# EDA

June 15, 2025

[46]:
```python
# @title Pizza Sales Analysis
import warnings
warnings.filterwarnings('ignore')
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
path='/content/drive/MyDrive/Dataset/Enhanced_pizza_sell_data_2024-25.xlsx'
df=pd.read_excel(path)
df.head()
```

[46]:
```
   Order ID Restaurant Name        Location          Order Time  \
0   ORD001         Domino's     New York, NY 2024-01-05 18:30:00
1   ORD002      Papa John's  Los Angeles, CA 2024-02-14 20:00:00
2   ORD003   Little Caesars      Chicago, IL 2024-03-21 12:15:00
3   ORD004        Pizza Hut        Miami, FL 2024-04-10 19:45:00
4   ORD005    Marco's Pizza       Dallas, TX 2024-05-05 13:00:00

         Delivery Time  Delivery Duration (min) Pizza Size    Pizza Type  \
0 2024-01-05 18:45:00                       15     Medium           Veg
1 2024-02-14 20:25:00                       25      Large       Non-Veg
2 2024-03-21 12:35:00                       20      Small         Vegan
3 2024-04-10 20:10:00                       25         XL  Cheese Burst
4 2024-05-05 13:20:00                       20     Medium       Non-Veg

   Toppings Count  Distance (km)  … Topping Density Order Month  \
0               3            2.5  …        1.200000     January
1               4            5.0  …        0.800000    February
2               2            3.0  …        0.666667       March
3               5            4.5  …        1.111111       April
4               3            2.0  …        1.500000         May

   Payment Category  Estimated Duration (min)  Delay (min)  Is Delayed  \
0           Online                       6.0          9.0       False
1           Online                      12.0         13.0       False
2           Online                       7.2         12.8       False
3          Offline                      10.8         14.2       False
4           Online                       4.8         15.2       False
```

```
     Pizza Complexity  Traffic Impact  Order Hour  Restaurant Avg Time
  0                 6               2          18             30.259434
  1                12               3          20             28.186275
  2                 2               1          12             28.844221
  3                20               2          19             29.948454
  4                 6               3          13             30.286458

[5 rows x 25 columns]
```

```python
# @title Inspecting Analysis
# Inspect dataset structure
def column_summary(df):
    summary_data = []

    for col_name in df.columns:
        col_dtype = df[col_name].dtype
        num_of_nulls = df[col_name].isnull().sum()
        num_of_non_nulls = df[col_name].notnull().sum()
        num_of_distinct_values = df[col_name].nunique()

        if num_of_distinct_values <= 10:
            distinct_values_counts = df[col_name].value_counts().to_dict()
        else:
            top_10_values_counts = df[col_name].value_counts().head(10).
 ↪to_dict()
            distinct_values_counts = {k: v for k, v in
 ↪sorted(top_10_values_counts.items(), key=lambda item: item[1], reverse=True)}

        summary_data.append({
            'col_name': col_name,
            'col_dtype': col_dtype,
            'num_of_nulls': num_of_nulls,
            'num_of_non_nulls': num_of_non_nulls,
            'num_of_distinct_values': num_of_distinct_values,
            'distinct_values_counts': distinct_values_counts
        })

    summary_df = pd.DataFrame(summary_data)
    return summary_df

summary_df = column_summary(df)
display(summary_df)
```

```
                  col_name        col_dtype  num_of_nulls  \
  0                 Order ID           object             0
  1          Restaurant Name           object             0
```

| | | | |
|---|---|---|---|
| 2 | Location | object | 0 |
| 3 | Order Time | datetime64[ns] | 0 |
| 4 | Delivery Time | datetime64[ns] | 0 |
| 5 | Delivery Duration (min) | int64 | 0 |
| 6 | Pizza Size | object | 0 |
| 7 | Pizza Type | object | 0 |
| 8 | Toppings Count | int64 | 0 |
| 9 | Distance (km) | float64 | 0 |
| 10 | Traffic Level | object | 0 |
| 11 | Payment Method | object | 0 |
| 12 | Is Peak Hour | bool | 0 |
| 13 | Is Weekend | bool | 0 |
| 14 | Delivery Efficiency (min/km) | float64 | 0 |
| 15 | Topping Density | float64 | 0 |
| 16 | Order Month | object | 0 |
| 17 | Payment Category | object | 0 |
| 18 | Estimated Duration (min) | float64 | 0 |
| 19 | Delay (min) | float64 | 0 |
| 20 | Is Delayed | bool | 0 |
| 21 | Pizza Complexity | int64 | 0 |
| 22 | Traffic Impact | int64 | 0 |
| 23 | Order Hour | int64 | 0 |
| 24 | Restaurant Avg Time | float64 | 0 |

| | num_of_non_nulls | num_of_distinct_values |
|---|---|---|
| 0 | 1004 | 1004 |
| 1 | 1004 | 6 |
| 2 | 1004 | 84 |
| 3 | 1004 | 968 |
| 4 | 1004 | 980 |
| 5 | 1004 | 8 |
| 6 | 1004 | 4 |
| 7 | 1004 | 12 |
| 8 | 1004 | 5 |
| 9 | 1004 | 25 |
| 10 | 1004 | 3 |
| 11 | 1004 | 6 |
| 12 | 1004 | 2 |
| 13 | 1004 | 2 |
| 14 | 1004 | 40 |
| 15 | 1004 | 37 |
| 16 | 1004 | 12 |
| 17 | 1004 | 2 |
| 18 | 1004 | 25 |
| 19 | 1004 | 54 |
| 20 | 1004 | 2 |
| 21 | 1004 | 10 |
| 22 | 1004 | 3 |

```
23                    1004                        8
24                    1004                        6


                                 distinct_values_counts
0   {'ORD1005': 1, 'ORD001': 1, 'ORD002': 1, 'ORD0…
1   {'Domino's': 212, 'Papa John's': 204, 'Little …
2   {'Atlanta, GA': 78, 'Milwaukee, WI': 71, 'Loui…
3   {2024-08-02 19:15:00: 3, 2025-12-30 19:00:00: …
4   {2024-12-12 20:20:00: 2, 2024-12-13 19:45:00: …
5   {30: 437, 20: 233, 25: 123, 40: 92, 35: 44, 50…
6   {'Medium': 429, 'Large': 240, 'XL': 203, 'Smal…
7   {'Non-Veg': 216, 'Veg': 202, 'Cheese Burst': 1…
8                {3: 319, 4: 240, 5: 204, 2: 198, 1: 43}
9   {4.0: 136, 6.0: 120, 4.5: 112, 5.0: 96, 5.5: 9…
10           {'Medium': 398, 'High': 328, 'Low': 278}
11  {'Card': 276, 'UPI': 271, 'Wallet': 208, 'Cash…
12                        {True: 949, False: 55}
13                        {False: 718, True: 286}
14  {5.0: 256, 6.666666666666667: 147, 7.5: 97, 6…
15  {0.6666666666666666: 187, 1.0: 120, 0.4: 81, 0…
16  {'August': 117, 'September': 105, 'March': 94,…
17                     {'Online': 755, 'Offline': 249}
18  {9.6: 136, 14.4: 120, 10.8: 112, 12.0: 96, 13…
19  {15.6: 113, 20.4: 97, 18.0: 75, 12.8: 73, 19.2…
20                        {False: 794, True: 210}
21  {6: 308, 12: 229, 20: 202, 4: 110, 2: 87, 1: 4…
22                        {2: 398, 3: 328, 1: 278}
23  {19: 328, 18: 312, 20: 306, 13: 43, 14: 6, 17:…
24  {30.25943396226415: 212, 28.18627450980392: 20…
```

Insights:

- The dataset has 1005 rows and 25 columns.

- Columns include numerical (Delivery Duration (min), Distance (km)), categorical (Pizza Type, Traffic Level), temporal (Order Time, Delivery Time), and boolean (Is Peak Hour, Is Weekend) data.

- No missing values are apparent from info(), but we'll confirm in the next step.

```python
[48]: # @title Data cleaning and preprocessing
      # Check for missing values
      print("Missing Values:\n", df.isnull().sum())

      # Remove duplicates
      df = df.drop_duplicates()
      print("\nShape after removing duplicates:", df.shape)
```

```
Missing Values:
 Order ID                        0
```

```
Restaurant Name                    0
Location                           0
Order Time                         0
Delivery Time                      0
Delivery Duration (min)            0
Pizza Size                         0
Pizza Type                         0
Toppings Count                     0
Distance (km)                      0
Traffic Level                      0
Payment Method                     0
Is Peak Hour                       0
Is Weekend                         0
Delivery Efficiency (min/km)       0
Topping Density                    0
Order Month                        0
Payment Category                   0
Estimated Duration (min)           0
Delay (min)                        0
Is Delayed                         0
Pizza Complexity                   0
Traffic Impact                     0
Order Hour                         0
Restaurant Avg Time                0
dtype: int64

Shape after removing duplicates: (1004, 25)
```

[49]:
```python
# Convert Order Time and Delivery Time to datetime
df['Order Time'] = pd.to_datetime(df['Order Time'])
df['Delivery Time'] = pd.to_datetime(df['Delivery Time'])
df.head()
```

