




## **Yulu: Transforming Urban Mobility with Smart, Sustainable Rides.**


Yulu, India's leading micro-mobility service provider, is redefining urban transportation by offering eco-friendly, tech-enabled mobility solutions. With a vision to make cities less congested and more sustainable, Yulu provides shared electric two-wheelers that seamlessly integrate into daily commutes. However, recent fluctuations in demand have prompted an in-depth analysis to uncover key factors influencing Yulu's growth and adoption across different cities in India. By leveraging data-driven insights, Yulu aims to optimize its operations, enhance user experience, and drive sustainable urban mobility forward.


 **Website:** [www.yulu.bike](http://www.yulu.bike)


---





### **The Yulu Advantage**

 **Dockless E-Mobility:** Unlock, ride, and park anywhere within designated zones using the Yulu app—no docking hassles!

 **Green and Clean:** Every Yulu ride contributes to reducing carbon emissions and easing urban traffic congestion.

 **Expanding Presence:** Operational in major cities like Bengaluru, Mumbai, and Delhi, with plans to scale further.

 **Efficient Charging Network:** AI-powered battery-swapping infrastructure ensures uninterrupted rides.


 **Data-Driven Growth:** Advanced analytics help optimize fleet management, pricing strategies, and rider engagement.

---



### **Why This Case Study?**

#### **From Yulu's Perspective:**


 **Strategic Expansion:** Yulu's decision to scale its operations in India is a calculated move to solidify its market presence. Understanding demand dynamics is crucial for tailoring services and strategies accordingly.


 **Revenue Recovery:** A recent decline in revenue calls for a deeper analysis of the factors




affecting demand. Data-driven insights will help Yulu refine pricing models, operational strategies, and customer engagement to regain profitability.

## From a Learner's Perspective:

 **Real-World Problem-Solving:** This case study presents an opportunity to apply machine learning and data analysis techniques to a practical business problem.

 **Market Insights:** Analyzing demand trends in the Indian market enhances market research skills, applicable across various industries.

 **Consulting Skills:** Learners gain experience in acting as data-driven consultants, extracting insights, and providing strategic recommendations.

---



## Business Problem Statement

### Key Questions to Address:

- 1 Which variables significantly influence the demand for shared electric cycles in the Indian market?
- 2 How well do these variables explain and predict fluctuations in demand?

By exploring these questions through data-driven analysis, Yulu can refine its operational strategy, improve user engagement, and drive sustainable mobility solutions tailored to the Indian urban landscape.

---



## Features of the Dataset

The dataset used in this analysis contains various features that impact Yulu bike demand. Below is a detailed column profiling:

Feature	Description
<b>datetime</b>	Timestamp of the ride data
<b>season</b>	Season (1: Spring, 2: Summer, 3: Fall, 4: Winter)
<b>holiday</b>	Whether the day is a holiday (extracted from <a href="#">DCHR Holiday Schedule</a> )
<b>workingday</b>	1 if the day is neither a weekend nor a holiday, otherwise 0



<b>temp</b>	Temperature in Celsius
<b>atemp</b>	"Feels like" temperature in Celsius
<b>humidity</b>	Humidity percentage
<b>windspeed</b>	Wind speed
<b>casual</b>	Count of casual/unregistered users
<b>registered</b>	Count of registered users
<b>count</b>	Total count of rental bikes (casual + registered users)

## Weather Categories

Category	Details
1	Clear, Few clouds, Partly cloudy
2	Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist
3	Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds
4	Heavy Rain + Ice Pellets + Thunderstorm + Mist, Snow + Fog






---

## The Road Ahead

Yulu envisions a future where cities are **less polluted, streets are more accessible, and mobility is effortless**. By harnessing the power of data science, AI, and smart mobility solutions, Yulu aims to make sustainable transportation the default choice for urban commuters.

 **Yulu = Smarter Rides. Greener Cities. A Better Tomorrow.**

### Content:

-  **Exploratory Data Analysis.**
-  **Test of Normality.**
-  **Hypothesis Testing.**
-  **Time Series Analysis.**
-  **Insights & Recommendations.**



## Exploratory Data Analysis:

🔌 Importing :

```
import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns

from scipy.stats import norm, zscore, boxcox, probplot

from statsmodels.stats import weightstats as tests

from statsmodels.stats.proportion import proportions_ztest

from scipy.stats import ttest_ind, ttest_rel, ttest_1samp, mannwhitneyu

from scipy.stats import chisquare, chi2, chi2_contingency

from scipy.stats import f_oneway, kruskal, shapiro, levene, kstest

from scipy.stats import pearsonr, spearmanr

import statsmodels.api as sm

from statsmodels.formula.api import ols

from sklearn.preprocessing import StandardScaler

from sklearn.preprocessing import MinMaxScaler

import warnings

warnings.filterwarnings('ignore')

url =

"https://d2bei9khq929f0.cloudfront.net/public_assets/assets/000/001/428/or

iginal/bike_sharing.csv?1642089089"

df_raw = pd.read_csv(url)

print(df_raw.shape)

df_raw.head()

df_raw.rename(columns={'count': 'total_rides'}, inplace=True)
```



(10886, 12)

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	total_rides
0	2011-01-01 00:00:00	1	0	0	1	9.84	14.395	81	0.0	3	13	16
1	2011-01-01 01:00:00	1	0	0	1	9.02	13.635	80	0.0	8	32	40
2	2011-01-01 02:00:00	1	0	0	1	9.02	13.635	80	0.0	5	27	32
3	2011-01-01 03:00:00	1	0	0	1	9.84	14.395	75	0.0	3	10	13
4	2011-01-01 04:00:00	1	0	0	1	9.84	14.395	75	0.0	0	1	1





```
df_raw.info()
```

```
>>> <class 'pandas.core.frame.DataFrame'>
RangeIndex: 10886 entries, 0 to 10885
Data columns (total 12 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   datetime        10886 non-null  object
 1   season          10886 non-null  int64
 2   holiday         10886 non-null  int64
 3   workingday      10886 non-null  int64
 4   weather         10886 non-null  int64
 5   temp           10886 non-null  float64
 6   atemp          10886 non-null  float64
 7   humidity        10886 non-null  int64
 8   windspeed       10886 non-null  float64
 9   casual          10886 non-null  int64
10  registered      10886 non-null  int64
11  total_rides     10886 non-null  int64
dtypes: float64(3), int64(8), object(1)
memory usage: 1020.7+ KB
```

```
df_raw.describe()
```

```
>>>
```

	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	total_rides
count	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000
mean	2.506614	0.028569	0.680875	1.418427	20.23086	23.655084	61.886460	12.799395	36.021955	155.552177	191.574132
std	1.116174	0.166599	0.466159	0.633839	7.79159	8.474601	19.245033	8.164537	49.960477	151.039033	181.144454
min	1.000000	0.000000	0.000000	1.000000	0.82000	0.760000	0.000000	0.000000	0.000000	0.000000	1.000000
25%	2.000000	0.000000	0.000000	1.000000	13.94000	16.665000	47.000000	7.001500	4.000000	36.000000	42.000000
50%	3.000000	0.000000	1.000000	1.000000	20.50000	24.240000	62.000000	12.998000	17.000000	118.000000	145.000000
75%	4.000000	0.000000	1.000000	2.000000	26.24000	31.060000	77.000000	16.997900	49.000000	222.000000	284.000000
max	4.000000	1.000000	1.000000	4.000000	41.00000	45.455000	100.000000	56.996900	367.000000	886.000000	977.000000

```
# Converting Necessary Columns
```

```
df_raw['datetime'] = pd.to_datetime(df_raw['datetime'])
catagorical_columns = ['season', 'holiday', 'workingday', 'weather']
for col in catagorical_columns:
    df_raw[col] = df_raw[col].astype('object')
```



```
# Checking data after converting
```

```
df_raw.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10886 entries, 0 to 10885
Data columns (total 12 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   datetime        10886 non-null  datetime64[ns]
 1   season          10886 non-null  object
 2   holiday          10886 non-null  object
 3   workingday       10886 non-null  object
 4   weather         10886 non-null  object
 5   temp            10886 non-null  float64
 6   atemp           10886 non-null  float64
 7   humidity        10886 non-null  int64
 8   windspeed       10886 non-null  float64
 9   casual          10886 non-null  int64
10  registered      10886 non-null  int64
11  total_rides     10886 non-null  int64
dtypes: datetime64[ns](1), float64(3), int64(4), object(4)
memory usage: 1020.7+ KB
```

```
# Checking null values in the data
```

```
df_raw.isnull().sum()
```

```
0
datetime 0
season    0
holiday   0
workingday 0
weather   0
temp      0
atemp     0
humidity  0
windspeed 0
casual    0
registered 0
total_rides 0
```

```
dtype: int64
```



```
min_date = df_raw['datetime'].min()
max_date = df_raw['datetime'].max()
print(f"The earliest datetime in the dataset is: {min_date}")
print(f"The latest datetime in the dataset is: {max_date}")
```



```
The earliest datetime in the dataset is: 2011-01-01 00:00:00
The latest datetime in the dataset is: 2012-12-19 23:00:00
```

```
# Iterate through each categorical column and display the counts
for col in catagorical_columns:
    print(f"\nColumn: {col}")
    print(df_raw[col].value_counts())
    print("-" * 30)
```



```
Column: season
season
4    2734
2    2733
3    2733
1    2686
Name: count, dtype: int64
-----

Column: holiday
holiday
0    10575
1     311
Name: count, dtype: int64
-----

Column: workingday
workingday
1    7412
0    3474
Name: count, dtype: int64
-----

Column: weather
weather
1    7192
2    2834
3     859
4         1
Name: count, dtype: int64
-----
```



```
df_raw.describe(include='all').T
```

	count	unique	top	freq	mean	min	25%	50%	75%	max	std
datetime	10886	NaN	NaN	NaN	2011-12-27 05:56:22.399411968	2011-01-01 00:00:00	2011-07-02 07:15:00	2012-01-01 20:30:00	2012-07-01 12:45:00	2012-12-19 23:00:00	NaN
season	10886.0	4.0	4.0	2734.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN
holiday	10886.0	2.0	0.0	10575.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN
workingday	10886.0	2.0	1.0	7412.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN
weather	10886.0	4.0	1.0	7192.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN
temp	10886.0	NaN	NaN	NaN	20.23086	0.82	13.94	20.5	26.24	41.0	7.79159
atemp	10886.0	NaN	NaN	NaN	23.655084	0.76	16.665	24.24	31.06	45.455	8.474601
humidity	10886.0	NaN	NaN	NaN	61.88646	0.0	47.0	62.0	77.0	100.0	19.245033
windspeed	10886.0	NaN	NaN	NaN	12.799395	0.0	7.0015	12.998	16.9979	56.9969	8.164537
casual_uesrs	10886.0	NaN	NaN	NaN	36.021955	0.0	4.0	17.0	49.0	367.0	49.960477
registerd_users	10886.0	NaN	NaN	NaN	155.552177	0.0	36.0	118.0	222.0	886.0	151.039033
total_rides	10886.0	NaN	NaN	NaN	191.574132	1.0	42.0	145.0	284.0	977.0	181.144454
year	10886.0	NaN	NaN	NaN	2011.501929	2011.0	2011.0	2012.0	2012.0	2012.0	0.500019
month	10886	12	May	912	NaN	NaN	NaN	NaN	NaN	NaN	NaN
hour	10886.0	NaN	NaN	NaN	11.541613	0.0	6.0	12.0	18.0	23.0	6.915838
day	10886.0	NaN	NaN	NaN	9.992559	1.0	5.0	10.0	15.0	19.0	5.476608
day_of_week	10886	7	Saturday	1584	NaN	NaN	NaN	NaN	NaN	NaN	NaN
seasons	10886	4	Winter	2734	NaN	NaN	NaN	NaN	NaN	NaN	NaN
weather_type	10886	4	Clear	7192	NaN	NaN	NaN	NaN	NaN	NaN	NaN
season_type	10886	4	Winter	2734	NaN	NaN	NaN	NaN	NaN	NaN	NaN

```
# Check for duplicates
```

```
duplicates = df_raw.duplicated()
```

```
print(df_raw[duplicates])
```

```
df_raw.duplicated().sum()
```



0





## Feature Engineering:

```
df_raw['year'] = df_raw['datetime'].dt.year
df_raw['month'] = df_raw['datetime'].dt.month
df_raw['day'] = df_raw['datetime'].dt.day
df_raw['hour'] = df_raw['datetime'].dt.hour
df_raw['month'] = df_raw['month'].replace({1: 'January',
                                           2: 'February',
                                           3: 'March',
                                           4: 'April',
                                           5: 'May',
                                           6: 'June',
                                           7: 'July',
                                           8: 'August',
                                           9: 'September',
                                           10: 'October',
                                           11: 'November',
                                           12: 'December'})

df_raw['day_of_week'] = df_raw['datetime'].dt.day_name()
df_raw['season_type'] = df_raw['season'].replace({1: 'Spring',
                                                  2: 'Summer',
                                                  3: 'Fall',
                                                  4: 'Winter'})

df_raw['weather_type'] = df_raw['weather'].replace({1: 'Clear',
                                                     2: 'Mist',
                                                     3: 'Light Snow',
                                                     4: 'Heavy Rain'})

df_raw.rename(columns={'casual': 'casual_uesrs'}, inplace=True)
df_raw.rename(columns={'registered': 'registerd_users'}, inplace=True)
```

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual_uesrs	registerd_users	total_rides	year	month	hour	day	day_of_week	seasons	weather_type	season_type
0	2011-01-01 00:00:00	1	0	0	1	9.84	14.395	81	0.0000	3	13	16	2011	January	0	1	Saturday	Spring	Clear	Spring
1	2011-01-01 01:00:00	1	0	0	1	9.02	13.635	80	0.0000	8	32	40	2011	January	1	1	Saturday	Spring	Clear	Spring
2	2011-01-01 02:00:00	1	0	0	1	9.02	13.635	80	0.0000	5	27	32	2011	January	2	1	Saturday	Spring	Clear	Spring
3	2011-01-01 03:00:00	1	0	0	1	9.84	14.395	75	0.0000	3	10	13	2011	January	3	1	Saturday	Spring	Clear	Spring
4	2011-01-01 04:00:00	1	0	0	1	9.84	14.395	75	0.0000	0	1	1	2011	January	4	1	Saturday	Spring	Clear	Spring
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
10881	2012-12-19 19:00:00	4	0	1	1	15.58	19.695	50	26.0027	7	329	336	2012	December	19	19	Wednesday	Winter	Clear	Winter
10882	2012-12-19 20:00:00	4	0	1	1	14.76	17.425	57	15.0013	10	231	241	2012	December	20	19	Wednesday	Winter	Clear	Winter
10883	2012-12-19 21:00:00	4	0	1	1	13.94	15.910	61	15.0013	4	164	168	2012	December	21	19	Wednesday	Winter	Clear	Winter
10884	2012-12-19 22:00:00	4	0	1	1	13.94	17.425	61	6.0032	12	117	129	2012	December	22	19	Wednesday	Winter	Clear	Winter
10885	2012-12-19 23:00:00	4	0	1	1	13.12	16.665	66	8.9981	4	84	88	2012	December	23	19	Wednesday	Winter	Clear	Winter

10886 rows × 20 columns



## Key Insights.

### Dataset Overview

- ✓ **Total Records & Features:** The dataset consists of **10,886 rows** and **12 columns**.
- ✓ **Data Integrity:** The dataset is **clean**—no **missing** values or **duplicate** entries.
- ✓ **Time Span:** The data covers a period from **January 1, 2011**, to **December 19, 2012**, spanning **23 months** and **719 days**.



