<p style="color:orange; text-align:center">Business Case</p>

# Yulu - Hypothesis Testing

### Suman Debnath

## Introduction

Yulu is India's leading micro-mobility service provider, which offers unique vehicles for the daily commute. Starting off as a mission to eliminate traffic congestion in India, Yulu provides the safest commute solution through a user-friendly mobile app to enable shared, solo and sustainable commuting.

Yulu zones are located at all the appropriate locations (including metro stations, bus stands, office spaces, residential areas, corporate offices, etc) to make those first and last miles smooth, affordable, and convenient!

Yulu has recently suffered considerable dips in its revenues. They have contracted a consulting company to understand the factors on which the demand for these shared electric cycles depends. Specifically, they want to understand the factors affecting the demand for these shared electric cycles in the Indian market.

**Business Problem**

The company wants to know:

- Which variables are significant in predicting the demand for shared electric cycles in the Indian market?
- How well those variables describe the electric cycle demands

**Dataset**

Dataset link: yulu_data.csv

The dataset have the following fields:

`datetime` : datetime

`season` : season (1: spring, 2: summer, 3: fall, 4: winter)

`holiday` : whether day is a holiday or not (extracted from http://dchr.dc.gov/page/holiday-schedule)

`workingday` : if day is neither weekend nor holiday is 1, otherwise is 0.

`weather` :

```
- Clear, Few clouds, partly cloudy, partly cloudy
- Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist
- Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds
- Heavy Rain + Ice Pallets + Thunderstorm + Mist, Snow + Fog
```

`temp` : temperature in Celsius

`atemp` : feeling temperature in Celsius

`humidity` : humidity

`windspeed` : wind speed

`casual` : count of casual users

`registered` : count of registered users

`count` : count of total rental bikes including both casual and registered

## Summary

- The data is given from Timestamp('2011-01-01 00:00:00') to Timestamp('2012-12-19 23:00:00'). The total time period for which the data is given is '718 days 23:00:00'.
- Out of every 100 users, around 19 are casual users and 81 are registered users.
- The mean total hourly count of rental bikes is 144 for the year 2011 and 239 for the year 2012. An annual growth rate of 65.41 % can be seen in the demand of electric vehicles on an hourly basis.
- There is a seasonal pattern in the count of rental bikes, with higher demand during the spring and summer months, a slight decline in the fall, and a further decrease in the winter months.
    - The average hourly count of rental bikes is the lowest in the month of January followed by February and March.
- There is a distinct fluctuation in count throughout the day, with low counts during early morning hours, a sudden increase in the morning, a peak count in the afternoon, and a gradual decline in the evening and nighttime.
- More than 80 % of the time, the temperature is less than 28 degrees celcius.
- More than 80 % of the time, the humidity value is greater than 40. Thus for most of the time, humidity level varies from optimum to too moist.
- More than 85 % of the total, windspeed data has a value of less than 20.
- The hourly count of total rental bikes is the highest in the clear and cloudy weather, followed by the misty weather and rainy weather. There are very few records for extreme weather conditions.
- The mean hourly count of the total rental bikes is statistically similar for both working and non- working days.
- There is statistically significant dependency of weather and season based on the hourly total number of bikes rented.
- The hourly total number of rental bikes is statistically different for different weathers.
- There is no statistically significant dependency of weather 1, 2, 3 on season based on the average hourly total number of bikes rented.
- The hourly total number of rental bikes is statistically different for different seasons.

## Recommendation

- **Seasonal Marketing**: Since there is a clear seasonal pattern in the count of rental bikes, Yulu can adjust its marketing strategies accordingly. Focus on promoting bike rentals during the spring and summer months when there is higher demand. Offer seasonal discounts or special packages to attract more customers during these periods.

- **Time-based Pricing**: Take advantage of the hourly fluctuation in bike rental counts throughout the day. Consider implementing time-based pricing where rental rates are lower during off-peak hours and higher during peak hours. This can encourage customers to rent bikes during less busy times, balancing out the demand and optimizing the resources.

- **Weather-based Promotions**: Recognize the impact of weather on bike rentals. Create weather-based promotions that target customers during clear and cloudy weather, as these conditions show the highest rental counts. Yulu can offer weather-specific discounts to attract more customers during these favorable weather conditions.

- **User Segmentation**: Given that around 81% of users are registered, and the remaining 19% are casual, Yulu can tailor its marketing and communication strategies accordingly. Provide loyalty programs, exclusive offers, or personalized recommendations for registered users to encourage repeat business. For casual users, focus on providing a seamless rental experience and promoting the benefits of bike rentals for occasional use.

- **Optimize Inventory**: Analyze the demand patterns during different months and adjust the inventory accordingly. During months with lower rental counts such as January, February, and March, Yulu can optimize its inventory levels to avoid excess bikes. On the other hand, during peak months, ensure having sufficient bikes available to meet the higher demand.

- **Improve Weather Data Collection**: Given the lack of records for extreme weather conditions, consider improving the data collection process for such scenarios. Having more data on extreme weather conditions can help to understand customer behavior and adjust the operations accordingly, such as offering specialized bike models for different weather conditions or implementing safety measures during extreme weather.

- **Customer Comfort**: Since humidity levels are generally high and temperature is often below 28 degrees Celsius, consider providing amenities like umbrellas, rain jackets, or water bottles to enhance the comfort and convenience of the customers. These small touches can contribute to a positive customer experience and encourage repeat business.

- **Collaborations with Weather Services**: Consider collaborating with weather services to provide real-time weather updates and forecasts to potential customers. Incorporate weather information into your marketing campaigns or rental app to showcase the ideal biking conditions and attract users who prefer certain weather conditions.