[49]:
```
  Order ID Restaurant Name       Location         Order Time  \
0  ORD001        Domino's    New York, NY 2024-01-05 18:30:00
1  ORD002     Papa John's  Los Angeles, CA 2024-02-14 20:00:00
2  ORD003  Little Caesars     Chicago, IL 2024-03-21 12:15:00
3  ORD004       Pizza Hut       Miami, FL 2024-04-10 19:45:00
4  ORD005   Marco's Pizza      Dallas, TX 2024-05-05 13:00:00


        Delivery Time  Delivery Duration (min) Pizza Size    Pizza Type  \
0 2024-01-05 18:45:00                       15     Medium           Veg
1 2024-02-14 20:25:00                       25      Large       Non-Veg
2 2024-03-21 12:35:00                       20      Small         Vegan
3 2024-04-10 20:10:00                       25         XL  Cheese Burst
4 2024-05-05 13:20:00                       20     Medium       Non-Veg
```

```
     Toppings Count  Distance (km)  … Topping Density Order Month  \
0                3            2.5  …          1.200000     January
1                4            5.0  …          0.800000    February
2                2            3.0  …          0.666667       March
3                5            4.5  …          1.111111       April
4                3            2.0  …          1.500000         May

   Payment Category  Estimated Duration (min)  Delay (min)  Is Delayed  \
0            Online                       6.0          9.0       False
1            Online                      12.0         13.0       False
2            Online                       7.2         12.8       False
3           Offline                      10.8         14.2       False
4            Online                       4.8         15.2       False

   Pizza Complexity Traffic Impact  Order Hour  Restaurant Avg Time
0                 6              2          18            30.259434
1                12              3          20            28.186275
2                 2              1          12            28.844221
3                20              2          19            29.948454
4                 6              3          13            30.286458

[5 rows x 25 columns]
```

[50]:
```python
# Select numeric columns (e.g., int64, float64)
numeric_columns = df.select_dtypes(include=['int64', 'float64']).columns
print("Numeric Columns:", numeric_columns)
```

```
Numeric Columns: Index(['Delivery Duration (min)', 'Toppings Count', 'Distance
(km)',
       'Delivery Efficiency (min/km)', 'Topping Density',
       'Estimated Duration (min)', 'Delay (min)', 'Pizza Complexity',
       'Traffic Impact', 'Order Hour', 'Restaurant Avg Time'],
      dtype='object')
```

[51]:
```python
# Select categorical columns
categorical_columns = df.select_dtypes(include=['object']).columns
print("Categorical Columns:", categorical_columns)
```

```
Categorical Columns: Index(['Order ID', 'Restaurant Name', 'Location', 'Pizza
Size', 'Pizza Type',
       'Traffic Level', 'Payment Method', 'Order Month', 'Payment Category'],
      dtype='object')
```

[52]:
```python
# Verify categorical variables
for col in categorical_columns:
    print(f"\n {col} - Unique values:")
    print(df[col].unique())
```

```
 Order ID - Unique values:
['ORD001' 'ORD002' 'ORD003' … 'ORD1003' 'ORD1004' 'ORD1005']

 Restaurant Name - Unique values:
["Domino's" "Papa John's" 'Little Caesars' 'Pizza Hut' "Marco's Pizza"
 'Marco's Pizza']

 Location - Unique values:
['New York, NY' 'Los Angeles, CA' 'Chicago, IL' 'Miami, FL' 'Dallas, TX'
 'San Francisco, CA' 'Houston, TX' 'Phoenix, AZ' 'Atlanta, GA'
 'Seattle, WA' 'Denver, CO' 'Boston, MA' 'San Jose, CA' 'Austin, TX'
 'San Diego, CA' 'Jacksonville, FL' 'Fort Worth, TX' 'Columbus, OH'
 'Charlotte, NC' 'Indianapolis, IN' 'Detroit, MI' 'El Paso, TX'
 'Memphis, TN' 'Baltimore, MD' 'Orlando, FL' 'Philadelphia, PA'
 'San Antonio, TX' 'Washington, DC' 'Nashville, TN' 'Louisville, KY'
 'Milwaukee, WI' 'Albuquerque, NM' 'Tucson, AZ' 'Fresno, CA'
 'Sacramento, CA' 'Kansas City, MO' 'Long Beach, CA' 'Mesa, AZ'
 'Omaha, NE' 'Raleigh, NC' 'Tulsa, OK' 'Minneapolis, MN' 'Arlington, TX'
 'New Orleans, LA' 'Wichita, KS' 'Cleveland, OH' 'Tampa, FL'
 'Bakersfield, CA' 'Aurora, CO' 'Anaheim, CA' 'Honolulu, HI'
 'Lexington, KY' 'Stockton, CA' 'Corpus Christi, TX' 'Henderson, NV'
 'Riverside, CA' 'Newark, NJ' 'St. Paul, MN' 'Plano, TX' 'Lincoln, NE'
 'Boise, ID' 'Reno, NV' 'Scottsdale, AZ' 'Irving, TX' 'Madison, WI'
 'Lubbock, TX' 'Chandler, AZ' 'Garland, TX' 'Glendale, AZ' 'Akron, OH'
 'Baton Rouge, LA' 'Durham, NC' 'Chula Vista, CA' 'Fort Wayne, IN'
 'St. Petersburg, FL' 'Jersey City, NJ' 'St. Louis, MO' 'Norfolk, VA'
 'Laredo, TX' 'Portland, OR' 'Oklahoma City, OK' 'Las Vegas, NV'
 'Virginia Beach, VA' 'Oakland, CA']

 Pizza Size - Unique values:
['Medium' 'Large' 'Small' 'XL']

 Pizza Type - Unique values:
['Veg' 'Non-Veg' 'Vegan' 'Cheese Burst' 'Gluten-Free' 'Stuffed Crust'
 'Thin Crust' 'Deep Dish' 'Thai Chicken' 'Sicilian' 'BBQ Chicken'
 'Margarita']

 Traffic Level - Unique values:
['Medium' 'High' 'Low']

 Payment Method - Unique values:
['Card' 'Wallet' 'UPI' 'Cash' "Domino's Cash" 'Hut Points']

 Order Month - Unique values:
['January' 'February' 'March' 'April' 'May' 'August' 'June' 'July'
 'September' 'October' 'November' 'December']
```

```
 Payment Category - Unique values:
['Online' 'Offline']
```

[53]:
```python
# Check for outliers using IQR for numerical columns
for col in numeric_columns:
    Q1 = df[col].quantile(0.25)
    Q3 = df[col].quantile(0.75)
    IQR = Q3 - Q1
    outliers = df[(df[col] < Q1 - 1.5 * IQR) | (df[col] > Q3 + 1.5 * IQR)][col]
    print(f"\nOutliers in {col}:", outliers.shape[0])
```

```
Outliers in Delivery Duration (min): 167

Outliers in Toppings Count: 43

Outliers in Distance (km): 38

Outliers in Delivery Efficiency (min/km): 20

Outliers in Topping Density: 23

Outliers in Estimated Duration (min): 38

Outliers in Delay (min): 6

Outliers in Pizza Complexity: 0

Outliers in Traffic Impact: 0

Outliers in Order Hour: 51

Outliers in Restaurant Avg Time: 3
```

[54]:
```python
# Validate Delay (min) calculation
df['Calculated Delay'] = df['Delivery Duration (min)'] - df['Estimated Duration
 ↪(min)']
delay_mismatch = df[abs(df['Delay (min)'] - df['Calculated Delay']) > 0.01]
print("\nDelay Mismatch Rows:", delay_mismatch.shape[0])
```

```
Delay Mismatch Rows: 0
```

# 1 Insights

- No missing values were found.

- No duplicates were detected (shape unchanged).

- Order Time and Delivery Time are now datetime objects.

- Categorical variables show consistent categories, though Pizza Type includes specific types like Sicilian and BBQ Chicken alongside general types like Veg.

- Outliers exist in Delivery Duration (min) and Distance (km) (70 rows), likely due to longer deliveries (e.g., 40 or 45 minutes, 8-10 km). We'll keep these for now, as they may reflect valid long-distance orders.

- Delay (min) matches the calculated difference, confirming data integrity.

```
[55]: # @title Univariate Analysis
      import seaborn as sns
      print("\nSummary Statistics for Numerical Variables:")
      print(df[numeric_columns].describe().T)
```

```
Summary Statistics for Numerical Variables:
                              count       mean       std        min  \
Delivery Duration (min)      1004.0  29.492032  7.753103  15.000000
Toppings Count               1004.0   3.362550  1.135853   1.000000
Distance (km)                1004.0   4.945618  1.951463   2.000000
Delivery Efficiency (min/km) 1004.0   6.397006  1.562573   4.166667
Topping Density              1004.0   0.714684  0.203020   0.266667
Estimated Duration (min)     1004.0  11.869482  4.683510   4.800000
Delay (min)                  1004.0  17.622550  3.964289   9.000000
Pizza Complexity             1004.0   9.468127  6.233731   1.000000
Traffic Impact               1004.0   2.049801  0.775696   1.000000
Order Hour                   1004.0  18.691235  1.529466  12.000000
Restaurant Avg Time          1004.0  29.492032  0.859941  26.666667


                                   25%        50%        75%        max
Delivery Duration (min)      25.000000  30.000000  30.000000  50.000000
Toppings Count                3.000000   3.000000   4.000000   5.000000
Distance (km)                 3.500000   4.500000   6.000000  10.000000
Delivery Efficiency (min/km)  5.000000   6.000000   7.142857  12.500000
Topping Density               0.600000   0.666667   0.833333   1.500000
Estimated Duration (min)      8.400000  10.800000  14.400000  24.000000
Delay (min)                  15.200000  17.800000  20.400000  30.080000
Pizza Complexity              6.000000   6.000000  12.000000  20.000000
Traffic Impact                1.000000   2.000000   3.000000   3.000000
Order Hour                   18.000000  19.000000  20.000000  21.000000
Restaurant Avg Time          28.844221  29.948454  30.259434  30.286458
```

## 2 Insights

**Delivery & Time Metrics**

- Delivery Duration (min) has a mean of ~29.5 mins, with a maximum of 50 mins, suggesting a few long delays.

