 1   user_id          43431 non-null  int64   
 2   vehicle_model_id 43431 non-null  int64   
 3   travel_type_id   43431 non-null  int64   
 4   from_area_id     43343 non-null  float64 
 5   to_area_id       34293 non-null  float64 
 6   from_date        43431 non-null  object  
 7   online_booking   43431 non-null  int64   
 8   mobile_site_booking 43431 non-null  int64   
 9   booking_created  43431 non-null  object  
 10  from_lat         43338 non-null  float64 
 11  from_long        43338 non-null  float64 
 12  to_lat           34293 non-null  float64 
 13  to_long          34293 non-null  float64 
 14  Car_Cancellation 43431 non-null  int64   
dtypes: float64(6), int64(7), object(2)
memory usage: 5.0+ MB
```

Handling Missing Values:

In [15]: 1 mean_from_area_id = df.from_area_id.mean()
2 mean_from_area_id
3
4 import math
5 mean_from_area_id = math.floor(df['from_area_id'].mean())
6 mean_from_area_id
7
8 df['from_area_id'].fillna(mean_from_area_id, inplace = True)
9

In [16]: 1 mean_from_area_id = df.from_area_id.mean()
2 mean_from_area_id
3
4 import math
5 mean_to_area_id = math.floor(df['to_area_id'].mean())
6 mean_to_area_id
7
8 df['to_area_id'].fillna(mean_from_area_id, inplace = True)

```
In [17]: 1 median_from_lat = df.from_lat.median()
2 median_from_lat
3
4 import math
5 median_from_lat = math.floor(df['from_lat'].median())
6 median_from_lat
7
8 df['from_lat'].fillna(median_from_lat, inplace = True)
```

```
In [18]: 1 median_from_long = df.from_long.median()
2 median_from_long
3
4 import math
5 median_from_long = math.floor(df['from_long'].median())
6 median_from_long
7
8 df['from_long'].fillna(median_from_long, inplace = True)
```

```
In [19]: 1 median_to_lat = df.to_lat.median()
2 median_to_lat
3
4 import math
5 median_to_lat = math.floor(df['to_lat'].median())
6 median_to_lat
7
8 df['to_lat'].fillna(median_to_lat, inplace = True)
```

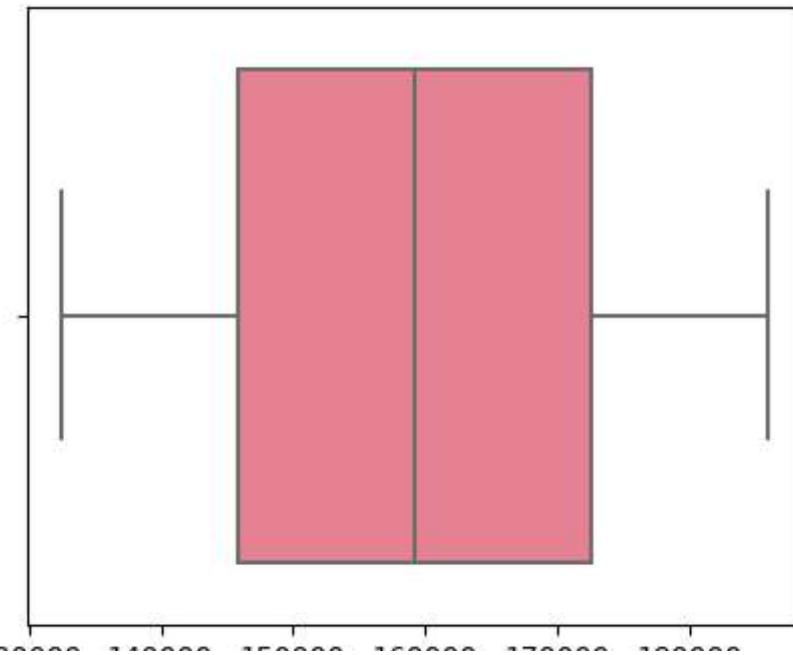
```
In [20]: 1 median_to_long = df.to_long.median()
2 median_to_long
3
4 import math
5 median_to_long = math.floor(df['to_long'].median())
6 median_to_long
7
8 df['to_long'].fillna(median_to_long, inplace = True)
```

```
In [21]: 1 df.isnull().sum()
```

```
Out[21]: id          0  
user_id      0  
vehicle_model_id 0  
travel_type_id 0  
from_area_id   0  
to_area_id     0  
from_date      0  
online_booking  0  
mobile_site_booking 0  
booking_created 0  
from_lat        0  
from_long       0  
to_lat          0  
to_long         0  
Car_Cancellation 0  
dtype: int64
```

Feature Scaling:

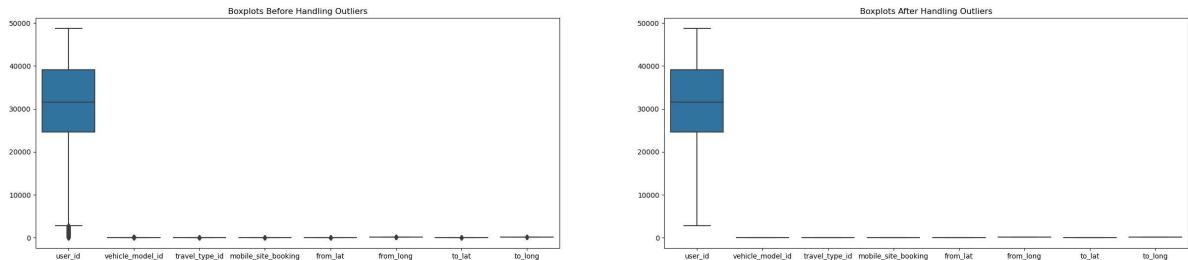
```
In [22]: 1 def boxplots(col):  
2     plt.figure(figsize=(5,4))  
3     sns.boxplot(df,x=col,palette='husl')  
4     plt.show()  
5  
6 # Assuming 'df' is your DataFrame  
7 for i in list(df.select_dtypes(exclude=['object']).columns)[0:]:  
8     boxplots(i)
```



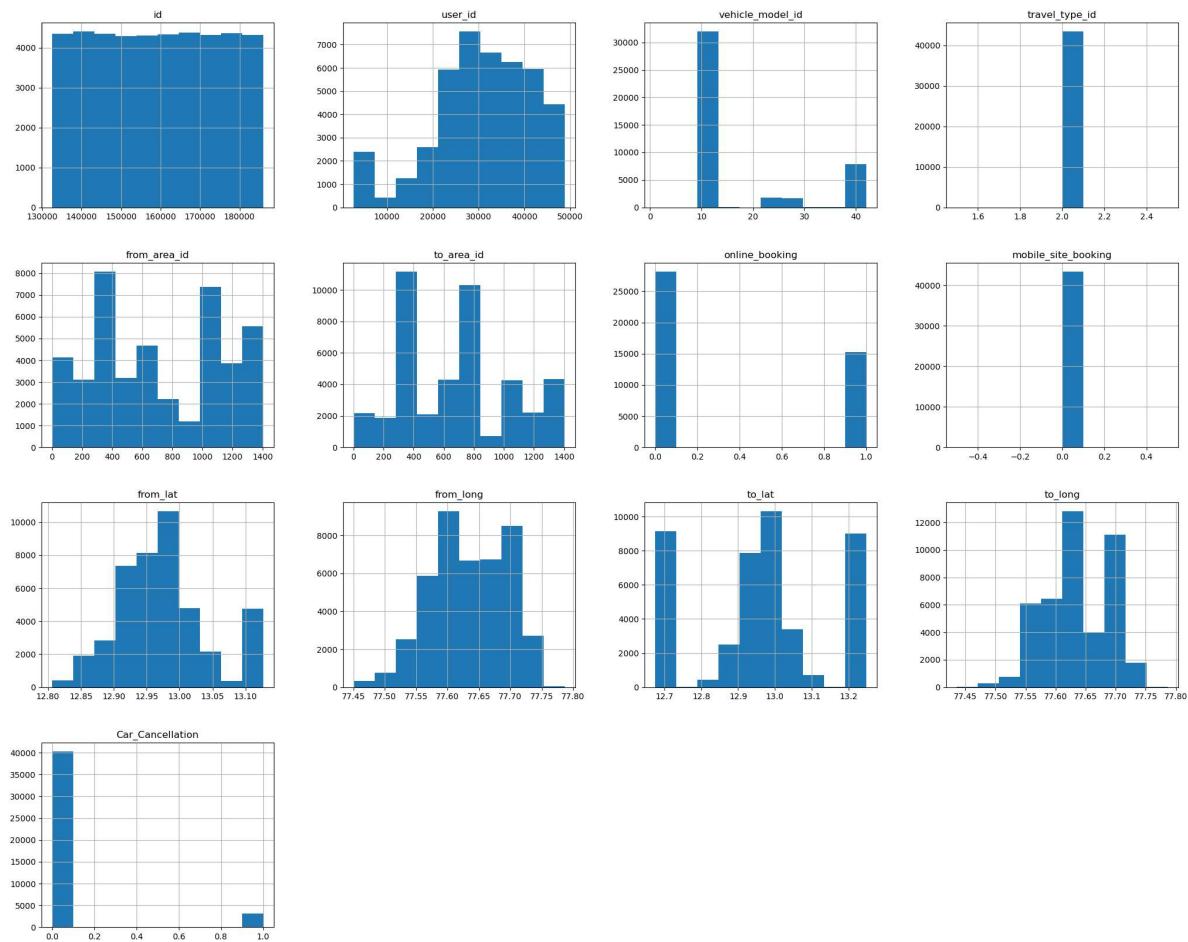
Handling Outliers:

In [23]:

```
1 import pandas as pd
2 import numpy as np
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5
6 # Load your own dataset
7 # Replace 'your_dataset.csv' with the actual file path or URL of your data
8 df = df
9
10 # Function to handle outliers using clip
11 def handle_outliers(dataframe, col):
12     q3 = dataframe[col].quantile(0.75)
13     q1 = dataframe[col].quantile(0.25)
14     IQR = q3 - q1
15     Lower = q1 - 1.5 * IQR
16     Upper = q3 + 1.5 * IQR
17
18     # Clip the values to be within the Lower and Upper bounds
19     dataframe[col] = dataframe[col].clip(lower=Lower, upper=Upper)
20
21 # Specify the columns to handle outliers
22 columns_to_handle_outliers = [
23     'user_id',
24     'vehicle_model_id',
25     'travel_type_id',
26     'mobile_site_booking' ,
27     'from_lat' ,
28     'from_long' ,
29     'to_lat' ,
30     'to_long' ,
31 ]
32
33 # Plot boxplots before handling outliers
34 plt.figure(figsize=(30, 6))
35 plt.subplot(1, 2, 1)
36 sns.boxplot(data=df[columns_to_handle_outliers].dropna())
37 plt.title('Boxplots Before Handling Outliers')
38
39 # Handle outliers for each specified column
40 for col in columns_to_handle_outliers:
41     handle_outliers(df, col)
42
43 # Plot boxplots after handling outliers
44 plt.subplot(1, 2, 2)
45 sns.boxplot(data=df[columns_to_handle_outliers].dropna())
46 plt.title('Boxplots After Handling Outliers')
47
48 plt.show()
49
```

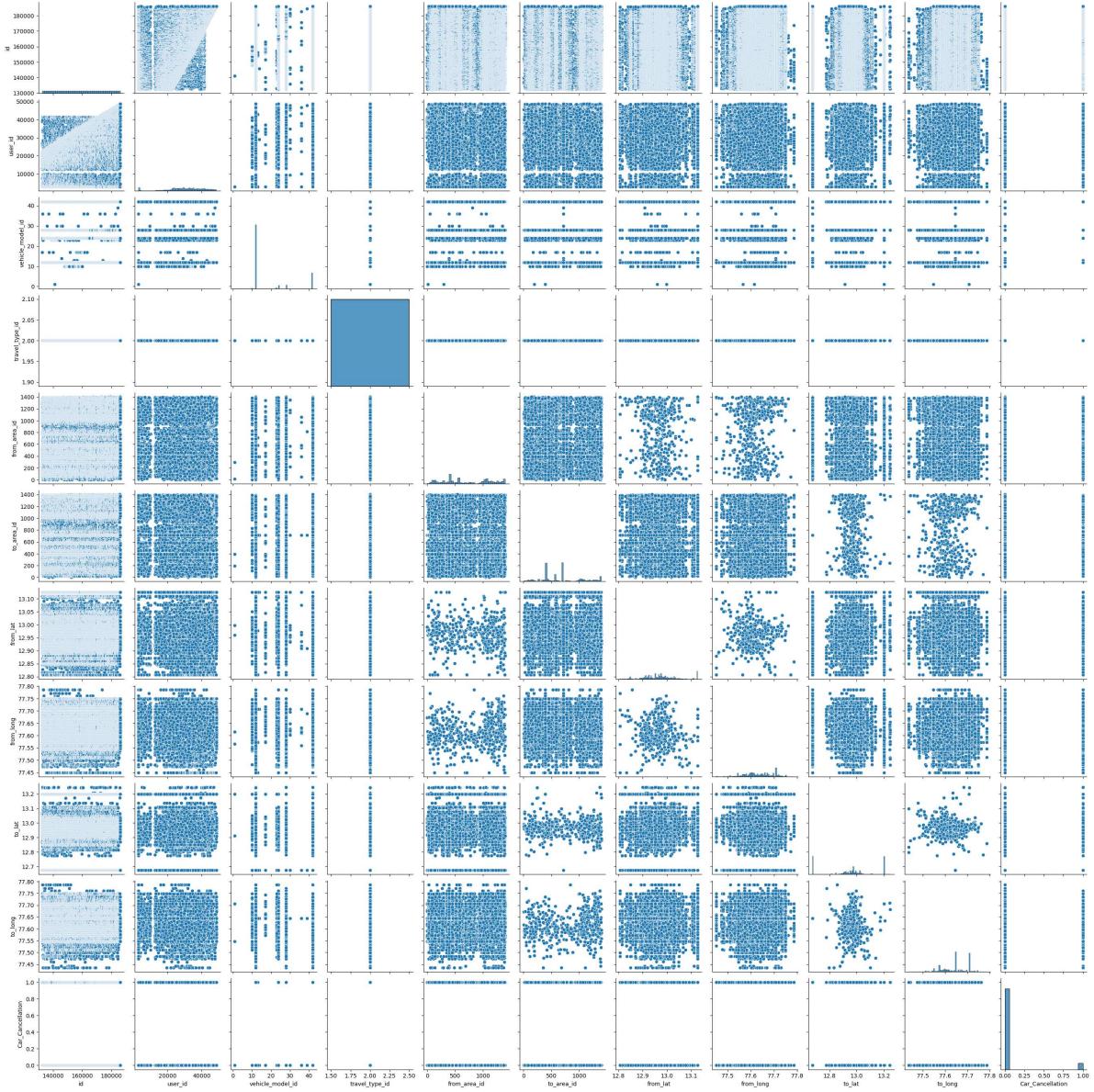


```
In [24]: 1 df.hist(figsize = (25,20))
2 plt.show()
```



In [79]:

```
1 sns.pairplot(df[['id',
2     'user_id',
3     'vehicle_model_id',
4     'travel_type_id',
5     'from_area_id',
6     'to_area_id',
7     'from_lat',
8     'from_long',
9     'to_lat',
10    'to_long',
11    'Car_Cancellation']])
12 plt.show()
```

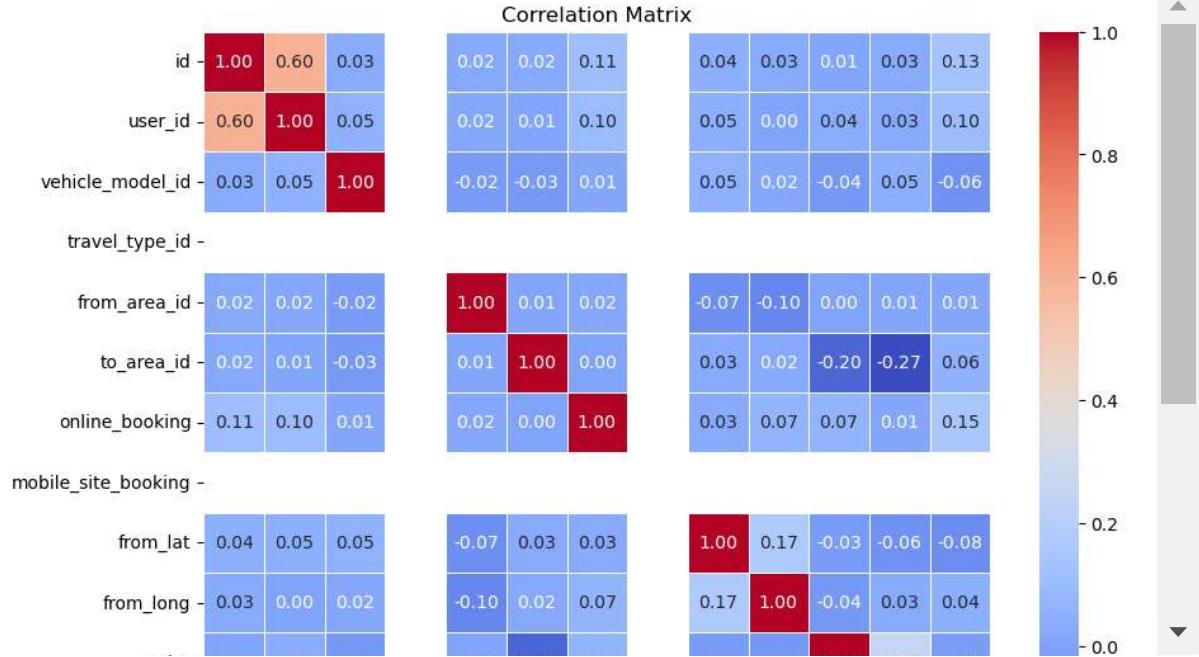


In [26]:

```

1 # Plot a heatmap of the correlation matrix
2 plt.figure(figsize=(10, 8))
3 sns.heatmap(df.corr(), annot=True, cmap='coolwarm', fmt='.2f', linewidths=0.5)
4 plt.title('Correlation Matrix')
5 plt.show()

```



Create a deep copy of dataset:

In [27]:

```

1 import copy
2
3 # Create a deep copy
4 deep_copy_df = df.copy(deep=True)

```

```
In [28]: 1 # Create a deep copy
          2 deep_copy_df.head()
```

Out[28]:

	id	user_id	vehicle_model_id	travel_type_id	from_area_id	to_area_id	from_date	online	
0	132512	22177.0		28		2	83.0	448.0	01-01-2013 02:00
1	132513	21413.0		12		2	1010.0	540.0	01-01-2013 09:00
2	132514	22178.0		12		2	1301.0	1034.0	01-01-2013 03:30
3	132515	13034.0		12		2	768.0	398.0	01-01-2013 05:45
4	132517	22180.0		12		2	1365.0	849.0	01-01-2013 09:00

Model Building :

- Linear Regression
- Decision Tree
- Random Forest
- SVM
- k-fold cross-validation

Linear Regression :

In this dataset, linear regression is used to model the relationship between the selected features (X) and the target variable "Car_Cancellation" (y). Specifically, the code I provided earlier splits the data into training and testing sets, then trains a linear regression model on the training set and evaluates its performance on the testing set.

Here's a summary of what happens for both the training and testing sets:

1. Training Set (X_train, y_train):

- The linear regression model is trained using the training features (X_train) and the corresponding target variable (y_train).
- The model learns the coefficients and intercept that minimize the difference between the predicted and actual values of "Car_Cancellation" in the training set.

2. Testing Set (X_test, y_test):

- The trained model is then used to make predictions on the testing features (X_{test}).
- The predictions (y_{pred}) are compared with the actual values of "Car_Cancellation" in the testing set (y_{test}).
- Evaluation metrics, such as Mean Squared Error (MSE) and R-squared, are calculated to assess how well the model performs on the testing set.

Idea is to check how well the model generalizes to new, unseen data (testing set) after being trained on the training set. The coefficients and intercept learned during training are then applied to the testing set to make predictions and evaluate the model's performance.

Here's a simplified summary of the process:

- Training:

- Train the linear regression model using X_{train} and y_{train} .
- Learn coefficients and intercept.

- Testing:

- Use the trained model to predict "Car_Cancellation" on X_{test} .
- Compare predictions (y_{pred}) with actual values (y_{test}).
- Evaluate the model's performance using metrics like MSE and R-squared on the testing set.

This process allows you to understand how well the linear regression model generalizes to new data and whether it can accurately predict "Car_Cancellation" based on the

It appears that the dataset is already arranged in tabular form with columns representing different attributes. Each row corresponds to a specific entry in the dataset. The columns include:

- 1. id : A unique identifier for each entry.
- 2. user_id : User identifier.
- 3. vehicle_model_id : Identifier for the vehicle model.
- 4. travel_type_id : Identifier for the travel type.
- 5. from_area_id : Identifier for the starting area.
- 6. to_area_id : Identifier for the destination area.
- 7. from_date : Date and time of the journey starting.
- 8. online_booking : Binary indicator (0 or 1) for online booking.
- 9. mobile_site_booking : Binary indicator (0 or 1) for mobile site booking.
- 10. booking_created : Date and time when the booking was created.
- 11. from_lat : Latitude of the starting location.
- 12. from_long : Longitude of the starting location.
- 13. to_lat : Latitude of the destination location.
- 14. to_long : Longitude of the destination location.
- 15. Car_Cancellation : Binary indicator (0 or 1) for car cancellation.

Library for ModelBuilding

In [29]:

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error, r2_score
4 from sklearn.preprocessing import StandardScaler
5 from datetime import datetime

```

split data in the independent and dependent variable

◀ ▶

In [30]:

```

1 # # Convert date-time strings to numerical representation
2 # df['from_date'] = pd.to_datetime(df['from_date']).astype(int) / 10**9 #
3 # df['booking_created'] = pd.to_datetime(df['booking_created']).astype(int)

```

In [31]:

```

1 # Select features (independent variables) and the target variable (y)
2 # (both code X(features), Y(target) are working the same)
3 # X = df[['from_area_id', 'to_area_id', 'from_lat', 'from_long', 'to_lat',
4 #          'to_long']]
5
6
7 x = df.drop(['Car_Cancellation', 'from_date', 'booking_created'], axis = 1)

```

In [32]:

```
1 x.head()
```

Out[32]:

	id	user_id	vehicle_model_id	travel_type_id	from_area_id	to_area_id	online_booking	Car_Cancellation
0	132512	22177.0		28	2	83.0	448.0	0
1	132513	21413.0		12	2	1010.0	540.0	0
2	132514	22178.0		12	2	1301.0	1034.0	0
3	132515	13034.0		12	2	768.0	398.0	0
4	132517	22180.0		12	2	1365.0	849.0	0

◀ ▶

In [33]:

```
1 y = df[['Car_Cancellation']]
```

In [34]: 1 y.head()

Out[34]: **Car_Cancellation**

0	0
1	0
2	0
3	0
4	0

In [35]: 1 # Split the data into training and testing sets

2 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, r

- The random_state parameter in the train_test_split function is used to ensure reproducibility.
- When you split your dataset into a training set and a testing set,
- the function randomly shuffles the data before the split.
- The random_state parameter allows you to set a seed for the random number generator,
- ensuring that the split is the same every time you run the code with the same random_state value
- The test_size parameter in the train_test_split function determines the proportion of the dataset that should be allocated to the testing set. It is expressed as a fraction between 0 and 1.