### Data Types & Structure

#### Categorical Features:

- Columns like **season**, **holiday**, **workingday**, and **weather** are stored as integers but should be converted to **categorical (object) format** for better interpretation.

#### Numerical Features:


- **temp**, **atemp**, and **windspeed** are in float format, representing continuous weather-related values.
- **humidity**, **casual**, **registered**, and **total\_riders** are in integer format, capturing count-based metrics.


#### Datetime Handling:


- The **datetime** column is currently in object format and should be converted to **datetime format** for time-based analysis.
- Further breakdown of the **datetime** column can reveal trends based on **year**, **month**, **day**, and **hour** to extract deeper insights.

### Actionable Steps for Better Analysis

 **Convert datetime** column to proper datetime format for accurate time-based insights.

 **Reclassify categorical variables** (season, holiday, workingday, weather) to object format for better visualization and analysis.

 **Extract time-based trends** by analyzing patterns across **months**, **years**, **weekdays**, and **hours** to understand demand fluctuations.

By structuring the data properly, we pave the way for **exploratory data analysis (EDA)** and predictive modeling! 



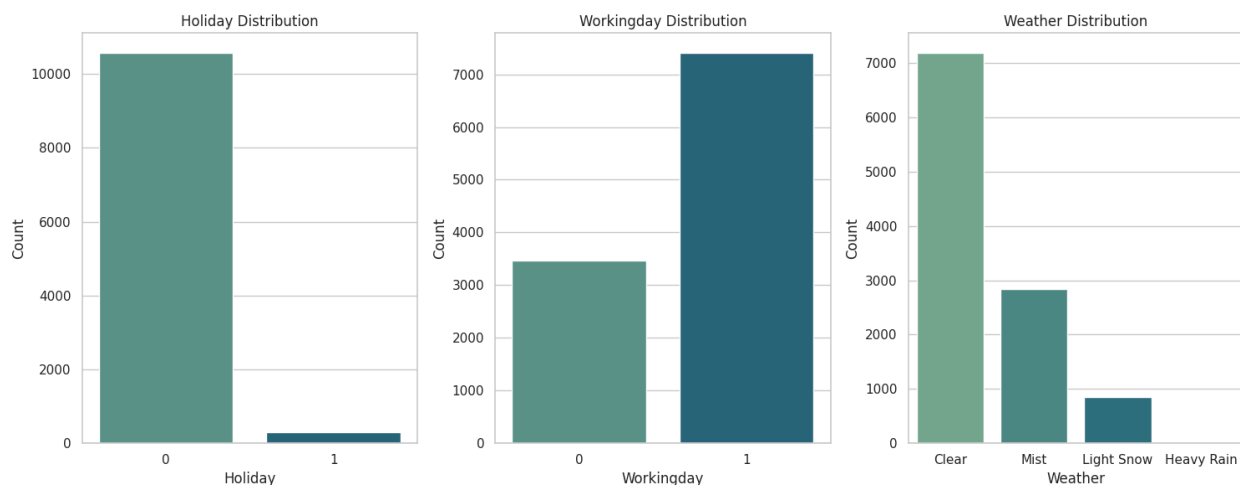
## Univariate & Bivariate Analysis.

```
plt.figure(figsize=(25, 6))
sns.set(style="whitegrid")

# Count Plot for 'holiday'
plt.subplot(1, 5, 1)
sns.countplot(data=df_raw, x='holiday', palette='crest')
plt.title('Holiday Distribution')
plt.xlabel('Holiday')
plt.ylabel('Count')

#Count Plot for 'workingday'
plt.subplot(1, 5, 2)
sns.countplot(data=df_raw, x='workingday', palette='crest')
plt.title('Workingday Distribution')
plt.xlabel('Workingday')
plt.ylabel('Count')

#Count Plot for 'weather_type'
plt.subplot(1, 5, 3)
sns.countplot(data=df_raw, x='weather_type', palette='crest')
plt.title('Weather Distribution')
plt.xlabel('Weather')
plt.ylabel('Count')
plt.tight_layout()
plt.show()
```

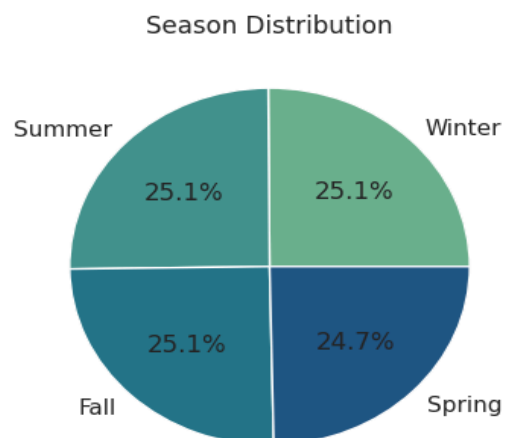
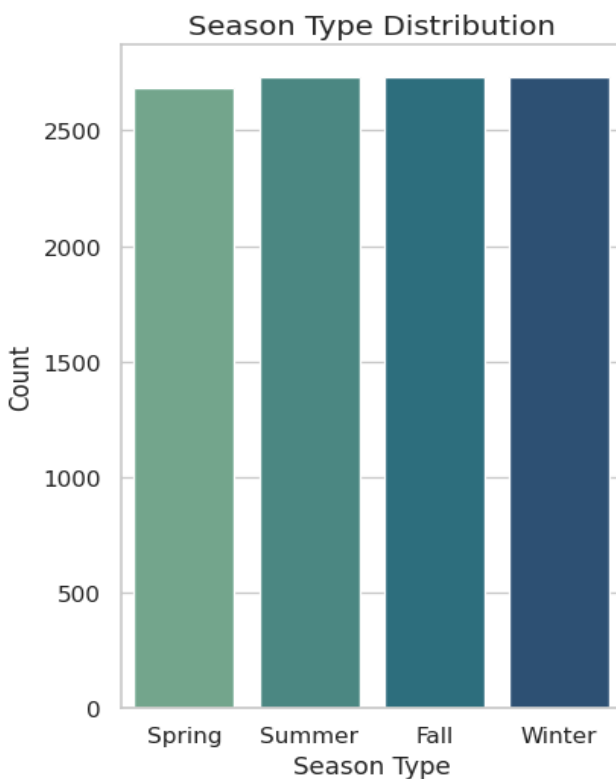




```
plt.figure(figsize=(10, 6))

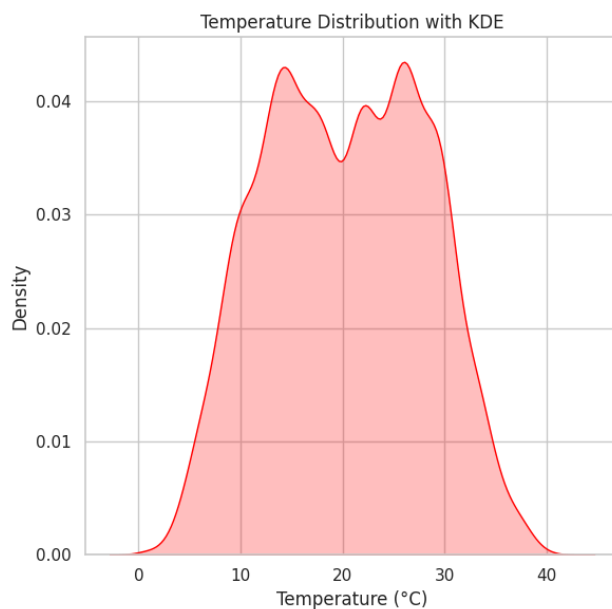
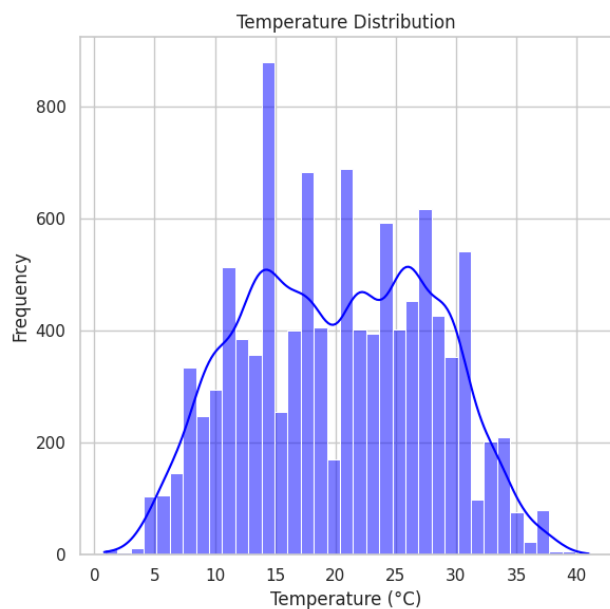
# Count Plot for 'season_type'
plt.subplot(1, 2, 1)
ax = sns.countplot(data=df_raw, x='season_type', palette='crest')
plt.title('Season Type Distribution', fontsize=14)
plt.xlabel('Season Type', fontsize=12)
plt.ylabel('Count', fontsize=12)

# Pie Chart for 'seasons'
plt.subplot(1, 2, 2)
season_counts = df_raw['season_type'].value_counts()
season_counts.plot(kind='pie', autopct='%1.1f%%', figsize=(8, 6),
colors=sns.color_palette("crest", len(season_counts)))
plt.title('Season Distribution')
plt.ylabel('')
plt.tight_layout()
plt.show()
```





```
plt.figure(figsize=(12, 6))
# Histogram for 'temp'
plt.subplot(1, 2, 1)
sns.histplot(data=df_raw, x='temp', kde=True, color='blue')
plt.title('Temperature Distribution')
plt.xlabel('Temperature (°C)')
plt.ylabel('Frequency')
# KDE Plot for 'temp'
plt.subplot(1, 2, 2)
sns.kdeplot(data=df_raw, x='temp', fill=True, color='red')
plt.title('Temperature Distribution with KDE')
plt.xlabel('Temperature (°C)')
plt.ylabel('Density')
plt.tight_layout()
plt.show()
```



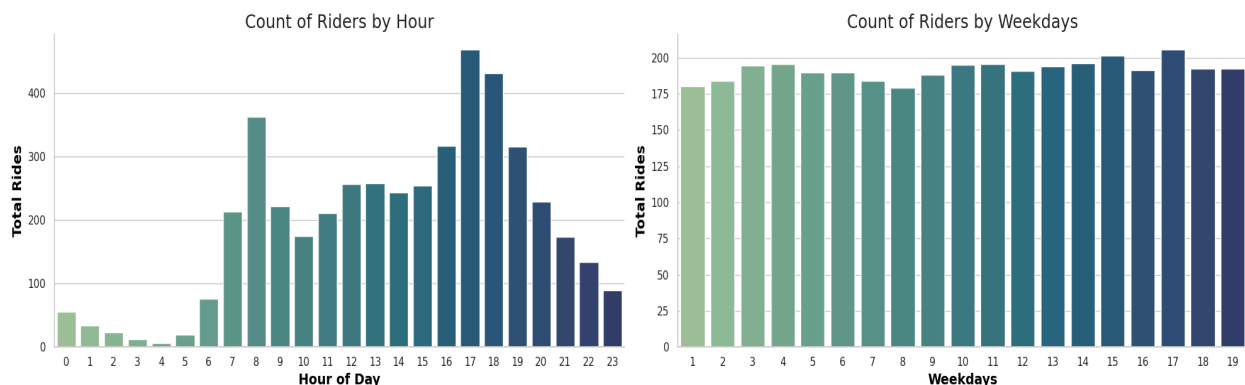


```
plt.figure(figsize=(20, 6))
sns.set(style="whitegrid")
plt.suptitle('Count of Riders', fontsize=24)

# Count of Riders by hour
plt.subplot(121)
b = sns.barplot(data=df_raw, x="hour", y="total_rides", palette='crest',
ci=None)
b.bar_label(b.containers[0], fmt='%d', color='white', fontsize=12)
plt.title('Count of Riders by Hour', fontsize=18)
plt.xlabel('Hour of Day', fontsize=14, fontweight='bold', color='black')
plt.ylabel('Total Rides', fontsize=14, fontweight='bold', color='black')

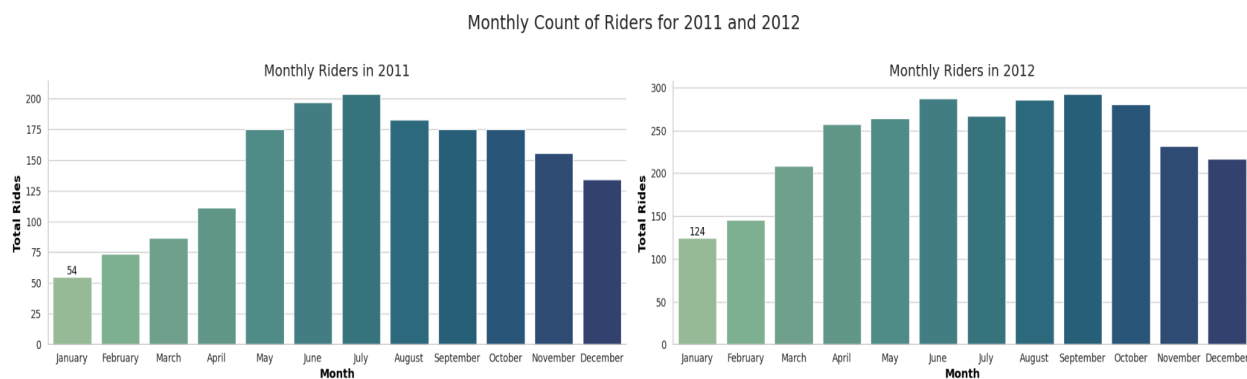
# Count of Riders by weekdays
plt.subplot(122)
b = sns.barplot(data=df_raw, x="day", y="total_rides", palette='crest',
ci=None)
b.bar_label(b.containers[0], label_type='edge', fmt='%d', color='white',
fontsize=12)
plt.title('Count of Riders by Weekdays', fontsize=18)
plt.xlabel('Weekdays', fontsize=14, fontweight='bold', color='black')
plt.ylabel('Total Rides', fontsize=14, fontweight='bold', color='black')
sns.despine()
plt.tight_layout(rect=[0, 0, 1, 0.96])
plt.show()
```

