

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import plotly.express as px
import pdfplumber

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score

from google.colab import drive

from google.colab import drive
drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

```
df = pd.read_csv('/content/drive/MyDrive/DS project/traffic.csv')
```

```
df.head()
```

	DateTime	Junction	Vehicles	ID
0	2015-11-01 00:00:00	1	15	20151101001
1	2015-11-01 01:00:00	1	13	20151101011
2	2015-11-01 02:00:00	1	10	20151101021
3	2015-11-01 03:00:00	1	7	20151101031
4	2015-11-01 04:00:00	1	9	20151101041

Next steps:

Generate code with df

View recommended plots

New interactive sheet

```
df
```

	DateTime	Junction	Vehicles	ID
0	2015-11-01 00:00:00	1	15	20151101001
1	2015-11-01 01:00:00	1	13	20151101011
2	2015-11-01 02:00:00	1	10	20151101021
3	2015-11-01 03:00:00	1	7	20151101031
4	2015-11-01 04:00:00	1	9	20151101041
...
48115	2017-06-30 19:00:00	4	11	20170630194
48116	2017-06-30 20:00:00	4	30	20170630204
48117	2017-06-30 21:00:00	4	16	20170630214
48118	2017-06-30 22:00:00	4	22	20170630224
48119	2017-06-30 23:00:00	4	12	20170630234

48120 rows × 4 columns

Next steps:

Generate code with df

View recommended plots

New interactive sheet

```
df.columns
```

Index(['DateTime', 'Junction', 'Vehicles', 'ID'], dtype='object')

```
#basic overview
print("Shape:", df.shape)
print("\nData Types:\n", df.dtypes)

print("\nMissing Values:\n", df.isnull().sum())

print("\nColumns:\n", df.columns)

df.describe()
```

↻ Shape: (48120, 4)

Data Types:

DateTime	object
Junction	int64
Vehicles	int64
ID	int64

dtype: object

Missing Values:

DateTime	0
Junction	0
Vehicles	0
ID	0

dtype: int64

Columns:

Index(['DateTime', 'Junction', 'Vehicles', 'ID'], dtype='object')

	Junction	Vehicles	ID
count	48120.000000	48120.000000	4.812000e+04
mean	2.180549	22.791334	2.016330e+10
std	0.966955	20.750063	5.944854e+06
min	1.000000	1.000000	2.015110e+10
25%	1.000000	9.000000	2.016042e+10
50%	2.000000	15.000000	2.016093e+10
75%	3.000000	29.000000	2.017023e+10
max	4.000000	180.000000	2.017063e+10

df.info()

↻ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 48120 entries, 0 to 48119
Data columns (total 7 columns):

#	Column	Non-Null Count	Dtype
0	DateTime	48120 non-null	datetime64[ns]
1	Junction	48120 non-null	int64
2	Vehicles	48120 non-null	int64
3	ID	48120 non-null	int64
4	hour	48120 non-null	int32
5	day_of_week	48120 non-null	int32
6	month	48120 non-null	int32

dtypes: datetime64[ns](1), int32(3), int64(3)
memory usage: 2.0 MB

```
#checking unique values
sns.set_style('whitegrid')
```

```
print("Unique values per column:")
print(df[['Junction', 'hour', 'day_of_week', 'month']].nunique())
```

```
print("\nUnique IDs:")
print(df['ID'].nunique())
```

