# Data Harmonization and Insights Extraction

Discover how harmonizing data unlocks valuable insights for improved decision-making and competitive advantage.

by Priyanka Mathur

# Data Harmonization using tools and techniques

Techniques: o

Merging datasets. o Imputation of

Outlier detection and handling. o
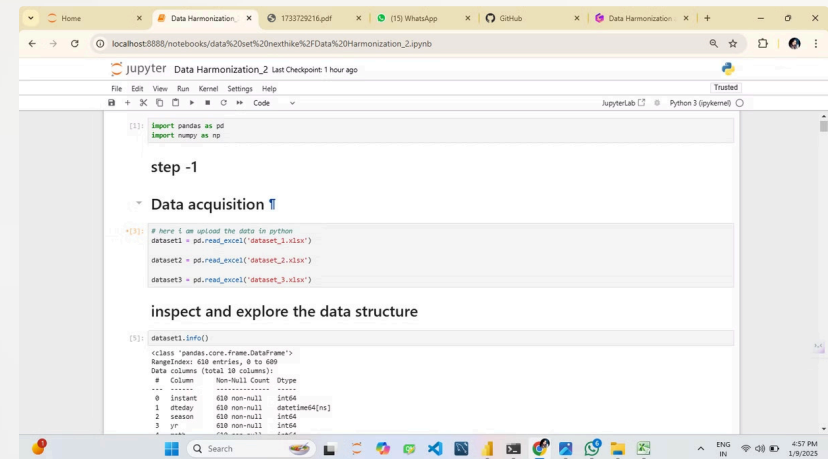
Data validation and exploration.

Tools: o

 Pandas for data manipulation.

NumPy for numerical computations. •.

# 1: Data Acquisition

• Load datasets into Python.

• Inspect and explore data structures using .info() and .head().

• Document initial observations about data quality and completeness.