Count of Riders





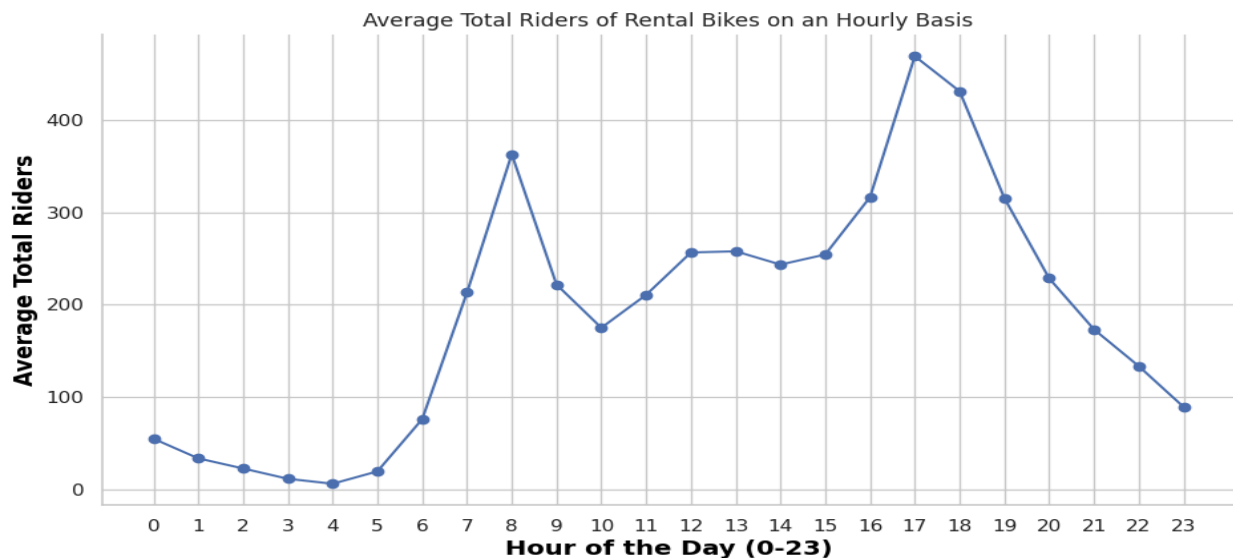
```
plt.figure(figsize=(25, 6))
sns.set_style("whitegrid")
plt.suptitle('Monthly Count of Riders for 2011 and 2012', fontsize=22)
# Count of Riders by Month for 2011
plt.subplot(121) # 1st subplot (1 row, 2 columns, 1st plot)
b = sns.barplot(data=df_raw[df_raw['year'] == 2011], x="month",
y="total_rides", palette='crest', ci=None)
b.bar_label(b.containers[0], label_type='edge', fmt='%d', color='black',
fontsize=12)
plt.title('Monthly Riders in 2011', fontsize=18)
plt.xlabel('Month', fontsize=14, fontweight='bold', color='black')
plt.ylabel('Total Rides', fontsize=14, fontweight='bold', color='black')
plt.xticks(fontsize=12)
plt.yticks(fontsize=12)
# Count of Riders by Month for 2012
plt.subplot(122) # 2nd subplot (1 row, 2 columns, 2nd plot)
b = sns.barplot(data=df_raw[df_raw['year'] == 2012], x="month",
y="total_rides", palette='crest', ci=None)
b.bar_label(b.containers[0], label_type='edge', fmt='%d', color='black',
fontsize=12)
plt.title('Monthly Riders in 2012', fontsize=18)
plt.xlabel('Month', fontsize=14, fontweight='bold', color='black')
plt.ylabel('Total Rides', fontsize=14, fontweight='bold', color='black')
plt.xticks(fontsize=12)
plt.yticks(fontsize=12)
sns.despine()
plt.tight_layout(rect=[0, 0, 1, 0.96])
plt.show()
```



```
# Count of riders by year
```



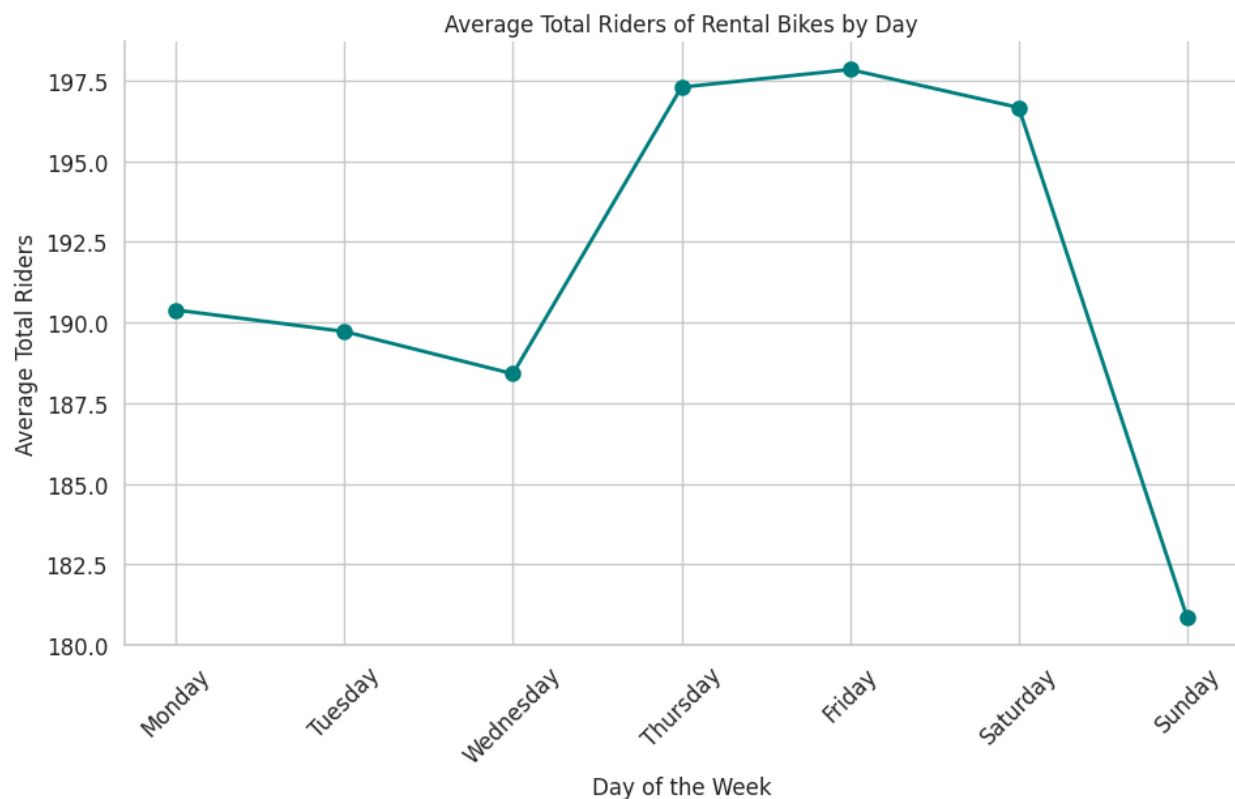
```
hourly_avg = df_raw.groupby('hour')['total_rides'].mean().reset_index()
plt.figure(figsize=(11, 6))
plt.plot(hourly_avg['hour'], hourly_avg['total_rides'], marker='o')
plt.title("Average Total Riders of Rental Bikes on an Hourly Basis")
plt.xlabel("Hour of the Day (0-23)", fontsize=14, fontweight='bold',
color='black')
plt.ylabel("Average Total Riders", fontsize=14, fontweight='bold',
color='black')
plt.xticks(np.arange(0, 24, 1))
plt.gca().spines['top'].set_visible(False)
plt.gca().spines['right'].set_visible(False)
plt.show()
```







```
day_order = ["Monday", "Tuesday", "Wednesday", "Thursday", "Friday",
             "Saturday", "Sunday"]
daily_avg =
df_raw.groupby('day_of_week')['total_rides'].mean().reset_index()
daily_avg['day_of_week'] = pd.Categorical(daily_avg['day_of_week'],
categories=day_order, ordered=True)
daily_avg = daily_avg.sort_values('day_of_week')
plt.figure(figsize=(11, 6))
plt.plot(daily_avg['day_of_week'], daily_avg['total_rides'], marker='o',
linestyle='-', color='teal', linewidth=2, markersize=8)
plt.title("Average Total Riders of Rental Bikes by Day")
plt.xlabel("Day of the Week")
plt.ylabel("Average Total Riders")
plt.xticks(fontsize=12, rotation=45)
plt.yticks(fontsize=12)
plt.gca().spines['top'].set_visible(False)
plt.gca().spines['right'].set_visible(False)
plt.show()
```



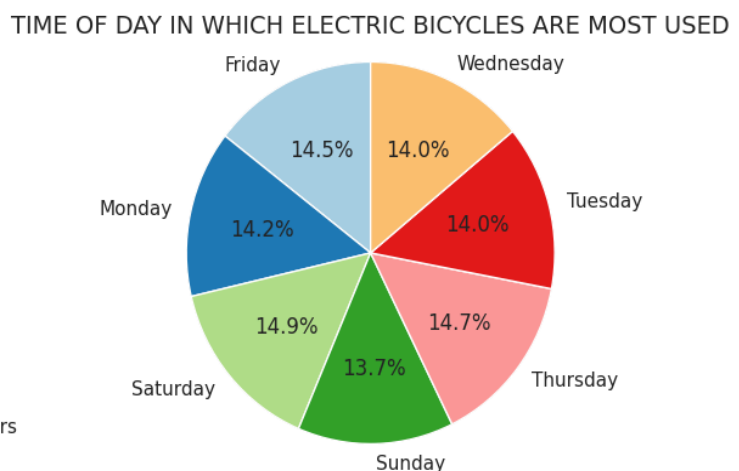
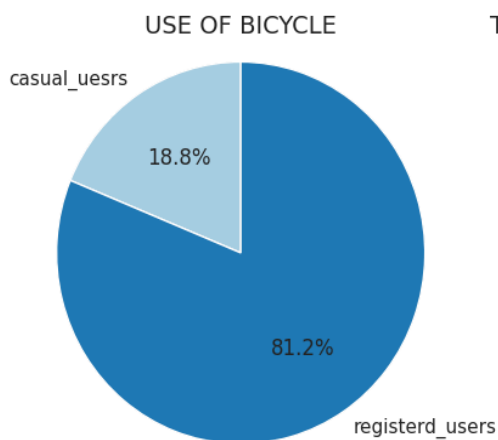


```
casual_reg = pd.DataFrame({
    'Type': ['casual_uesrs', 'registerd_users'],
    'Total': [df_raw['casual_uesrs'].sum(),
df_raw['registerd_users'].sum()]
})

timeofday_pie = df_raw.groupby('day_of_week')['total_rides'].sum()
fig, axes = plt.subplots(1, 2, figsize=(10, 4))

# First pie chart: Use of Bicycle
axes[0].pie(casual_reg['Total'], labels=casual_reg['Type'],
autopct='%1.1f%%', startangle=90, colors=plt.cm.Paired.colors)
axes[0].set_title('USE OF BICYCLE', fontsize=14)
axes[0].axis('equal')

# Second pie chart: Time of Day in which Electric Bicycles are Most Used
axes[1].pie(timeofday_pie, labels=timeofday_pie.index, autopct='%1.1f%%',
startangle=90, colors=plt.cm.Paired.colors)
axes[1].set_title('TIME OF DAY IN WHICH ELECTRIC BICYCLES ARE MOST USED',
fontsize=14)
axes[1].axis('equal')
plt.tight_layout()
plt.show()
```





```
df_raw.drop(['datetime', 'month', 'day_of_week', 'season_type',  
'weather_type'], axis=1).corr().style.background_gradient(cmap='coolwarm')
```