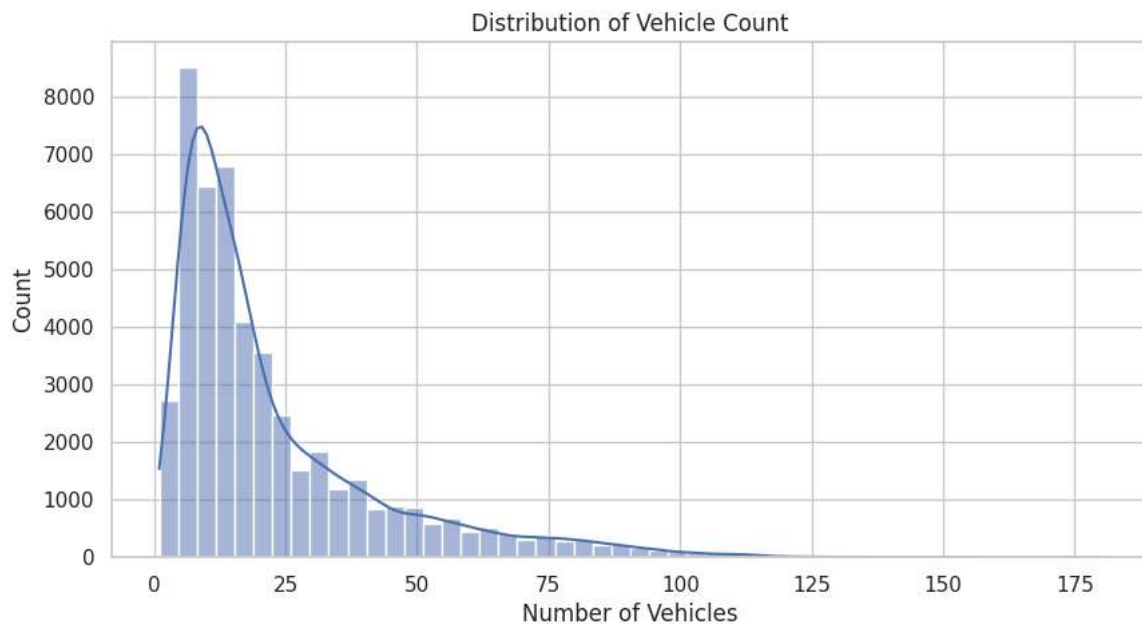
↻ Unique values per column:

Junction	4
hour	24
day_of_week	7
month	12

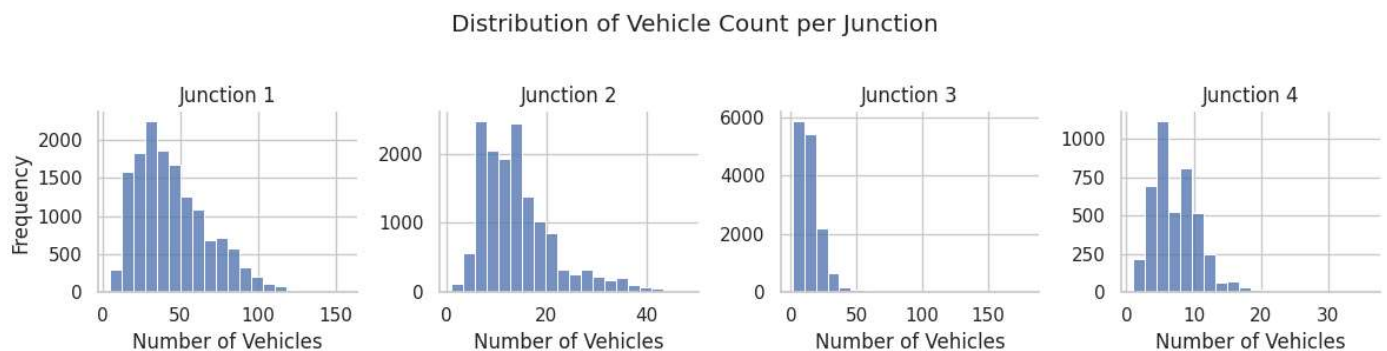
dtype: int64

Unique IDs:
48120

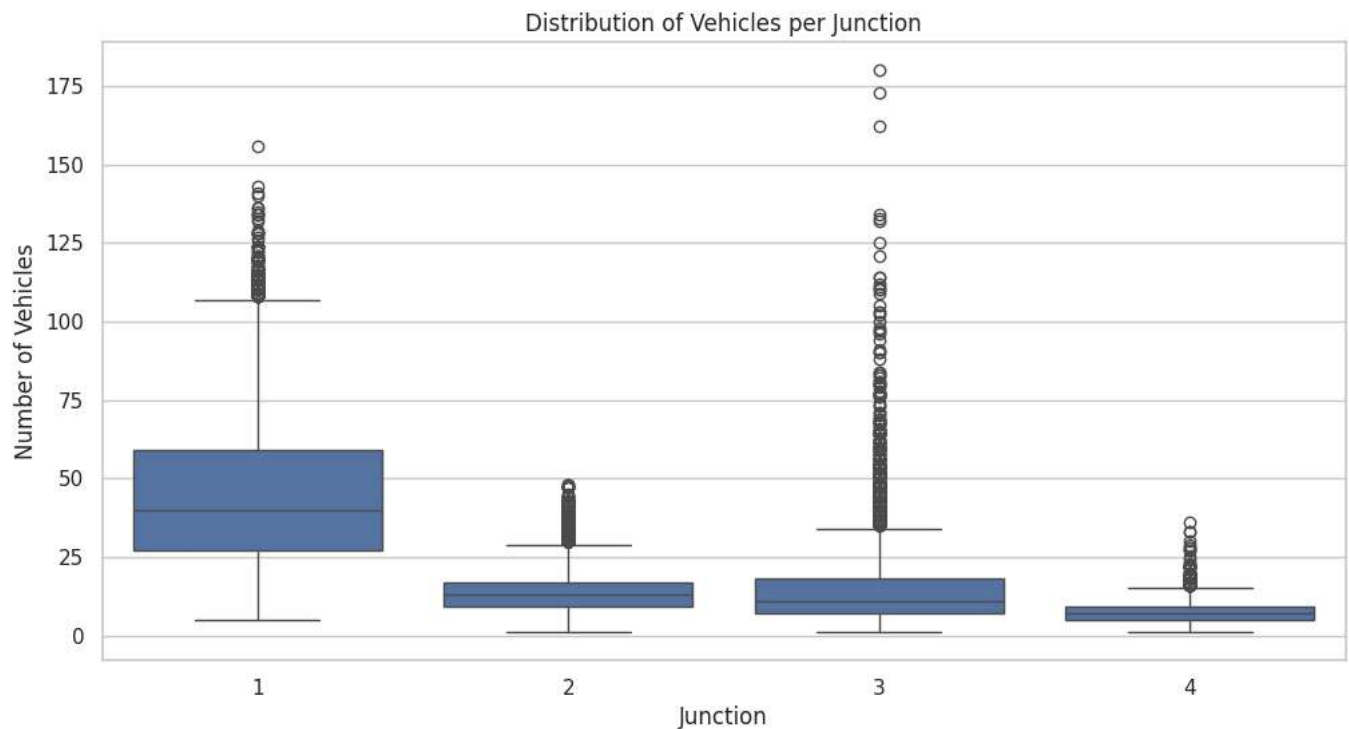
```
# Distribution curve of vehicle count
plt.figure(figsize=(10, 5))
sns.histplot(df['Vehicles'], kde=True, bins=50) # Using histplot instead of just hist
plt.title('Distribution of Vehicle Count')
plt.xlabel('Number of Vehicles')
plt.show()
```