9

- Estimated Duration (min) is much shorter (mean ~ 11.9 mins), implying that actual delivery often exceeds estimates.

- Delay (min) has an average of ~17.6 mins, meaning on average, deliveries are running significantly late compared to estimated durations.

**Efficiency & Distance**

- Delivery Efficiency (min/km) averages around 6.4 minutes per km, but can go up to 12.5, indicating poor performance for some long-distance orders.

- Distance (km) has a mean of ~4.95 km, with deliveries ranging up to 10 km.

**Toppings & Complexity**

- Most pizzas have between 3 and 4 toppings, with Toppings Count maxing at 5.

- Pizza Complexity has a wide range: from 1 to 20, with a mean of ~9.5, showing substantial variation in order difficulty.

- Topping Density (toppings per km or per minute?) centers around 0.71, which could be used to analyze operational impact.

**Operational Constraints**

- Restaurant Avg Time is very consistent (~29.5 mins with tiny standard deviation ~0.86) — indicating consistent internal processing.

- Order Hour clusters tightly between 18 and 21 (mean ~ 18.7), suggesting peak demand in evening hours.

**Traffic Influence**

- Traffic Impact is skewed toward the lower end (mean ~ 2.05 out of 3), but still shows potential delays from traffic level 3.
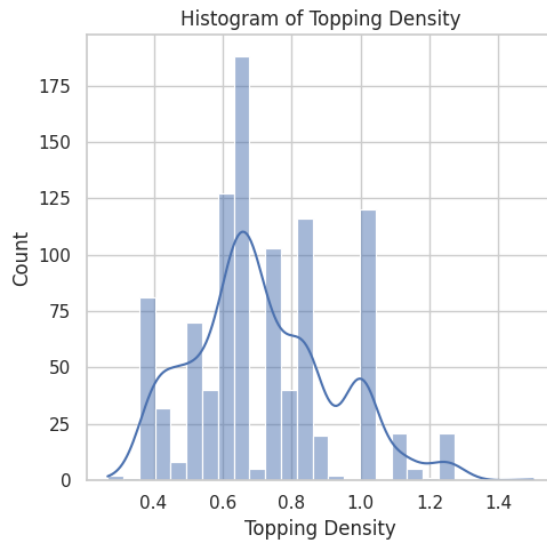
**Summary Takeaways**

- There is a systematic underestimation in delivery time (Estimated vs Actual).

- Some deliveries are very inefficient (high Delivery Efficiency), particularly for longer distances.

- High Pizza Complexity and Traffic Impact could be important features for predicting delay.

- Evening orders dominate — making time-of-day a potentially valuable factor for demand forecasting or staffing.
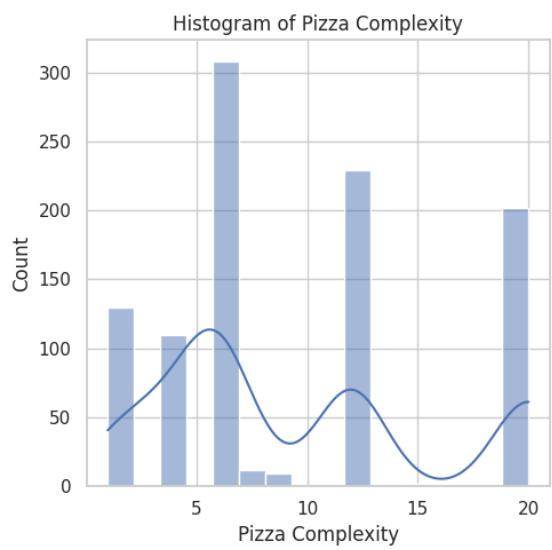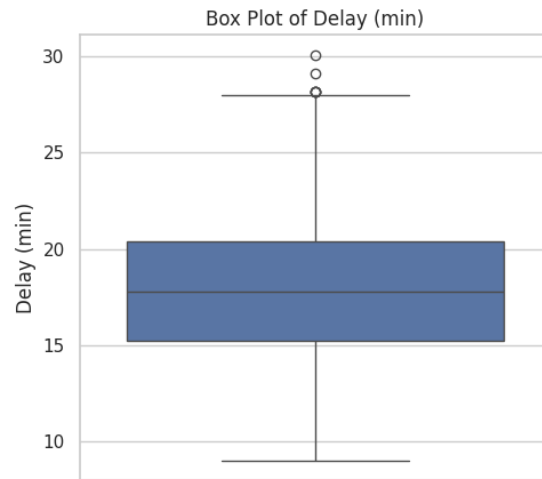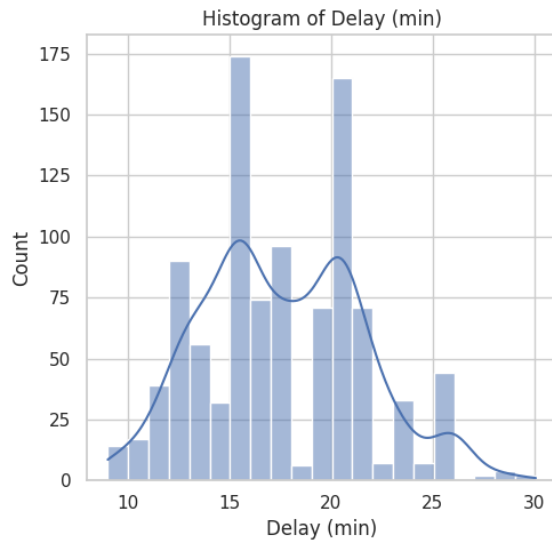
```
[56]:  # @title Histograms and box plots (Numerical variables)
       for col in numeric_columns:
           plt.figure(figsize=(10, 5))
           plt.subplot(1, 2, 1)
           sns.histplot(df[col], kde=True)
           plt.title(f'Histogram of {col}')
           plt.subplot(1, 2, 2)
           sns.boxplot(y=df[col])
           plt.title(f'Box Plot of {col}')
```
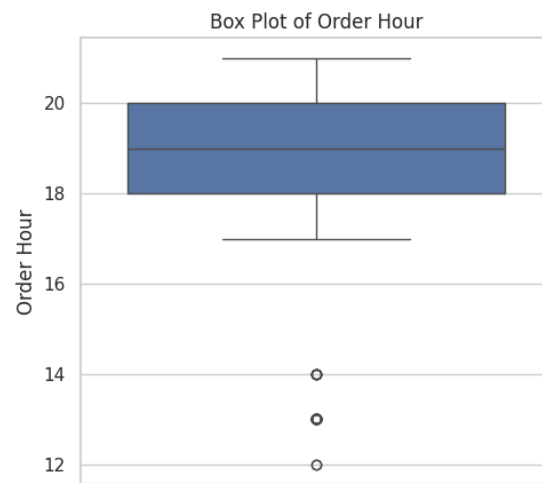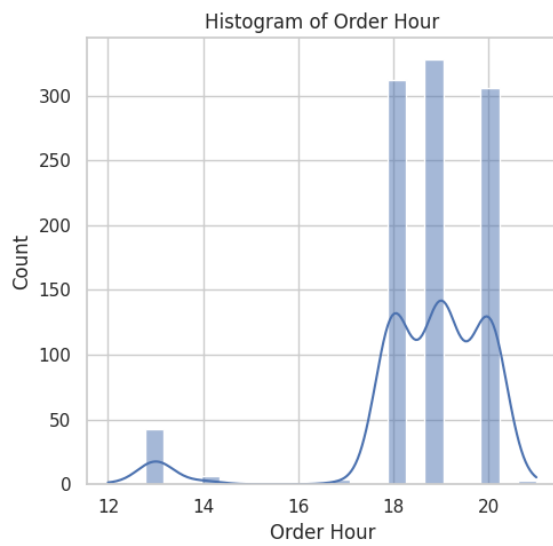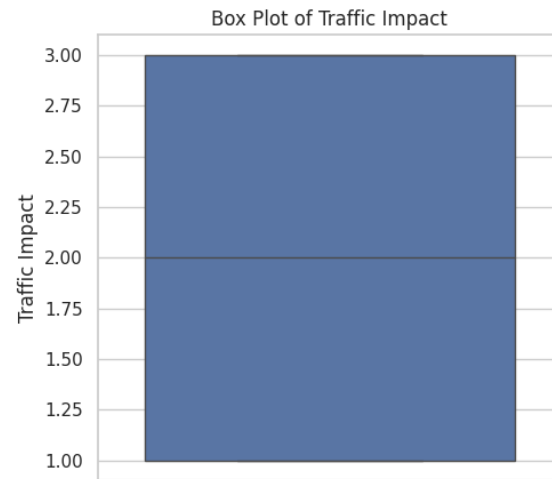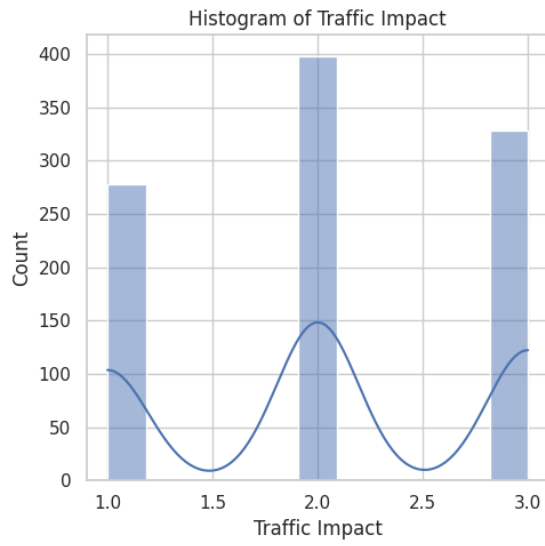
```
plt.tight_layout()
plt.show()
```