In [36]: 1 # Standardize the features

2
3 scaler = StandardScaler()
4 x_train_scaled = scaler.fit_transform(x_train)
5 x_test_scaled = scaler.transform(x_test)

In [37]: 1 # Train the Linear regression model

2 model = LinearRegression()
3 model.fit(x_train_scaled, y_train)

Out[37]: LinearRegression()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [38]: 1 # Print coefficients and intercept
2 print('Coefficients:', model.coef_)
3 print('Intercept:', model.intercept_)
```

```
Coefficients: [[ 2.60066861e-02  1.05901660e-02 -1.49229009e-02  6.93889390e-18
   1.54232220e-03  6.64009455e-03  3.68165398e-02  3.46944695e-18
   -2.49852487e-02  1.35334266e-02 -1.98421158e-04 -2.59994250e-02]]
Intercept: [0.07186852]
```

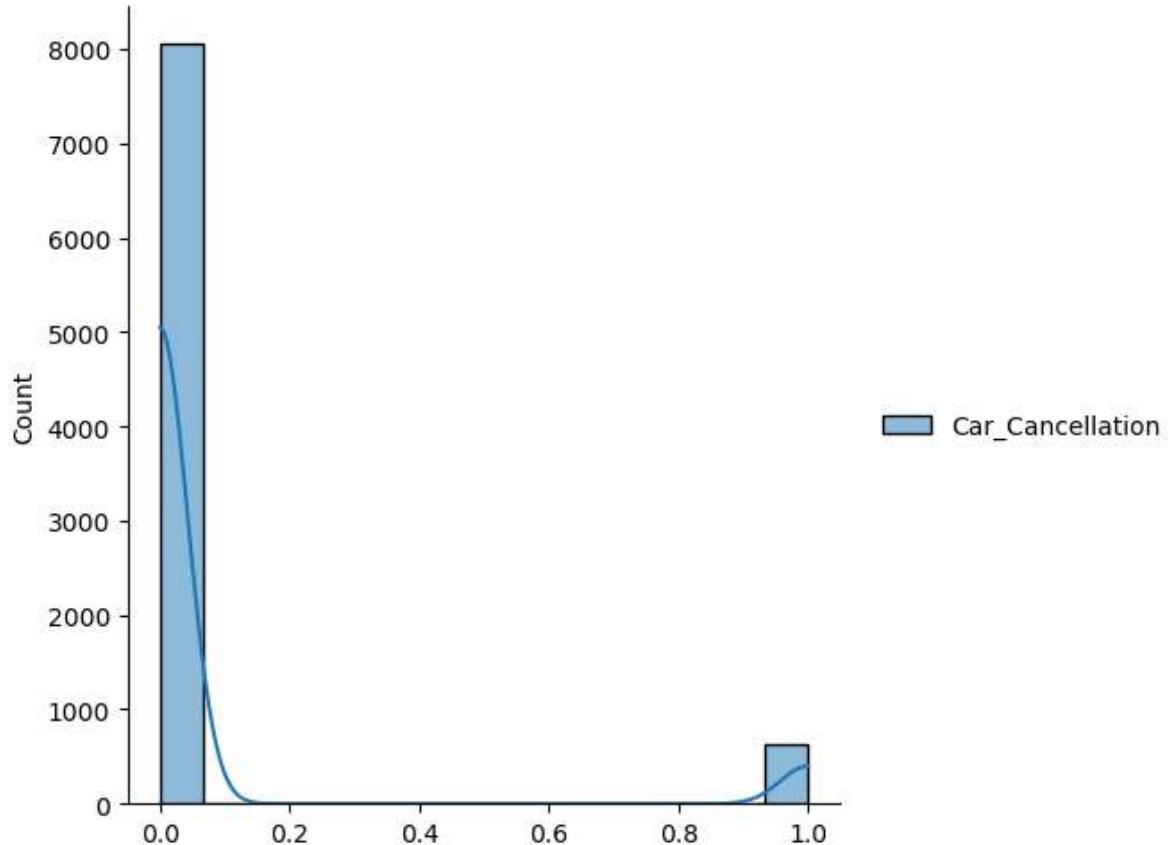
```
In [39]: 1 # Make predictions on the test set
2 y_pred = model.predict(x_test)
3 y_pred_train = model.predict(x_train)
```

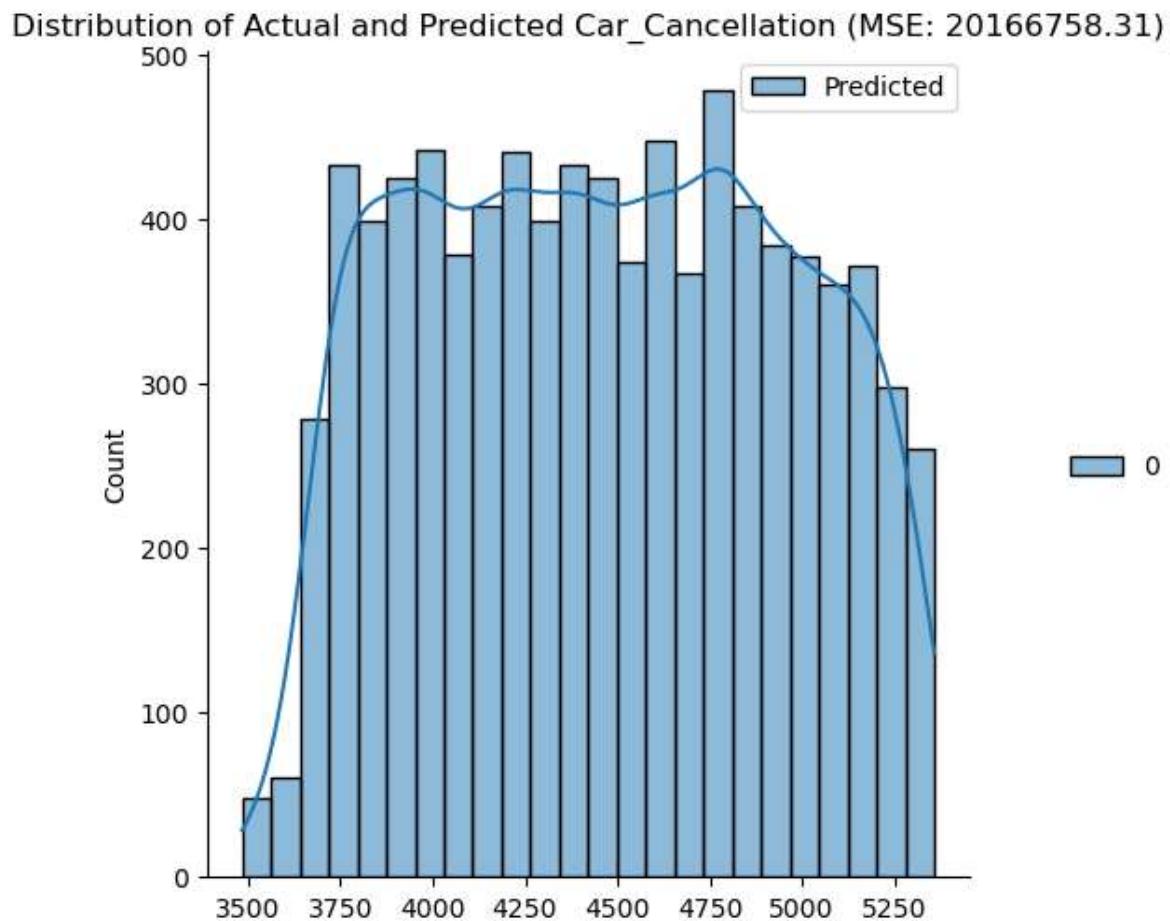
```
In [40]: 1 # Evaluate the model
2 mse = mean_squared_error(y_test, y_pred)
3 r2 = r2_score(y_test, y_pred)
4
5 print(f'Mean Squared Error: {mse}')
6 print(f'R-squared: {r2}')
```

