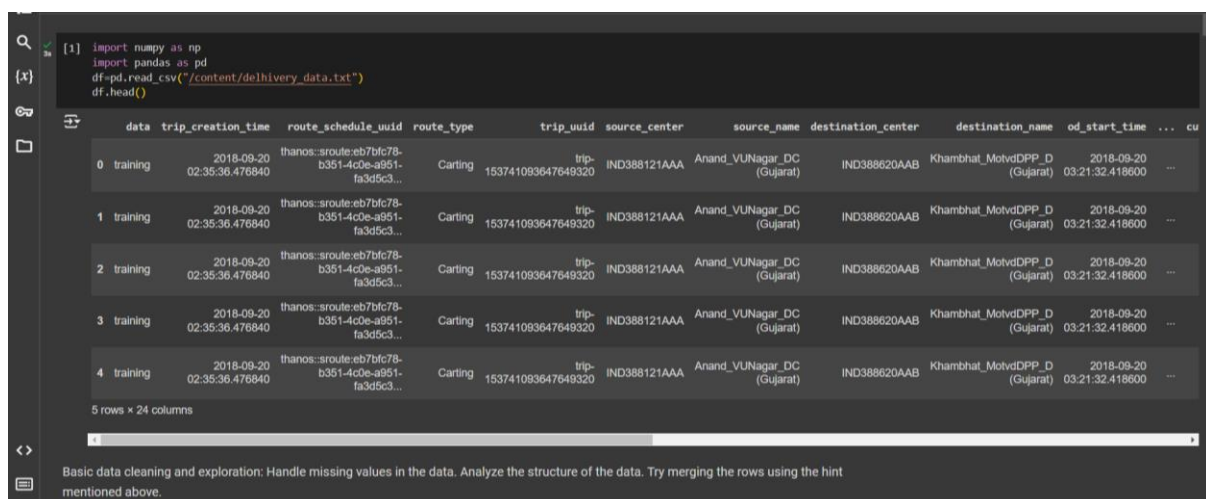## 1. Define the Problem Statement

- **Objective**: The goal is to analyze trip data to understand patterns in delivery times, distances, and other route-related attributes. The insights will be used to improve route efficiency, identify key factors impacting delivery performance, and provide actionable recommendations for business improvements.

- **Key Aspects**:

  - Identify the main sources of delays or extended delivery times.

  - Explore the relationship between distance, time, and route type.

  - Detect patterns in trip creation times to optimize resource allocation.

---

## 2. Perform Exploratory Data Analysis (EDA)

- **Data Overview**:

  - **Shape of Data**: Look at the number of rows and columns.

  -



```
[1] import numpy as np
    import pandas as pd
    df=pd.read_csv("/content/delhivery_data.txt")
    df.head()
```

| | data | trip_creation_time | route_schedule_uuid | route_type | trip_uuid | source_center | source_name | destination_center | destination_name | od_start_time | ... | cu |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | trip-153741093647649320 | IND388121AAA | Anand_VUNagar_DC (Gujarat) | IND388620AAB | Khambhat_MotvdDPP_D (Gujarat) | 2018-09-20 03:21:32.418600 | ... | |
| 1 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | trip-153741093647649320 | IND388121AAA | Anand_VUNagar_DC (Gujarat) | IND388620AAB | Khambhat_MotvdDPP_D (Gujarat) | 2018-09-20 03:21:32.418600 | ... | |
| 2 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | trip-153741093647649320 | IND388121AAA | Anand_VUNagar_DC (Gujarat) | IND388620AAB | Khambhat_MotvdDPP_D (Gujarat) | 2018-09-20 03:21:32.418600 | ... | |
| 3 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | trip-153741093647649320 | IND388121AAA | Anand_VUNagar_DC (Gujarat) | IND388620AAB | Khambhat_MotvdDPP_D (Gujarat) | 2018-09-20 03:21:32.418600 | ... | |
| 4 | training | 2018-09-20 02:35:36.476840 | thanos::sroute:eb7bfc78-b351-4c0e-a951-fa3d5c3... | Carting | trip-153741093647649320 | IND388121AAA | Anand_VUNagar_DC (Gujarat) | IND388620AAB | Khambhat_MotvdDPP_D (Gujarat) | 2018-09-20 03:21:32.418600 | ... | |

5 rows × 24 columns

Basic data cleaning and exploration: Handle missing values in the data. Analyze the structure of the data. Try merging the rows using the hint mentioned above.