- **Seasonal Bike Maintenance**: Allocate resources for seasonal bike maintenance. Before the peak seasons, conduct thorough maintenance checks on the bike fleet to ensure they are in top condition. Regularly inspect and service bikes throughout the year to prevent breakdowns and maximize customer satisfaction.

- **Customer Feedback and Reviews**: Encourage customers to provide feedback and reviews on their biking experience. Collecting feedback can help identify areas for improvement, understand customer preferences, and tailor the services to better meet customer expectations.

- **Social Media Marketing**: Leverage social media platforms to promote the electric bike rental services. Share captivating visuals of biking experiences in different weather conditions, highlight customer testimonials, and engage with potential customers through interactive posts and contests. Utilize targeted advertising campaigns to reach specific customer segments and drive more bookings.

- **Special Occasion Discounts**: Since Yulu focusses on providing a sustainable solution for vehicular pollution, it should give special discounts on the occassions like Zero Emissions Day (21st September), Earth day (22nd April), World Environment Day (5th June) etc in order to attract new users.

## Detailed Analysis

### Importing all the `libs`

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import datetime as dt
import scipy.stats as stats
```

### Loading the `data`

```python
# data_set = 'https://d2beiqkhq929f0.cloudfront.net/public_assets/assets/000/001/428/original/bike_sharing.csv'
data_set = 'bike_sharing.csv'
```

### Exploratory Data Exploration (EDA)

In [3]:
```python
df = pd.read_csv(data_set)
```

In [4]:
```python
df.shape
```
Out[4]: `(10886, 12)`

In [5]:
```python
df.head()
```
Out[5]:

|   | datetime | season | holiday | workingday | weather | temp | atemp | humidity | windspeed | casual | registered | count |
|---|----------|--------|---------|------------|---------|------|-------|----------|-----------|--------|------------|-------|
| 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 |

In [6]:
```python
df.dtypes
```
Out[6]:
```
datetime       object
season          int64
holiday         int64
workingday      int64
weather         int64
temp          float64
atemp         float64
humidity        int64
windspeed     float64
casual          int64
registered      int64
count           int64
dtype: object
```

In [7]:
```python
df.columns
```
Out[7]:
```
Index(['datetime', 'season', 'holiday', 'workingday', 'weather', 'temp',
       'atemp', 'humidity', 'windspeed', 'casual', 'registered', 'count'],
      dtype='object')
```

### Check for `null` values

In [8]:
```python
np.any(df.isna())
```
Out[8]: `False`

### Check for `duplicate` values

In [9]:
```python
np.any(df.duplicated())
```
Out[9]: `False`

```
In [10]: df.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  count       10886 non-null  int64
dtypes: float64(3), int64(8), object(1)
memory usage: 1020.7+ KB
```

### Converting the datatype of datetime column from object to datetime

```
In [11]: df['datetime'] = pd.to_datetime(df['datetime'])
```

```
In [12]: df.dtypes
```

```
Out[12]: datetime      datetime64[ns]
season                int64
holiday               int64
workingday            int64
weather               int64
temp                float64
atemp               float64
humidity              int64
windspeed           float64
casual                int64
registered            int64
count                 int64
dtype: object
```

```
In [13]: df['datetime'].min(), df['datetime'].max()
```

```
Out[13]: (Timestamp('2011-01-01 00:00:00'), Timestamp('2012-12-19 23:00:00'))
```

### Setting the `datetime` column as the index of the DataFrame `df`

```
In [14]: df.set_index('datetime', inplace = True)
```
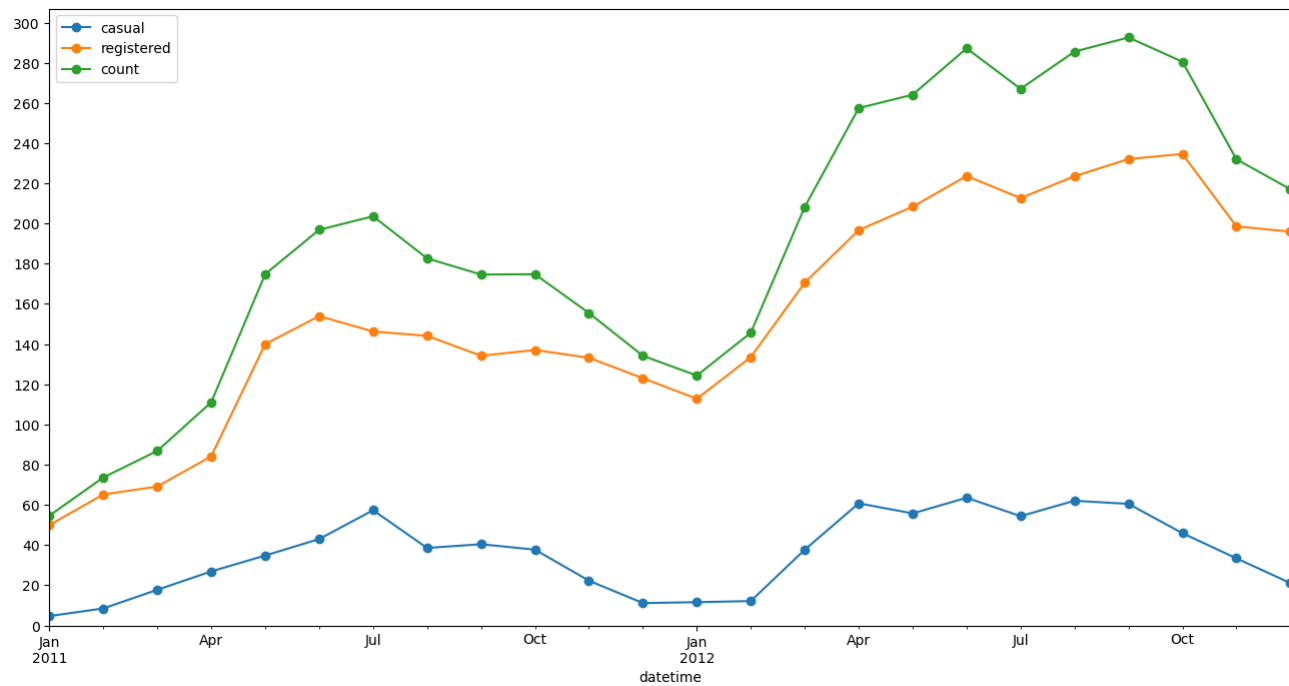
```
In [15]: df.head()
```

Out[15]:

| datetime | season | holiday | workingday | weather | temp | atemp | humidity | windspeed | casual | registered | count |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 |

## Slicing Data by Time