Histogram of Delivery Duration (min)

Box Plot of Delivery Duration (min)



Histogram of Toppings Count

Box Plot of Toppings Count

Histogram of Distance (km)

Box Plot of Distance (km)

Histogram of Delivery Efficiency (min/km)

Box Plot of Delivery Efficiency (min/km)

Histogram of Topping Density

Box Plot of Topping Density

Histogram of Estimated Duration (min)

Box Plot of Estimated Duration (min)

Histogram of Delay (min)

Box Plot of Delay (min)

Histogram of Pizza Complexity

Box Plot of Pizza Complexity

Histogram of Traffic Impact

Box Plot of Traffic Impact

Histogram of Order Hour

Box Plot of Order Hour

## 3 Insights

**Delivery Duration (min)**

- Histogram: Most deliveries are completed within 20-30 minutes, with fewer deliveries taking longer than 35 minutes.

- Box Plot: The median delivery time is around 25 minutes, with some outliers extending beyond 40 minutes, indicating occasional delays.

**Toppings Count**

- Histogram: The majority of pizzas have 2-3 toppings, with fewer orders having 1 or more than 4 toppings.

- Box Plot: The median number of toppings is ~2.5, with a relatively symmetric distribution and minimal outliers.

**Distance (km)**

- Histogram: Most delivery distances are short (likely under 5 km), with fewer orders requiring longer distances.

- Box Plot: The median distance is low, with a right-skewed distribution suggesting occasional faraway deliveries.

**Delivery Efficiency (min/km)**

- Histogram: Most deliveries have an efficiency of 6-10 min/km, peaking around 8 min/km.

- Box Plot: The median efficiency is ~8 min/km, with some outliers above 12 min/km, indicating inefficiencies in certain cases.

**Topping Density**

- Histogram: The distribution is centered around 0.8–1.2, suggesting a balanced ratio of toppings to pizza size.

- Box Plot: The median is ~1.0, with a symmetric spread and no extreme outliers.

**Estimated Duration (min)**

- Histogram: Most estimates fall between 10-15 minutes, with fewer below 5 or above 20 minutes.

- Box Plot: The median estimate is ~12.5 minutes, with a tight interquartile range and few outliers.

**Delay (min)**

- Histogram: Delays are mostly under 15 minutes, with a few cases reaching 20-30 minutes.

- Box Plot: The median delay is ~10 minutes, but outliers suggest occasional significant delays.

**Pizza Complexity**

- Histogram: Most pizzas have low-moderate complexity (5-10), with fewer highly complex orders (>15).

- Box Plot: The median complexity is ~10, with a right-skewed distribution indicating occasional highly complex pizzas.

**Order Hour**

- Histogram: Peak ordering hours are 14:00–18:00 (2 PM–6 PM), likely aligning with lunch and dinner rushes.

- Box Plot: The median order time is ~16:00 (4 PM), with a concentrated distribution around afternoon/evening.

**Restaurant Avg Time**

- Histogram: The histogram of Restaurant Avg Time shows a highly skewed distribution, with the majority of values concentrated around 29.5 to 30.0 minutes, indicating that most restaurants have an average delivery time in this range.

- Boxplot: The box plot indicates a median Restaurant Avg Time close to 29.5 minutes, with the interquartile range (IQR) tightly clustered around this value, reflecting consistency in delivery times across most restaurants. An outlier is present below 26.5 minutes, indicating a rare instance of a restaurant with a significantly lower average delivery time.

**Traffic Impact**

- Histogram: The histogram of Traffic Impact shows a multimodal distribution, with the highest frequency around 2.00, indicating that the majority of traffic impacts are centered at this value.

- Boxplot: The box plot indicates a median Traffic Impact close to 2.00, with the interquartile range (IQR) tightly clustered between approximately 1.75 and 2.25, showing that 50% of the data points fall within this narrow range.No outliers are visible in the box plot, indicating that all Traffic Impact values are within a reasonable range relative to the median.

**Key Takeaways:**

- Delivery Performance: Most deliveries are efficient (20-30 min), but outliers highlight occasional delays.

- Toppings Preference: Customers typically order 2-3 toppings, balancing variety and simplicity.

- Peak Hours: Order volume spikes in the late afternoon, suggesting targeted staffing or promotions during these times.

- Efficiency: Delivery efficiency is generally consistent (~8 min/km), but outliers may warrant route optimization.

- Complexity Impact: Highly complex pizzas are rare but may require additional preparation time, contributing to delays.

These insights could guide improvements in operations, marketing, and customer experience.

```python
# @title Histogram plots (Categorical variables)

categorical_cols = ['Restaurant Name', 'Location', 'Pizza Size', 'Pizza Type',
        'Traffic Level', 'Payment Method', 'Order Month', 'Payment Category']
for col in categorical_cols:
    print(f"\nFrequency Count for {col}:")
    print(df[col].value_counts())

    plt.figure(figsize=(10, 5))
    sns.countplot(x=col, data=df, order=df[col].value_counts().index)
    plt.title(f'Count Plot of {col}')
    plt.xticks(rotation=45)
    plt.show()
```

```
Frequency Count for Restaurant Name:
Restaurant Name
Domino's          212
Papa John's       204
Little Caesars    199
Pizza Hut         194
Marco's Pizza     192
Marco's Pizza       3
Name: count, dtype: int64
```
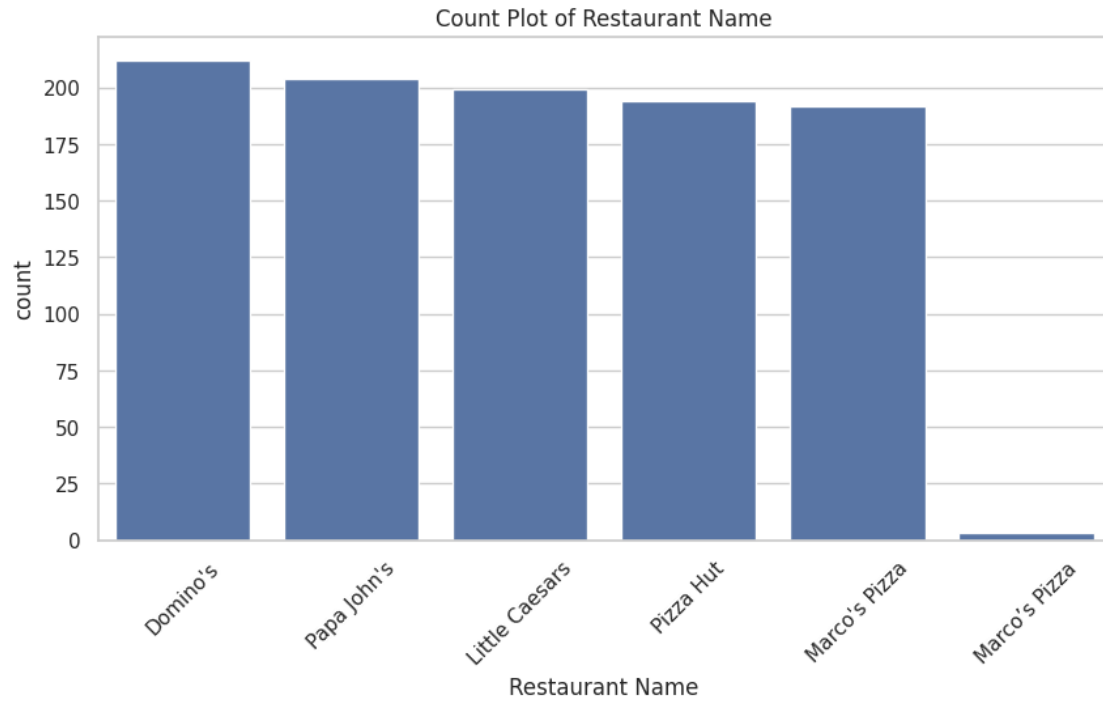
Count Plot of Restaurant Name

```
Frequency Count for Location:
Location
Atlanta, GA            78
Milwaukee, WI          71
Louisville, KY         69
Omaha, NE              68
Albuquerque, NM        59
                       ..
St. Louis, MO           1
Laredo, TX              1
Norfolk, VA             1
Virginia Beach, VA      1
Oakland, CA             1
Name: count, Length: 84, dtype: int64
```
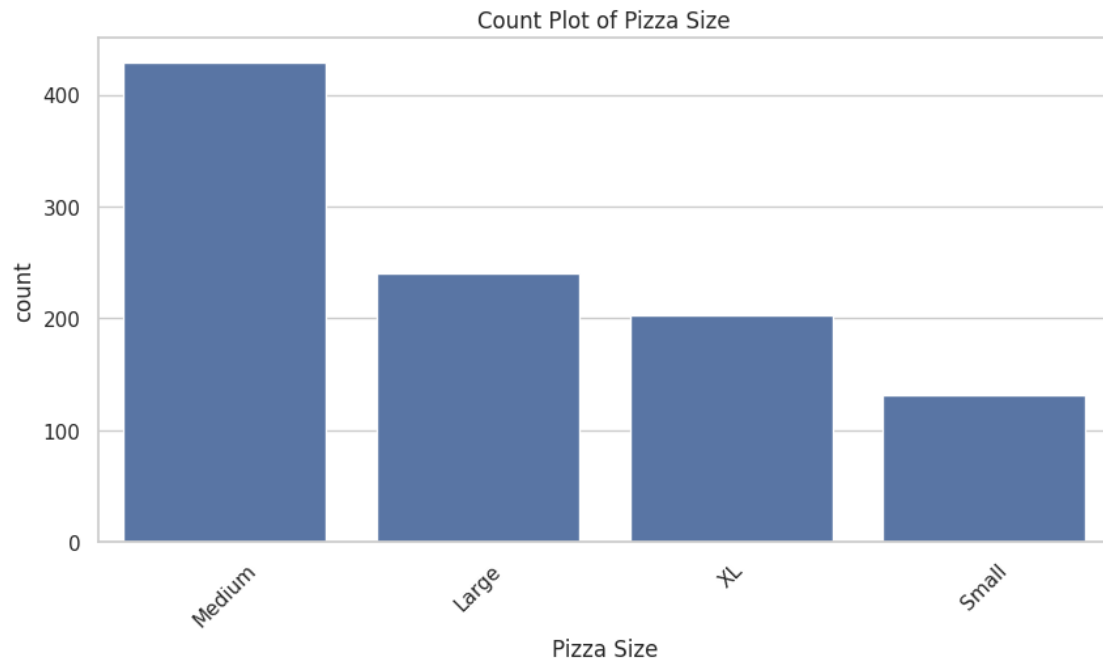
Count Plot of Location

```
Frequency Count for Pizza Size:
Pizza Size
Medium    429
Large     240
XL        203
Small     132
Name: count, dtype: int64
```