- o **Convert Categorical Attributes to 'Category'**:

- o **Data Types**: List data types of all columns using df.dtypes.

```
for col in ['route_type', 'source_name', 'destination_name']:
    df[col] = df[col].astype('category')
```

```
df.dtypes
```

|  | 0 |
|---|---|
| data | object |
| trip_creation_time | object |
| route_schedule_uuid | object |
| route_type | category |
| trip_uuid | object |
| source_center | object |
| source_name | category |
| destination_center | object |
| destination_name | category |
| od_start_time | object |
| od_end_time | object |
| start_scan_to_end_scan | float64 |
| is_cutoff | bool |
| cutoff_factor | int64 |
| cutoff_timestamp | object |
| actual_distance_to_destination | float64 |

- **Missing Values**:

  - o Detect missing values with final_df.isnull().sum().

```
df.isnull().sum()
```

|  | 0 |
|---|---|
| data | 0 |
| trip_creation_time | 0 |
| route_schedule_uuid | 0 |
| route_type | 0 |
| trip_uuid | 0 |
| source_center | 0 |
| source_name | 0 |
| destination_center | 0 |
| destination_name | 0 |
| od_start_time | 0 |
| od_end_time | 0 |
| start_scan_to_end_scan | 0 |
| is_cutoff | 0 |
| cutoff_factor | 0 |
| cutoff_timestamp | 0 |
| actual_distance_to_destination | 0 |
| actual_time | 0 |
| osrm_time | 0 |
| osrm_distance | 0 |

✓ 0s    completed at 2:35 PM

- **Statistical Summary**:

  o View the summary statistics for numerical columns using final_df.describe().



  o For categorical columns, use final_df.describe(include='category').



- **Visual Analysis**:

  o **Continuous Variables**:

    ▪ Plot distributions of each continuous variable to observe spread, skewness, and potential outliers.
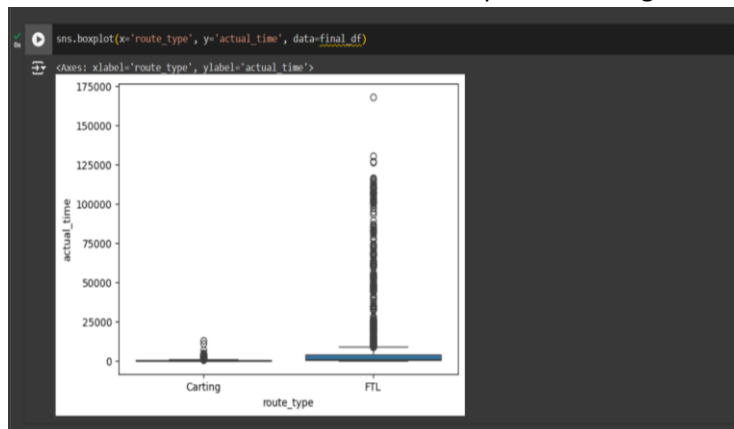
- **Categorical Variables**:

  - Use box plots to understand the distribution and spread of categorical



  variables.

- **Insights from EDA**:

  - Identify outliers in attributes like actual_time and actual_distance_to_destination.

  - Comment on the range of each attribute, noting variables with high variance.

  - Observe relationships between variables using correlation heatmaps and pairplots.

  - Insights into potential relationships, e.g., actual_time vs. osrm_time or actual_distance_to_destination vs. osrm_distance.

---

## 3. Feature Creation

- **Extract features from columns**:

  - Split destination_name and source_name into city, place code, and state.

  - Extract year, month, and day from trip_creation_time.