	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual_uesrs	registerd_users	total_rides	year	day	hour
season	1.000000	0.029368	-0.008126	0.008879	0.258689	0.264744	0.190610	-0.147121	0.096758	0.164011	0.163439	-0.004797	0.001729	-0.006546
holiday	0.029368	1.000000	-0.250491	-0.007074	0.000295	-0.005215	0.001929	0.008409	0.043799	-0.020956	-0.005393	0.012021	-0.015877	-0.000354
workingday	-0.008126	-0.250491	1.000000	0.033772	0.029966	0.024660	-0.010880	0.013373	-0.319111	0.119460	0.011594	-0.002482	0.009829	0.002780
weather	0.008879	-0.007074	0.033772	1.000000	-0.055035	-0.055376	0.406244	0.007261	-0.135918	-0.109340	-0.128655	-0.012548	-0.007890	-0.022740
temp	0.258689	0.000295	0.029966	-0.055035	1.000000	0.984948	-0.064949	-0.017852	0.467097	0.318571	0.394454	0.061226	0.015551	0.145430
atemp	0.264744	-0.005215	0.024660	-0.055376	0.984948	1.000000	-0.043536	-0.057473	0.462067	0.314635	0.389784	0.058540	0.011866	0.140343
humidity	0.190610	0.001929	-0.010880	0.406244	-0.064949	-0.043536	1.000000	-0.318607	-0.348187	-0.265458	-0.317371	-0.078606	-0.011335	-0.278011
windspeed	-0.147121	0.008409	0.013373	0.007261	-0.017852	-0.057473	-0.318607	1.000000	0.092276	0.091052	0.101369	-0.015221	0.036157	0.146631
casual_uesrs	0.096758	0.043799	-0.319111	-0.135918	0.467097	0.462067	-0.348187	0.092276	1.000000	0.497250	0.690414	0.145241	0.014109	0.302045
registerd_users	0.164011	-0.020956	0.119460	-0.109340	0.318571	0.314635	-0.265458	0.091052	0.497250	1.000000	0.970948	0.264265	0.019111	0.380540
total_rides	0.163439	-0.005393	0.011594	-0.128655	0.394454	0.389784	-0.317371	0.101369	0.690414	0.970948	1.000000	0.260403	0.019826	0.400601
year	-0.004797	0.012021	-0.002482	-0.012548	0.061226	0.058540	-0.078606	-0.015221	0.145241	0.264265	0.260403	1.000000	0.001800	-0.004234
day	0.001729	-0.015877	0.009829	-0.007890	0.015551	0.011866	-0.011335	0.036157	0.014109	0.019111	0.019826	0.001800	1.000000	0.001132
hour	-0.006546	-0.000354	0.002780	-0.022740	0.145430	0.140343	-0.278011	0.146631	0.302045	0.380540	0.400601	-0.004234	0.001132	1.000000

## Key Insights.

### Rider Trends & Time-Based Patterns.

#### ✓ Hourly Fluctuations:

- **Low rider count** in early morning hours.
- **Sharp increase in the morning** as people commute.
- **Peak demand between 4 PM - 7 PM** (evening rush hours).
- **Gradual decline at night** as activity slows.

#### ✓ Monthly & Yearly Growth:

- The **total rider count has increased year-over-year** (2011 vs. 2012).
- **Seasonality effect observed**—rider count is **highest from April to October**, with **June and July** being peak months.

#### ✓ Weekly Trends:

- **Thursday, Friday, and Saturday** have the highest average rentals.
- **Sunday** sees the **lowest** rental activity.

### Weather & Seasonal Impact on Demand



#### ✓ Weather Distribution:

- **Clear weather** sees the **most rides**, followed by **mist and light snow**.
- **Heavy rain** leads to the **lowest demand**, suggesting **adverse weather affects ridership**.

#### ✓ Temperature Analysis:

- **Normally distributed**, indicating that rides occur mostly in **moderate weather conditions**.



#### User Type Breakdown

##### ✓ Registered vs. Casual Riders:

- **Registered users dominate the dataset**, accounting for **81.2%** of total riders.
- **Casual users make up only 18.8%**, suggesting that most users rely on the service regularly.



#### Collinearity Check

##### ✓ Heatmap Analysis:

- **Correlation analysis performed** to check relationships between variables.
- Certain variables show **high collinearity**, which may impact model performance.



#### Key Takeaways:

📌 **Peak demand occurs between 4 PM - 7 PM**, highlighting an opportunity to optimize bike availability during these hours.

📌 **Registered users (81.2%) drive the majority of demand**, emphasizing the need for loyalty programs or membership perks.

📌 **Weather significantly affects demand**, with **clear weather driving the highest rides** and **heavy rain reducing usage**.

📌 **Seasonal demand fluctuations** suggest **higher fleet allocation from April to October**, especially in **June and July**.

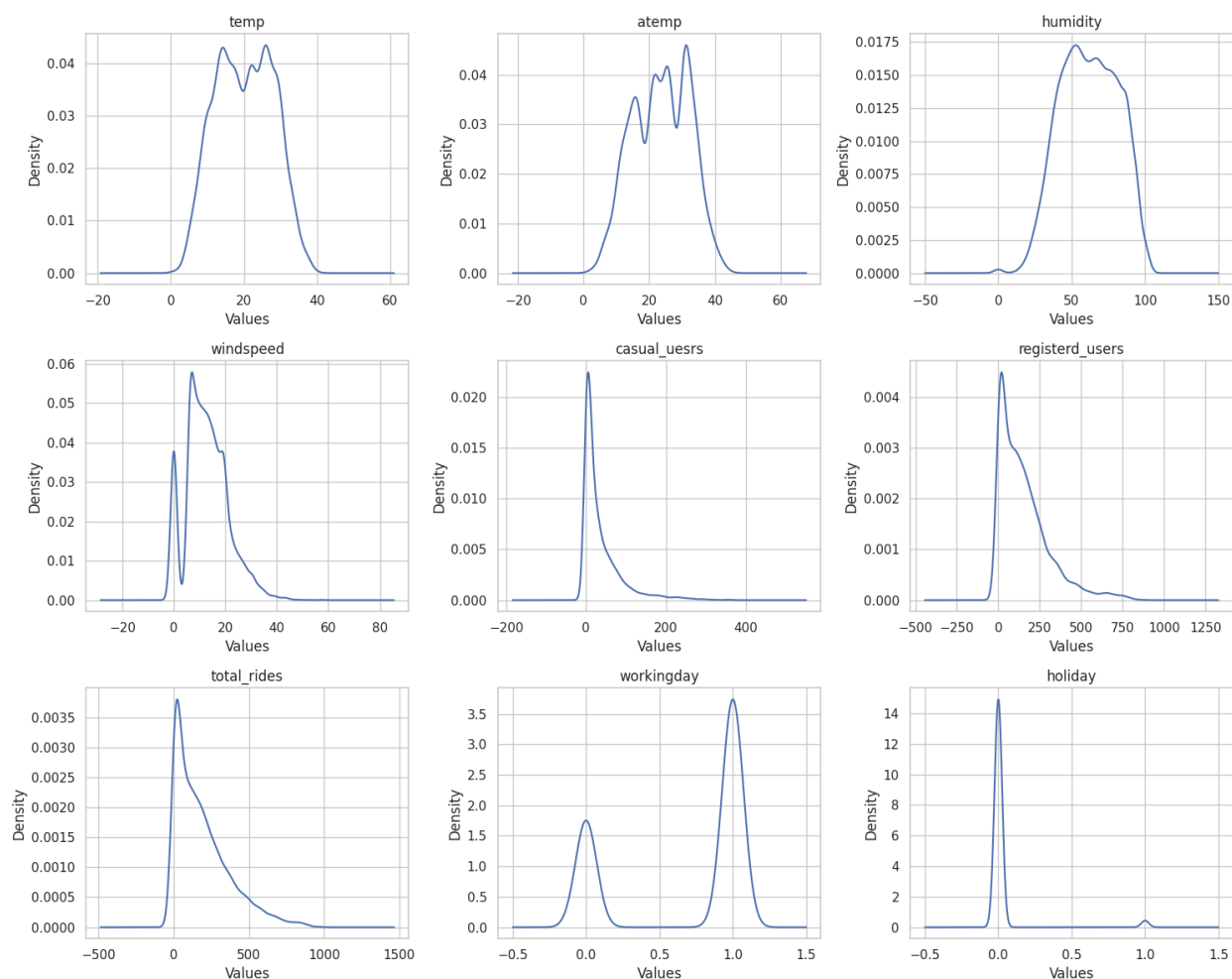
📌 **Thursday-Saturday see the highest rentals**, while **Sundays have the least demand**.

🚀 **Overall, the data suggests strong seasonal and weather-based influences, with a dominant share of registered users and predictable daily peaks. Optimizing fleet availability, pricing strategies, and weather-based promotions can enhance overall service efficiency!**



Now let's check whether the features are normally distributed or not.

```
columns =  
df_raw[['temp', 'atemp', 'humidity', 'windspeed', 'casual_uesrs', 'registered_us  
ers', 'total_rides']]  
fig, axes = plt.subplots(4, 3, figsize=(15, 12))  
axes = axes.flatten()  
for i, col in enumerate(columns.columns):  
    columns[col].plot(kind='density', ax=axes[i], title=col)  
    axes[i].set_xlabel("Values")  
    axes[i].set_ylabel("Density")  
for j in range(len(columns.columns), len(axes)):  
    axes[j].axis('off')  
plt.tight_layout()  
plt.show()
```





## Key Insights.

### Observations on Distributions of Continuous Variables

#### ✓ Right-Skewed Distributions:

- The variables **Casual**, **Registered**, and **Total Count** are **positively skewed**, indicating that **most days** have low to moderate rentals, but there are occasional **high-demand days**.

#### ✓ Binomial Distribution in Windspeed:

- **Windspeed** follows a **binomial-like distribution**, with some days recording **zero windspeed**, while others experience **low to moderate levels**.

#### ✓ Normality Trends in Temp, Atemp & Humidity:

- The distributions of **Temp**, **Atemp**, and **Humidity** appear **approximately normal**, as most values cluster **around the mean**.
- To **confirm normality**, we will apply the **Shapiro-Wilk test** for statistical validation.



#### Key Takeaways:

 **Bike rentals exhibit a right-skewed pattern**, with occasional high-demand spikes.

 **Windspeed shows a distinct pattern**, with many zero-windspeed days.

 **Temperature, Atemp, and Humidity suggest a normal distribution**, but further testing is required.



**Next Steps:** Conduct the **Shapiro-Wilk test** to statistically validate normality assumptions.



## Shapiro-Wilk Test Results:

The **Shapiro-Wilk test** was applied to check for normality in various continuous variables. The test results indicate that **none of the distributions are normal**. Below are the details:

```
for column_name in columns:
    print()
    test_statistic, p_value = shapiro(columns[column_name])
    print(f"The test-statistic for {column_name} is {test_statistic} with p-value {p_value}")
    if p_value > 0.05:
        print(f"The distribution of {column_name} is normal")
    else:
        print(f"The distribution of {column_name} is not normal")
```

1. **Temperature (Temp)**
  - **Test Statistic:** 0.9804
  - **P-value:** 4.44e-36
  - **Insight:** The distribution of **Temp** is **not normal**.
2. **Apparent Temperature (Atemp)**
  - **Test Statistic:** 0.9815
  - **P-value:** 3.22e-35
  - **Insight:** The distribution of **Atemp** is **not normal**.
3. **Humidity**
  - **Test Statistic:** 0.9823
  - **P-value:** 1.22e-34
  - **Insight:** The distribution of **Humidity** is **not normal**.
4. **Windspeed**
  - **Test Statistic:** 0.9587
  - **P-value:** 7.59e-48
  - **Insight:** The distribution of **Windspeed** is **not normal**.
5. **Casual Users**
  - **Test Statistic:** 0.7056
  - **P-value:** 3.54e-87
  - **Insight:** The distribution of **Casual Users** is **not normal**.
6. **Registered Users**
  - **Test Statistic:** 0.8563
  - **P-value:** 1.97e-71
  - **Insight:** The distribution of **Registered Users** is **not normal**.
7. **Total Rides**
  - **Test Statistic:** 0.8784
  - **P-value:** 5.37e-68
  - **Insight:** The distribution of **Total Rides** is **not normal**.



## 🎯 T test:

Before applying the **t-test**, I need to check whether my data follows a **normal distribution**. Since my data contains **outliers** and is **not normally distributed**, I will first check the normality of the data. If it is not normally distributed, I will apply a **log transformation** to see if it helps normalize the distribution and reduce the impact of outliers.

To ensure the transformation is effective, I will **compare the original and transformed data side by side** using:

- ✓ **Statistical Tests** – Shapiro-Wilk test or Kolmogorov-Smirnov test
- ✓ **Visualizations** – Histograms and Q-Q plots

- ♦ **If the transformed data follows a normal distribution**, I will proceed with the **t-test**, as the assumption of normality is met.
- ♦ **If the data remains non-normal even after transformation**, I will use a **non-parametric alternative**, such as the **Mann-Whitney U test**, which does not require normality.