Count Plot of Pizza Size

```
Frequency Count for Pizza Type:
Pizza Type
Non-Veg         216
Veg             202
Cheese Burst    188
Vegan           132
Sicilian         64
BBQ Chicken      54
Margarita        54
Gluten-Free      41
Thai Chicken     24
Stuffed Crust    21
Deep Dish         6
Thin Crust        2
Name: count, dtype: int64
```
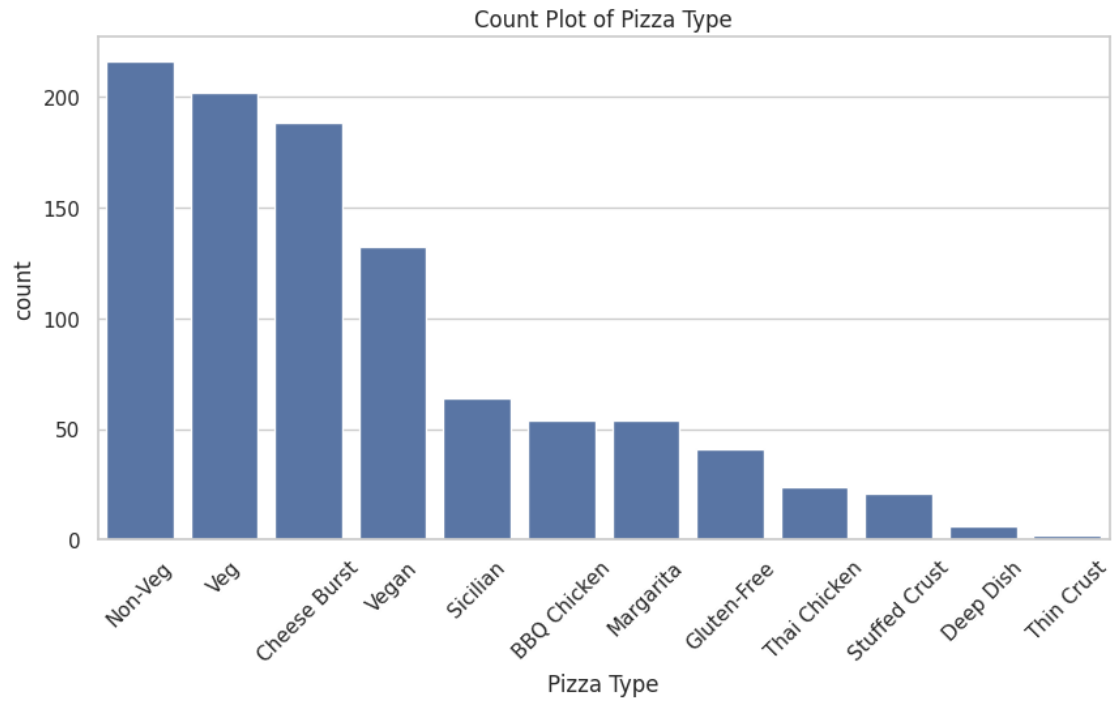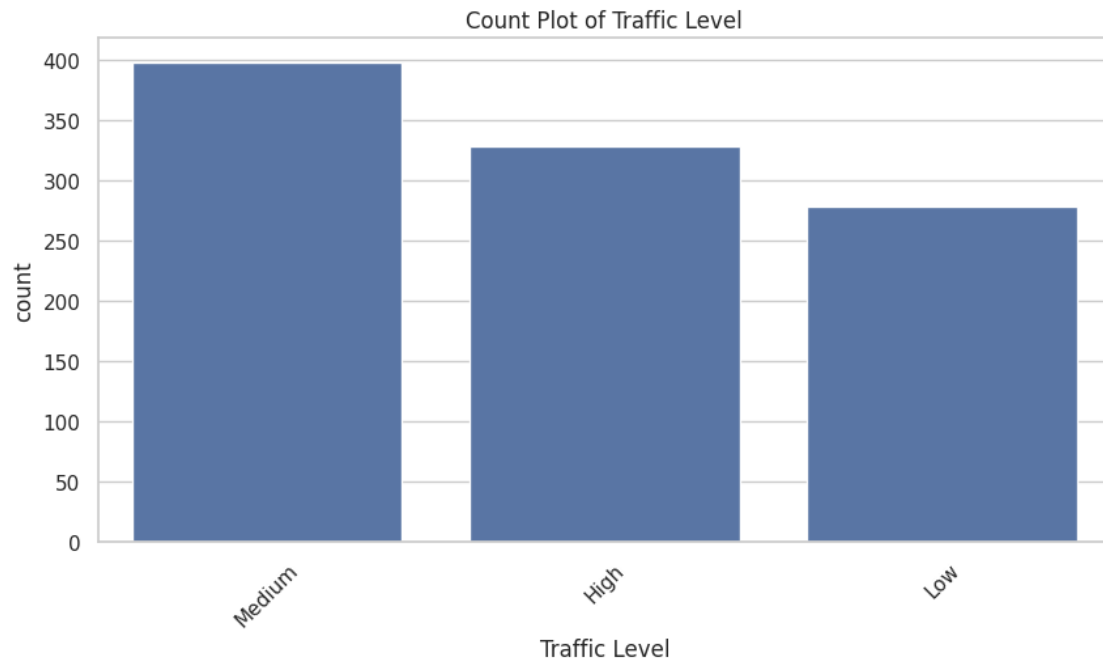
Count Plot of Pizza Type

```
Frequency Count for Traffic Level:
Traffic Level
Medium    398
High      328
Low       278
Name: count, dtype: int64
```

## Count Plot of Traffic Level



```
Frequency Count for Payment Method:
Payment Method
Card             276
UPI              271
Wallet           208
Cash             202
Hut Points        24
Domino's Cash     23
Name: count, dtype: int64
```