  - Calculate additional metrics, e.g., time_taken_to_reach_destination.

```python
[19]  # 1. Split and extract features from 'destination_name'
      final_df['destination_city_place_code'] = final_df['destination_name'].str.extract(r'([A-Za-z_]+)')
      final_df['destination_state'] = final_df['destination_name'].str.extract(r'\(([^)]+)\)')

[20]  # 2. Split and extract features from 'source_name'
      final_df['source_city_place_code'] = final_df['source_name'].str.extract(r'([A-Za-z_]+)')
      final_df['source_state'] = final_df['source_name'].str.extract(r'\(([^)]+)\)')

[21]  # Convert 'trip_creation_time' to datetime format
      final_df['trip_creation_time'] = pd.to_datetime(final_df['trip_creation_time'])

      # 3. Extract features from 'trip_creation_time'
      final_df['trip_year'] = final_df['trip_creation_time'].dt.year
      final_df['trip_month'] = final_df['trip_creation_time'].dt.month
      final_df['trip_day_of_week'] = final_df['trip_creation_time'].dt.day_name()
      final_df['trip_day'] = final_df['trip_creation_time'].dt.day
      final_df['trip_hour'] = final_df['trip_creation_time'].dt.hour
```

In-depth analysis and feature engineering: Calculate the time taken between od_start_time and od_end_time and keep it as a feature. Drop the

```python
[20]  final_df['source_city_place_code'] = final_df['source_name'].str.extract(r'([A-Za-z_]+)')
      final_df['source_state'] = final_df['source_name'].str.extract(r'\(([^)]+)\)')

[21]  # Convert 'trip_creation_time' to datetime format
      final_df['trip_creation_time'] = pd.to_datetime(final_df['trip_creation_time'])

[22]  # 3. Extract features from 'trip_creation_time'
      final_df['trip_year'] = final_df['trip_creation_time'].dt.year
      final_df['trip_month'] = final_df['trip_creation_time'].dt.month
      final_df['trip_day_of_week'] = final_df['trip_creation_time'].dt.day_name()
      final_df['trip_day'] = final_df['trip_creation_time'].dt.day
      final_df['trip_hour'] = final_df['trip_creation_time'].dt.hour
```

```python
final_df.head()
```

| ata | ... | start_scan_to_end_scan | destination_city_place_code | destination_state | source_city_place_code | source_state | trip_year | trip_month | trip_day_of_week | trip_day | trip_hour |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ing | ... | 98.0 | Delhi_Bhogal | Delhi | Delhi_Lajpat_IP | Delhi | 2018 | 9 | Wednesday | 12 | 0 |
| ing | ... | 38.0 | Faridabad | Haryana | FBD_Balabhgarh_DPC | Haryana | 2018 | 9 | Wednesday | 12 | 0 |
| ing | ... | 780.0 | MAA_Poonamallee_HB | Tamil Nadu | Chennai_Chrompet_L | Tamil Nadu | 2018 | 9 | Wednesday | 12 | 1 |
| ing | ... | 150.0 | Silvassa_Samrvmi_D | Dadra and Nagar Haveli | Vapi_IndEstat_I | Gujarat | 2018 | 9 | Wednesday | 12 | 1 |
| ing | ... | 228.0 | Dhoraji_JmnvadRd_DC | Gujarat | Jetpur_DC | Gujarat | 2018 | 9 | Wednesday | 12 | 1 |