```
In [16]: plt.figure(figsize = (16, 8))

# plotting a lineplot by resampling the data on a monthly basis, and calculating the mean value
    # of 'casual', 'registered' and 'count' users for each month
df.resample('M')['casual'].mean().plot(kind = 'line', legend = 'casual', marker = 'o')
df.resample('M')['registered'].mean().plot(kind = 'line', legend = 'registered', marker = 'o')
df.resample('M')['count'].mean().plot(kind = 'line', legend = 'count', marker = 'o')

plt.yticks(np.arange(0, 301, 20))
plt.ylim(0,)
plt.show()
```

Check if there is an increase in the average hourly count of rental bikes from the year 2011 to 2012 ?

```
In [17]: # resampling the DataFrame by the year
         df1 = df.resample('Y')['count'].mean().to_frame().reset_index()

         # Create a new column 'prev_count' by shifting the 'count' column one position up
         # to compare the previous year's count with the current year's count
         df1['prev_count'] = df1['count'].shift(1)

         # Calculating the growth percentage of 'count' with respect to the 'count' of previous year
         df1['growth_percent'] = np.round((df1['count'] - df1['prev_count']) * 100 / df1['prev_count'], 2)
         df1
```

Out[17]:

| | datetime | count | prev_count | growth_percent |
|---|---|---|---|---|
| 0 | 2011-12-31 | 144.223349 | NaN | NaN |
| 1 | 2012-12-31 | 238.560944 | 144.223349 | 65.41 |

# Observation

- This data suggests that there was substantial growth in the count of the variable over the course of one year.
- The mean total hourly count of rental bikes is 144 for the year 2011 and 239 for the year 2012. An annual growth rate of 65.41 % can be seen in the demand of electric vehicles on an hourly basis.

*It indicates positive growth and potentially a successful outcome or increasing demand for the variable being measured.*

How does the average hourly count of rental bikes varies for different month ?

```
In [18]: df.head()
```

Out[18]:

| | season | holiday | workingday | weather | temp | atemp | humidity | windspeed | casual | registered | count |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **datetime** | | | | | | | | | | | |
| **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 |

```
In [19]: df.reset_index(inplace=True)
```

```
In [20]: # Grouping the DataFrame by the month
         df1 = df.groupby(by = df['datetime'].dt.month)['count'].mean().reset_index()
         df1.rename(columns = {'datetime' : 'month'}, inplace = True)

         # Create a new column 'prev_count' by shifting the 'count' column one position up
             # to compare the previous month's count with the current month's count
         df1['prev_count'] = np.round(df1['count'].shift(1), 2)

         # Calculating the growth percentage of 'count' with respect to the 'count' of previous month
         df1['growth_percent'] = np.round((df1['count'] - df1['prev_count']) * 100 / df1['prev_count'], 2)
         df1.set_index('month', inplace = True)
         df1
```

Out[20]:

| month | count | prev_count | growth_percent |
|---|---|---|---|
| 1 | 90.366516 | NaN | NaN |
| 2 | 110.003330 | 90.37 | 21.73 |
| 3 | 148.169811 | 110.00 | 34.70 |
| 4 | 184.160616 | 148.17 | 24.29 |
| 5 | 219.459430 | 184.16 | 19.17 |
| 6 | 242.031798 | 219.46 | 10.29 |
| 7 | 235.325658 | 242.03 | -2.77 |
| 8 | 234.118421 | 235.33 | -0.51 |
| 9 | 233.805281 | 234.12 | -0.13 |
| 10 | 227.699232 | 233.81 | -2.61 |
| 11 | 193.677278 | 227.70 | -14.94 |
| 12 | 175.614035 | 193.68 | -9.33 |

# Observation

- The count of rental bikes shows an increasing trend from January to March, with a significant growth rate of 34.70% between February and March.
- The growth rate starts to stabilize from April to June, with a relatively smaller growth rate.
- From July to September, there is a slight decrease in the count of rental bikes, with negative growth rates.
- The count further declines from October to December, with the largest drop observed between October and November (~14.94%).

In [21]:
```python
# Setting the figure size for the plot
plt.figure(figsize = (15, 6))

# Setting the title for the plot
plt.title("The average hourly distribution of count of rental bikes across different months")

# Grouping the DataFrame by the month and calculating the mean of the 'count' column for each month.
# Ploting the line graph using markers ('o') to represent the average count per month.
df.groupby(by = df['datetime'].dt.month)['count'].mean().plot(kind = 'line', color = 'green', marker = 'o')

plt.ylim(0,)
plt.xticks(np.arange(1, 13))
plt.legend('count')
plt.yticks(np.arange(0, 400, 50))
plt.grid(axis = 'both', linestyle = '--')
plt.ylabel('count')
plt.xlabel('month')
plt.plot()      # Displaing the plot.
```

Out[21]: []



# Observation

- The average hourly count of rental bikes is the highest in the month of June followed by July and August.
- The average hourly count of rental bikes is the lowest in the month of January followed by February and March.
- Overall, these trends suggest a seasonal pattern in the count of rental bikes, with higher demand during the spring and summer months, a slight decline in the fall, and a further decrease in the winter months. It could be useful for the rental bike company to consider these patterns for resource allocation, marketing strategies, and operational planning throughout the year.

### What is the distribution of average count of rental bikes on an hourly basis in a single day ?

In [22]:
```python
# Grouping the DataFrame by the hour
df1 = df.groupby(by = df['datetime'].dt.hour)['count'].mean().reset_index()
df1.rename(columns = {'datetime' : 'hour'}, inplace = True)