```
Mean Squared Error: 20166758.3090341
R-squared: -297644761.5597869
```

In [41]:

```
1 # Create a displot for actual and predicted values
2 sns.displot(y_test, kde=True, label='Actual')
3 sns.displot(y_pred, kde=True, label='Predicted')
4 plt.title(f'Distribution of Actual and Predicted Car_Cancellation (MSE: {r})
5 plt.legend()
6 plt.show()
```





In [42]: 1 df.head()

Out[42]:

	id	user_id	vehicle_model_id	travel_type_id	from_area_id	to_area_id	from_date	online
0	132512	22177.0		28	2	83.0	448.0	01-01-2013 02:00
1	132513	21413.0		12	2	1010.0	540.0	01-01-2013 09:00
2	132514	22178.0		12	2	1301.0	1034.0	01-01-2013 03:30
3	132515	13034.0		12	2	768.0	398.0	01-01-2013 05:45
4	132517	22180.0		12	2	1365.0	849.0	01-01-2013 09:00

Model - Decision Tree

In [43]:

```

1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.tree import DecisionTreeClassifier, plot_tree
4 from sklearn.preprocessing import LabelEncoder
5 from sklearn.metrics import accuracy_score, confusion_matrix
6 import matplotlib.pyplot as plt
7
8 # Load the dataset
9 # Assuming the dataset is stored in a variable named 'df'
10 # If not, replace 'df' with the actual variable name you are using

```

1. Data Loading and Preprocessing:

In [44]:

```
1 df = df.drop(['id', 'user_id', 'from_date', 'booking_created'], axis=1)
```

- This line drops unnecessary columns ('id', 'user_id', 'from_date', 'booking_created') from the dataset.
- The 'from_date' column is removed to simplify the example, and if you need it for your analysis, you should keep it.

In [45]:

```

1 label_encoder = LabelEncoder()
2 df['from_area_id'] = label_encoder.fit_transform(df['from_area_id'])
3 df['to_area_id'] = label_encoder.fit_transform(df['to_area_id'])
4

```

- These lines use Label Encoding to convert categorical variables ('from_area_id' and 'to_area_id') into numerical format.
- Label Encoding assigns a unique integer to each category.

2. Train-Test Split:

In [46]:

```

1 X = df.drop('Car_Cancellation', axis=1)
2 y = df['Car_Cancellation']
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
5

```

- The dataset is split into features (X) and the target variable (y).
- Then, it is further split into training and testing sets using 80% of the data for training and 20% for testing.

3. Decision Tree Model with Pre-Pruning:

In [47]:

```

1 classifier = DecisionTreeClassifier(
2     random_state=42,
3     max_depth=5,
4     min_samples_split=10,
5     min_samples_leaf=5
6 )
7

```

- A Decision Tree Classifier is initialized with pre-pruning parameters: max_depth limits the maximum depth of the tree, min_samples_split sets the minimum number of samples required to split an internal node,
- and min_samples_leaf sets the minimum number of samples required to be at a leaf node.

The model is trained on the training data.

In [48]:

```
1 classifier.fit(X_train, y_train)
```

Out[48]: DecisionTreeClassifier(max_depth=5, min_samples_leaf=5, min_samples_split=10, random_state=42)

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with [nbviewer.org](#).

4. Model Evaluation:

In [49]:

```

1 y_pred = classifier.predict(X_test)
2
3 accuracy = accuracy_score(y_test, y_pred)
4 conf_matrix = confusion_matrix(y_test, y_pred)
5
6 print(f"Accuracy: {accuracy}")
7 print("Confusion Matrix:")
8 print(conf_matrix)
9

```

Accuracy: 0.9266720386784851

Confusion Matrix:

8047	5
632	3

- The model is used to make predictions on the test set, and its accuracy and confusion matrix are printed.

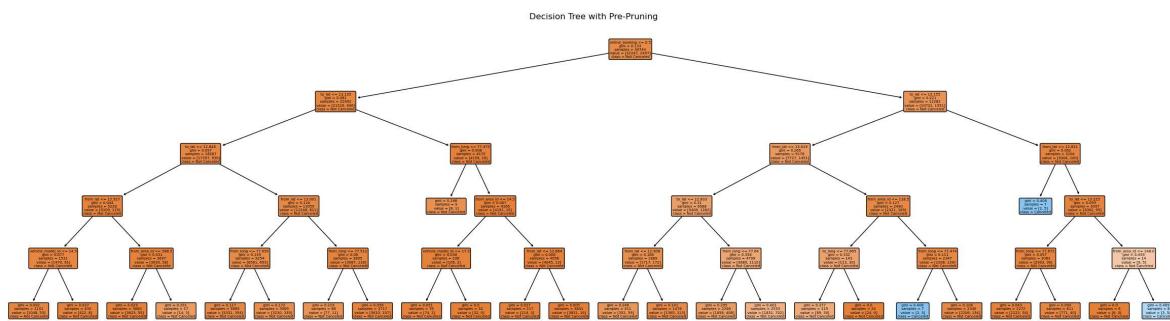
5. Visualization of the Decision Tree:

In [50]:

```

1 plt.figure(figsize=(30, 8))
2 plot_tree(classifier, filled=True, feature_names=X.columns, class_names=[]
3 plt.title("Decision Tree with Pre-Pruning")
4 plt.show()
5

```



The Decision Tree is visualized using the `plot_tree` function.

- It shows the structure of the tree, the decisions made at each node,
- and the leaf nodes with the predicted classes. The tree is displayed with filled nodes,
- feature names, class names, and rounded corners for better readability.

This code essentially performs the following tasks: data preprocessing, splitting the dataset,

- training a Decision Tree model with pre-pruning, evaluating the model, and visualizing the resulting tree.
- The pre-pruning parameters are used to control the size of the tree and prevent overfitting.
- Adjusting these parameters might be necessary based on your specific dataset and goals.

SVM (Support Vector Maching)

Certainly! Let's break it down even further:

Support Vector Machine (SVM):

- Imagine you have two types of things, like cats and dogs, and you want to draw a line on the ground to separate them.
- SVM helps you find the best line to do this.

Key Concepts:

1. Line in the Sand:

- SVM is like drawing a line in the sand to separate different things.

2. Best Line:

- It's not just any line; SVM helps you find the best line that puts as much space as possible between the different things.

3. Closest Friends:

- The things (data points) that are right at the line are like your closest friends. SVM cares a lot about them.

4. Safe Distance:

- The best line is the one that keeps your closest friends at a safe distance from it.
- This distance is called the "margin."

5. Handling Messy Situations:

- Sometimes, things are mixed up and not easy to separate.
- SVM is okay with a few mistakes and finds the best line even when it's a bit messy.

How SVM Works:

1. Finding the Line:

- You look at your data points (cats and dogs) and figure out the best line to put between them.

2. Keeping Friends Safe:

- You want the line to be far enough from your closest friends (data points) on both sides.

3. Messy Ground:

- If the ground is messy and a straight line won't work, SVM is clever.
- It can use tricks to handle the mess and still find a good line.

In Simple Terms:

- SVM is like drawing a smart line in the sand to keep your closest friends (data points) safe