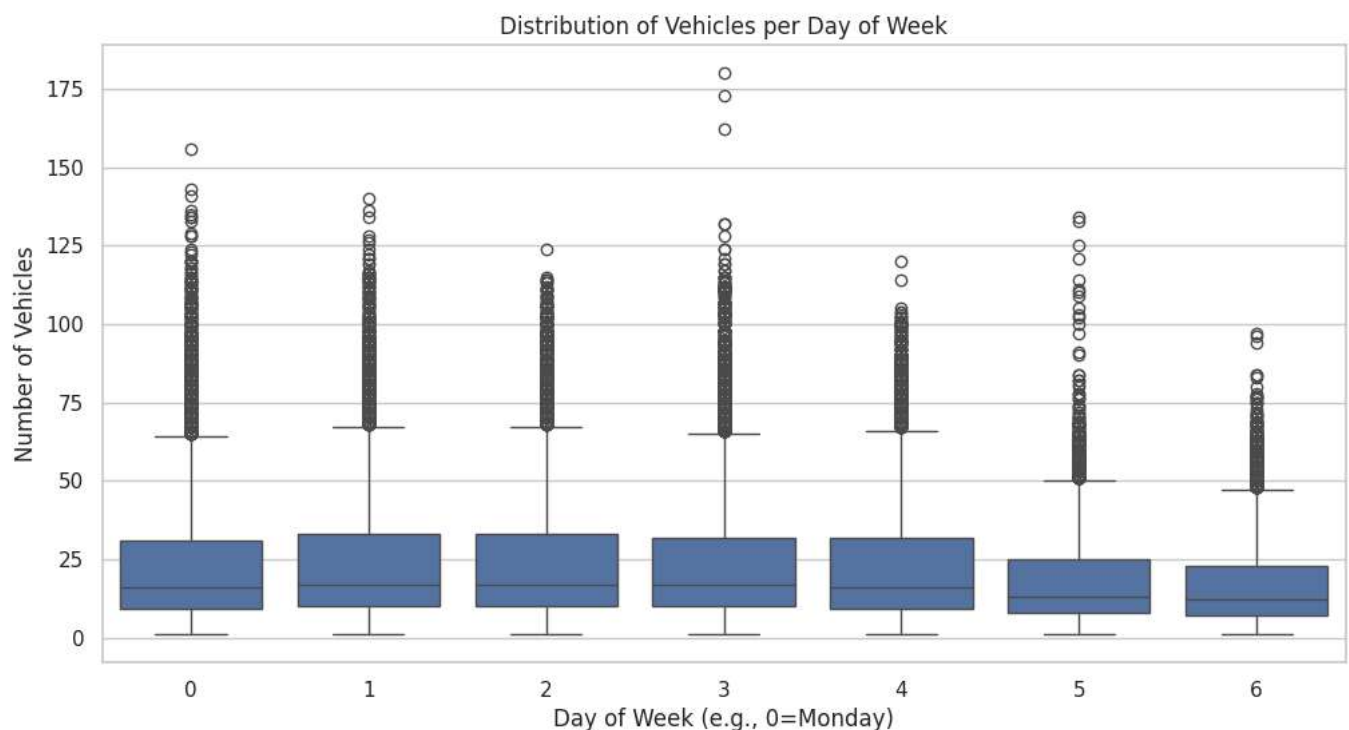
```
# FacetGrid visualization
g = sns.FacetGrid(df, col="Junction", col_wrap=4, sharex=False, sharey=False, height=3)
g.map(sns.histplot, "Vehicles", bins=20, kde=False)
g.fig.suptitle('Distribution of Vehicle Count per Junction', y=1.02)
g.set_axis_labels("Number of Vehicles", "Frequency")
g.set_titles("Junction {col_name}")
plt.tight_layout()
plt.show()
```



```
#Box and whiskers plot for each junction
plt.figure(figsize=(12, 6))
sns.boxplot(data=df, x='Junction', y='Vehicles')
plt.title('Distribution of Vehicles per Junction')
plt.xlabel('Junction')
plt.ylabel('Number of Vehicles')
plt.show()
```



```
#Box and whiskers plot for each day
plt.figure(figsize=(12, 6))
sns.boxplot(data=df, x='day_of_week', y='Vehicles')
plt.title('Distribution of Vehicles per Day of Week')
plt.xlabel('Day of Week (e.g., 0=Monday)')
plt.ylabel('Number of Vehicles')
plt.show()
```