# Create a new column 'prev_count' by shifting the 'count' column one position up
    # to compare the previous hour's count with the current hour's count
df1['prev_count'] = df1['count'].shift(1)

# Calculating the growth percentage of 'count' with respect to the 'count' of previous hour
df1['growth_percent'] = (df1['count'] - df1['prev_count']) * 100 / df1['prev_count']
df1.set_index('hour', inplace = True)
df1
```
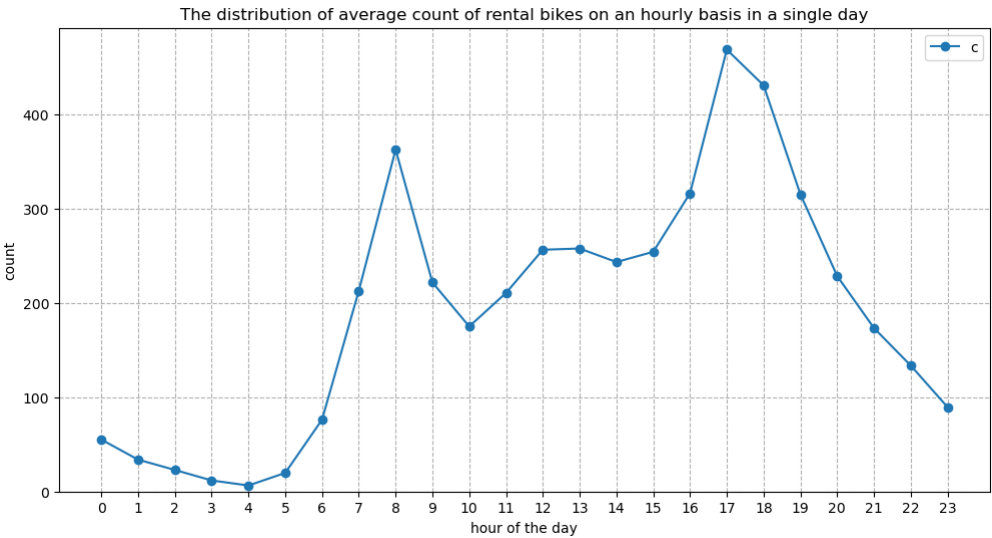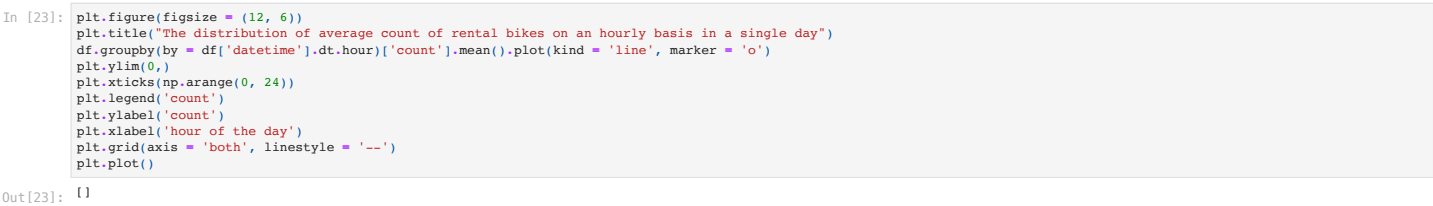
Out[22]:

| hour | count | prev_count | growth_percent |
|---|---|---|---|
| 0 | 55.138462 | NaN | NaN |
| 1 | 33.859031 | 55.138462 | -38.592718 |
| 2 | 22.899554 | 33.859031 | -32.367959 |
| 3 | 11.757506 | 22.899554 | -48.656179 |
| 4 | 6.407240 | 11.757506 | -45.505110 |
| 5 | 19.767699 | 6.407240 | 208.521293 |
| 6 | 76.259341 | 19.767699 | 285.777526 |
| 7 | 213.116484 | 76.259341 | 179.462793 |
| 8 | 362.769231 | 213.116484 | 70.221104 |
| 9 | 221.780220 | 362.769231 | -38.864655 |
| 10 | 175.092308 | 221.780220 | -21.051432 |
| 11 | 210.674725 | 175.092308 | 20.322091 |
| 12 | 256.508772 | 210.674725 | 21.755835 |
| 13 | 257.787281 | 256.508772 | 0.498427 |
| 14 | 243.442982 | 257.787281 | -5.564393 |
| 15 | 254.298246 | 243.442982 | 4.459058 |
| 16 | 316.372807 | 254.298246 | 24.410141 |
| 17 | 468.765351 | 316.372807 | 48.168661 |
| 18 | 430.859649 | 468.765351 | -8.086285 |
| 19 | 315.278509 | 430.859649 | -26.825705 |
| 20 | 228.517544 | 315.278509 | -27.518833 |
| 21 | 173.370614 | 228.517544 | -24.132471 |
| 22 | 133.576754 | 173.370614 | -22.953059 |
| 23 | 89.508772 | 133.576754 | -32.990757 |

## Observation

- During the early morning hours (hours 0 to 5), there is a significant decrease in the count, with negative growth percentages ranging from -38.59% to -48.66%.
- However, starting from hour 5, there is a sudden increase in count, with a sharp positive growth percentage of 208.52% observed from hour 4 to hour 5.
- The count continues to rise significantly until reaching its peak at hour 17, with a growth percentage of 48.17% compared to the previous hour.
- After hour 17, there is a gradual decrease in count, with negative growth percentages ranging from -8.08% to -32.99% during the late evening and nighttime hours.

In [23]:
```python
plt.figure(figsize = (12, 6))
plt.title("The distribution of average count of rental bikes on an hourly basis in a single day")
df.groupby(by = df['datetime'].dt.hour)['count'].mean().plot(kind = 'line', marker = 'o')
plt.ylim(0,)
plt.xticks(np.arange(0, 24))
plt.legend('count')
plt.ylabel('count')
plt.xlabel('hour of the day')
plt.grid(axis = 'both', linestyle = '--')
plt.plot()
```

Out[23]: []



The distribution of average count of rental bikes on an hourly basis in a single day

## Observation

- The average count of rental bikes is the highest at 5 PM followed by 6 PM and 8 AM of the day.
- The average count of rental bikes is the lowest at 4 AM followed by 3 AM and 5 AM of the day.
- These patterns indicate that there is a distinct fluctuation in count throughout the day, with low counts during early morning hours, a sudden increase in the morning, a peak count in the afternoon, and a gradual decline in the evening and nighttime.*

In [24]:
```python
# 1: spring, 2: summer, 3: fall, 4: winter
def season_category(x):
    if x == 1:
        return 'spring'
    elif x == 2:
        return 'summer'
    elif x == 3:
        return 'fall'
```