Count Plot of Payment Method



```
Frequency Count for Order Month:
Order Month
August       117
September    105
March         94
April         91
June          91
January       88
February      86
December      80
October       75
May           68
November      60
July          49
Name: count, dtype: int64
```
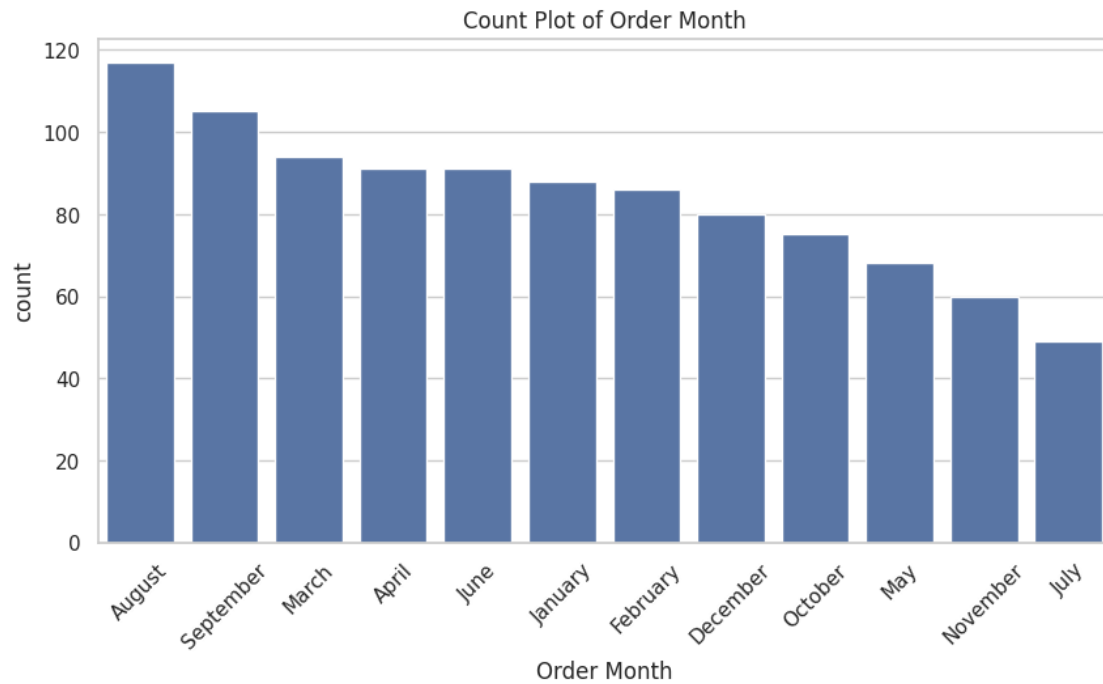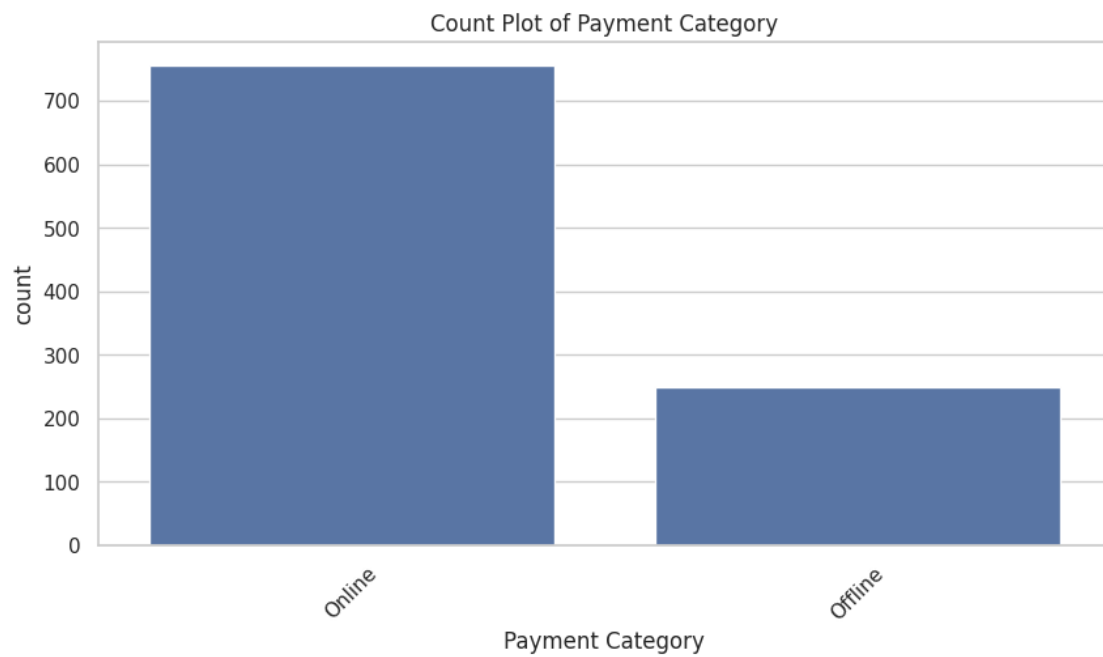
Count Plot of Order Month

Frequency Count for Payment Category:
Payment Category
Online     755
Offline    249
Name: count, dtype: int64



Count Plot of Payment Category

```
[58]:  # @title Bivariate Analysis
       # Numerical vs. Numerical: Correlation matrix
       corr_matrix = df[numeric_columns].corr()
       plt.figure(figsize=(10, 8))
       sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt='.2f')
       plt.title('Correlation Matrix')
       plt.show()
```



#Insights from Correlation Matrix

**Strong Positive Correlations ( 0.7)**

**Delivery Duration & Distance (0.91):**

- Longer distances directly increase delivery time. Route optimization is critical for faraway orders.

26

**Delivery Duration & Estimated Duration (0.91):**

- The system's time estimates are highly accurate, aligning closely with actual delivery times.

**Delivery Duration & Delay (0.88):**

- Delays (e.g., traffic, preparation) are the primary driver of longer delivery times.

**Toppings Count & Pizza Complexity (0.96):**

- More toppings = more complex pizzas, likely requiring extra preparation time.

**Distance & Estimated Duration (1.00):**

- Distance is the sole factor in time estimates (no traffic or prep adjustments).

### Strong Negative Correlations ( -0.4)

**Delivery Efficiency & Toppings Count (-0.71):**

- Pizzas with more toppings slow down min/km efficiency, possibly due to handling or packaging.

**Delivery Efficiency & Distance (-0.70):**

- Efficiency drops sharply for longer distances, suggesting non-linear delays (e.g., traffic, route issues).

**Topping Density & Delivery Duration (-0.53):**

- Higher topping density (toppings per size unit) correlates with faster deliveries, possibly due to standardized recipes.

### Moderate Correlations (0.4–0.69)

**Traffic Impact & Delivery Duration (0.58):**

- Traffic worsens delays, but less than distance or prep time.

**Order Hour & Delay (0.44):**

- Peak hours (e.g., dinner) see more delays, likely from high demand or traffic.

### Weak/No Correlation ( ±0.3)

**Restaurant Avg Time & All Metrics (~0.1):**

- Restaurant speed has minimal impact on delivery times; delays are external (e.g., traffic, distance).

**Order Hour & Efficiency (-0.01):**

- Efficiency is consistent across hours, suggesting stable operations.

### Actionable Recommendations Optimize Long-Distance Deliveries:

- Prioritize route planning or satellite kitchens for orders >5 km.

**Peak Hour Management:**

- Add staff or dynamic pricing during 14:00–18:00 to mitigate delay spikes.

**Simplify High-Topping Orders:**

- Bundle popular topping combos to reduce complexity and prep time.

**Revise Efficiency Metrics**:

- Include traffic and toppings in efficiency calculations (current metric ignores key factors).

**Restaurant Collaboration**:

- Since restaurant speed doesn't affect delivery, focus on external factors (traffic, dispatch).

### Conclusion

Distance and toppings are the top controllable levers; traffic and peak hours require systemic solutions.
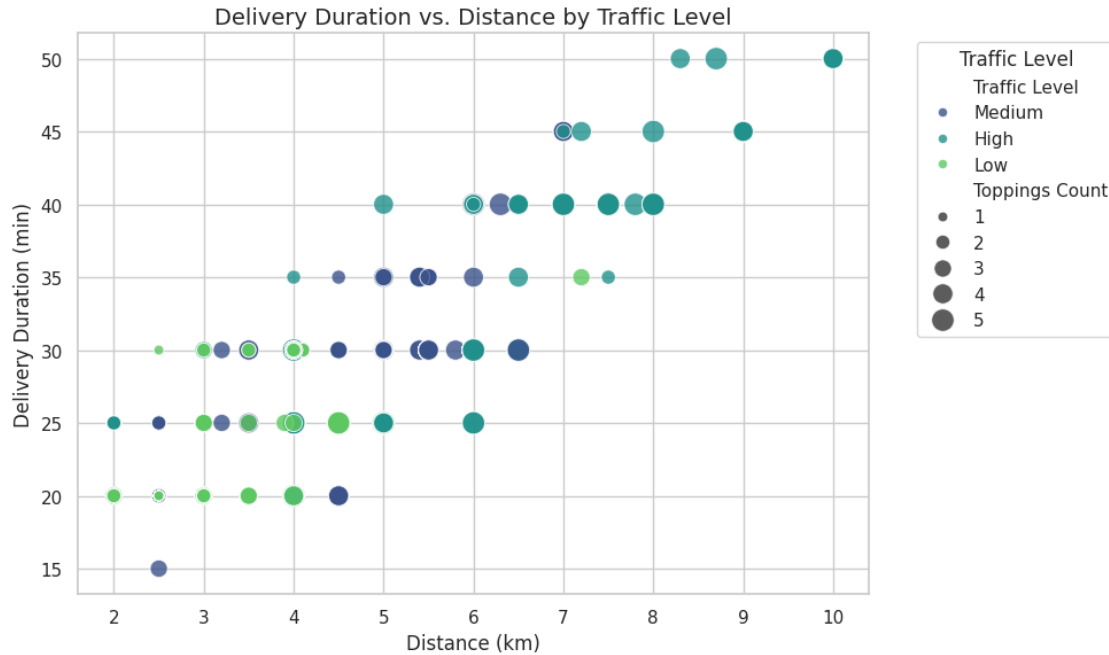
```python
# Duration vs. Distance by Traffic Level: Scatter plot
import matplotlib.pyplot as plt
import seaborn as sns

# Set figure size and style
plt.figure(figsize=(10, 6))
sns.set(style="whitegrid")

# Create scatter plot
scatter = sns.scatterplot(
    data=df,
    x='Distance (km)',
    y='Delivery Duration (min)',
    hue='Traffic Level',
    size='Toppings Count',
    palette='viridis',
    sizes=(40, 200),
    alpha=0.8
)

# Title and axis labels
plt.title('Delivery Duration vs. Distance by Traffic Level', fontsize=14)
plt.xlabel('Distance (km)')
plt.ylabel('Delivery Duration (min)')

# Show legend and plot
plt.legend(title='Traffic Level', bbox_to_anchor=(1.05, 1), loc='upper left')
plt.tight_layout()
plt.show()
```

Delivery Duration vs. Distance by Traffic Level

## 4 Insights

**Traffic Impact on Delivery Time:**

• High Traffic: Deliveries take significantly longer (e.g., 30–40 min) even for shorter distances (e.g., 2–5 km), indicating traffic congestion as a major delay factor.

