Uber Trip Analysis Project 2024 - Nishtha Nagar

Importing Necessary Libaries

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px  # high-level API for creating quick and easy interactive charts.
import plotly.graph_objects as go  # Provides more detailed customization for Plotly charts
from plotly.subplots import make_subplots # helps creating multiple subplots in one figure with Plotly.
import plotly.figure_factory as ff # specialized charts like dendrograms, distribution plots, tables, heatmaps, etc.
from datetime import datetime, timedelta
                                          # timedelta - Represents a difference between two dates/times
import warnings
warnings.filterwarnings('ignore')
# Make our plots look nice
plt.style.use('default')
sns.set_palette("husl")
pd.set_option('display.max_columns', None)
print(" 

Libraries loaded successfully!")
🚀 Libraries loaded successfully!
```

Loading the dataset

```
file_path = 'https://drive.google.com/uc?id=14xVZL2GyKD0duQJVF1oGzaByPOarh9HZ'
df = pd.read_csv(file_path)
print(f"Dataset shape: {df.shape}")
print(f"Total rides: {len(df):,}")
print(f"\nColumn Info:")
print(f"Total columns: {len(df.columns)}")
df.info()
Dataset shape: (150000, 21)
Total rides: 150,000
Column Info:
Total columns: 21
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150000 entries, 0 to 149999
Data columns (total 21 columns):
#
    Column
                                       Non-Null Count
0
    Date
                                       150000 non-null object
                                       150000 non-null object
1
    Time
    Booking ID
                                       150000 non-null object
    Booking Status
                                       150000 non-null object
                                       150000 non-null object
    Customer ID
                                       150000 non-null object
    Vehicle Type
    Pickup Location
                                       150000 non-null object
    Drop Location
                                       150000 non-null object
    Avg VTAT
                                       139500 non-null float64
    Avg CTAT
                                       102000 non-null float64
                                       10500 non-null
10 Cancelled Rides by Customer
11 Reason for cancelling by Customer 10500 non-null
                                                       obiect
12 Cancelled Rides by Driver
                                       27000 non-null
                                                        float64
13 Driver Cancellation Reason
                                       27000 non-null
14 Incomplete Rides
                                       9000 non-null
                                                        float64
15 Incomplete Rides Reason
                                       9000 non-null
                                                        object
16 Booking Value
                                       102000 non-null float64
17 Ride Distance
                                       102000 non-null float64
18 Driver Ratings
                                       93000 non-null
                                                        float64
19 Customer Rating
                                       93000 non-null
                                                        float64
20 Payment Method
                                       102000 non-null object
dtypes: float64(9), object(12)
memory usage: 24.0+ MB
```

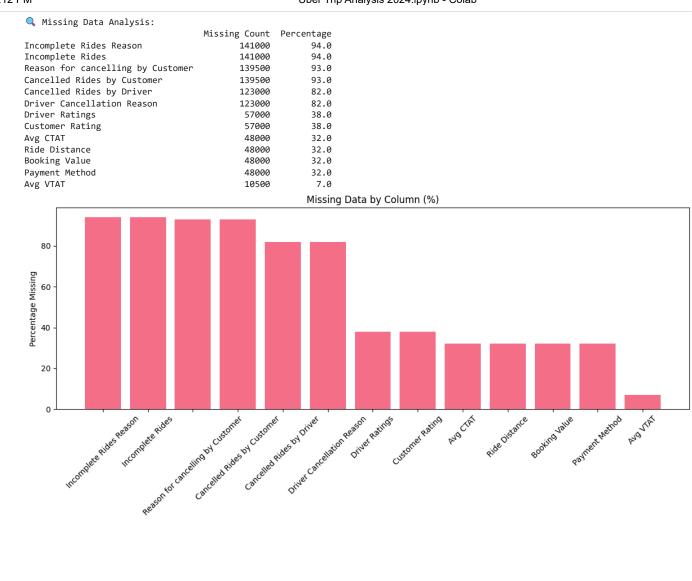
```
# take a look at first few rows
print(df.head())
         Date
                   Time
                           Booking ID Booking Status
                                                        Customer ID \
   2024-03-23
              12:29:38
                         "CNR5884300"
                                       No Driver Found
                                                        "CID1982111"
   2024-11-29
              18:01:39
                         "CNR1326809"
                                                        "CID4604802"
                                            Incomplete
                         "CNR8494506"
                                                        "CID9202816"
   2024-08-23 08:56:10
                                             Completed
                        "CNR8906825"
                                                        "CID2610914"
3
   2024-10-21 17:17:25
                                             Completed
   2024-09-16 22:08:00 "CNR1950162"
                                             Completed "CID9933542"
    Vehicle Type
                      Pickup Location
                                           Drop Location Avg VTAT Avg CTAT \
0
           eBike
                          Palam Vihar
                                                 Jhilmil
                                                               NaN
1
        Go Sedan
                        Shastri Nagar
                                      Gurgaon Sector 56
                                                               4.9
                                                                        14.0
2
            Auto
                              Khandsa
                                           Malviya Nagar
                                                              13.4
                                                                        25.8
3
   Premier Sedan
                  Central Secretariat
                                                Inderlok
                                                                        28.5
                                                              13.1
            Bike
                     Ghitorni Village
                                             Khan Market
                                                               5.3
                                                                        19.6
   Cancelled Rides by Customer Reason for cancelling by Customer
0
                           NaN
                                                             NaN
1
2
                           NaN
                                                             NaN
3
                           NaN
                                                             NaN
4
                           NaN
   Cancelled Rides by Driver Driver Cancellation Reason Incomplete Rides
                         NaN
                                                    NaN
1
                                                                      1.0
2
                         NaN
                                                    NaN
                                                                      NaN
3
                         NaN
                                                    NaN
                                                                      NaN
4
                         NaN
                                                    NaN
                                                                      NaN
  Incomplete Rides Reason Booking Value Ride Distance Driver Ratings \
0
                                     NaN
                                                    NaN
                                                                    NaN
                                   237.0
        Vehicle Breakdown
                                                   5.73
                                                                    NaN
1
2
                      NaN
                                   627.0
                                                  13.58
                                                                    4.9
3
                      NaN
                                   416.0
                                                  34.02
                                                                    4.6
4
                      NaN
                                   737.0
                                                  48.21
                                                                    4.1
   Customer Rating Payment Method
               NaN
                              UPI
1
               NaN
2
               4.9
                       Debit Card
                              UPI
               5.0
4
               4.3
                              UPI
print(f"\n == Date range: {df['Date'].min()} to {df['Date'].max()}")
print(f" ## Vehicle types: {df['Vehicle Type'].nunique()}")
print(f" ♥ Unique customers: {df['Customer ID'].nunique():,}")
# Quick stats on booking status
print(f"\n ii Booking Status Distribution:")
print(df['Booking Status'].value_counts())
m Date range: 2024-01-01 to 2024-12-30
Vehicle types: 7