```
        else:
            return 'winter'
df['season'] = df['season'].apply(season_category)
```

## Optimizing Memory Usage of the Dataframe

```
In [25]:  # Updating dtype of season column

          print('Memory usage of season column : ', df['season'].memory_usage())
          # Since the dtype of season column is object, we can convert the dtype to category to save memory
          df['season'] = df['season'].astype('category')
          print('Updated Memory usage of season column : ', df['season'].memory_usage())

          Memory usage of season column :  87216
          Updated Memory usage of season column :  11218
```

```
In [26]:  # Updating dtype of holiday column

          print('Max value entry in holiday column : ', df['holiday'].max())
          print('Memory usage of holiday column : ', df['holiday'].memory_usage())
          # Since the maximum entry in holiday column is 1 and the dtype is int64, we can convert the dtype to category to save memory
          df['holiday'] = df['holiday'].astype('category')
          print('Updated Memory usage of holiday column : ', df['holiday'].memory_usage())

          Max value entry in holiday column :  1
          Memory usage of holiday column :  87216
          Updated Memory usage of holiday column :  11138
```

## Basic Description of the dataset

```
In [27]:  df.describe()
```

Out[27]:

|  | workingday | weather | temp | atemp | humidity | windspeed | casual | registered | count |
|---|---|---|---|---|---|---|---|---|---|
| count | 10886.000000 | 10886.000000 | 10886.00000 | 10886.000000 | 10886.000000 | 10886.000000 | 10886.000000 | 10886.000000 | 10886.000000 |
| mean | 0.680875 | 1.418427 | 20.23086 | 23.655084 | 61.886460 | 12.799395 | 36.021955 | 155.552177 | 191.574132 |
| std | 0.466159 | 0.633839 | 7.79159 | 8.474601 | 19.245033 | 8.164537 | 49.960477 | 151.039033 | 181.144454 |
| min | 0.000000 | 1.000000 | 0.82000 | 0.760000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 1.000000 |
| 25% | 0.000000 | 1.000000 | 13.94000 | 16.665000 | 47.000000 | 7.001500 | 4.000000 | 36.000000 | 42.000000 |
| 50% | 1.000000 | 1.000000 | 20.50000 | 24.240000 | 62.000000 | 12.998000 | 17.000000 | 118.000000 | 145.000000 |
| 75% | 1.000000 | 2.000000 | 26.24000 | 31.060000 | 77.000000 | 16.997900 | 49.000000 | 222.000000 | 284.000000 |
| max | 1.000000 | 4.000000 | 41.00000 | 45.455000 | 100.000000 | 56.996900 | 367.000000 | 886.000000 | 977.000000 |

```
In [28]:  np.round(df['season'].value_counts(normalize = True) * 100, 2)
```

Out[28]:
```
winter    25.11
fall      25.11
summer    25.11
spring    24.67
Name: season, dtype: float64
```

```
In [29]:  np.round(df['holiday'].value_counts(normalize = True) * 100, 2)
```

Out[29]:
```
0    97.14
1     2.86
Name: holiday, dtype: float64
```

```
In [30]:  np.round(df['workingday'].value_counts(normalize = True) * 100, 2)
```

Out[30]:
```
1    68.09
0    31.91
Name: workingday, dtype: float64
```

```
In [31]:  np.round(df['weather'].value_counts(normalize = True) * 100, 2)
```

Out[31]:
```
1    66.07
2    26.03
3     7.89
4     0.01
Name: weather, dtype: float64
```

### Distribution of Season

```
In [32]:  plt.figure(figsize = (6, 6))        # setting the figure size to 6*6

          # setting the title of the plot
          plt.title('Distribution of season', fontdict = {'fontsize' : 18})

          df_season = np.round(df['season'].value_counts(normalize = True) * 100, 2).to_frame()

          # Creating the color palette
          colors = ['skyblue', 'yellowgreen', 'coral', 'gold']

          # Creating the pie-chart
          plt.pie(x = df_season['season'],
                  explode = [0.025, 0.025, 0.025, 0.025],
                  labels = df_season.index,
                  colors = colors,   # use custom color palette
                  autopct = '%.2f%%',
                  startangle = 140,  # change start angle
                  shadow = True,   # add shadow for 3D effect
                  textprops = {'fontsize' : 14})