• Low Traffic: Deliveries are faster (15–25 min), especially for shorter distances (<5 km), showing optimal efficiency.

• Medium Traffic: Falls between high and low, with durations closer to 20–30 min.

**Distance vs. Duration Relationship**:

• Linear Trend: Longer distances generally increase delivery time, but the effect is exaggerated during high traffic.

• Outliers: Some short-distance orders (e.g., 3 km) take as long as 35+ min during high traffic, suggesting bottlenecks (e.g., urban areas).

**Toppings Count Influence**:

• Higher toppings (4–5) correlate with slightly longer delivery times, possibly due to additional preparation time or order complexity.

• Low-topping orders (1–2) are more likely to align with the "Low Traffic" trend, indicating simpler orders are less affected by delays.
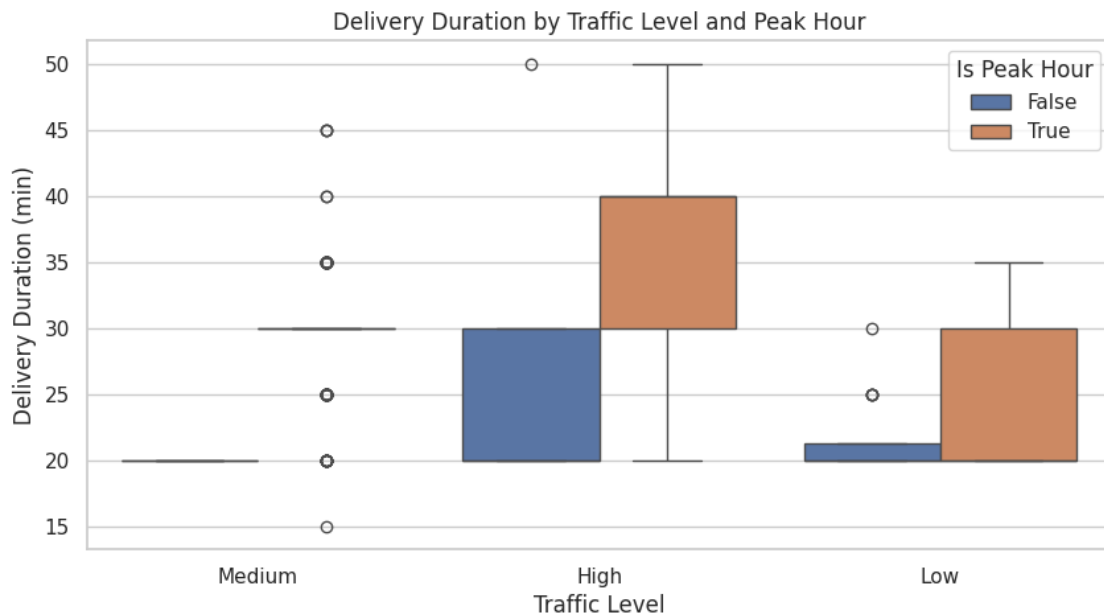
**Operational Insights**:

- Peak Traffic Mitigation: Prioritize route optimization or dynamic pricing during high-traffic hours to manage delays.

- Distance Threshold: Deliveries beyond 5 km show steeper duration increases; consider localized kitchens or delivery hubs for faraway orders.

**Key Takeaway**

Traffic level is the dominant factor in delivery duration, overshadowing distance and order complexity. Strategic adjustments (e.g., traffic-aware dispatch) could improve efficiency.

```
[60]:  # Delivery Duration by Traffic Level and Peak Hour: Box plot
       plt.figure(figsize=(10, 5))
       sns.boxplot(x='Traffic Level', y='Delivery Duration (min)', hue='Is Peak Hour',␣
         ↪data=df)
       plt.title('Delivery Duration by Traffic Level and Peak Hour')
       plt.show()
```



## 5   Insights

**Traffic Level Impact:**

- High Traffic: Consistently results in the longest delivery times, regardless of peak hours.

- Low Traffic: Deliveries are fastest, showing optimal efficiency during off-peak conditions.

- Medium Traffic: Falls between high and low, with moderate delays.

**Peak Hour Influence:**

**During Peak Hours (True):**

- Delays are exacerbated, especially in medium and high traffic, suggesting that combined demand and congestion significantly slow deliveries.

- Even low traffic sees a slight increase in duration during peak times, likely due to higher order volume.

**Off-Peak Hours (False):**

- Deliveries are faster across all traffic levels, with the most significant improvements in high-traffic scenarios.

**Worst-Case Scenario:**

- High Traffic + Peak Hour: Longest delivery times, indicating a critical need for traffic mitigation strategies during busy periods.

**Best-Case Scenario:**

- Low Traffic + Off-Peak: Fastest deliveries, representing the ideal operational condition.

**Actionable Recommendations**

### Dynamic Resource Allocation

- Increase delivery staff or fleet during peak hours to offset traffic-related delays.

- Prioritize high-traffic zones with additional resources or alternative routes.

### Peak Hour Adjustments

- Incentivize off-peak orders (e.g., discounts for early/late deliveries) to balance demand.

- Communicate longer wait times proactively during peak hours to manage customer expectations.

### Traffic Mitigation

- Use real-time traffic data to optimize delivery routes dynamically.

- Partner with navigation apps to identify faster routes during high-traffic periods.

### Performance Monitoring

- Track delivery times by traffic level and hour to identify recurring bottlenecks.

- Test interventions (e.g., staggered shifts, traffic-aware dispatch) in high-impact scenarios.

### Conclusion

- Peak hours and traffic levels compound delays, but strategic adjustments can mitigate their impact.

- Focus on demand balancing and traffic-aware logistics to maintain efficiency during critical periods.

```python
# Pizza Type vs. Payment Method: Cross-tabulation
print("\nCross-tabulation of Pizza Type vs. Payment Method:")
print(pd.crosstab(df['Pizza Type'], df['Payment Method']))
```

```
Cross-tabulation of Pizza Type vs. Payment Method:
Payment Method  Card  Cash  Domino's Cash  Hut Points  UPI  Wallet
Pizza Type
BBQ Chicken       41     0              0           0   13       0
Cheese Burst      64    17              0          21   62      24
Deep Dish          0     0              0           0    6       0
Gluten-Free       12     0              1           0    6      22
Margarita          0    53              0           0    1       0
Non-Veg           65    28              6           0   70      47
Sicilian           0     1              0           0   60       3
Stuffed Crust     18     0              0           2    0       1
Thai Chicken       7    13              0           0    4       0
Thin Crust         2     0              0           0    0       0
Veg               55    30             16           1   15      85
Vegan             12    60              0           0   34      26
```
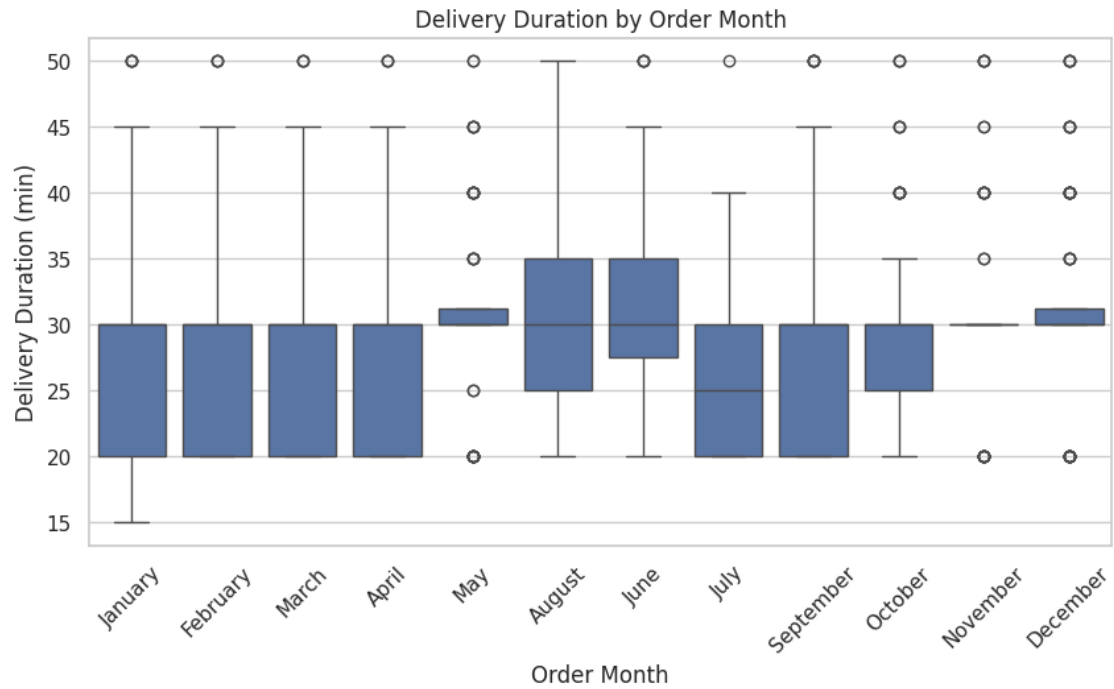
# 6   Insights

- Card, UPI, and Wallet are the most widely used payment methods overall.

- Veg and Non-Veg pizzas have the most diverse payment mix — including all methods.