This approach ensures that I select the most appropriate statistical test based on the data characteristics.

```
# Separate data into two groups: working day and non-working day
```

```
working_day = df_raw[df_raw['workingday'] == 1]['total_rides']
```

```
non_working_day = df_raw[df_raw['workingday'] == 0]['total_rides']
```

```
# Apply log transformation
```

```
log_working_day = np.log(working_day)
```

```
log_non_working_day = np.log(non_working_day)
```

```
# 1. Histogram to check skewness
```

```
plt.figure(figsize=(12, 5))
```

```
# Original Data
```

```
plt.subplot(1, 2, 1)
```





```
sns.histplot(working_day, kde=True, color="blue", label="Working Day",
bins=30)

sns.histplot(non_working_day, kde=True, color="red", label="Non-Working
Day", bins=30)

plt.legend()

plt.title("Distribution of Bike Rentals on Working vs Non-Working Days")

# Log-transformed Data

plt.subplot(1, 2, 2)

sns.histplot(log_working_day, kde=True, color="blue", label="Log-Working
Day", bins=30)

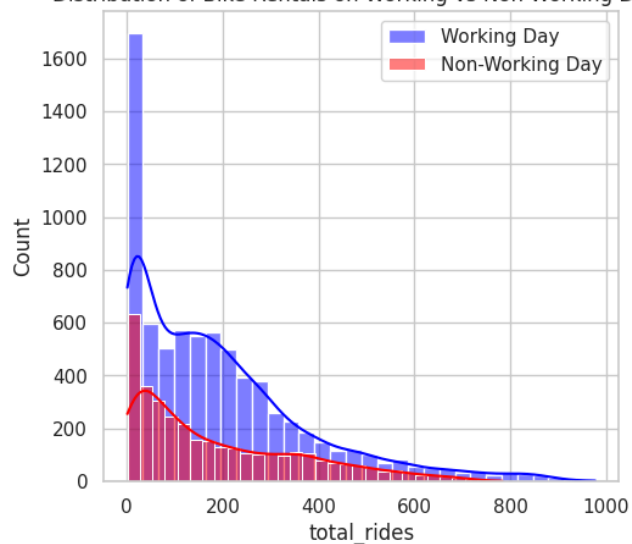
sns.histplot(log_non_working_day, kde=True, color="red",
label="Log-Non-Working Day", bins=30)

plt.legend()

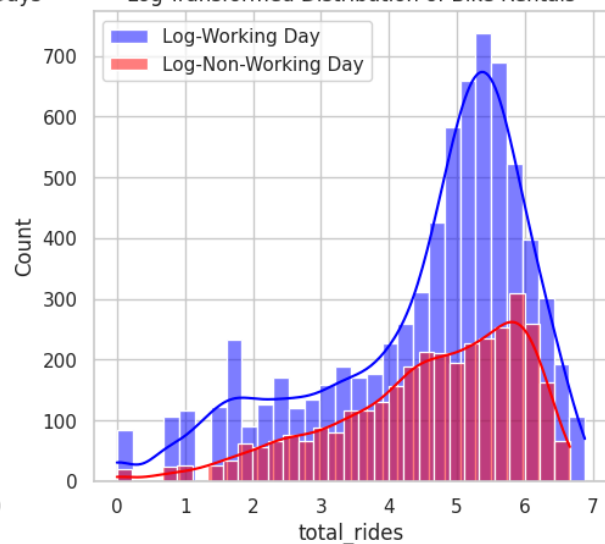
plt.title("Log-Transformed Distribution of Bike Rentals")

plt.show()
```

Distribution of Bike Rentals on Working vs Non-Working Days



Log-Transformed Distribution of Bike Rentals

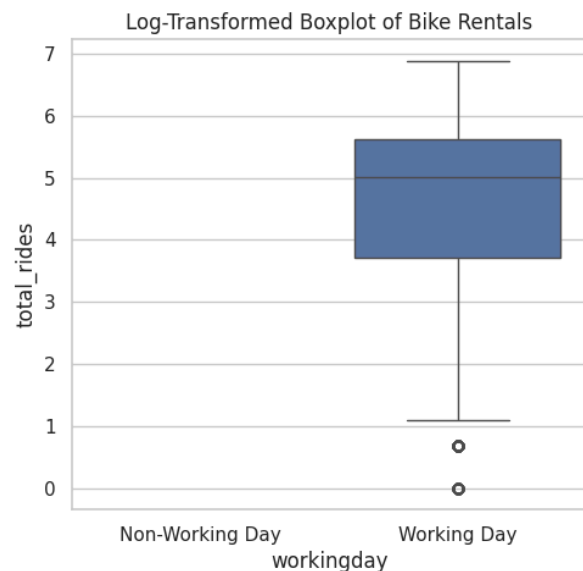
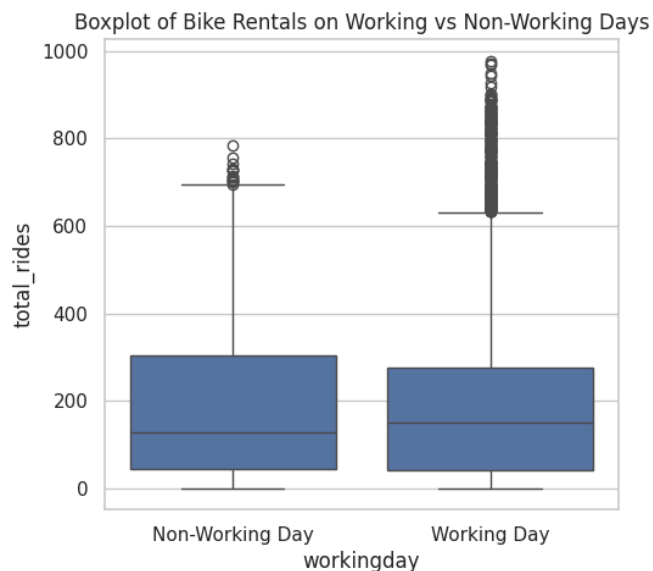




```
# 2. Boxplot to check outliers
plt.figure(figsize=(12, 5))

# Original Data
plt.subplot(1, 2, 1)
sns.boxplot(x=df_raw['workingday'], y=df_raw['total_rides'])
plt.xticks([0, 1], ['Non-Working Day', 'Working Day'])
plt.title("Boxplot of Bike Rentals on Working vs Non-Working Days")

# Log-transformed Data
plt.subplot(1, 2, 2)
sns.boxplot(x=df_raw['workingday'], y=log_working_day)
plt.xticks([0, 1], ['Non-Working Day', 'Working Day'])
plt.title("Log-Transformed Boxplot of Bike Rentals")
plt.show()
```

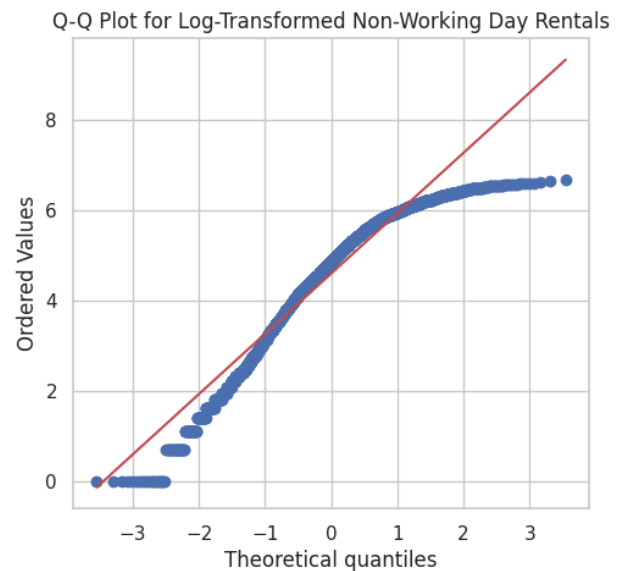
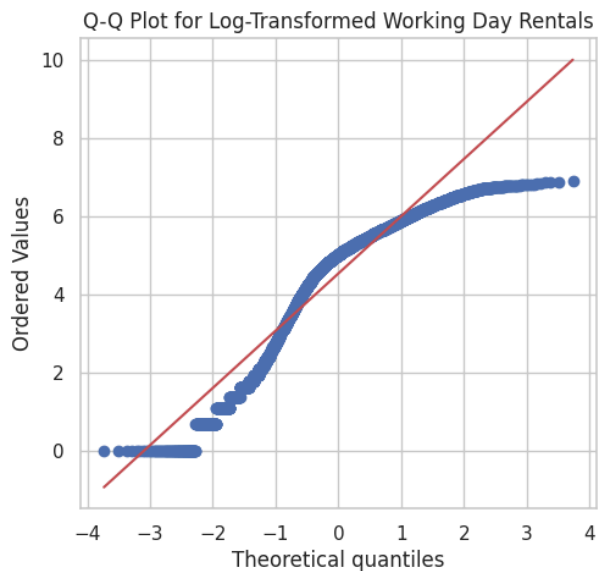
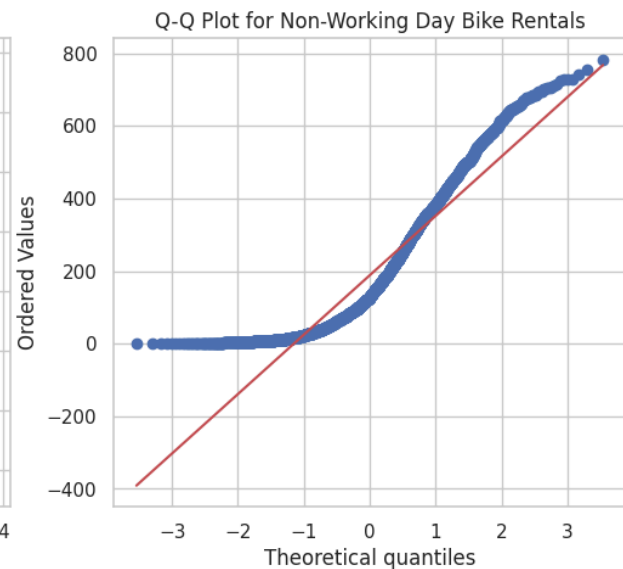
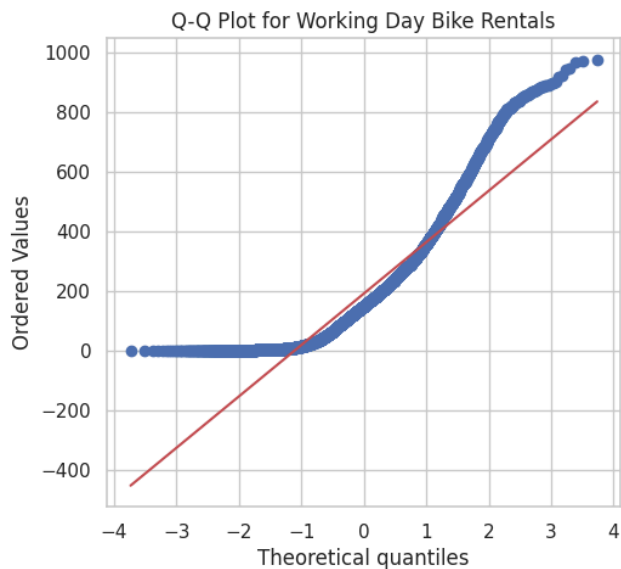


```
import scipy.stats as stats

# 3. Q-Q Plot to check normality
plt.figure(figsize=(12, 5))
# Q-Q Plot for original data
plt.subplot(1, 2, 1)
stats.probplot(working_day, dist="norm", plot=plt)
plt.title("Q-Q Plot for Working Day Bike Rentals")
plt.subplot(1, 2, 2)
stats.probplot(non_working_day, dist="norm", plot=plt)
```



```
plt.title("Q-Q Plot for Non-Working Day Bike Rentals")
plt.show()
# Q-Q Plot for log-transformed data
plt.figure(figsize=(12, 5))
# Q-Q Plot for log-transformed data
plt.subplot(1, 2, 1)
stats.probplot(log_working_day, dist="norm", plot=plt)
plt.title("Q-Q Plot for Log-Transformed Working Day Rentals")
plt.subplot(1, 2, 2)
stats.probplot(log_non_working_day, dist="norm", plot=plt)
plt.title("Q-Q Plot for Log-Transformed Non-Working Day Rentals")
plt.show()
```





```
# 4. Check skewness numerically for both original and log-transformed data
skew_working = stats.skew(working_day)
skew_non_working = stats.skew(non_working_day)

skew_log_working = stats.skew(log_working_day)
skew_log_non_working = stats.skew(log_non_working_day)

print(f"Skewness of Working Day Rentals (Original): {skew_working}")
print(f"Skewness of Non-Working Day Rentals (Original):
{skew_non_working}")
print(f"Skewness of Working Day Rentals (Log-Transformed):
{skew_log_working}")
print(f"Skewness of Non-Working Day Rentals (Log-Transformed):
{skew_log_non_working}")
```



```
Skewness of Working Day Rentals (Original): 1.3464498800708273
Skewness of Non-Working Day Rentals (Original): 0.963621480100446
Skewness of Working Day Rentals (Log-Transformed): -0.997898828621721
Skewness of Non-Working Day Rentals (Log-Transformed): -0.8606344088721859
```

```
# 5. Perform Shapiro-Wilk Test for normality
shapiro_working = stats.shapiro(working_day)
shapiro_non_working = stats.shapiro(non_working_day)

shapiro_log_working = stats.shapiro(log_working_day)
shapiro_log_non_working = stats.shapiro(log_non_working_day)

print(f"Shapiro-Wilk Test for Working Day (Original): p-value =
{shapiro_working.pvalue}")
print(f"Shapiro-Wilk Test for Non-Working Day (Original): p-value =
{shapiro_non_working.pvalue}")
print(f"Shapiro-Wilk Test for Log-Transformed Working Day: p-value =
{shapiro_log_working.pvalue}")
print(f"Shapiro-Wilk Test for Log-Transformed Non-Working Day: p-value =
{shapiro_log_non_working.pvalue}")

# Interpretation
alpha = 0.05
if shapiro_working.pvalue < alpha:
```



```
print("The working day rentals are NOT normally distributed
(Original).")
else:
    print("The working day rentals are normally distributed (Original).")

if shapiro_non_working.pvalue < alpha:
    print("The non-working day rentals are NOT normally distributed
(Original).")
else:
    print("The non-working day rentals are normally distributed
(Original).")

if shapiro_log_working.pvalue < alpha:
    print("The log-transformed working day rentals are NOT normally
distributed.")
else:
    print("The log-transformed working day rentals are normally
distributed.")





if shapiro_log_non_working.pvalue < alpha:
    print("The log-transformed non-working day rentals are NOT normally
distributed.")
else:
    print("The log-transformed non-working day rentals are normally
distributed.")
```

⇒ Shapiro-Wilk Test for Working Day (Original): p-value = 2.2521124830019574e-61  
Shapiro-Wilk Test for Non-Working Day (Original): p-value = 4.4728547627911074e-45  
Shapiro-Wilk Test for Log-Transformed Working Day: p-value = 1.3876984374347798e-55  
Shapiro-Wilk Test for Log-Transformed Non-Working Day: p-value = 9.752105323817043e-37  
The working day rentals are NOT normally distributed (Original).  
The non-working day rentals are NOT normally distributed (Original).  
The log-transformed working day rentals are NOT normally distributed.  
The log-transformed non-working day rentals are NOT normally distributed.