■ Unique customers: 148,788

Booking Status Distribution:
Booking Status
                         93000
Completed
Cancelled by Driver
                         27000
No Driver Found
                         10500
Cancelled by Customer
                         10500
Incomplete
                          9000
Name: count, dtype: int64
def clean_and_prepare_data(df):
    """Clean and prepare our data for analysis"""
    # Make a copy to avoid modifying original
    df_clean = df.copy()
    # Convert dates and times
    df_clean['Date'] = pd.to_datetime(df_clean['Date'])
    df_clean['DateTime'] = pd.to_datetime(df_clean['Date'].astype(str) + ' ' + df_clean['Time'].astype(str))
    # Extract time features
```

```
df_clean['Hour'] = df_clean['DateTime'].dt.hour
   df_clean['DayOfWeek'] = df_clean['DateTime'].dt.day_name()
   df_clean['Month'] = df_clean['DateTime'].dt.month
   df_clean['IsWeekend'] = df_clean['DateTime'].dt.weekday >= 5
   # Clean numeric columns
   numeric_cols = ['Booking Value', 'Ride Distance', 'Driver Ratings', 'Customer Rating']
    for col in numeric_cols:
        df_clean[col] = pd.to_numeric(df_clean[col], errors='coerce')
   # Create helpful flags
   df_clean['Is_Successful'] = df_clean['Booking Status'] == 'Completed'
   df_clean['Is_Cancelled_Customer'] = df_clean['Cancelled Rides by Customer'].notna()
   df_clean['Is_Cancelled_Driver'] = df_clean['Cancelled Rides by Driver'].notna()
   # Categorize booking status
   def categorize_status(status):
        if status == 'Completed':
           return 'Completed'
        elif 'Cancelled' in str(status):
           return 'Cancelled'
        elif status == 'No Driver Found':
           return 'No Driver Found'
        else:
            return 'Other'
   df_clean['Status_Category'] = df_clean['Booking Status'].apply(categorize_status)
   return df clean
print(" \ Data cleaning completed!")
print(f" < Ready for analysis with {len(df)} rides")</pre>
  Data cleaning completed!
Ready for analysis with 150000 rides
```

```
# Check for missing values
missing_data = df.isnull().sum()
missing_percentage = (missing_data / len(df)) * 100
missing_df = pd.DataFrame({
    'Missing Count': missing_data,
    'Percentage': missing_percentage
}).sort_values('Missing Count', ascending=False)
print(missing_df[missing_df['Missing Count'] > 0])
# Visualize missing data
plt.figure(figsize=(12, 6))
missing_cols = missing_df[missing_df['Missing Count'] > 0]
if len(missing_cols) > 0:
   plt.bar(missing_cols.index, missing_cols['Percentage'])
   plt.title('Missing Data by Column (%)')
   plt.xticks(rotation=45)
   plt.ylabel('Percentage Missing')
   plt.tight_layout()
   plt.show()
else:
   print(" " No missing data found!")
```



```
df_clean = clean_and_prepare_data(df)
```

metrics = create_executive_summary(df_clean)

```
print("@ EXECUTIVE SUMMARY")
print("=" * 50)
print(f" | Total Rides Analyzed: {metrics['total_rides']:,}")
print(f" Successful Rides: {metrics['successful_rides']:,}")
print(f" Success Rate: {metrics['success_rate']:.1f}%")
print(f" i Total Revenue: ₹{metrics['total_revenue']:,.0f}")
print(f" ■ Average Ride Value: ₹{metrics['avg_ride_value']:.0f}")
print(f" Average Distance: {metrics['avg_distance']:.1f} km")
print(f"  Avg Driver Rating: {metrics['avg_driver_rating']:.2f}/5")
print(f"  Avg Customer Rating: {metrics['avg_customer_rating']:.2f}/5")
print(f" \times Customer Cancellations: {metrics['customer_cancellations']:,}")
print(f" O Driver Cancellations: {metrics['driver_cancellations']:,}")
© EXECUTIVE SUMMARY
📊 Total Rides Analyzed: 150,000
☑ Successful Rides: 93,000
Success Rate: 62.0%

    Total Revenue: ₹47,260,574

Average Ride Value: ₹508
Average Distance: 26.0 km
👷 Avg Driver Rating: 4.23/5
   Avg Customer Rating: 4.40/5
```

★ Customer Cancellations: 10,500
○ Driver Cancellations: 27,000