          # Add a legend
          plt.legend(df_season.index, loc="upper right")
          plt.show()
```
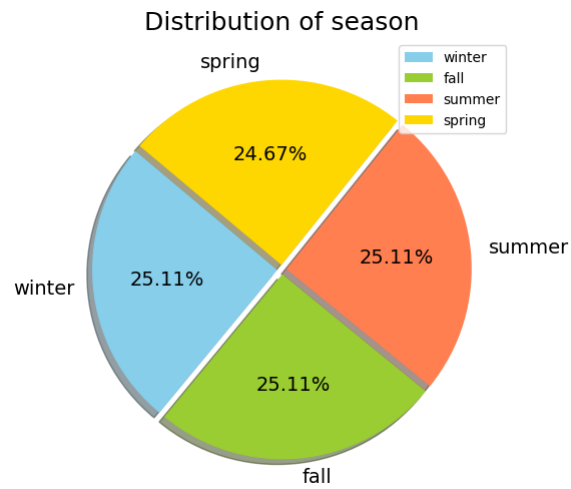
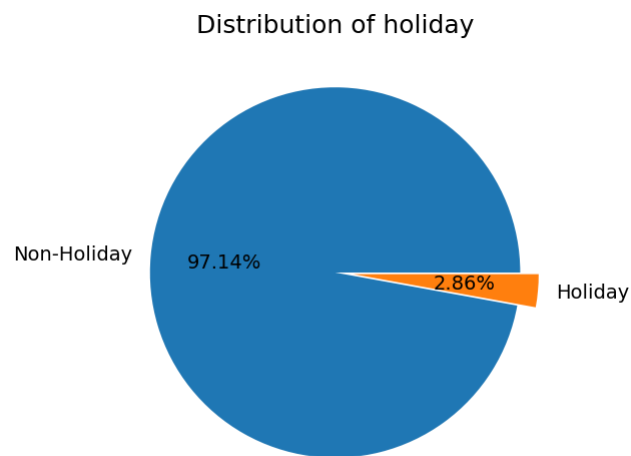## Distribution of season



```
In [33]:  plt.figure(figsize = (6, 6))      # setting the figure size to 6*6

          # setting the title of the plot
          plt.title('Distribution of holiday', fontdict = {'fontsize' : 18})

          df_holiday = np.round(df['holiday'].value_counts(normalize = True) * 100, 2).to_frame()

          # Creating the pie-chart
          plt.pie(x = df_holiday['holiday'],
                  explode = [0, 0.1],
                  labels = ['Non-Holiday', 'Holiday'],
                  autopct = '%.2f%%',
                  textprops = {'fontsize' : 14})

          plt.show()
```
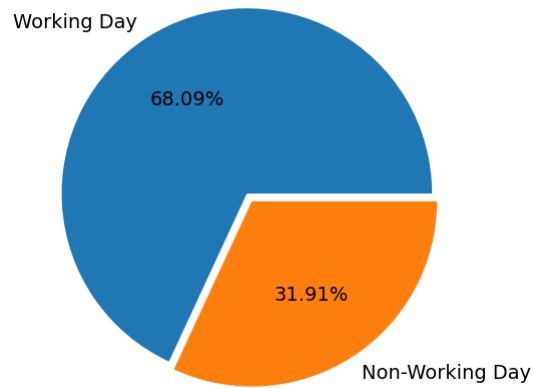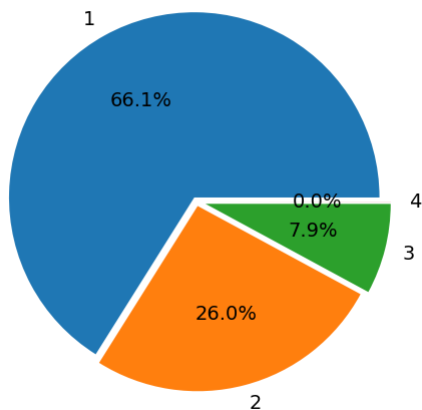
## Distribution of holiday



```
In [34]:  plt.figure(figsize = (6, 6))      # setting the figure size to 6*6

          # setting the title of the plot
          plt.title('Distribution of workingday', fontdict = {'fontsize' : 18})

          df_workingday = np.round(df['workingday'].value_counts(normalize = True) * 100, 2).to_frame()

          # Creating the pie-chart
          plt.pie(x = df_workingday['workingday'],
                  explode = [0, 0.05],
                  labels = ['Working Day', 'Non-Working Day'],
                  autopct = '%.2f%%',
                  textprops = {'fontsize' : 14})

          plt.show()
```

## Distribution of workingday



```
In [35]:  plt.figure(figsize = (6, 6))        # setting the figure size to 6*6

          # setting the title of the plot
          plt.title('Distribution of weather', fontdict = {'fontsize' : 18})

          df_weather = np.round(df['weather'].value_counts(normalize = True) * 100, 2).to_frame()

          # Creating the pie-chart
          plt.pie(x = df_weather['weather'],
                  explode = [0.025, 0.025, 0.05, 0.05],
                  labels = df_weather.index,
                  autopct = '%.1f%%',
                  textprops = {'fontsize' : 14})

          plt.plot()
```

Out[35]:  []

## Distribution of weather



## Univariate Analysis