## Key Insights.

### Skewness Analysis:

- **Original Data:**
  - **Working Day Rentals: Skewness = 1.35** (Right-skewed  positively skewed)
  - **Non-Working Day Rentals: Skewness = 0.96** (Right-skewed  moderately positively skewed)
- **Log-Transformed Data:**
  - **Working Day Rentals: Skewness = -0.99** (Left-skewed  slightly negatively skewed)
  - **Non-Working Day Rentals: Skewness = -0.86** (Left-skewed  slightly negatively skewed)

### Shapiro-Wilk Test Results

- **Original Data:**
  - **Working Day Rentals: p-value = 2.25e-61** (Highly significant, **not normally distributed** )
  - **Non-Working Day Rentals: p-value = 4.47e-45** (Highly significant, **not normally distributed** )
- **Log-Transformed Data:**
  - **Working Day Rentals: p-value = 1.39e-55** (Highly significant, **not normally distributed** )
  - **Non-Working Day Rentals: p-value = 9.75e-37** (Highly significant, **not normally distributed** )

### Conclusion

Despite attempting a **log transformation**:

- Both the **original** and **log-transformed** data remain **non-normally distributed**.
- The **skewness** reduction was minimal, and normality was not achieved.

### Next Step :

Since the data does not meet the assumption of normality, I will proceed with the **Mann-Whitney U test** (a non-parametric test) to compare **Working Day Rentals** and **Non-Working Day Rentals**, ensuring the results are valid without assuming a normal distribution.



## 💡 Mann-Whitney U Test:

```
u_stat, p_value = stats.mannwhitneyu(working_day, non_working_day,
alternative="two-sided")
print(f"Mann-Whitney U Test: U-statistic = {u_stat}, p-value = {p_value}")
if p_value < 0.05:
    print("There is a significant difference in bike rentals between
working and non-working days.")
else:
    print("Fail to reject the null hypothesis: There is no significant
difference in bike rentals between working and non-working days.")
```

➡ Mann-Whitney U Test: U-statistic = 12868495.5, p-value = 0.9679139953914079  
Fail to reject the null hypothesis: There is no significant difference in bike rentals between working and non-working days.

## Mann-Whitney U Test Results 📊

- U-statistic: 12868495.5
- p-value: 0.9679

## Interpretation 🔍:

Since the **p-value** is **0.9679** (which is much greater than the significance level of 0.05), we **fail to reject the null hypothesis**.

## Conclusion 🚫:

There is **no significant difference** in bike rentals between **working** and **non-working days**. The data suggests that the rentals do not vary significantly based on whether it is a working day or a non-working day.



## ANNOVA:

```
# Visualize the distribution of total_rides across different weather conditions
```

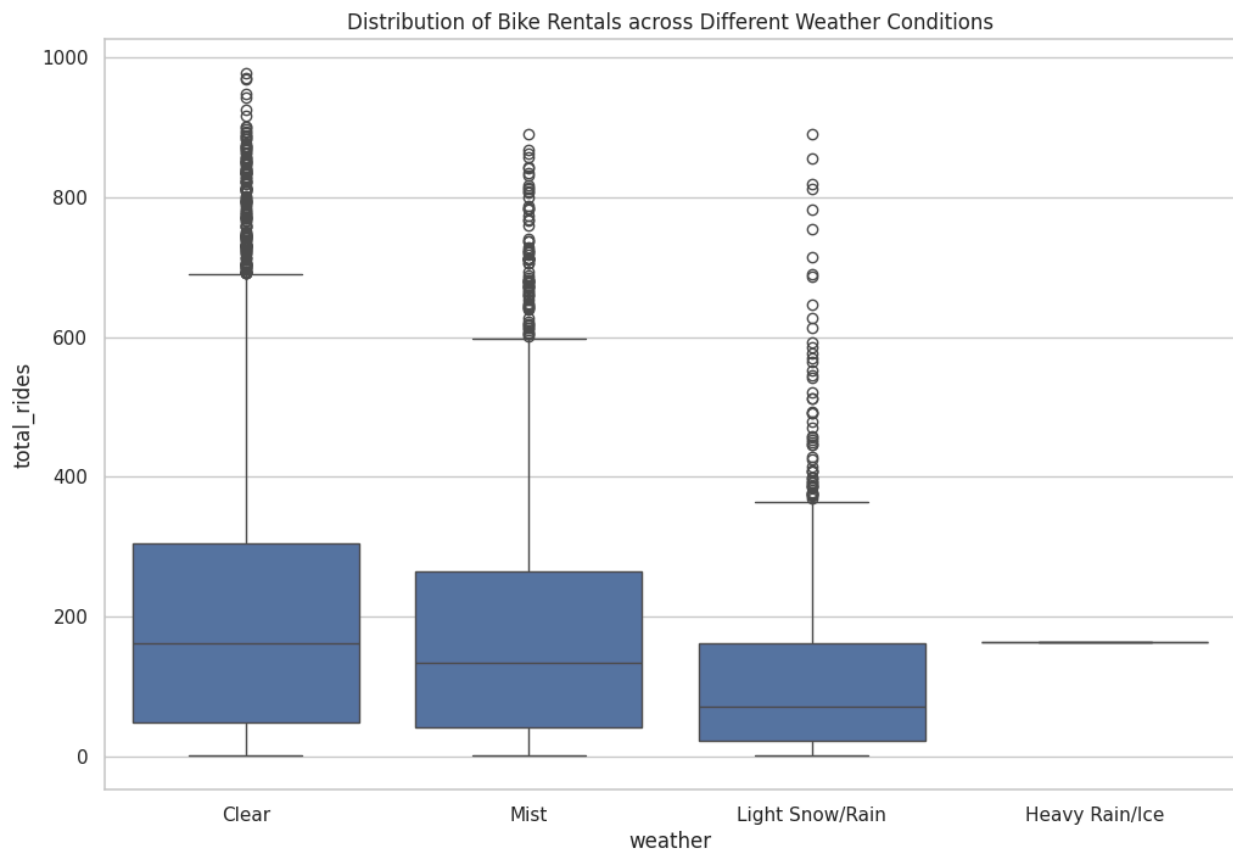
```
plt.figure(figsize=(12, 8))
```

```
sns.boxplot(x='weather', y='total_rides', data=df_raw)
```

```
plt.title('Distribution of Bike Rentals across Different Weather Conditions')
```

```
plt.xticks([0, 1, 2, 3], ['Clear', 'Mist', 'Light Snow/Rain', 'Heavy Rain/Ice'])
```

```
plt.show()
```







```
weather_conditions = df_raw['weather'].unique()

# Creating a grid of subplots
fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(12, 8))

# Flattening axes array to iterate over
axes = axes.flatten()

# Loop over weather conditions and plot Q-Q plot
for i, weather_condition in enumerate(weather_conditions):
    data_weather = df_raw[df_raw['weather'] ==
weather_condition]['total_rides']

    # Q-Q Plot
    stats.probplot(data_weather, dist="norm", plot=axes[i])
    axes[i].set_title(f"Q-Q Plot for Weather {weather_condition}")
    axes[i].grid(True)

plt.tight_layout()
plt.show()

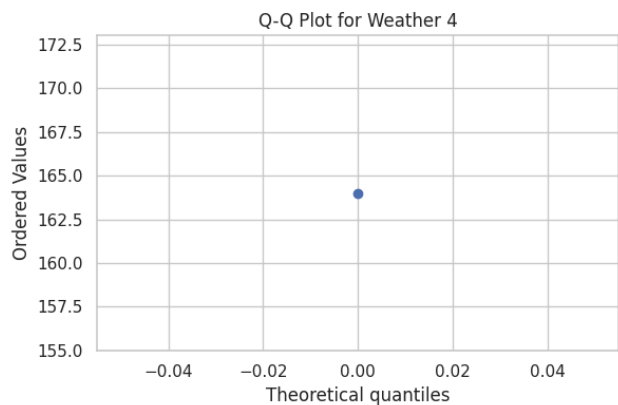
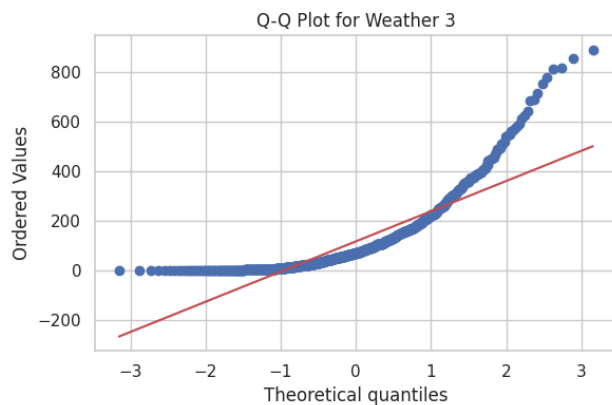
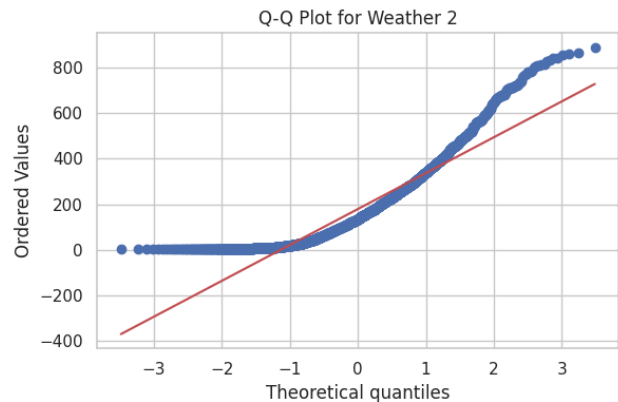
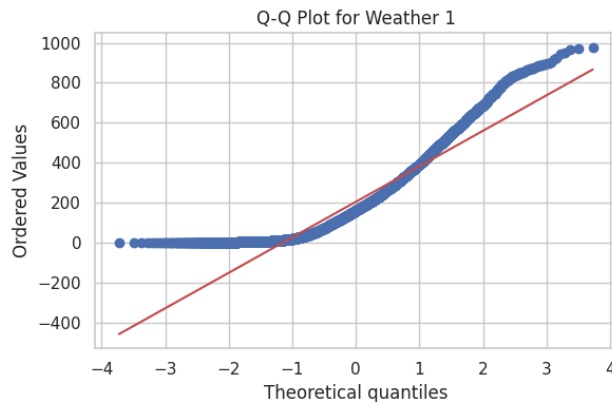
# Levene's Test for Equal Variance across weather conditions
levene_stat, levene_p = stats.levene(
    df_raw[df_raw['weather'] == 1]['total_rides'], # Clear
    df_raw[df_raw['weather'] == 2]['total_rides'], # Mist
    df_raw[df_raw['weather'] == 3]['total_rides'], # Light Snow/Rain
    df_raw[df_raw['weather'] == 4]['total_rides']   # Heavy Rain/Ice
)
print(f"Levene's Test p-value: {levene_p}")

alpha = 0.05
if levene_p < alpha:
    print("The variances across weather conditions are significantly
different.")
else:
    print("The variances across weather conditions are not significantly
different.")

# One-way ANOVA to check if total_rides differ across weather conditions
```



```
anova_stat, anova_p = stats.f_oneway(  
    df_raw[df_raw['weather'] == 1]['total_rides'], # Clear  
    df_raw[df_raw['weather'] == 2]['total_rides'], # Mist  
    df_raw[df_raw['weather'] == 3]['total_rides'], # Light Snow/Rain  
    df_raw[df_raw['weather'] == 4]['total_rides']  # Heavy Rain/Ice  
)  
print(f"ANOVA Test p-value: {anova_p}")  
if anova_p < alpha:  
    print("There is a significant difference in the number of cycles  
rented across different weather conditions.")  
else:  
    print("There is no significant difference in the number of cycles  
rented across different weather conditions.")
```



Levene's Test p-value: 3.504937946833238e-35

The variances across weather conditions are significantly different.

ANOVA Test p-value: 5.482069475935669e-42

There is a significant difference in the number of cycles rented across different weather conditions.



## Key Insights.

### Levene's Test for Homogeneity of Variance

- **Levene's Test p-value: 3.5049e-35**
- This p-value is **extremely small**, indicating that the assumption of **equal variances** across the weather conditions is **violated**. The variances across the weather categories are significantly different, meaning the assumption of homogeneity of variance for ANOVA is **not met**.

### ANOVA Test Results

- **ANOVA Test p-value: 5.4821e-42**
- This **p-value is very small**, suggesting a **significant difference** in the number of **cycles rented** across the different weather conditions. While this indicates that weather does impact rental behavior, the violation of variance homogeneity makes these results less reliable.

## Conclusion

Since the assumption of **equal variances** is not met, the results from ANOVA might not be valid. Therefore, I have decided to proceed with the **Kruskal-Wallis test**, which is a **non-parametric test** that does not require the assumption of equal variances. This test will allow me to assess whether there is a significant difference in bike rentals across different weather conditions without relying on ANOVA's assumptions.