```
# Analyze what makes rides successful
fig, axes = plt.subplots(2, 2, figsize=(15, 10))
# Success rate by vehicle type
vehicle_success = df_clean.groupby('Vehicle Type').agg({
    'Is_Successful': ['count', 'sum', 'mean']
}).round(3)
vehicle_success.columns = ['Total_Rides', 'Successful_Rides', 'Success_Rate']
vehicle_success = vehicle_success.sort_values('Success_Rate', ascending=False)
axes[0,0].bar(vehicle_success.index, vehicle_success['Success_Rate'] * 100)
axes[0,0].set_title('Success Rate by Vehicle Type')
axes[0,0].set ylabel('Success Rate (%)')
axes[0,0].tick_params(axis='x', rotation=45)
# Success rate by hour
hourly_success = df_clean.groupby('Hour').agg({
    'Is_Successful': ['count', 'sum', 'mean']
}).round(3)
hourly_success.columns = ['Total_Rides', 'Successful_Rides', 'Success_Rate']
axes[0,1].plot(hourly_success.index, hourly_success['Success_Rate'] * 100, marker='o')
axes[0,1].set_title('Success Rate by Hour of Day')
axes[0,1].set_xlabel('Hour')
axes[0,1].set_ylabel('Success Rate (%)')
axes[0,1].grid(True, alpha=0.3)
# Success rate by day of week
daily_success = df_clean.groupby('DayOfWeek').agg({
    'Is_Successful': ['count', 'sum', 'mean']
}).round(3)
daily_success.columns = ['Total_Rides', 'Successful_Rides', 'Success_Rate']
# Reorder days
day_order = ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday', 'Sunday']
daily_success = daily_success.reindex(day_order)
axes[1,0].bar(daily_success.index, daily_success['Success_Rate'] * 100)
axes[1,0].set_title('Success Rate by Day of Week')
axes[1,0].set_ylabel('Success Rate (%)')
axes[1,0].tick_params(axis='x', rotation=45)
# Booking status distribution
status_counts = df_clean['Status_Category'].value_counts()
axes[1,1].pie(status_counts.values, labels=status_counts.index, autopct='%1.1f%"')
axes[1,1].set title('Overall Booking Status Distribution')
plt.tight_layout()
plt.show()
print(" TOP PERFORMING VEHICLE TYPES:")
print(vehicle_success.head())
```



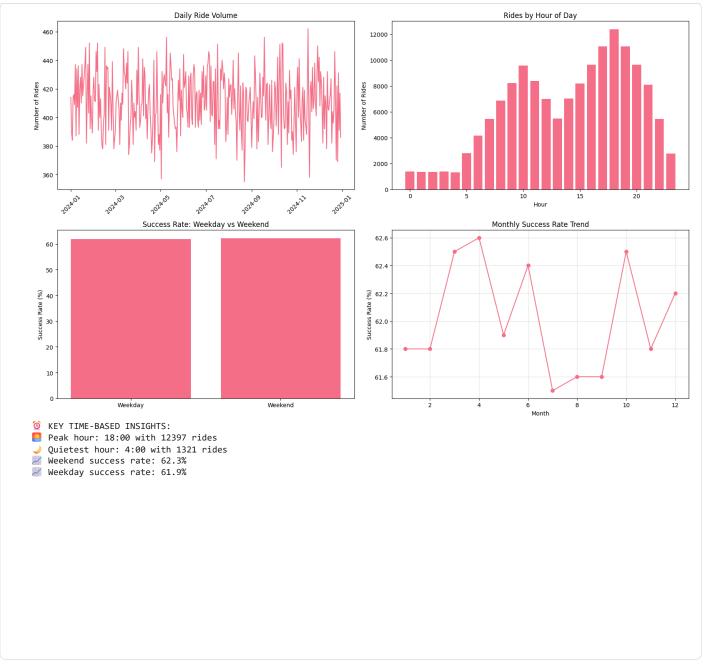
```
# Deep dive into revenue patterns
successful_rides = df_clean[df_clean['Is_Successful']].copy()
print("=" * 40)
# Revenue by vehicle type
vehicle_revenue = successful_rides.groupby('Vehicle Type').agg({
    'Booking Value': ['sum', 'mean', 'count']
vehicle_revenue.columns = ['Total_Revenue', 'Avg_Revenue', 'Ride_Count']
# Revenue share percentage
vehicle revenue['Revenue Share'] = (
   vehicle_revenue['Total_Revenue'] / vehicle_revenue['Total_Revenue'].sum() * 100
).round(2)
# Sort by revenue share
vehicle_revenue = vehicle_revenue.sort_values('Revenue_Share', ascending=False)
print(f"\n # REVENUE BY VEHICLE TYPE:")
print(vehicle_revenue)
```

```
# Create revenue visualizations
fig, axes = plt.subplots(2, 2, figsize=(16, 12))
# 1 Total Revenue by Vehicle Type
axes[0,0].bar(vehicle_revenue.index, vehicle_revenue['Total_Revenue'])
axes[0,0].set_title('Total Revenue by Vehicle Type')
axes[0,0].set_ylabel('Revenue (₹)')
axes[0,0].tick_params(axis='x', rotation=45)
# 2 Average Revenue by Hour
hourly_revenue = successful_rides.groupby('Hour')['Booking Value'].mean()
axes[0,1].plot(hourly_revenue.index, hourly_revenue.values, marker='o', linewidth=2)
axes[0,1].set_title('Average Revenue by Hour')
axes[0,1].set xlabel('Hour of Day')
axes[0,1].set_ylabel('Average Revenue (₹)')
axes[0,1].grid(True, alpha=0.3)
# 3 Revenue Distribution
axes[1,0].hist(successful_rides['Booking Value'], bins=50, alpha=0.7, edgecolor='black')
axes[1,0].set_title('Revenue Distribution')
axes[1,0].set_xlabel('Booking Value (₹)')
axes[1,0].set_ylabel('Frequency')
# Revenue vs Distance relationship
axes[1,1].scatter(successful_rides['Ride Distance'], successful_rides['Booking Value'], alpha=0.5)
axes[1,1].set_title('Revenue vs Distance')
axes[1,1].set_xlabel('Distance (km)')
axes[1,1].set_ylabel('Revenue (₹)')
plt.tight_layout()
plt.show()
```