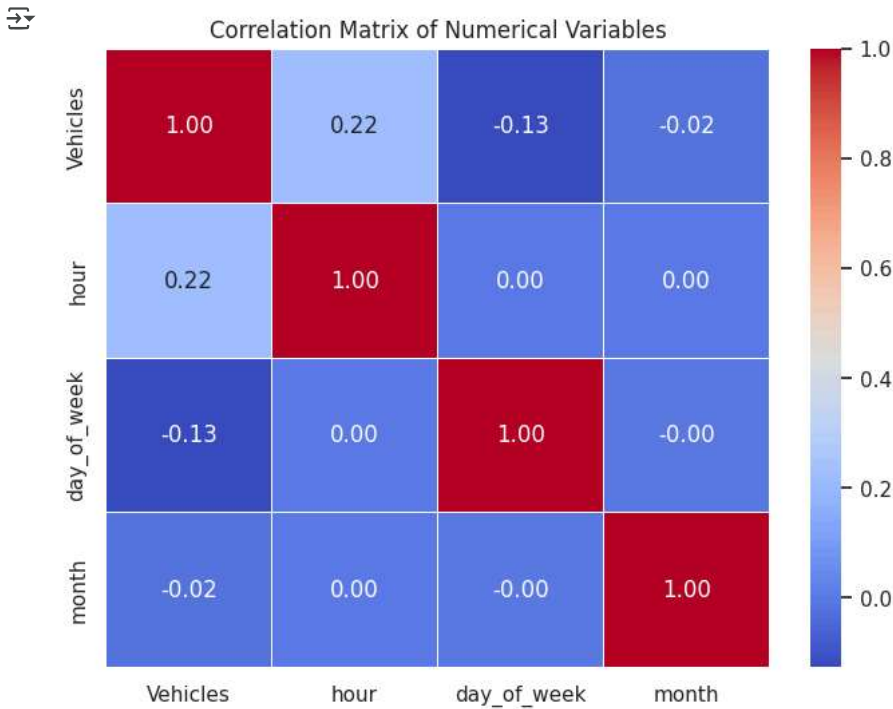


```
#heatmap
numerical_cols = ['Vehicles', 'hour', 'day_of_week', 'month']
correlation_matrix = df[numerical_cols].corr()

plt.figure(figsize=(8, 6))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=.5)

plt.title('Correlation Matrix of Numerical Variables')
plt.show()

print("\nCorrelation Matrix:")
print(correlation_matrix)
```



```
Correlation Matrix:
      Vehicles      hour  day_of_week      month
Vehicles  1.000000  2.199377e-01 -1.260265e-01 -2.272345e-02
hour      0.219938  1.000000e+00  6.442124e-18  4.967201e-16
day_of_week -0.126027  6.442124e-18  1.000000e+00 -3.208219e-03
month      -0.022723  4.967201e-16 -3.208219e-03  1.000000e+00
```

```
data = {
    'City': ['Mumbai', 'Pune', 'Nagpur', 'Nashik', 'Thane'],
    'Road_Length_km': [2200, 19391, 14734, 19769, 4479],
    'Vehicle_Count_lakhs': [42, 32, 18, 15, 25],
    'Population_lakhs': [124, 95, 46, 35, 75],
    'Congestion_Index_2023': [52, 46, 38, 33, 40]
}
```

```
city_df = pd.DataFrame(data)
city_df
```

	City	Road_Length_km	Vehicle_Count_lakhs	Population_lakhs	Congestion_Index_2023
0	Mumbai	2200	42	124	52
1	Pune	19391	32	95	46
2	Nagpur	14734	18	46	38
3	Nashik	19769	15	35	33
4	Thane	4479	25	75	40

```
df = pd.read_csv('/content/drive/MyDrive/DS project/city_traffic_data_1000.csv')

df.head()
```

	City	Road_Length_km	Vehicle_Count_lakhs	Population_lakhs	Congestion_Index_2023
0	City_1	15895	43.31	144.38	30.67
1	City_2	960	46.23	76.96	50.88
2	City_3	5490	23.82	135.41	68.50
3	City_4	12064	24.56	55.11	35.23
4	City_5	11384	46.00	27.53	34.75

Next steps:

[Generate code with df](#)

[View recommended plots](#)

[New interactive sheet](#)

```
# Convert lakh-based values to actual numbers
df['Vehicle_Count'] = df['Vehicle_Count_lakhs'] * 100000
df['Population'] = df['Population_lakhs'] * 100000

df['Vehicles_per_km'] = (df['Vehicle_Count'] / df['Road_Length_km']).round(2)
df['Population_per_km'] = (df['Population'] / df['Road_Length_km']).round(2)
df['Road_km_per_lakh_population'] = (df['Road_Length_km'] / df['Population_lakhs']).round(2)
df['Congestion_per_km'] = (df['Congestion_Index_2023'] / df['Road_Length_km']).round(4)

df.head()
```