## Kruskal-Wallis Test:

```
# Kruskal-Wallis Test to check if total_rides differ across weather
conditions

kruskal_stat, kruskal_p = stats.kruskal(
    df_raw[df_raw['weather'] == 1]['total_rides'], # Clear
    df_raw[df_raw['weather'] == 2]['total_rides'], # Mist
    df_raw[df_raw['weather'] == 3]['total_rides'], # Light Snow/Rain
    df_raw[df_raw['weather'] == 4]['total_rides']  # Heavy Rain/Ice
)

# Display Kruskal-Wallis result
print(f"Kruskal-Wallis Test p-value: {kruskal_p}")

# Interpretation based on Kruskal-Wallis Test p-value
if kruskal_p < alpha:
    print("There is a significant difference in the number of cycles
    rented across different weather conditions.")
else:
    print("There is no significant difference in the number of cycles
    rented across different weather conditions.")
```

### Key Insights.

#### Kruskal-Wallis Test Results

- **Kruskal-Wallis Test p-value: 3.5016e-44**
- The p-value is **extremely small**, indicating a **significant difference** in the number of **cycles rented** across different **weather conditions**.

#### Conclusion

Since the assumptions of ANOVA were not met (due to unequal variances across weather conditions), I proceeded with the **Kruskal-Wallis test**, which is a non-parametric test suitable for comparing more than two groups when the data doesn't meet ANOVA assumptions. The test results confirm that there is a **statistically significant difference** in bike rentals between different weather conditions.

This outcome suggests that the **weather** does indeed have a significant impact on the number of **cycles rented**, which could be valuable for understanding rental patterns under various weather conditions.



## Chi-Square Test:

```
# Create a contingency table for 'weather' and 'season'
contingency_table = pd.crosstab(df_raw['weather_type'],
df_raw['season_type'], margins=True)

print("\nContingency Table")
print(contingency_table.to_string(index=True, header=True,
justify='center'))

# Check for missing values in the contingency table
if contingency_table.isnull().values.any():
    print("\nWarning: The contingency table contains missing
values, which might affect the Chi-Square test.")
else:
    print("\nNo missing values in the contingency table.")

# Perform the Chi-Square Test of Independence
# Exclude the 'All' row/column (margins) for the test
chi2_stat, p_value, dof, expected =
stats.chi2_contingency(contingency_table.iloc[:-1, :-1])

# Output the results with rounded values for better readability
print(f"\nChi-Square Test Results:")
print(f"{'Chi-Square Statistic':<25}: {chi2_stat:.4f}")
print(f"{'Degrees of Freedom':<25}: {dof}")
print(f"{'P-Value':<25}: {p_value:.4f}")

print(f"\n{'Expected Frequencies':<20}")
expected_df = pd.DataFrame(expected,
columns=contingency_table.columns[:-1],
index=contingency_table.index[:-1])
print(expected_df.to_string(index=True, header=True,
justify='center'))

# Interpretation
```



```
alpha = 0.05
if p_value < alpha:
    print("\nConclusion: There is a significant relationship
between 'Weather' and 'Season' (Reject H0).")
else:
    print("\nConclusion: 'Weather' and 'Season' are independent
(Fail to reject H0).")

# Check if the expected frequency is too low (less than 5) in any
cell, which can affect the validity of the test
if (expected < 5).any():
    print("\nWarning: Some expected frequencies are less than 5,
which may affect the reliability of the Chi-Square test.")
```



Contingency Table

season_type	Fall	Spring	Summer	Winter	All
weather_type					
Clear	1930	1759	1801	1702	7192
Heavy Rain	0	1	0	0	1
Light Snow	199	211	224	225	859
Mist	604	715	708	807	2834
All	2733	2686	2733	2734	10886

No missing values in the contingency table.

Chi-Square Test Results:

Chi-Square Statistic : 49.1587  
Degrees of Freedom : 9  
P-Value : 0.0000

Expected Frequencies

season_type	Fall	Spring	Summer	Winter
weather_type				
Clear	1805.597648	1774.546390	1805.597648	1806.258313
Heavy Rain	0.251056	0.246739	0.251056	0.251148
Light Snow	215.657450	211.948742	215.657450	215.736359
Mist	711.493845	699.258130	711.493845	711.754180

Conclusion: There is a significant relationship between 'Weather' and 'Season' (Reject H0).

Warning: Some expected frequencies are less than 5, which may affect the reliability of the Chi-Square test.



## Key Insights.

### Chi-Square Test Results 🧑🔬:

- **Chi-Square Statistic: 49.1587**
- **Degrees of Freedom: 9**
- **p-value: 0.0000**

The **p-value** is **extremely small**, leading us to **reject the null hypothesis (H0)**, which states that there is **no relationship** between 'Weather' and 'Season'. This result indicates that there is a **significant relationship** between weather conditions and the season of bike rentals.

### Conclusion ✅:

- The **Chi-Square test results** reveal a **significant relationship** between **Weather** and **Season**.
- We **reject the null hypothesis (H0)**, which means there is enough evidence to conclude that the **distribution of weather conditions varies across seasons**.
- However, be mindful that **low expected frequencies** for some cells may limit the test's reliability. A **Fisher's Exact Test** might be more appropriate in such cases if further refinement is needed.

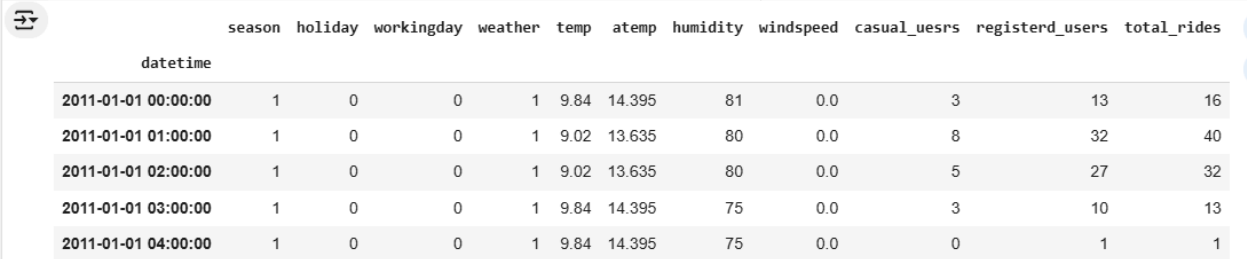


## Time Series Analysis:

```
df_raw.set_index('datetime', inplace=True) # Set 'datetime' as the
index
```

```
df_raw.drop(['day_of_week', 'month', 'year', 'day', 'hour',
'season_type', 'weather_type'], axis=1, inplace=True) # Drop
unnecessary columns
```

```
df_raw.head(5)
```



datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual_uesrs	registerd_users	total_rides
2011-01-01 00:00:00	1	0	0	1	9.84	14.395	81	0.0	3	13	16
2011-01-01 01:00:00	1	0	0	1	9.02	13.635	80	0.0	8	32	40
2011-01-01 02:00:00	1	0	0	1	9.02	13.635	80	0.0	5	27	32
2011-01-01 03:00:00	1	0	0	1	9.84	14.395	75	0.0	3	10	13
2011-01-01 04:00:00	1	0	0	1	9.84	14.395	75	0.0	0	1	1

```
# Plot the total rides to visualize the time series
```

```
plt.figure(figsize=(55, 10))
```

```
plt.plot(df_raw.index, df_raw['total_rides'], marker='o',
linestyle='-')
```

```
plt.title('Total Rides Over Time')
```

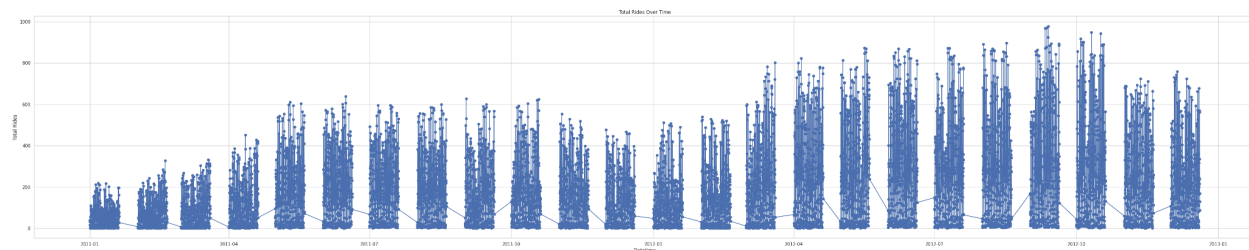
```
plt.xlabel('Datetime')
```

```
plt.ylabel('Total Rides')
```

```
plt.grid(True)
```

```
plt.show()
```





```
df_daily = df_raw['total_rides'].resample('D').sum()
```

```
# Plot daily data
```

```
plt.figure(figsize=(35, 5))
```

```
plt.plot(df_daily, label='Daily Total Rides')
```

```
plt.title('Daily Total Rides')
```

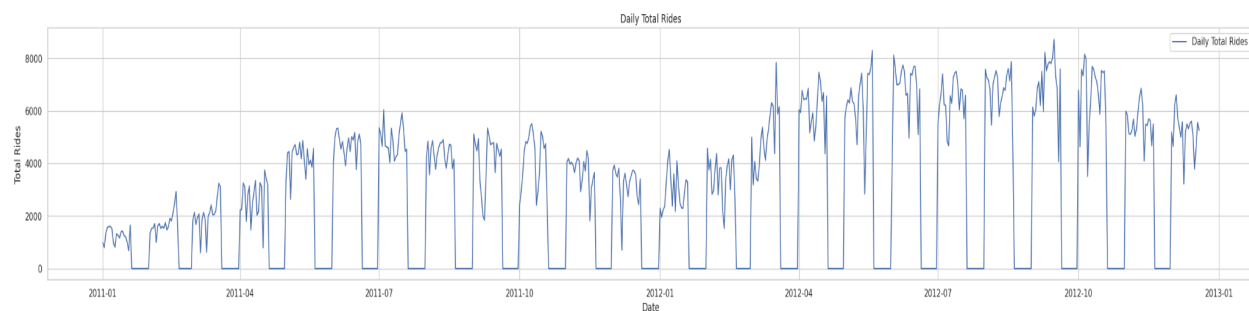
```
plt.xlabel('Date')
```

```
plt.ylabel('Total Rides')
```

```
plt.grid(True)
```

```
plt.legend()
```

```
plt.show()
```





```
# Aggregate to monthly data
```

```
df_monthly = df_raw['total_rides'].resample('M').sum()
```

```
# Plot monthly data
```

```
plt.figure(figsize=(20, 6))
```

```
plt.plot(df_monthly, marker='o', label='Monthly Total Rides')
```

```
plt.title('Monthly Total Rides')
```

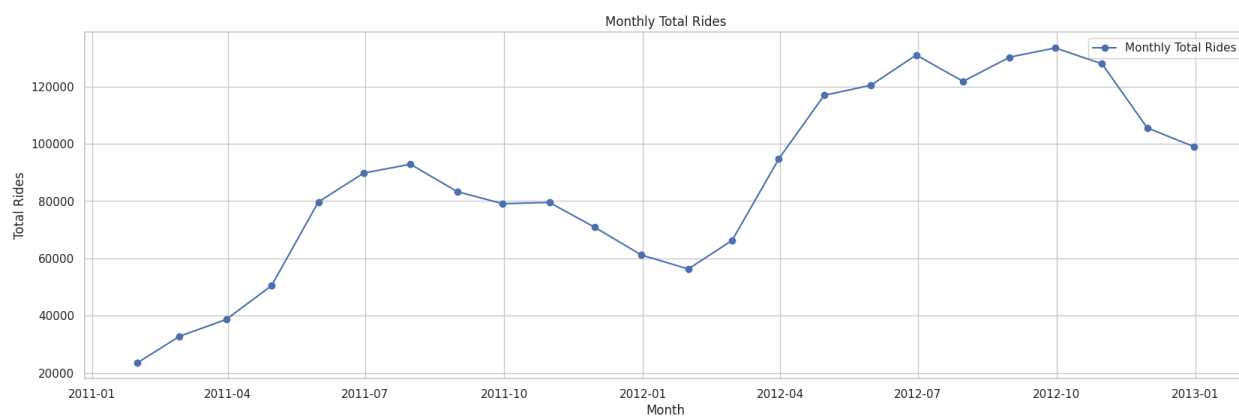
```
plt.xlabel('Month')
```

```
plt.ylabel('Total Rides')
```

```
plt.grid(True)
```

```
plt.legend()
```

```
plt.show()
```





```
from statsmodels.tsa.seasonal import seasonal_decompose

import matplotlib.pyplot as plt

# Resample to daily data and fill missing values

df_daily =
df_raw['total_rides'].resample('D').sum().interpolate(method='time')

# Decompose the data (weekly seasonality)

decompose_result = seasonal_decompose(df_daily, model='additive',
period=7)

# Custom plot with larger size

plt.figure(figsize=(18, 12))

# Plot each component separately


plt.subplot(4, 1, 1)

plt.plot(df_daily, label='Original')

plt.legend(loc='upper left')

plt.title('Original Time Series')


plt.subplot(4, 1, 2)

plt.plot(decompose_result.trend, label='Trend')

plt.legend(loc='upper left')

plt.title('Trend Component')


plt.subplot(4, 1, 3)

plt.plot(decompose_result.seasonal, label='Seasonal')
```



```
plt.legend(loc='upper left')

plt.title('Seasonal Component')

plt.subplot(4, 1, 4)

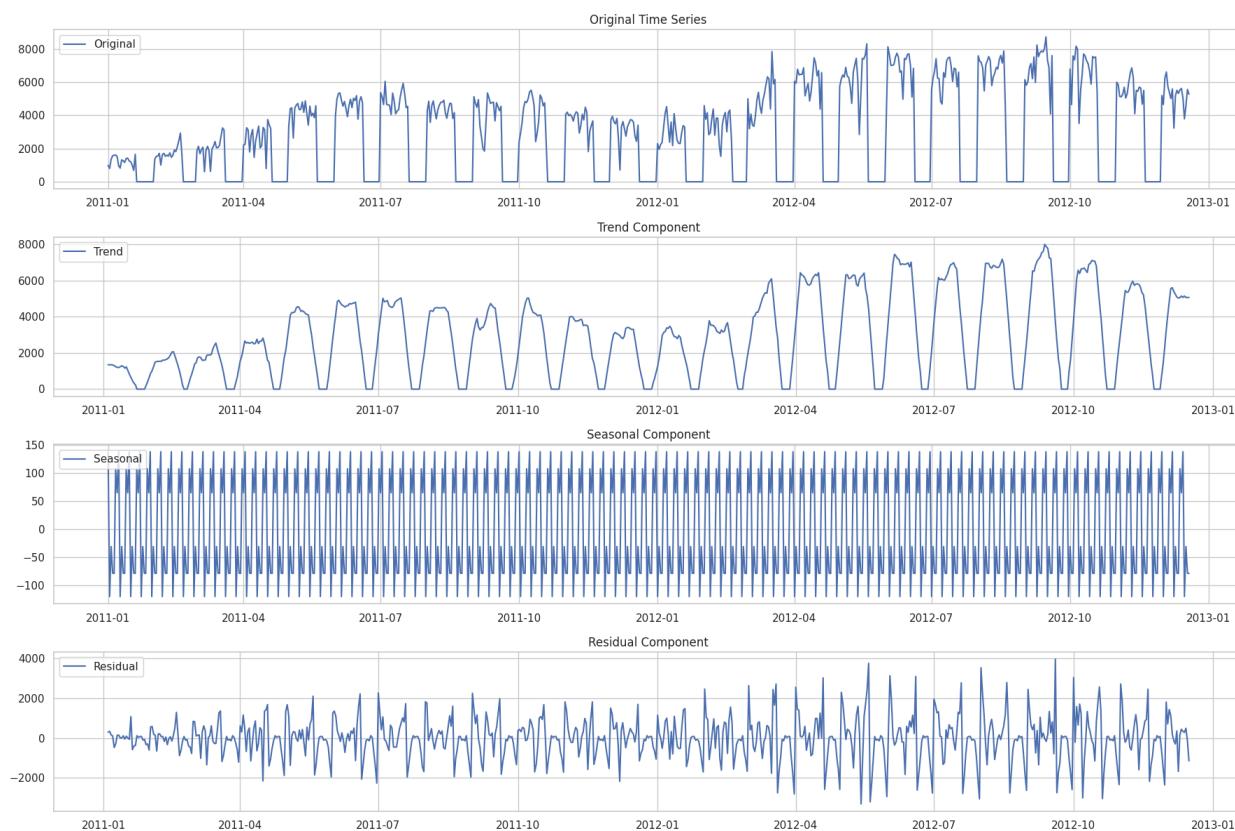
plt.plot(decompose_result.resid, label='Residual')

plt.legend(loc='upper left')

plt.title('Residual Component')

plt.tight_layout()

plt.show()
```