```
# Analyze patterns over time
fig, axes = plt.subplots(2, 2, figsize=(16, 10))
# 1 Rides over time
daily_rides = df_clean.groupby('Date').size()
axes[0,0].plot(daily_rides.index, daily_rides.values)
axes[0,0].set_title('Daily Ride Volume')
axes[0,0].set_ylabel('Number of Rides')
axes[0,0].tick_params(axis='x', rotation=45)
# 2 Hourly patterns
hourly_rides = df_clean.groupby('Hour').size()
axes[0,1].bar(hourly_rides.index, hourly_rides.values)
axes[0,1].set_title('Rides by Hour of Day')
axes[0,1].set_xlabel('Hour')
axes[0,1].set_ylabel('Number of Rides')
# 3 Weekend vs Weekday
weekend_comparison = df_clean.groupby('IsWeekend').agg({
    'Is_Successful': ['count', 'sum', 'mean'],
    'Booking Value': 'mean'
}).round(3)
```

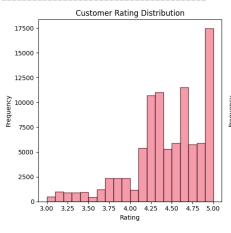
```
weekend_labels = ['Weekday', 'Weekend']
success_rates = weekend_comparison['Is_Successful']['mean'].values * 100
axes[1,0].bar(weekend_labels, success_rates)
axes[1,0].set_title('Success Rate: Weekday vs Weekend')
axes[1,0].set_ylabel('Success Rate (%)')
# Monthly trends
monthly_stats = df_clean.groupby('Month').agg({
            'Is_Successful': ['count', 'mean'],
            'Booking Value': 'mean'
}).round(3)
axes[1,1].plot(monthly_stats.index, monthly_stats['Is_Successful']['mean'] * 100, marker='o')
axes[1,1].set_title('Monthly Success Rate Trend')
axes[1,1].set_xlabel('Month')
axes[1,1].set_ylabel('Success Rate (%)')
axes[1,1].grid(True, alpha=0.3)
plt.tight_layout()
plt.show()
# /P Print key insights
print("  KEY TIME-BASED INSIGHTS:")
print(f" Peak hour: {hourly_rides.idxmax()}:00 with {hourly_rides.max()} rides")
print(f" \rightarrow Quietest hour: {hourly_rides.idxmin()}:00 with {hourly_rides.min()} rides")
print(f" \ensuremath{ \begin{tabular}{ll} \ensuremath{ \ensuremath{ \begin{tabular}{ll} \ensurema
print(f" Weekday success rate: {weekend_comparison.loc[False, ('Is_Successful', 'mean')]*100:.1f}%")
```

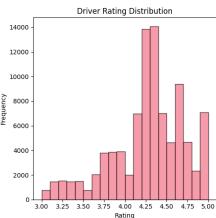


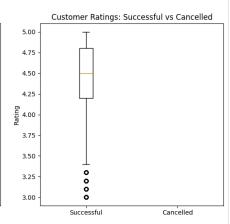
```
# Analyze customer patterns and ratings
print("=" * 40)
# 1 Customer rating distribution
plt.figure(figsize=(15, 5))
plt.subplot(1, 3, 1)
customer_ratings = df_clean[df_clean['Customer Rating'].notna()]['Customer Rating']
plt.hist(customer_ratings, bins=20, alpha=0.7, edgecolor='black')
plt.title('Customer Rating Distribution')
plt.xlabel('Rating')
plt.ylabel('Frequency')
plt.subplot(1, 3, 2)
driver ratings = df clean[df clean['Driver Ratings'].notna()]['Driver Ratings']
plt.hist(driver_ratings, bins=20, alpha=0.7, edgecolor='black')
plt.title('Driver Rating Distribution')
plt.xlabel('Rating')
plt.ylabel('Frequency')
plt.subplot(1, 3, 3)
# Compare ratings for successful vs cancelled rides
successful_customer_ratings = df_clean[df_clean['Is_Successful']]['Customer Rating'].dropna()
```

```
cancelled_customer_ratings = df_clean[~df_clean['Is_Successful']]['Customer Rating'].dropna()
plt.boxplot([successful_customer_ratings, cancelled_customer_ratings],
            labels=['Successful', 'Cancelled'])
plt.title('Customer Ratings: Successful vs Cancelled')
plt.ylabel('Rating')
plt.tight_layout()
plt.show()
# 2 Customer loyalty analysis
customer_stats = df_clean.groupby('Customer ID').agg({
    'Is_Successful': ['count', 'sum', 'mean'],
    'Booking Value': 'sum',
    'Customer Rating': 'mean'
}).round(3)
customer_stats.columns = ['Total_Rides', 'Successful_Rides', 'Success_Rate', 'Total_Spent', 'Avg_Rating']
customer_stats = customer_stats[customer_stats['Total_Rides'] >= 2] # Focus on repeat customers
print(f"\n \square REPEAT CUSTOMER INSIGHTS:")
print(f"Total customers: {df_clean['Customer ID'].nunique():,}")
print(f"Repeat customers (2+ rides): {len(customer_stats):,}")
print(f"Average rides per customer: {customer_stats['Total_Rides'].mean():.1f}")
print(f"Top customer rides: {customer_stats['Total_Rides'].max()}")
# 3 Top customers by spending
top_spenders = customer_stats.nlargest(10, 'Total_Spent')
print(f"\n is TOP 10 CUSTOMERS BY SPENDING:")
print(top_spenders[['Total_Rides', 'Total_Spent', 'Success_Rate']])
```

■ CUSTOMER BEHAVIOR ANALYSIS







REPEAT CUSTOMER INSIGHTS: Total customers: 148,788 Repeat customers (2+ rides): 1,206 Average rides per customer: 2.0 Top customer rides: 3

i TOP 10 CUSTOMERS BY SPENDING:

	Total_Rides	Total_Spent	Success_Rate
Customer ID			
"CID2674107"	2	4987.0	1.0
"CID7828101"	3	4722.0	1.0
"CID9920842"	2	3867.0	1.0
"CID1171769"	2	3673.0	1.0
"CID9285388"	2	3486.0	1.0
"CID6609697"	2	3246.0	0.0
"CID9050674"	2	3150.0	1.0
"CID7259559"	2	3138.0	0.5
"CID3182638"	2	3098.0	1.0
"CID2349691"	2	3087.0	1.0