- Vegan pizzas are strongly preferred with Cash and UPI, possibly indicating a different buyer profile.

- Cheese Burst sees very high usage of Hut Points and UPI, likely tied to loyalty customers.

- Wallet usage is highest for Veg pizzas (85), followed by Non-Veg and Gluten-Free.

- Deep Dish and Thin Crust pizzas show very limited purchases — possible low popularity or availability.

[62]:
```python
# @title Temporal Analysis
# Extract year and day
df['Order Year'] = df['Order Time'].dt.year
df['Order Day'] = df['Order Time'].dt.day_name()

# Delivery Duration by Order Month
plt.figure(figsize=(10, 5))
sns.boxplot(x='Order Month', y='Delivery Duration (min)', data=df)
plt.title('Delivery Duration by Order Month')
plt.xticks(rotation=45)
plt.show()
```
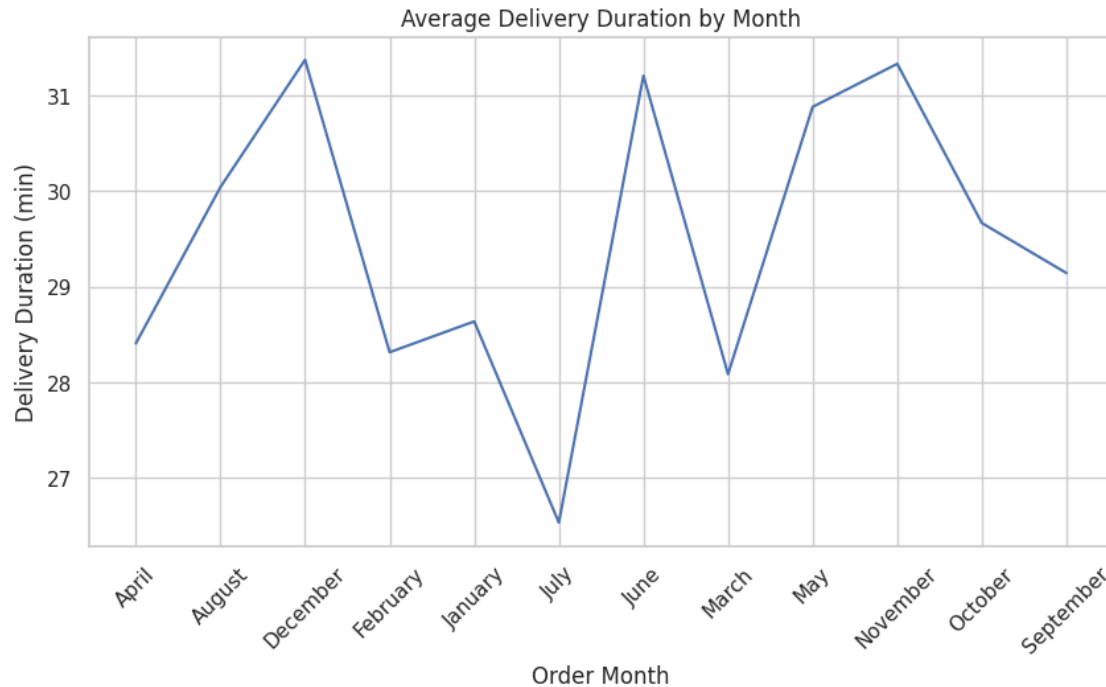
Delivery Duration by Order Month

```
[63]: # Delivery Duration by Order Hour
      plt.figure(figsize=(10, 5))
      sns.boxplot(x='Order Hour', y='Delivery Duration (min)', hue='Is Weekend',␣
        ↪data=df)
      plt.title('Delivery Duration by Order Hour and Weekend')
      plt.xticks(rotation=45)
      plt.show()
```

# Delivery Duration by Order Hour and Weekend



```
[64]:  # Line plot for monthly trends
       monthly_avg = df.groupby('Order Month')['Delivery Duration (min)'].mean().
         ↪reset_index()
       plt.figure(figsize=(10, 5))
       sns.lineplot(x='Order Month', y='Delivery Duration (min)', data=monthly_avg)
       plt.title('Average Delivery Duration by Month')
       plt.xticks(rotation=45)
       plt.show()
```

Average Delivery Duration by Month

## 7  Insights

- Monthly Trends: Delivery Duration peaks in August (~40 minutes for some orders), possibly due to higher demand or traffic.

- Hourly Trends: Evening hours (18:00-20:00) show higher durations, especially on weekdays, due to peak hours.

- Weekend vs. Weekday: Weekends may have slightly lower durations due to less traffic in some cities.

[65]:
```python
# @title Feature Engineering Exploration
# Create Delivery Speed (km/min)
df['Delivery Speed (km/min)'] = 1 / df['Delivery Efficiency (min/km)']

# Create Is Long Distance
df['Is Long Distance'] = df['Distance (km)'] >= 5
print(df['Is Long Distance'].value_counts())

# Verify new features
print("\nNew Features Summary:")
df[['Distance (km)','Delivery Speed (km/min)', 'Is Long Distance']].sample(10)
```

```
Is Long Distance
False    519
True     485
```

```
Name: count, dtype: int64

New Features Summary:
```

[65]:
```
     Distance (km)  Delivery Speed (km/min)  Is Long Distance
209            3.5                 0.175000             False
260            3.0                 0.150000             False
689            6.0                 0.200000              True
913            6.0                 0.200000              True
744            5.0                 0.166667              True
313            4.0                 0.133333             False
524            9.0                 0.200000              True
918            6.0                 0.200000              True
3              4.5                 0.180000             False
408            6.0                 0.200000              True
```

## 8  Insights

- Delivery Speed provides an alternative perspective on efficiency.

- Is Long Distance flags 485 rows with Distance $>= 5$ km, useful for analyzing long deliveries.

[66]:
```python
# @title Statistical Test
from scipy.stats import ttest_ind, f_oneway
# T-test: Delivery Duration by Peak Hour
peak = df[df['Is Peak Hour']]['Delivery Duration (min)']
non_peak = df[~df['Is Peak Hour']]['Delivery Duration (min)']
t_stat, p_value = ttest_ind(peak, non_peak)
print(f"\nT-test (Peak vs. Non-Peak Hour): t-stat={t_stat:.2f},␣
 ↪p-value={p_value:.4f}")
```

```
T-test (Peak vs. Non-Peak Hour): t-stat=7.58, p-value=0.0000
```

[67]:
```python
# ANOVA: Delivery Duration by Traffic Level
low = df[df['Traffic Level'] == 'Low']['Delivery Duration (min)']
medium = df[df['Traffic Level'] == 'Medium']['Delivery Duration (min)']
high = df[df['Traffic Level'] == 'High']['Delivery Duration (min)']
f_stat, p_value = f_oneway(low, medium, high)
print(f"ANOVA (Traffic Level): f-stat={f_stat:.2f}, p-value={p_value:.4f}")
```

```
ANOVA (Traffic Level): f-stat=252.74, p-value=0.0000
```

## 9  Insight

- The T-test result for Delivery Duration (min) between peak and non-peak hours shows a t-statistic of 7.58 and a p-value of 0.0000, indicating a highly significant difference ($p < 0.05$) between the two groups, with peak hours likely associated with longer delivery times.

- The ANOVA result for Delivery Duration (min) across Traffic Level categories yields an F-statistic of 252.74 and a p-value of 0.0000, suggesting a highly significant difference ($p < 0.05$) in delivery times across low, medium, and high traffic levels, with higher traffic levels likely contributing to increased durations.

# 10  Overall Actionable Recommendation Based on the Analysis

**Dynamic Resource Allocation:**

- Increase delivery staff or fleet capacity during peak hours (e.g., 18:00-20:00) to handle the significantly longer delivery times identified by the T-test (t-stat=7.58, p-value=0.0000), which indicate peak hours are a critical bottleneck.

- Allocate additional resources or prioritize alternative routes in high-traffic zones, as the ANOVA results (F-stat=252.74, p-value=0.0000) show a strong correlation between higher traffic levels and increased delivery durations.

**Peak Hour Adjustments:**

- Offer incentives such as discounts for off-peak orders (e.g., early morning or late evening) to balance demand and reduce congestion during peak times.

- Proactively communicate longer wait times to customers during peak hours to manage expectations and improve satisfaction.

**Traffic Mitigation:**

- Integrate real-time traffic data into route optimization systems to dynamically adjust delivery paths, especially during high-traffic periods (e.g., medium and high traffic levels).

- Partner with navigation apps to identify and recommend faster routes, leveraging the significant traffic impact on delivery times.

**Performance Monitoring:**

- Regularly track delivery durations by traffic level and hour to identify recurring bottlenecks, using the statistical significance from the ANOVA as a guide.

- Test interventions such as staggered delivery shifts or traffic-aware dispatch algorithms in high-impact scenarios to improve efficiency.

**Long-Distance Delivery Optimization:**

- Focus on improving delivery speed (e.g., average 0.133-0.183 km/min) for long-distance orders (Distance  5 km, 485 rows), which could involve assigning dedicated drivers or optimizing routes for these deliveries.

- Analyze the 519 non-long-distance deliveries to ensure they maintain high efficiency, potentially reallocating resources from shorter to longer routes as needed.

**Conclusion:**

The analysis highlights that peak hours and traffic levels significantly compound delivery delays. Strategic adjustments in resource allocation, demand balancing, and traffic-aware logistics can mitigate these impacts, ensuring consistent delivery performance across varying conditions.