### Seasonality Check:

```
from statsmodels.tsa.stattools import adfuller

# Perform the ADF test on the original data

adf_test = adfuller(df_raw['total_rides'])

print("ADF Statistic:", adf_test[0])

print("p-value:", adf_test[1])

# If p-value > 0.05, the data is non-stationary. Apply
differencing:

df_raw['total_rides_diff'] = df_raw['total_rides'].diff()

# Re-run the ADF test on the differenced data

adf_test_diff = adfuller(df_raw['total_rides_diff'].dropna())

print("ADF Statistic (Differenced):", adf_test_diff[0])

print("p-value (Differenced):", adf_test_diff[1])
```

```
➡ ADF Statistic: -6.419975656501509
  p-value: 1.80161952866896e-08
  ADF Statistic (Differenced): -23.300534023635578
  p-value (Differenced): 0.0
```



```
from statsmodels.tsa.statespace.sarimax import SARIMAX

# Fit SARIMA model on the original data

sarima_model = SARIMAX(df_raw['total_rides'], order=(1, 0, 1),
seasonal_order=(1, 0, 1, 24))

sarima_fit = sarima_model.fit(dispatch=False)

# Print model summary

print(sarima_fit.summary())
```



#### SARIMAX Results

```
=====
Dep. Variable:          total_rides    No. Observations:          10886
Model:                SARIMAX(1, 0, 1)x(1, 0, 1, 24)    Log Likelihood          -60722.849
Date:                  Wed, 05 Feb 2025    AIC                  121455.698
Time:                  04:15:24    BIC                  121492.175
Sample:                0    HQIC                  121467.993
                    - 10886
Covariance Type:                opg
=====
              coef    std err          z      P>|z|      [0.025    0.975]
-----
ar.L1          0.5648      0.007     85.421     0.000      0.552      0.578
ma.L1          0.4601      0.009     50.642     0.000      0.442      0.478
ar.S.L24       0.9903      0.001    1239.809     0.000      0.989      0.992
ma.S.L24      -0.8071      0.003    -237.699     0.000     -0.814     -0.800
sigma2       4081.4802     32.992     123.712     0.000    4016.818    4146.143
=====
Ljung-Box (L1) (Q):                7.11    Jarque-Bera (JB):                9672.85
Prob(Q):                          0.01    Prob(JB):                      0.00
Heteroskedasticity (H):            2.68    Skew:                          -0.13
Prob(H) (two-sided):              0.00    Kurtosis:                      7.61
=====
```

#### Warnings:

```
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
```



```
residuals = sarima_fit.resid

plt.figure(figsize=(14, 7))

plt.subplot(2, 1, 1)

plt.plot(residuals, label='Residuals')

plt.title('Residuals of SARIMA Model')

plt.legend()

plt.subplot(2, 1, 2)

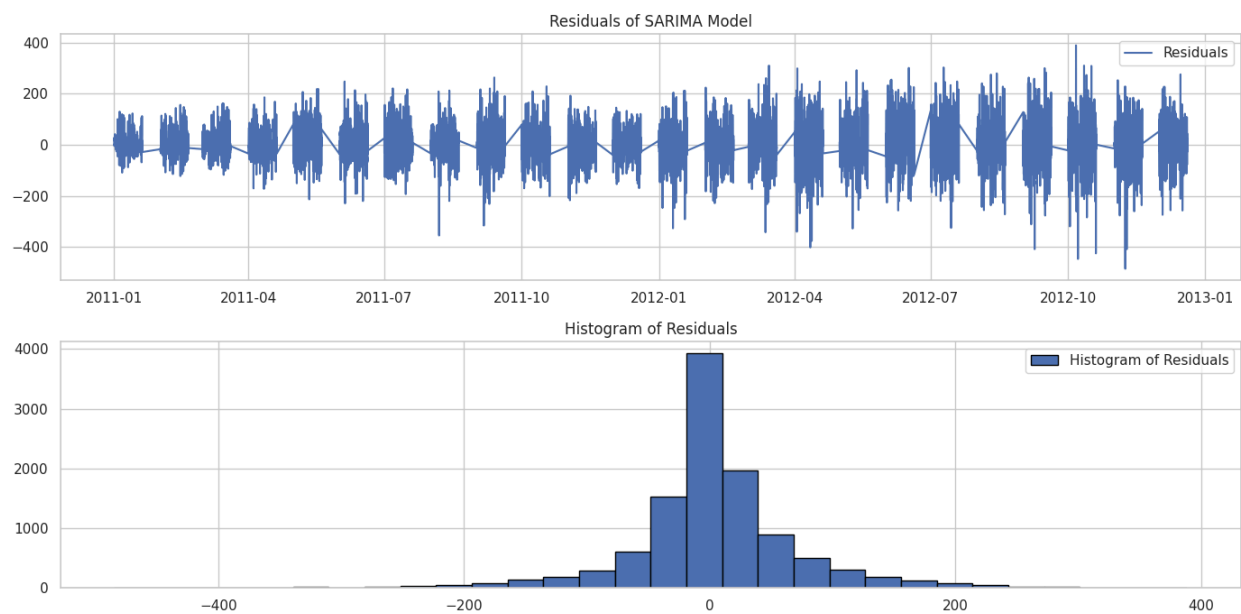
plt.hist(residuals, bins=30, edgecolor='black', label='Histogram of
Residuals')

plt.title('Histogram of Residuals')

plt.legend()

plt.tight_layout()

plt.show()
```





## Model Evaluation:

```
from sklearn.metrics import mean_squared_error
import numpy as np

# Calculate the RMSE on the fitted values
rmse = np.sqrt(mean_squared_error(df_raw['total_rides'],
sarima_fit.fittedvalues))
print(f"RMSE: {rmse:.2f}")
```

⇒ RMSE: 63.90

```
recent_data = df_raw[df_raw.index >= (df_raw.index[-1] -
pd.Timedelta(days=7))]

forecast = sarima_fit.forecast(steps=24)

plt.figure(figsize=(25, 5))

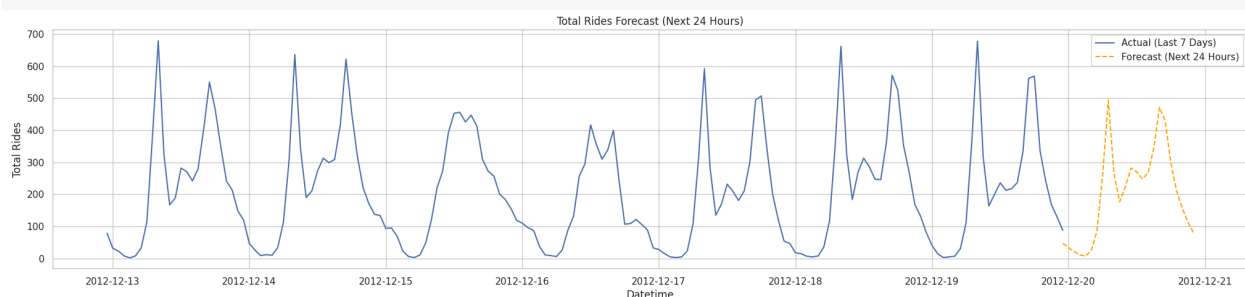
plt.plot(recent_data.index, recent_data['total_rides'], label='Actual
(Last 7 Days)')

plt.plot(pd.date_range(recent_data.index[-1], periods=24, freq='H'),
forecast, linestyle='--', color='orange', label='Forecast (Next 24
Hours)')

plt.title('Total Rides Forecast (Next 24 Hours)')

plt.xlabel('Datetime')

plt.ylabel('Total Rides')
```







## Overall Insights & Key Takeaways

### Rider Trends & Time-Based Patterns

- **Hourly Fluctuations:**
    - Low rider count during early morning hours (12 AM – 6 AM).
    - **Peak demand** from **4 PM to 7 PM**, driven by evening commuters.
    - Gradual decline in rentals at night.
  - **Monthly & Yearly Growth:**
    - Significant **year-over-year growth** from **2011 to 2012**.
    - **Seasonal effect** observed—rider count peaks between **April and October**, with **June and July** being the busiest months.
  - **Weekly Trends:**
    - Highest rental activity on **Thursday, Friday, and Saturday**.
    - **Sunday** experiences the lowest rentals, reflecting a relaxed weekend pattern.
- 

### Weather & Seasonal Impact on Demand

- **Chi-Square Test Insight:** Strong relationship between **weather conditions** and **season**.
  - **Clear weather** sees the highest number of rides, while **heavy rain and snow** reduce demand.
  - **Temperature:**
    - Most rides occur in **moderate temperatures** (~14°C).
    - Extreme cold or heat leads to decreased activity.
- 

### User Type Breakdown

- **Registered Users (81.2%)** dominate, indicating that most riders use the service regularly.
  - **Casual Users (18.8%)** are likely influenced by special events, good weather, or weekends.
- 

### Collinearity Check & Feature Relationships

- Correlation analysis shows that **temperature, weather conditions, and seasonal patterns** are highly interrelated, which may affect predictive model performance.
- 

### Time Series Analysis Insights: Total Rides Prediction



## Key Findings

- **Stationarity Test (ADF):** The series is stationary (p-value: 1.8e-08), meaning it's stable and ready for direct modeling.
- **SARIMA Model Summary:**
  - **Model:** SARIMA(1, 0, 1) × (1, 0, 1, 24), capturing the **24-hour seasonality** and short-term patterns.
  - **AIC:** 121455.70
  - **RMSE:** 63.90—indicating a reasonable level of prediction accuracy.
  - **Ljung-Box Test (Q-Stat):** Residuals are sufficiently random.



## Forecast for the Next 7 Days

- **Daily total rides** follow a predictable pattern with **24-hour seasonality**.
- **Weekday vs. Weekend:** Higher ride demand on weekdays; slight dips on weekends.
- **Weather Impact:** Clear days result in higher predicted rides, while adverse weather reduces activity.



## Impact of External Factors



- **Weather:** Clear weather encourages more rides, while bad weather (rain, snow) discourages them.
- **Working Days vs. Holidays:**
  - **Working Days:** Peak hours during commuting times—**8–10 AM** and **5–7 PM**.
  - **Holidays:** Irregular patterns with lower overall demand.



## Model Performance & Business Impact

- **Operational Efficiency:** Use forecast data to optimize resource allocation, especially during peak hours.
- **Weather-Driven Strategy:** Surge pricing and targeted promotions on clear days can boost demand.
- **Future Enhancements:** Incorporate real-time weather data and holiday schedules for adaptive forecasting and long-term planning.

## Key Takeaways & Recommendations

-  **Peak Hours:** Focus on fleet optimization during **4 PM - 7 PM** and **April to October**.
-  **Registered Users:** Since they account for **81.2% of demand**, consider loyalty programs or



subscription perks to retain and grow this base.

📌 **Weather-Based Promotions:** Boost demand with discounts on clear days and incentivize rides during low-demand periods.

📌 **Sunday Promotions:** Increase activity with targeted offers, as Sunday sees the lowest rentals.

📌 **Higher Fleet Allocation:** Especially in **June and July**, to meet peak seasonal demand.

---

## 🚀 Overall Strategy

The data highlights **seasonality, time-based demand patterns, weather sensitivity**, and **predictable daily peaks**. To maximize operational efficiency and customer satisfaction, focus on:

- **Dynamic pricing,**
- **Weather-based strategies,**
- **Targeted marketing,** and
- **Accurate demand forecasting.**

These efforts will help balance supply and demand, improve customer experience, and drive business growth.