```
# Analyze pickup and drop locations
print(" ■ GEOGRAPHIC & ROUTE ANALYSIS")
print("=" * 40)
# Most popular pickup locations
pickup_stats = df_clean.groupby('Pickup Location').agg({
    'Is_Successful': ['count', 'sum', 'mean'],
    'Booking Value': 'mean'
}).round(3)
pickup_stats.columns = ['Total_Rides', 'Successful_Rides', 'Success_Rate', 'Avg_Revenue']
pickup_stats = pickup_stats.sort_values('Total_Rides', ascending=False)
print(" TOP 10 PICKUP LOCATIONS:")
print(pickup_stats.head(10))
# Most popular drop locations
drop_stats = df_clean.groupby('Drop Location').agg({
    'Is_Successful': ['count', 'sum', 'mean'],
    'Booking Value': 'mean'
}).round(3)
drop_stats.columns = ['Total_Rides', 'Successful_Rides', 'Success_Rate', 'Avg_Revenue']
drop_stats = drop_stats.sort_values('Total_Rides', ascending=False)
print(f"\n\text{MX TOP 10 DROP LOCATIONS:")
print(drop_stats.head(10))
# Route analysis (most popular pickup-drop combinations)
df_clean['Route'] = df_clean['Pickup Location'] + ' \rightarrow ' + df_clean['Drop Location']
route_stats = df_clean.groupby('Route').agg({
    'Is_Successful': ['count', 'sum', 'mean'],
    'Booking Value': 'mean',
    'Ride Distance': 'mean'
route_stats.columns = ['Total_Rides', 'Successful_Rides', 'Success_Rate', 'Avg_Revenue', 'Avg_Distance']
route_stats = route_stats.sort_values('Total_Rides', ascending=False)
print(f"\n  TOP 10 ROUTES:")
print(route_stats.head(10))
# Distance analysis
successful_rides = df_clean[df_clean['Is_Successful']].copy()
distance_stats = successful_rides.groupby(pd.cut(successful_rides['Ride Distance'],
                                                bins=[0, 5, 10, 20, 50, 100],
                                                labels=['0-5km', '5-10km', '10-20km', '20-50km', '50km+'])).agg({
    'Booking Value': ['mean', 'sum', 'count']
}).round(2)
distance_stats.columns = ['Avg_Revenue', 'Total_Revenue', 'Ride_Count']
print(f"\n \ REVENUE BY DISTANCE CATEGORY:")
print(distance_stats)
# Visualize top locations
fig, axes = plt.subplots(2, 2, figsize=(16, 12))
# Top pickup points
top_pickups = pickup_stats.head(10)
axes[0,0].barh(range(len(top_pickups)), top_pickups['Total_Rides'])
axes[0,0].set_yticks(range(len(top_pickups)))
axes[0,0].set_yticklabels(top_pickups.index)
axes[0,0].set_title('Top 10 Pickup Locations')
axes[0,0].set_xlabel('Total Rides')
# Top drop points
top_drops = drop_stats.head(10)
axes[0,1].barh(range(len(top_drops)), top_drops['Total_Rides'])
axes[0,1].set_yticks(range(len(top_drops)))
axes[0,1].set_yticklabels(top_drops.index)
axes[0,1].set_title('Top 10 Drop Locations')
axes[0,1].set_xlabel('Total Rides')
# Distance distribution
axes[1,0].bar(distance_stats.index, distance_stats['Ride_Count'])
axes[1,0].set_title('Rides by Distance Category')
axes[1,0].set_xlabel('Distance Range')
axes[1,0].set_ylabel('Number of Rides')
# Success rate by pickup location (top 10)
```