## 4. Merging Rows and Aggregation

- **Grouping by Trip_uuid**:
  - Aggregate based on Trip_uuid, preserving the first and last values for selected fields.

```python
# Step 1: Aggregation based on Trip_uuid, Source ID, and Destination ID
grouped_df = df_new.groupby(['trip_uuid', 'source_center', 'destination_center']).agg({
    'actual_distance_to_destination': 'sum',   # Sum the distances
    'actual_time': 'sum',                      # Sum the times
    'osrm_time': 'sum',
    'osrm_distance': 'sum',
    'segment_actual_time': 'sum',
    'segment_osrm_time': 'sum',
    'segment_osrm_distance': 'sum',
    'segment_factor': 'sum',                   # Sum for segment_factor or other numeric fields
    # For categorical or non-summable fields, you can keep first or last values
    'data': 'first',                           # Keeping the first value
    'trip_creation_time': 'first',             # Keeping the first value
    'route_schedule_uuid': 'first',            # Keeping the first value
    'route_type': 'first',                     # Keeping the first value
    'source_name': 'first',                    # Keeping the first value
    'destination_name': 'first',               # Keeping the first value
    'od_start_time': 'first',                  # Keeping the first value
    'od_end_time': 'last',                     # Keeping the last value
    'start_scan_to_end_scan': 'sum',           # Sum this numeric field
}).reset_index()

[16]  # Step 2: Further aggregation based on Trip_uuid only
      final_df = grouped_df.groupby('trip_uuid').agg({
          'actual_distance_to_destination': 'sum',   # Further sum for the entire trip
          'actual_time': 'sum',
          'osrm_time': 'sum',
          'osrm_distance': 'sum',
          'segment_actual_time': 'sum',
          'segment_osrm_time': 'sum',
          'segment_osrm_distance': 'sum',
          'segment_factor': 'sum',
          # Keep the first or last values for non-numeric fields as appropriate
```

```
                    'segment_factor': 'sum',                    # Sum for segment_factor or other numeric fields
                    # For categorical or non-summable fields, you can keep first or last values
                    'data': 'first',                            # Keeping the first value
                    'trip_creation_time': 'first',              # Keeping the first value
                    'route_schedule_uuid': 'first',             # Keeping the first value
                    'route_type': 'first',                      # Keeping the first value
                    'source_name': 'first',                     # Keeping the first value
                    'destination_name': 'first',                # Keeping the first value
                    'od_start_time': 'first',                   # Keeping the first value
                    'od_end_time': 'last',                      # Keeping the last value
                    'start_scan_to_end_scan': 'sum',            # Sum this numeric field
                }).reset_index()


# Step 2: Further aggregation based on Trip_uuid only
final_df = grouped_df.groupby('trip_uuid').agg({
    'actual_distance_to_destination': 'sum',  # Further sum for the entire trip
    'actual_time': 'sum',
    'osrm_time': 'sum',
    'osrm_distance': 'sum',
    'segment_actual_time': 'sum',
    'segment_osrm_time': 'sum',
    'segment_osrm_distance': 'sum',
    'segment_factor': 'sum',
    # Keep the first or last values for non-numeric fields as appropriate
    'data': 'first',
    'trip_creation_time': 'first',
    'route_schedule_uuid': 'first',
    'route_type': 'first',
    'source_name': 'first',
    'destination_name': 'first',
    'od_start_time': 'first',
    'od_end_time': 'last',
    'start_scan_to_end_scan': 'sum'
}).reset_index()
```

○

## 5. Comparison and Visualization of Time and Distance Fields

- **Time Fields**:

    ○ Compare actual_time and osrm_time to check for discrepancies.

```python
from scipy.stats import ttest_ind
from scipy.stats import f_oneway
from scipy.stats import ttest_rel
import matplotlib.pyplot as plt
import seaborn as sns
# Paired t-test: Comparing actual_time with start_scan_to_end_scan
# Directly call ttest_rel instead of stats.ttest_rel
t_stat_actual, p_value_actual = ttest_rel(final_df['actual_time'], final_df['start_scan_to_end_scan'])

# Paired t-test: Comparing segment_actual_time with start_scan_to_end_scan
# Directly call ttest_rel instead of stats.ttest_rel
t_stat_segment, p_value_segment = ttest_rel(final_df['segment_actual_time'], final_df['start_scan_to_end_scan'])

# Results of hypothesis testing
print(f"Paired T-test (actual_time vs start_scan_to_end_scan): t-statistic = {t_stat_actual}, p-value = {p_value_actual}")
print(f"Paired T-test (segment_actual_time vs start_scan_to_end_scan): t-statistic = {t_stat_segment}, p-value = {p_value_segment}")

# Visual Analysis: Boxplot for the three variables
plt.figure(figsize=(10, 6))
sns.boxplot(data=[final_df['actual_time'], final_df['segment_actual_time'], final_df['start_scan_to_end_scan']], palette='Set3')
plt.xticks([0, 1, 2], ['Actual Time', 'Segment Actual Time', 'Start to End Scan'])
plt.title('Boxplot of Actual Time, Segment Actual Time, and Start to End Scan')
plt.show()

# ANOVA test if comparing all three together
# Directly call f_oneway instead of stats.f_oneway
anova_stat, anova_p_value = f_oneway(final_df['actual_time'], final_df['segment_actual_time'], final_df['start_scan_to_end_scan'])
print(f"ANOVA Test: statistic = {anova_stat}, p-value = {anova_p_value}")
```
```
Paired T-test (actual_time vs start_scan_to_end_scan): t-statistic = -34.77431075476514, p-value = 8.585866880917482e-255
Paired T-test (segment_actual_time vs start_scan_to_end_scan): t-statistic = -33.10951648563558, p-value = 5.79077673889952e-232
```
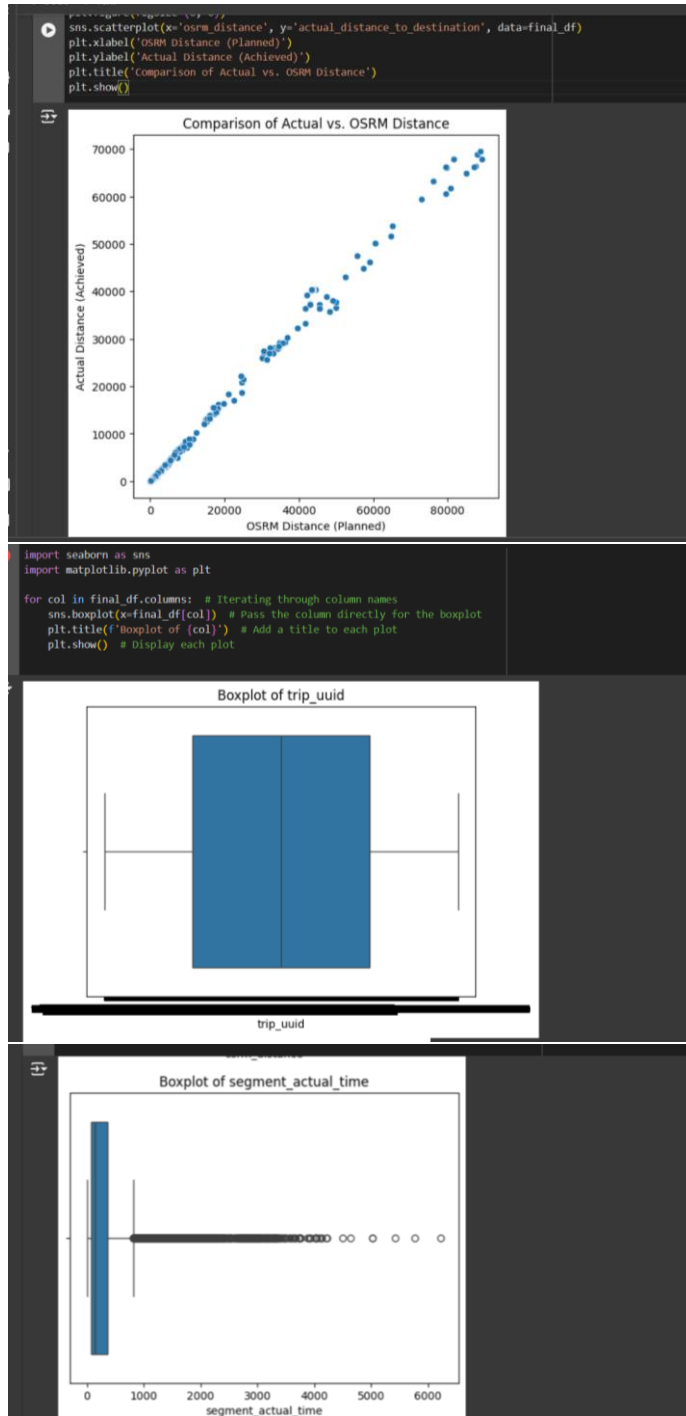
```python
ttest_rel(final_df['osrm_time'], final_df['actual_time'])
```
```
TtestResult(statistic=-32.4251460522982, pvalue=7.093080524136582e-223, df=14786)
```