```
In [51]: 1 from sklearn.svm import SVC
          2 from sklearn.metrics import accuracy_score, classification_report
```

```
In [52]: 1 model = SVC()
          2 model.fit(X_train, y_train)
```

Out[52]: SVC()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [53]: 1 scaler = StandardScaler()
          2 X_scaled = scaler.fit_transform(X)
```

```
In [54]: 1 X = df.drop('Car_Cancellation', axis=1)
          2 y = df['Car_Cancellation']
          3
          4 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=
          5
          6 # Scale the numerical features using StandardScaler to ensure they have th
```

Tune parameters

1. Regularization (C)

```
In [55]: 1 model_C = SVC(C=1)
          2 model_C.fit(X_train, y_train)
          3 model_C.score(X_test, y_test)
```

Out[55]: 0.9269022677564176

```
In [56]: 1 y_pred = model.predict(X_test)
          2 y_pred
```

Out[56]: array([0, 0, 0, ..., 0, 0, 0], dtype=int64)

```
In [57]: 1 accuracy = accuracy_score(y_test, y_pred)
          2 print(f"Accuracy: {accuracy:.2f}")
```

Accuracy: 0.93

2. Gamma

```
In [58]: 1 model_g = SVC(gamma=10)
2 model_g.fit(X_train, y_train)
3 model_g.score(X_test, y_test)
```

Out[58]: 0.9259813514446875

Using RBF kernel

```
In [59]: 1 from sklearn.svm import SVC
2 rbf_model = SVC(kernel='rbf')
```

```
In [60]: 1 len(X_train)
```

Out[60]: 34744

```
In [61]: 1 rbf_model.fit(X_train, y_train)
```

Out[61]: SVC()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [62]: 1 rbf_model.score(X_test,y_test)
```

Out[62]: 0.9269022677564176

Using Linear kernel

```
In [63]: 1 # Build and Train SVM Model
2 svm_model = SVC(kernel ='linear') # You can try other kernels like 'rbf'
3 svm_model.fit(X_train, y_train)
```

Out[63]: SVC(kernel='linear')

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [64]: 1 # Make predictions
2 y_pred = svm_model.predict(X_test)
```

```
In [65]: 1 # Evaluate the Model  
2 accuracy = accuracy_score(y_test, y_pred)  
3 print(f"Accuracy: {accuracy:.2f}")
```

Accuracy: 0.93

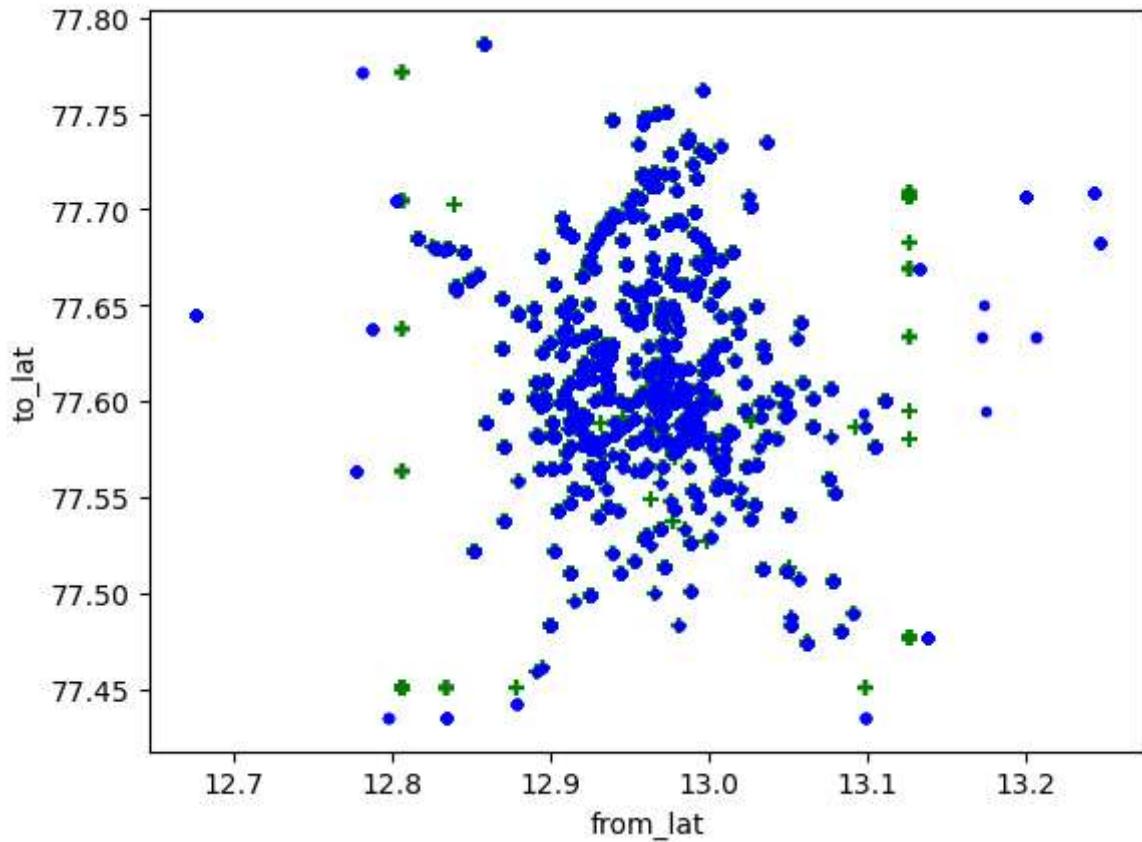
```
In [66]: 1 # Additional metrics if needed  
2 print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.93	1.00	0.96	8052
1	0.00	0.00	0.00	635
accuracy			0.93	8687
macro avg	0.46	0.50	0.48	8687
weighted avg	0.86	0.93	0.89	8687

- Evaluate the model's performance using accuracy and
- additional classification metrics such as precision, recall, and F1-score.

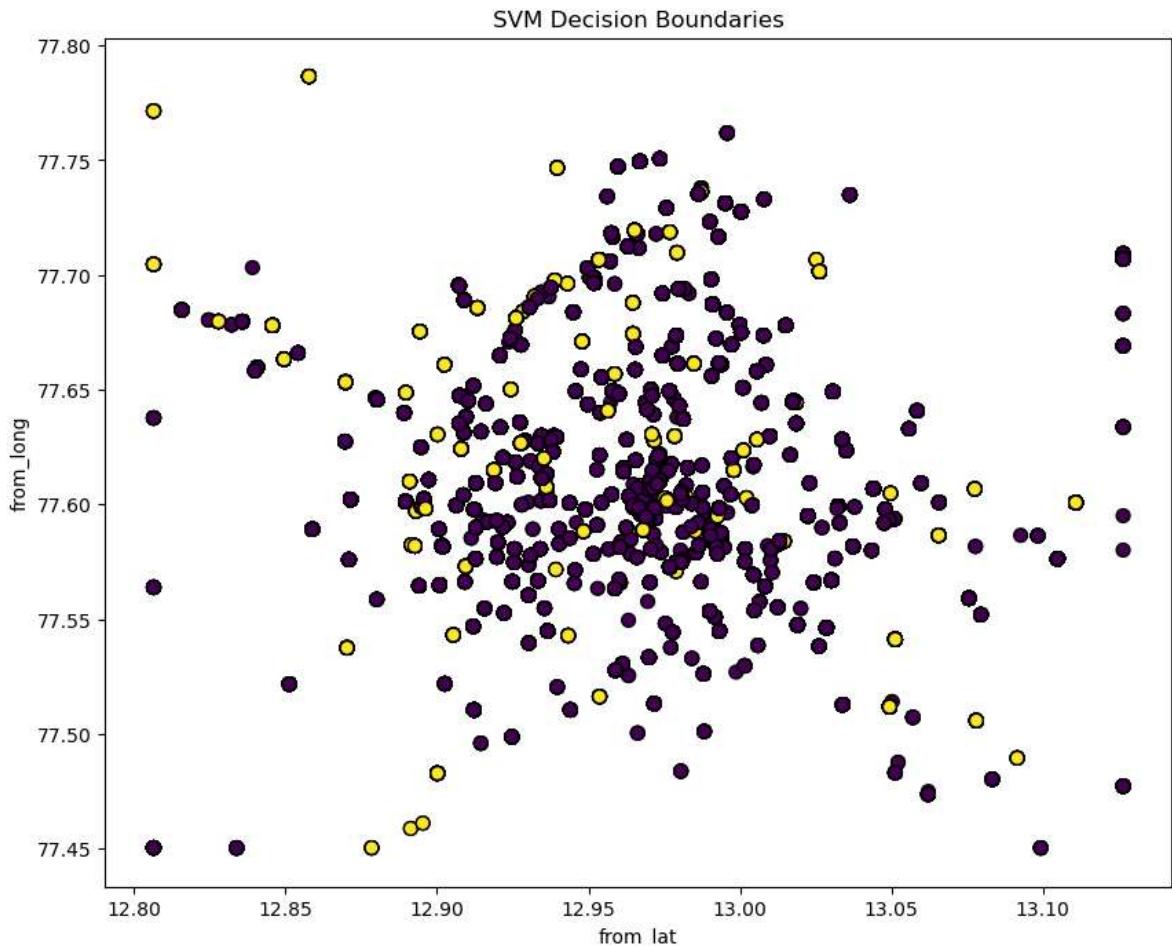
```
In [67]: 1 plt.xlabel('from_lat')
2 plt.ylabel('to_lat')
3 plt.scatter(df['from_lat'],df['from_long'],color="green",marker='+')
4 plt.scatter(df['to_lat'], df['to_long'],color="blue",marker='.' )
```

Out[67]: <matplotlib.collections.PathCollection at 0x1e1f64d7280>



In [68]:

```
1 import matplotlib.pyplot as plt
2 from sklearn.svm import SVC
3 from sklearn.preprocessing import StandardScaler
4 import pandas as pd
5
6 # Load the dataset
7
8 df = df
9
10 # Choose two features for visualization
11 feature1 = 'from_lat'
12 feature2 = 'from_long'
13
14 # Select only the relevant columns
15 data = df[[feature1, feature2, 'Car_Cancellation']]
16
17 # Feature Scaling
18 scaler = StandardScaler()
19 data_scaled = scaler.fit_transform(data[[feature1, feature2]])
20
21 # Train the SVM model
22 svm_model = SVC(kernel='linear')
23 svm_model.fit(data_scaled, data['Car_Cancellation'])
24
25 # Create a scatter plot
26 plt.figure(figsize=(10, 8))
27
28 # Plot data points
29 plt.scatter(data[feature1], data[feature2], c=data['Car_Cancellation'], cr
30
31 # Plot decision boundaries
32 ax = plt.gca()
33 xlim = ax.get_xlim()
34 ylim = ax.get_ylim()
35
36 xx, yy = np.meshgrid(np.linspace(xlim[0], xlim[1], 50), np.linspace(ylim[0],
37 Z = svm_model.decision_function(scaler.transform(np.c_[xx.ravel(), yy.ravel()]))
38 Z = Z.reshape(xx.shape)
39
40 plt.contour(xx, yy, Z, colors='k', levels=[0], alpha=0.5, linestyles=['-'])
41
42 # Plotting details
43 plt.title('SVM Decision Boundaries')
44 plt.xlabel(feature1)
45 plt.ylabel(feature2)
46
47 plt.show()
48
```



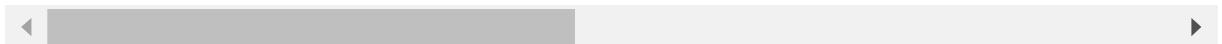
- Replace 'from_lat' and 'from_long' with the features you want to visualize.
- This code creates a scatter plot of the data points and overlays the decision boundaries from the SVM model.
- The decision boundaries are shown as a contour line.

Random Forest

In [69]: 1 deep_copy_df.head()

Out[69]:

	id	user_id	vehicle_model_id	travel_type_id	from_area_id	to_area_id	from_date	online	
0	132512	22177.0		28		2	83.0	448.0	01-01-2013 02:00
1	132513	21413.0		12		2	1010.0	540.0	01-01-2013 09:00
2	132514	22178.0		12		2	1301.0	1034.0	01-01-2013 03:30
3	132515	13034.0		12		2	768.0	398.0	01-01-2013 05:45
4	132517	22180.0		12		2	1365.0	849.0	01-01-2013 09:00



In [70]:

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
4
5 # Separate features (X) and target variable (y)
6 X = df.drop('Car_Cancellation', axis=1)
7 y = df['Car_Cancellation']
8
9 # Split the dataset into training and testing sets
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
11
12 # Check the dimensions of the splits
13 print(f"X_train shape: {X_train.shape}")
14 print(f"X_test shape: {X_test.shape}")
15 print(f"y_train shape: {y_train.shape}")
16 print(f"y_test shape: {y_test.shape}")
17 print("-----\n")
18
19 # Create a Random Forest classifier
20 rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)
21
22 # Train the classifier on the training data
23 rf_classifier.fit(X_train, y_train)
24
25 # Make predictions on the test set
26 y_pred = rf_classifier.predict(X_test)
27
28 # Evaluate the model
29 accuracy = accuracy_score(y_test, y_pred)
30 conf_matrix = confusion_matrix(y_test, y_pred)
31 classification_rep = classification_report(y_test, y_pred)
32
33 # Print the results
34 print(f"Accuracy: {accuracy}")
35 print("-----\n")
36
37 print("Confusion Matrix:")
38 print(conf_matrix)
39 print("-----\n")
40
41 print("Classification Report:")
42 print(classification_rep)
43
```

```
X_train shape: (34744, 10)
X_test shape: (8687, 10)
y_train shape: (34744,)
y_test shape: (8687,)

-----
Accuracy: 0.9201105099574076

-----
Confusion Matrix:
[[7930 122]
 [ 572  63]]


-----
Classification Report:
      precision    recall  f1-score   support

          0       0.93      0.98      0.96     8052
          1       0.34      0.10      0.15      635

   accuracy                           0.92     8687
  macro avg       0.64      0.54      0.56     8687
weighted avg       0.89      0.92      0.90     8687
```

In []:

1

How do you use K-fold cross validation?

- K-fold cross-validation is used to assess the performance of a machine learning model and to estimate its
- generalization ability. Here are the steps to utilize K-fold cross-validation:

1. Split the data:

- Divide your dataset into k equal-sized subsets (folds). Typically, k is chosen as 5 or 10, but you can
- adjust it based on your needs.

2. Train and validate:

- Iterate over the k folds. In each iteration, use k-1 folds for training the model and the remaining
- fold for validation. Train your model on the training folds and evaluate its performance on the validation fold.