# • Dataset_3 Integration
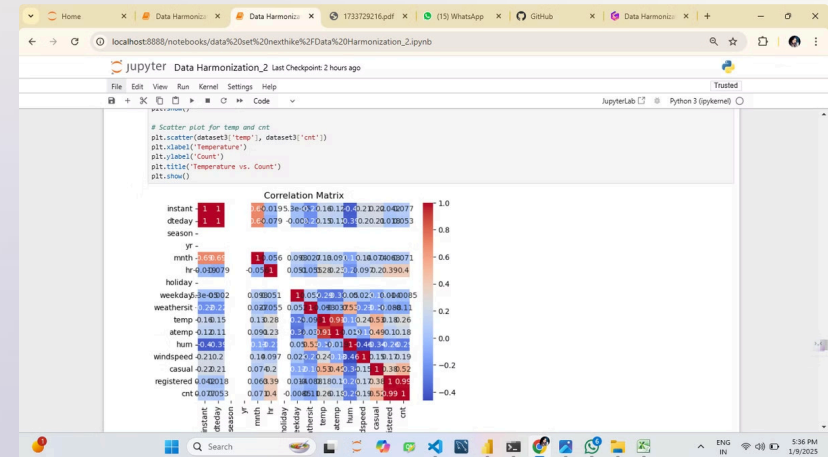


data set integration



handling missing value



handle outliers

# Compute correlations between attributes to identify relationships.
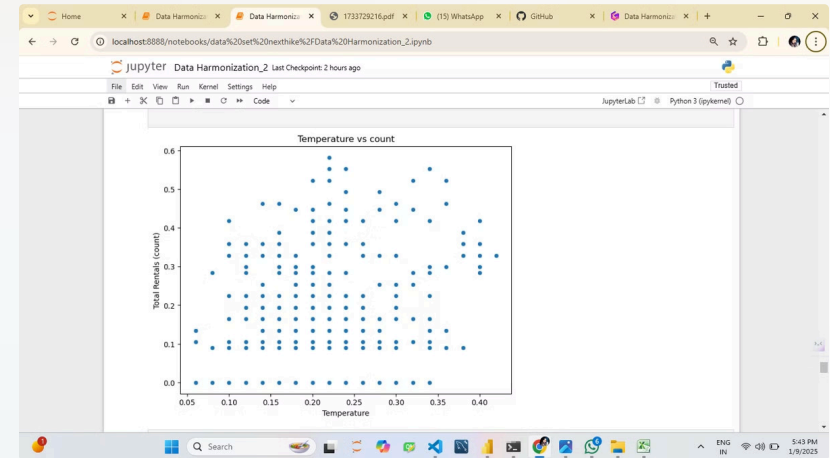


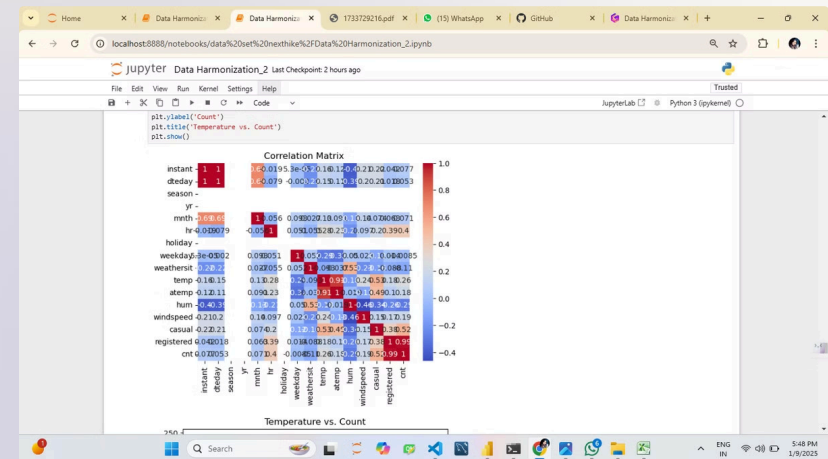using matplolip

using seaborn

# scatter plots
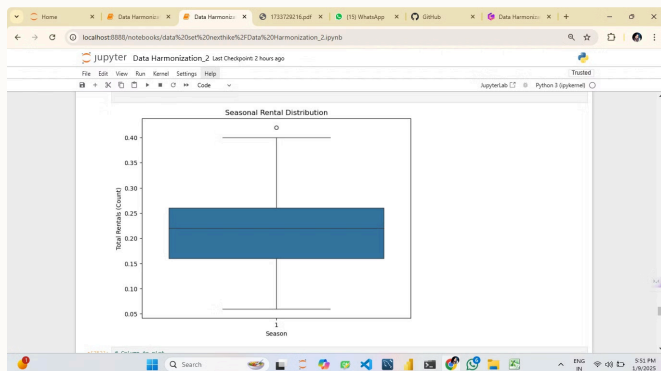
Visualizations

# Scatterplot of Temperature vs Count

plt.figure(figsize=(8, 6)) sns.scatterplot(dataset3, x='temp', y='windspeed') plt.title("Temperature vs count") plt.xlabel("Temperature") plt.ylabel("Total Rentals (count)") plt.show()
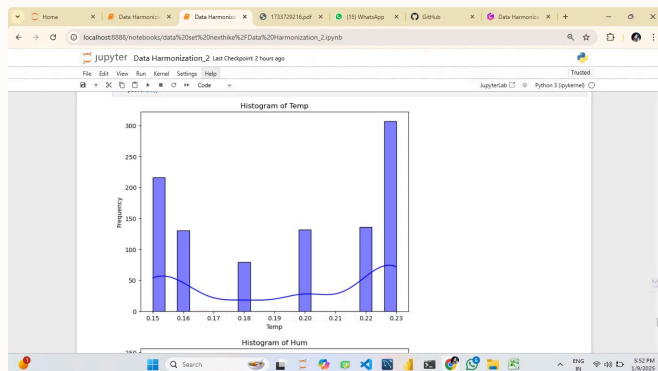
# heapmap

correlation_matrix = dataset3.corr() sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm') plt.title("Correlation Matrix") plt.show()
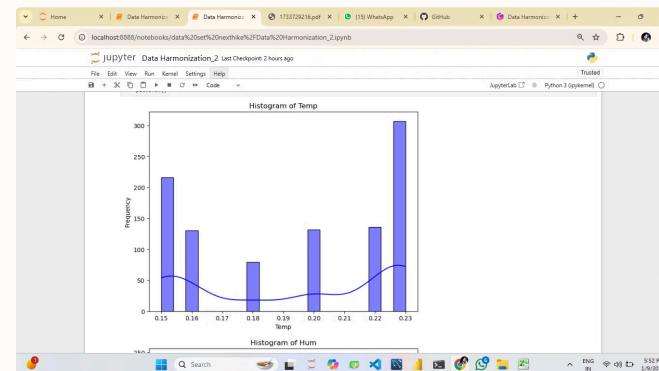
box plot



histogram



Histogram