	City	Road_Length_km	Vehicle_Count_lakhs	Population_lakhs	Congestion_Index_2023	Vehicle_Count	Population	Vehicles_per_km	Popu
0	City_1	15895	43.31	144.38	30.67	4331000.0	14438000.0	272.48	
1	City_2	960	46.23	76.96	50.88	4623000.0	7696000.0	4815.62	
2	City_3	5490	23.82	135.41	68.50	2382000.0	13541000.0	433.88	
3	City_4	12064	24.56	55.11	35.23	2456000.0	5511000.0	203.58	
4	City_5	11384	46.00	27.53	34.75	4600000.0	2753000.0	404.08	

Next steps:

[Generate code with df](#)

[View recommended plots](#)

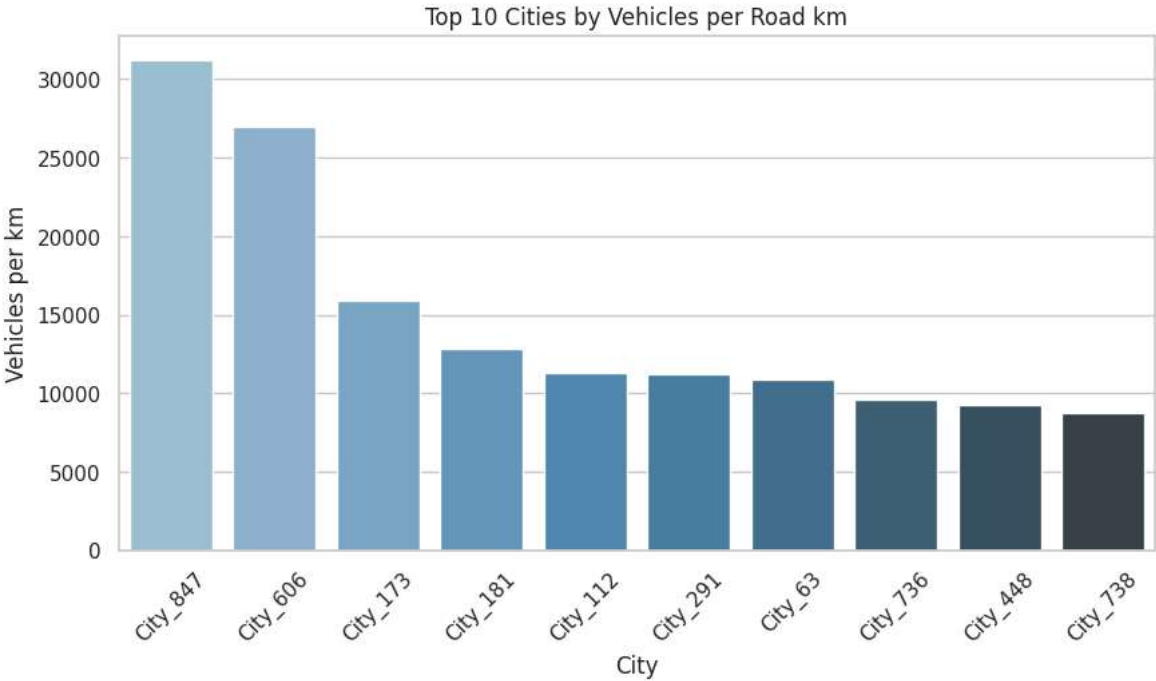
[New interactive sheet](#)

```
plt.figure(figsize=(10, 6))
sns.heatmap(df.corr(numeric_only=True), annot=True, cmap='coolwarm', fmt=".2f")
plt.title('Correlation Heatmap - Traffic & Infrastructure Factors')
plt.show()
```



```
top_vehicles_per_km = df.sort_values('Vehicles_per_km', ascending=False).head(10)

plt.figure(figsize=(10, 5))
sns.barplot(x='City', y='Vehicles_per_km', data=top_vehicles_per_km, palette='Blues_d')
plt.title('Top 10 Cities by Vehicles per Road km')
plt.xticks(rotation=45)
plt.ylabel('Vehicles per km')
plt.show()
```



Start coding or [generate](#) with AI.