3. Performance metrics:

- Calculate the performance metric(s) of interest (e.g., accuracy, precision, recall) for each fold.
- These metrics quantify how well the model generalizes to unseen data.

4. Average the results:

- Average the performance metrics obtained from the k folds to obtain a more robust estimate of
- the model's performance. This average value represents the overall performance of the model.

5. Model selection and tuning:

- Based on the cross-validation results, you can compare different models or
- hyperparameter settings and select the one that performs the best on average across the folds.

6. Final evaluation:

- After selecting the model or hyperparameters, you can retrain the model using the entire dataset
- and evaluate its performance on a separate test set to obtain a final performance estimation.
- K-fold cross-validation helps to mitigate the risk of overfitting and provides a more reliable assessment of how well the

Advantages of Cross Validation:

1,Overcoming Overfitting:

- Cross validation helps to prevent overfitting by providing a more robust estimate of the model's performance on unseen data.

2,Model Selection:

- Cross validation can be used to compare different models and select the one that performs the best on average.

3,Hyperparameter tuning:

- Cross validation can be used to optimize the hyperparameters of a model, such as the regularization parameter,
- by selecting the values that result in the best performance on the validation set.

4,Data Efficient:

- Cross validation allows the use of all the available data for both training and validation,
- making it a more data-efficient method compared to traditional validation techniques.

Disadvantages of Cross Validation:

1, Computationally Expensive:

- Cross validation can be computationally expensive, especially when the number of folds is large
- or when the model is complex and requires a long time to train.

2, Time-Consuming:

- Cross validation can be time-consuming, especially when there are many hyperparameters to tune or when multiple models need to be compared.

3, Bias-Variance Tradeoff:

- The choice of the number of folds in cross validation can impact the bias-variance tradeoff,
- i.e., too few folds may result in high variance, while too many folds may result in high bias.

Type *Markdown* and *LaTeX*: α^2

```
In [71]: 1 from sklearn.model_selection import KFold
```

```
In [72]: 1 # Load copy dataset
2 df = deep_copy_df
```

```
In [73]: 1 # Assuming 'Car_Cancellation' is the target variable
2 X = df.drop(columns=['Car_Cancellation'])
3 y = df['Car_Cancellation']
```

```
In [74]: 1 # Set the number of folds (k)
2 k = 5
```

```
In [75]: 1 # Initialize KFold
2 kf = KFold(n_splits=k, shuffle=True, random_state=42)
```

```
In [76]: 1 # Lists to store training and testing set sizes for each fold
2 train_sizes = []
3 test_sizes = []
```

In [77]:

```

1 # Enumerate through the folds
2 for fold, (train_index, test_index) in enumerate(kf.split(X)):
3     X_train, X_test = X.iloc[train_index], X.iloc[test_index]
4     y_train, y_test = y.iloc[train_index], y.iloc[test_index]
5
6     # Append sizes to the lists
7     train_sizes.append(len(X_train))
8     test_sizes.append(len(X_test))
9
10    # Print information (optional)
11    print(f"Fold {fold + 1}:")
12    print("Training Set Size:", len(X_train))
13    print("Testing Set Size:", len(X_test))
14    print("=" * 30)

```

Fold 1:

Training Set Size: 34744

Testing Set Size: 8687

=====

Fold 2:

Training Set Size: 34745

Testing Set Size: 8686

=====

Fold 3:

Training Set Size: 34745

Testing Set Size: 8686

=====

Fold 4:

Training Set Size: 34745

Testing Set Size: 8686

=====

Fold 5:

Training Set Size: 34745

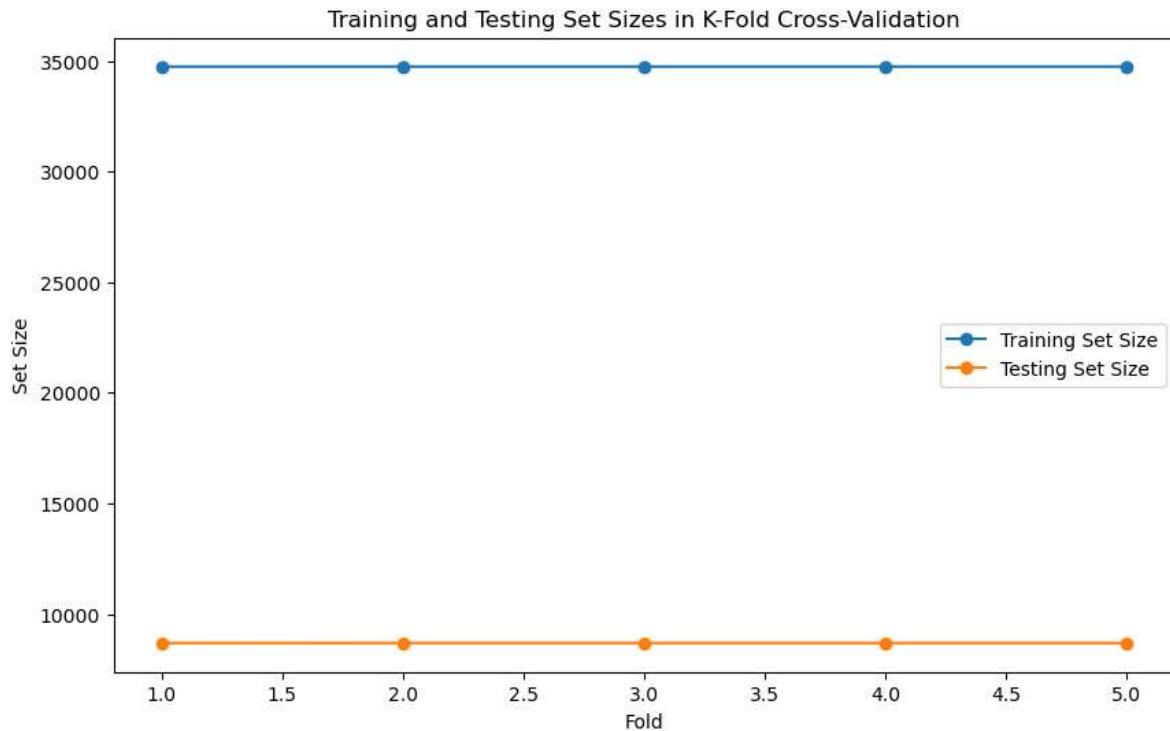
Testing Set Size: 8686

=====

- kf.split(X) generates pairs of indices for training and testing sets for each fold.
- enumerate is used to loop through these pairs and keep track of the fold number (fold).
- So, during each iteration of the loop, fold tells you which fold you're currently working with,
- and (train_index, test_index) gives you the indices of the training and testing sets for that fold.
- This information is helpful when you want to do something specific for each fold, like printing or plotting information.

In [78]:

```
1 # Plotting
2 plt.figure(figsize=(10, 6))
3 plt.plot(range(1, k + 1), train_sizes, label='Training Set Size', marker='.')
4 plt.plot(range(1, k + 1), test_sizes, label='Testing Set Size', marker='o')
5 plt.xlabel('Fold')
6 plt.ylabel('Set Size')
7 plt.title('Training and Testing Set Sizes in K-Fold Cross-Validation')
8 plt.legend()
9 plt.show()
```



conclusion

- the script is a basic framework for performing k-fold cross-validation on a dataset.
- The sizes of the training and testing sets for each fold are printed and optionally plotted.
- The primary purpose of k-fold cross-validation is to assess the performance
- and generalization ability of a machine learning model in a more robust way than a single train-test split.
- The script could be extended to include model training and evaluation within the loop for each fold if desired.

In []:

1