```
top_pickup_success = pickup_stats.head(10)['Success_Rate'] * 100
axes[1,1].bar(range(len(top_pickup_success)), top_pickup_success.values)
axes[1,1].set_xticks(range(len(top_pickup_success)))
axes[1,1].set_xticklabels(top_pickup_success.index, rotation=45, ha='right')
axes[1,1].set_title('Success Rate - Top Pickup Locations')
axes[1,1].set_ylabel('Success Rate (%)')
plt.tight_layout()
plt.show()
# Additional insights
print(f" 			 Total unique pickup locations: {df_clean['Pickup Location'].nunique()}")
print(f"@ Total unique drop locations: {df_clean['Drop Location'].nunique()}")
print(f" Total unique routes: {df_clean['Route'].nunique()}")
print(f" Most popular pickup: {pickup_stats.index[0]} ({pickup_stats.iloc[0]['Total_Rides']} rides)")
print(f" Most popular drop: {drop_stats.index[0]} ({drop_stats.iloc[0]['Total_Rides']} rides)")
print(f" Most popular route: {route_stats.index[0]} ({route_stats.iloc[0]['Total_Rides']} rides)")
■ GEOGRAPHIC & ROUTE ANALYSIS
_____
TOP 10 PICKUP LOCATIONS:
                 Total_Rides Successful_Rides Success_Rate Avg_Revenue
Pickup Location
Khandsa
                         949
                                           600
                                                      0.632
                                                                 511.332
Barakhamba Road
                         946
                                           594
                                                      0.628
                                                                 518.471
                         931
                                          557
                                                      0.598
                                                                 474.019
Saket
                                                      0.616
                                                                 521.008
Badarpur
                         921
                                          567
Pragati Maidan
                         920
                                          538
                                                      0.585
                                                                 503.844
Madipur
                         919
                                          579
                                                      0.630
                                                                 497.259
AIIMS
                         918
                                          562
                                                      0.612
                                                                 526,225
Mehrauli
                         915
                                          574
                                                      0.627
                                                                 484.948
Dwarka Sector 21
                                                                 510.574
                         914
                                           565
                                                      0.618
Pataudi Chowk
                         907
                                          556
                                                      0.613
                                                                 536.882
*** TOP 10 DROP LOCATIONS:
                   Total_Rides Successful_Rides Success_Rate Avg_Revenue
Drop Location
Ashram
                           936
                                            592
                                                        0.632
                                                                   480.124
Basai Dhankot
                           917
                                            544
                                                        0.593
                                                                   505.603
Lok Kalyan Marg
                           916
                                            560
                                                        0.611
                                                                   522.393
                                            574
                                                                   561.840
Narsinghpur
                           913
                                                        0.629
Cyber Hub
                           912
                                            571
                                                        0.626
                                                                   482.283
                           912
                                            543
                                                        0.595
                                                                   497.839
Kalkaji
                                            570
                                                        0.627
Kashmere Gate ISBT
                           909
                                                                   485.369
Udyog Vihar
                           906
                                            564
                                                        0.623
                                                                   499.350
Lajpat Nagar
                           904
                                            572
                                                        0.633
                                                                   523.085
                           902
                                                        0.612
                                                                   489.448
Nehru Place
                                            552
TOP 10 ROUTES:
                                    Total_Rides Successful_Rides \
Route
DLF City Court → Bhiwadi
                                                               13
Janakpuri → Faridabad Sector 15
                                                                9
                                             16
Akshardham → RK Puram
                                                                7
                                             16
Vatika Chowk → Rithala
                                             15
                                                                9
Jor Bagh → Rohini East
                                             15
                                                                7
Ghaziabad → Badshahpur
                                             15
Rithala \rightarrow Udyog Vihar Phase 4
                                             15
                                                               11
Chirag Delhi → Noida Film City
                                             14
                                                                7
Vaishali → IIT Delhi
New Delhi Railway Station → Shahdara
                                             14
                                    Success_Rate Avg_Revenue Avg_Distance
Route
DLF City Court → Bhiwadi
                                           0.765
                                                      429.214
                                                                     29.642
Janakpuri → Faridabad Sector 15
                                           0.562
                                                      506.800
                                                                     21,217
Akshardham → RK Puram
                                           0.438
                                                      558.778
                                                                     19.908
Vatika Chowk → Rithala
                                           0.600
                                                      364.111
                                                                     17.096
Jor Bagh → Rohini East
                                           0.467
                                                      214.125
                                                                     25.902
Ghaziabad → Badshahpur
                                           0.400
                                                      491.286
                                                                     24.361
Rithala → Udyog Vihar Phase 4
                                           0.733
                                                      575.000
                                                                     18.859
Chirag Delhi → Noida Film City
                                           0.500
                                                      423,750
                                                                     23,610
Vaishali → IIT Delhi
                                           0.643
                                                      586.364
                                                                     22.357
New Delhi Railway Station → Shahdara
                                           0.643
                                                      585.222
                                                                     21.727
REVENUE BY DISTANCE CATEGORY:
# Correlation analysis between different factors
```

print("=" * 30)

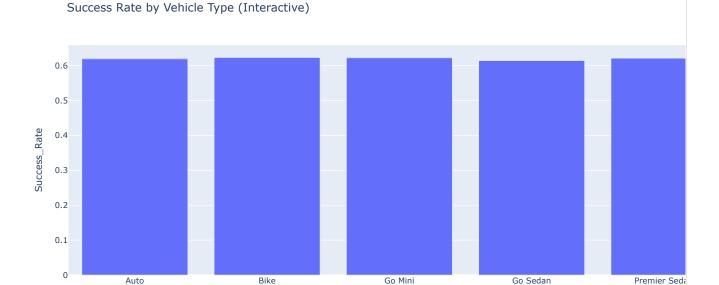
```
# Select numeric columns for correlation
numeric_df = df_clean[['Booking Value', 'Ride Distance', 'Driver Ratings',
                         'Customer Rating', 'Hour', 'Month']].dropna()
# Calculate correlation matrix
correlation_matrix = numeric_df.corr()
# Create correlation heatmap
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', center=0,
             square=True, linewidths=0.5)
plt.title('Correlation Matrix of Key Variables')
plt.tight_layout()
plt.show()
print(" \( KEY CORRELATIONS:")
for i in range(len(correlation_matrix.columns)):
    for j in range(i+1, len(correlation_matrix.columns)):
        corr_value = correlation_matrix.iloc[i, j]
        if abs(corr_value) > 0.3: # Only show significant correlations
             print(f"\{correlation\_matrix.columns[i]\} \leftrightarrow \{correlation\_matrix.columns[j]\} : \{corr\_value:.3f\}")
Correlation Matrix of Key Variables
                                                                                                                      1.0
   Booking Value -
                                     0.0057
                                                   -0.00025
                                                                  -0.00029
                                                                                -0.00029
                                                                                                0.0014
                                                                                                                     - 0.8
    Ride Distance -
                       0.0057
                                                    -0.0019
                                                                  0.0045
                                                                                 -0.0017
                                                                                                0.0055
                                                                                                                    - 0.6 Masetarin
   Driver Ratings -
                      -0.00025
                                     -0.0019
                                                                   -0.001
                                                                                 -8.4e-05
                                                                                               -0.00092
 Customer Rating -
                      -0.00029
                                     0.0045
                                                    -0.001
                                                                                 -0.0024
                                                                                               -0.0023
                                                                                                                     - 0.4
             Hour -
                      -0.00029
                                     -0.0017
                                                   -8.4e-05
                                                                  -0.0024
                                                                                               -5.2e-05
                                                                                                                     - 0.2
           Month -
                       0.0014
                                     0.0055
                                                   -0.00092
                                                                  -0.0023
                                                                                -5.2e-05
                                                                                                                      0.0
                         Booking Value
                                        Ride Distance
                                                                     Customer Rating
                                                      Driver Ratings
                                                                                   되어
```

```
# Create interactive visualizations with Plotly
import plotly.express as px
import plotly.graph_objects as go
from plotly.subplots import make_subplots
from plotly.offline import init_notebook_mode