- **Distance Fields**:

o  Plot the comparison between actual_distance_to_destination and osrm_distance.



```
sns.scatterplot(x='osrm_distance', y='actual_distance_to_destination', data=final_df)
plt.xlabel('OSRM Distance (Planned)')
plt.ylabel('Actual Distance (Achieved)')
plt.title('Comparison of Actual vs. OSRM Distance')
plt.show()
```



```
import seaborn as sns
import matplotlib.pyplot as plt

for col in final_df.columns:  # Iterating through column names
    sns.boxplot(x=final_df[col])  # Pass the column directly for the boxplot
    plt.title(f'Boxplot of {col}')  # Add a title to each plot
    plt.show()  # Display each plot
```





---

**6. Missing Values Treatment & Outlier Treatment**

- **Missing Value Treatment**:

  o  Use imputation methods like mean, median, or forward fill depending on the nature of the missing data.

- **Outlier Treatment**:

    o   Calculate IQR for each column, filtering out values beyond the upper and lower bounds.



---

## 7. Checking Relationship Between Aggregated Fields

- **Correlations**:

      o   Use a heatmap to explore correlations between aggregated fields like actual_time, actual_distance_to_destination, and other fields.





- **Visual Analysis**:

      o   Plot scatter plots to see how aggregated fields interact.

o Example: Relationship between cumulative distance and time can reveal insights about delivery efficiency.

---

**8. Handling Categorical Values**

- **One-Hot Encoding**:

    o Convert categorical columns like route_type to one-hot encoding to prepare them for machine learning models.

    for col in ['route_type', 'source_name', 'destination_name']:

    final_df[col] = final_df[col].astype('category')