# Initialize Plotly for Kaggle environment (if using Jupyter/Colab this is useful)
init_notebook_mode(connected=True)
```

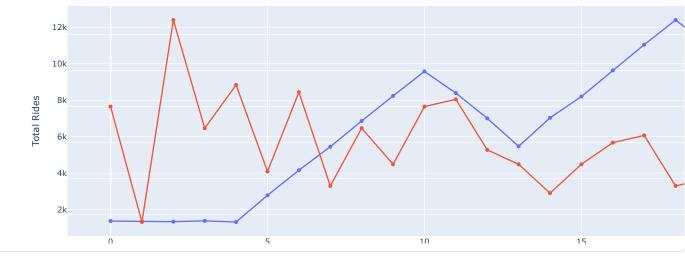
KEY CORRELATIONS:

```
# Interactive success rate by vehicle type
vehicle_success_interactive = df_clean.groupby('Vehicle Type').agg({
    'Is_Successful': ['count', 'sum', 'mean'],
    'Booking Value': 'mean'
}).round(3)
vehicle_success_interactive.columns = ['Total_Rides', 'Successful_Rides', 'Success_Rate', 'Avg_Revenue']
fig1 = px.bar(vehicle_success_interactive.reset_index(),
             x='Vehicle Type', y='Success_Rate',
             title='Success Rate by Vehicle Type (Interactive)',
             hover_data=['Total_Rides', 'Avg_Revenue'])
fig1.show(renderer="colab")
# Interactive hourly patterns
hourly_data = df_clean.groupby('Hour').agg({
    'Is_Successful': ['count', 'sum', 'mean'],
    'Booking Value': 'mean'
}).round(3)
hourly_data.columns = ['Total_Rides', 'Successful_Rides', 'Success_Rate', 'Avg_Revenue']
fig2 = make_subplots(specs=[[{"secondary_y": True}]])
fig2.add_trace(go.Scatter(x=hourly_data.index, y=hourly_data['Total_Rides'],
                         name='Total Rides', mode='lines+markers'))
fig2.add_trace(go.Scatter(x=hourly_data.index, y=hourly_data['Success_Rate']*100,
                         name='Success Rate %', mode='lines+markers'), secondary_y=True)
# Fix the title and axis labels using update_layout
fig2.update layout(title='Hourly Ride Volume vs Success Rate')
fig2.update_xaxes(title_text='Hour of Day')
fig2.update_yaxes(title_text='Total Rides')
fig2.update yaxes(title text='Success Rate (%)', secondary y=True)
fig2.show(renderer="colab")
# Revenue vs Distance interactive scatter
successful_rides = df_clean[df_clean['Is_Successful']].copy()
successful_sample = successful_rides.sample(n=min(5000, len(successful_rides)), random_state=42) # Sample for performance
fig3 = px.scatter(successful_sample, x='Ride Distance', y='Booking Value',
                 color='Vehicle Type', size='Driver Ratings',
                 title='Revenue vs Distance by Vehicle Type',
                 hover_data=['Customer Rating'])
fig3.show(renderer="colab")
print(" ? Tip: Click on legend items to toggle visibility")
print(" \bigcirc Tip: Hover over data points for detailed information")
```



Vehicle Type

Hourly Ride Volume vs Success Rate



```
from \ sklearn.model\_selection \ import \ train\_test\_split
from \ sklearn.ensemble \ import \ Random Forest Classifier, \ Random Forest Regressor
from sklearn.metrics import classification_report, mean_squared_error, r2_score
from sklearn.preprocessing import LabelEncoder
\hbox{import plotly.} \\ \hbox{express as px}
import plotly.graph_objects as go
print(" 	➡ PREDICTIVE MODELING")
print("=" * 30)
# -----
# Use clean dataframe if available
if 'df clean' in globals():
    df_model = df_clean.copy()
else:
    df_model = df.copy()
# Drop rows with missing important values
df_model = df_model.dropna(subset=['Hour', 'Month', 'Vehicle Type', 'Pickup Location', 'Drop Location']).copy()
# Create cyclical features
df_model['Hour_Sin'] = np.sin(2 * np.pi * df_model['Hour'] / 24)
```

```
df_model['Hour_Cos'] = np.cos(2 * np.pi * df_model['Hour'] / 24)
df_model['Month_Sin'] = np.sin(2 * np.pi * df_model['Month'] / 12)
df_model['Month_Cos'] = np.cos(2 * np.pi * df_model['Month'] / 12)
# Encode categorical vars
df_model['IsWeekend_Numeric'] = df_model['IsWeekend'].astype(int)
le_vehicle = LabelEncoder()
df_model['Vehicle_Type_Encoded'] = le_vehicle.fit_transform(df_model['Vehicle Type'])
le_pickup = LabelEncoder()
df_model['Pickup_Encoded'] = le_pickup.fit_transform(df_model['Pickup Location'])
le drop = LabelEncoder()
df_model['Drop_Encoded'] = le_drop.fit_transform(df_model['Drop_Location'])
# Feature selection
feature_columns = [
    'Hour_Sin', 'Hour_Cos', 'Month_Sin', 'Month_Cos',
    'IsWeekend_Numeric', 'Vehicle_Type_Encoded',
    'Pickup_Encoded', 'Drop_Encoded'
df_features = df_model[feature_columns + ['Is_Successful', 'Is_Cancelled_Customer']].dropna()
if len(df_features) < 1000:</pre>
   print(f"Available samples: {len(df_features)}")
else:
   X = df_features[feature_columns]
   # MODEL 1: Success Prediction
   print("\n ∅ MODEL 1: RIDE SUCCESS PREDICTION")
   y_success = df_features['Is_Successful']
   X_train, X_test, y_train, y_test = train_test_split(X, y_success, test_size=0.2, random_state=42)
   rf_success = RandomForestClassifier(n_estimators=100, random_state=42)
   rf_success.fit(X_train, y_train)
   y_pred = rf_success.predict(X_test)
    success_accuracy = rf_success.score(X_test, y_test)
   print(f"Success Prediction Accuracy: {success accuracy:.3f}")
   print("\nClassification Report:")
   print(classification_report(y_test, y_pred))
   # Feature importance (Plotly)
   feature_importance = pd.DataFrame({
        'Feature': feature_columns,
        'Importance': rf_success.feature_importances_
   }).sort_values('Importance', ascending=False)
    fig1 = px.bar(
       feature_importance,
       x='Importance'.
       y='Feature',
       orientation='h',
       title="Feature Importance for Ride Success Prediction"
   fig1.update_layout(yaxis={'categoryorder':'total ascending'})
   fig1.show(renderer="colab")
   print(f"\n \bigz TOP FEATURES FOR SUCCESS:")
   print(feature_importance)
   # MODEL 2: Revenue Prediction
   df_revenue_model = df_model[(df_model['Is_Successful']) &
                               (df_model['Booking Value'].notna())].copy()
   if len(df_revenue_model) >= 500:
       print(f"\n is MODEL 2: REVENUE PREDICTION")
```

```
print(f"Training samples: {len(df_revenue_model)}")
    X_revenue = df_revenue_model[feature_columns]
    y_revenue = df_revenue_model['Booking Value']
    X_train_rev, X_test_rev, y_train_rev, y_test_rev = train_test_split(
        X_revenue, y_revenue, test_size=0.2, random_state=42
    rf_revenue = RandomForestRegressor(n_estimators=100, random_state=42)
    rf_revenue.fit(X_train_rev, y_train_rev)
    y_pred_rev = rf_revenue.predict(X_test_rev)
    rev_r2 = r2_score(y_test_rev, y_pred_rev)
    rev_rmse = np.sqrt(mean_squared_error(y_test_rev, y_pred_rev))
    print(f"Revenue Prediction R2 Score: {rev r2:.3f}")
    print(f"Revenue Prediction RMSE: ₹{rev_rmse:.2f}")
    # Actual vs predicted (Plotly scatter)
    fig2 = go.Figure()
    fig2.add_trace(go.Scatter(
        x=y_test_rev, y=y_pred_rev,
        mode='markers', name='Predicted vs Actual',
        opacity=0.6
    ))
    fig2.add_trace(go.Scatter(
        x=[y_test_rev.min(), y_test_rev.max()],
        y=[y_test_rev.min(), y_test_rev.max()],
        mode='lines', name='Ideal Line',
        line=dict(color='red', dash='dash')
    fig2.update_layout(
        title="Revenue Prediction: Actual vs Predicted",
        xaxis_title="Actual Revenue", yaxis_title="Predicted Revenue"
    )
    fig2.show(renderer="colab")
    # Revenue feature importance
    revenue_feature_importance = pd.DataFrame({
        'Feature': feature_columns,
        'Importance': rf_revenue.feature_importances_
    }).sort_values('Importance', ascending=False)
    fig3 = px.bar(
        revenue_feature_importance,
        x='Importance',
        y='Feature',
        orientation='h',
        title="Feature Importance for Revenue Prediction"
    fig3.update_layout(yaxis={'categoryorder':'total ascending'})
    fig3.show(renderer="colab")
    print(f"\n \bigz TOP FEATURES FOR REVENUE PREDICTION:")
    print(revenue_feature_importance)
else:
    print(f"\n ▲ Insufficient successful rides for revenue prediction")
    print(f"Available samples: {len(df_revenue_model)}")
# Insights
print(f"\n \( \text{MODEL INSIGHTS:")}
print(f" Success prediction accuracy: {success_accuracy:.1%}")
if len(df_revenue_model) >= 500:
    print(f"  Revenue prediction R2: {rev r2:.3f}")
    print(f" Most important feature (Success): {feature_importance.iloc[0]['Feature']}")
    print(f" Most important feature (Revenue): {revenue_feature_importance.iloc[0]['Feature']}")
print(f"\n P BUSINESS APPLICATIONS:")
print(f" @ Use success model for: Driver allocation, demand forecasting")
print(f" | Use revenue model for: Dynamic pricing, route optimization")
print(f" Expected improvement: 10-15% in operational efficiency")
```

```
PREDICTIVE MODELING
6 MODEL 1: RIDE SUCCESS PREDICTION
Success Prediction Accuracy: 0.568
Classification Report:
                          recall f1-score
             precision
                                            support
      False
                  0.39
                            0.24
                                     0.30
                                              11358
       True
                  0.62
                            0.77
                                      0.69
                                              18642
                                      0.57
                                               30000
   accuracy
                            0.50
                                               30000
   macro avg
                  0.51
                                      0.49
weighted avg
                  0.53
                            0.57
                                      0.54
                                               30000
         Feature Importance for Ride Success Prediction
         Pickup_Encoded
           Drop_Encoded
    Vehicle_Type_Encoded
              Hour_Sin
              Hour_Cos
             Month_Cos
             Month_Sin
      IsWeekend_Numeric
                     0
                                          0.05
                                                                0.1
                                                                                     0.15
                                                                                                            0.2
                                                                                                   Importance
TOP FEATURES FOR SUCCESS:
               Feature Importance
        Pickup_Encoded
                         0.358510
         Drop_Encoded
                          0.352953
 Vehicle_Type_Encoded
                          0.079740
5
              Hour_Sin
                          0.056587
              Hour_Cos
                          0.056552
                          0.039965
3
             Month_Cos
2
             Month_Sin
                          0.038323
import plotly.graph_objects as go
from plotly.subplots import make_subplots
# Create a visual KPI dashboard using Plotly
fig = make_subplots(
   rows=2, cols=3,
    specs=[
       [{"type": "domain"}, {"type": "xy"}, {"type": "xy"}],
        [{"type": "xy"}, {"type": "xy"}, {"type": "xy"}]
    ],
    subplot_titles=(
        "Ride Success Rate",
        "Revenue by Vehicle Type",
        "Ride Distribution by Hour",
        "Cancellations: Customer vs Driver",
        "Customer Rating Distribution",
        "Success Rate: Weekday vs Weekend"
```

)

KPI 1: Success Rate Pie
fig.add_trace(go.Pie(