o

---

**9. Column Normalization / Column Standardization**

- Use **MinMaxScaler** for normalization or **StandardScaler** for standardization based on your analysis needs.

```python
from sklearn.preprocessing import StandardScaler, MinMaxScaler
numerical_cols = ['actual_distance_to_destination', 'actual_time', 'osrm_time',
        'osrm_distance', 'segment_actual_time', 'segment_osrm_time',
        'segment_osrm_distance', 'segment_factor', 'start_scan_to_end_scan',
        'trip_year', 'trip_month', 'trip_day', 'trip_hour',
        'time_taken_to_reach_destination']
scaler = StandardScaler()
# Fit the scaler to the data and transform
final_df[numerical_cols] = scaler.fit_transform(final_df[numerical_cols])

# Check the transformed data
print(final_df[numerical_cols].head())
```

```
   actual_distance_to_destination  actual_time  osrm_time  osrm_distance  \
0                        0.746056     0.761748   0.714944       0.723431
1                       -0.232835    -0.241639  -0.236403      -0.233759
2                        7.480025     7.100159   7.994873       8.046926
3                       -0.256872    -0.262451  -0.259757      -0.255838
4                       -0.232972    -0.231331  -0.236780      -0.234050

   segment_actual_time  segment_osrm_time  segment_osrm_distance  \
0             2.147833           2.629714               2.633597
1            -0.381163          -0.367090              -0.332307
2             5.311326           5.594737               5.571936
3            -0.528553          -0.522809              -0.486596
4            -0.023473          -0.208192              -0.182120

   segment_factor  start_scan_to_end_scan  trip_year  trip_month  trip_day  \
0        1.917776                1.015411        0.0   -0.369459 -0.808828
1        0.040583               -0.251954        0.0   -0.369459 -0.808828
2        4.633265                7.091575        0.0   -0.369459 -0.808828
3       -0.467573               -0.272882        0.0   -0.369459 -0.808828
4        0.063635               -0.231796        0.0   -0.369459 -0.808828

   trip_hour  time_taken_to_reach_destination
```

```python
scaler = MinMaxScaler()
# Fit the scaler to the data and transform
final_df[numerical_cols] = scaler.fit_transform(final_df[numerical_cols])

# Check the transformed data
print(final_df[numerical_cols].head())
```

```
   actual_distance_to_destination  actual_time  osrm_time  osrm_distance  \
0                        0.104014     0.093341   0.101122       0.103203
1                        0.002717     0.002323   0.002651       0.002542
2                        0.800858     0.668306   0.854640       0.873362
3                        0.000229     0.000435   0.000234       0.000220
4                        0.002703     0.003258   0.002612       0.002512

   segment_actual_time  segment_osrm_time  segment_osrm_distance  \
0             0.247388           0.391712               0.373134
1             0.021218           0.023065               0.021373
2             0.530301           0.756450               0.721625
3             0.008037           0.003909               0.003074
4             0.053207           0.042611               0.039185

   segment_factor  start_scan_to_end_scan  trip_year  trip_month  trip_day  \
0        0.129781                0.109969        0.0         0.0   0.37931
1        0.037828                0.002218        0.0         0.0   0.37931
2        0.262796                0.626566        0.0         0.0   0.37931
3        0.012937                0.000439        0.0         0.0   0.37931
4        0.038958                0.003932        0.0         0.0   0.37931

   trip_hour  time_taken_to_reach_destination
0        0.0                         0.176292
1        0.0                         0.176292
2        0.0                         0.176292
3        0.0                         0.186772
4        0.0                         0.200544
```
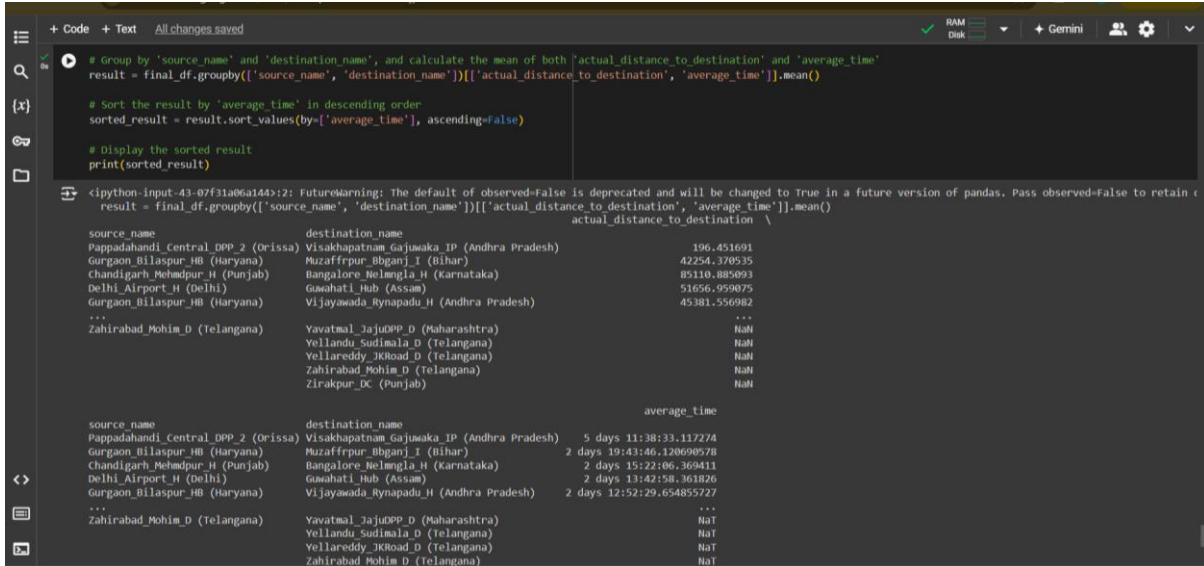
-

## 10. Business Insights

- **Pattern Analysis**:

  - Identify the busiest corridors, average time, and average distance between destinations.



  - Observe which states or regions are generating the highest volume of deliveries.

  - Insights into peak trip creation times and delays.



## 11. Recommendations

- **Route Optimization**:

  - Focus on optimizing the busiest routes by redistributing resources during peak times.

- **State-Specific Focus**:

  - Increase delivery resources in states with the highest order volume.

- **Improving Time Estimates**:

  - Update osrm_time estimates where there is a significant discrepancy between planned and actual times.