```
In [36]:  sns.countplot(data = df, x = 'season')
          plt.show()
```



```
In [37]:  sns.countplot(data = df, x = 'holiday')
          plt.show()
```

In [38]: 
```
sns.countplot(data = df, x = 'workingday')
plt.show()
```



In [39]: 
```
sns.countplot(data = df, x = 'weather')
plt.show()
```



In [40]: 
```
sns.histplot(data = df, x = 'temp', kde = True, bins = 40, color='green')
plt.xlabel('temperature')
plt.show()
```
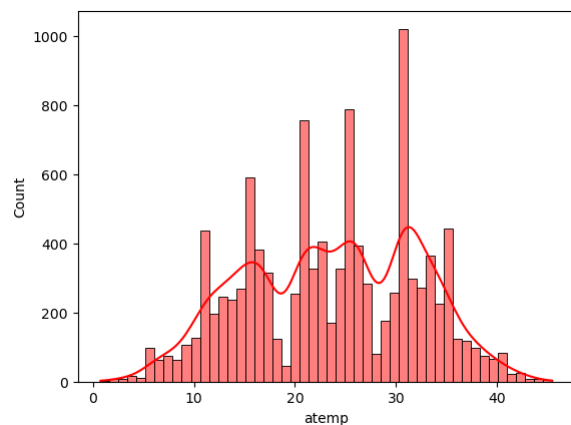
In [41]:
```python
temp_mean = np.round(df['temp'].mean(), 2)
temp_std = np.round(df['temp'].std(), 2)
temp_mean, temp_std
```
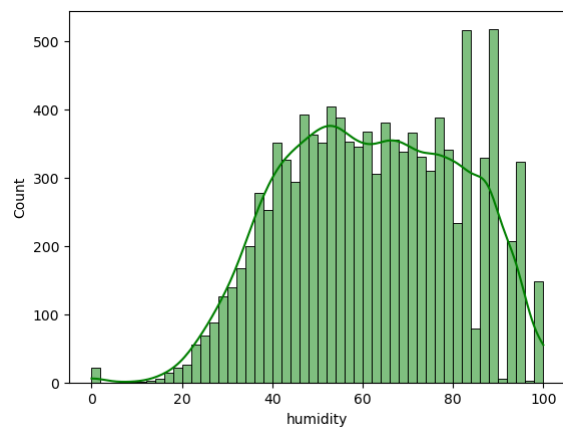
Out[41]: (20.23, 7.79)

In [42]:
```python
sns.histplot(data = df, x = 'temp', kde = True, cumulative = True, stat = 'percent', color='orange')
plt.grid(axis = 'y', linestyle = '--')
plt.yticks(np.arange(0, 101, 10))
plt.show()
```

In [43]:
```python
sns.histplot(data = df, x = 'atemp', kde = True, bins = 50, color='red')
plt.show()
```
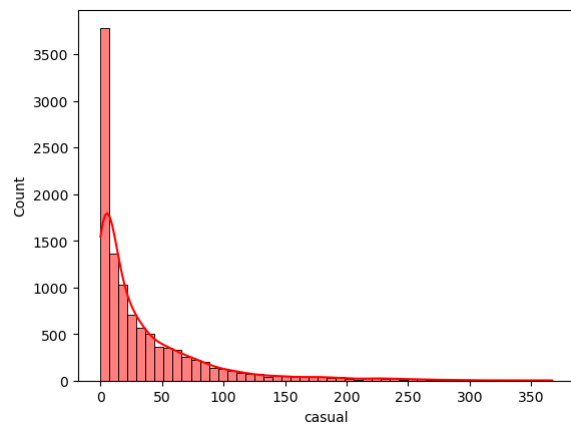
In [44]:
```python
sns.histplot(data = df, x = 'humidity', kde = True, bins = 50, color='green')
plt.show()
```

In [45]:
```python
humidity_mean = np.round(df['humidity'].mean(), 2)
humidity_std = np.round(df['humidity'].std(), 2)
humidity_mean, humidity_std
```
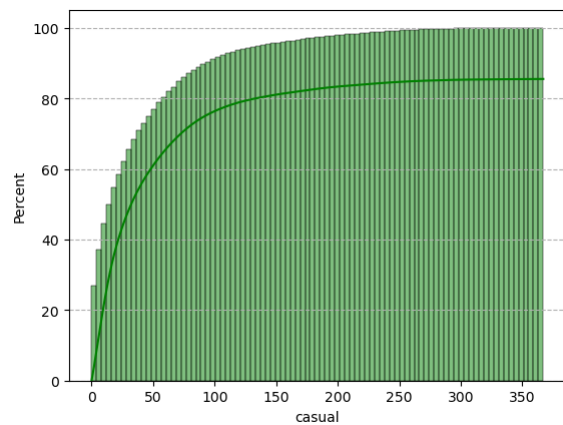
Out[45]: (61.89, 19.25)

In [46]:
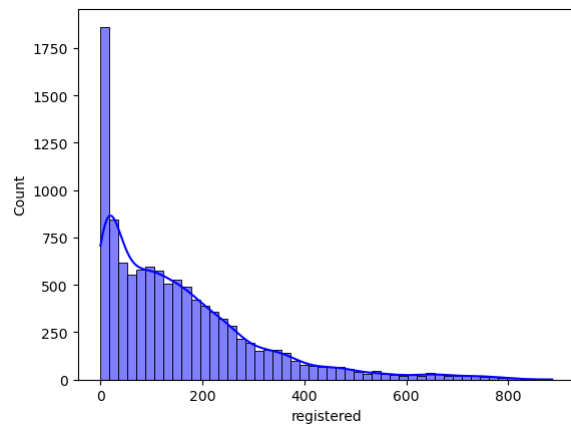```python
sns.histplot(data = df, x = 'casual', kde = True, bins = 50, color='red')
plt.show()
```
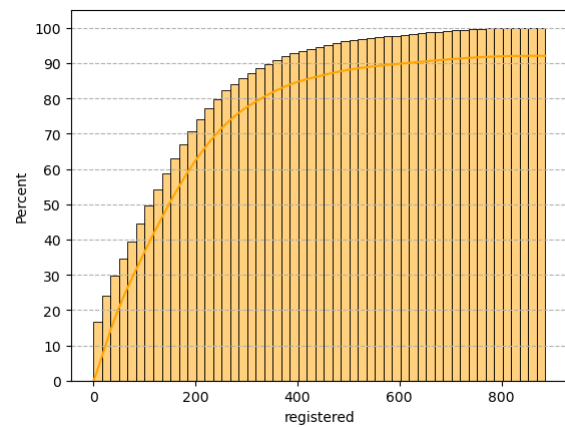
```
In [47]:  sns.histplot(data = df, x = 'casual', kde = True, cumulative = True, stat = 'percent', color='green')
          plt.grid(axis = 'y', linestyle = '--')
          plt.yticks(np.arange(0, 101, 20))
          plt.show()
```



```
In [48]:  sns.histplot(data = df, x = 'registered', kde = True, bins = 50, color='blue')
          plt.show()
```
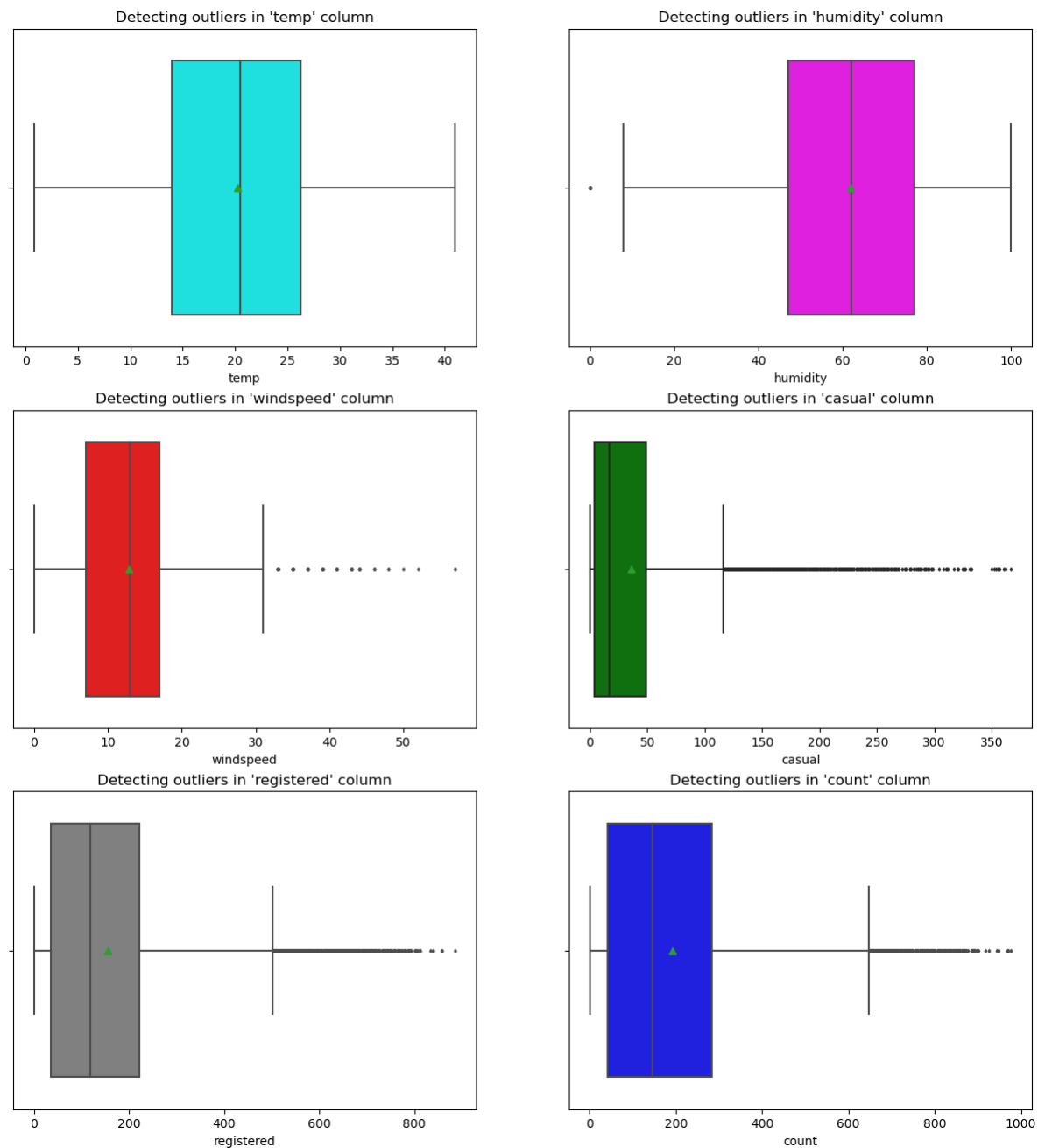


```
In [49]:  sns.histplot(data = df, x = 'registered', kde = True, cumulative = True, stat = 'percent', color='orange')
          plt.grid(axis = 'y', linestyle = '--')
          plt.yticks(np.arange(0, 101, 10))
          plt.show()
```

## Outliers Detection

```
In [50]: columns = ['temp', 'humidity', 'windspeed', 'casual', 'registered', 'count']
         colors = np.random.permutation(['red', 'blue', 'green', 'magenta', 'cyan', 'gray'])
         count = 1
         plt.figure(figsize = (15, 16))
         for i in columns:
             plt.subplot(3, 2, count)
             plt.title(f"Detecting outliers in '{i}' column")
             sns.boxplot(data = df, x = df[i], color = colors[count - 1], showmeans = True, fliersize = 2)
             plt.plot()
             count += 1
```
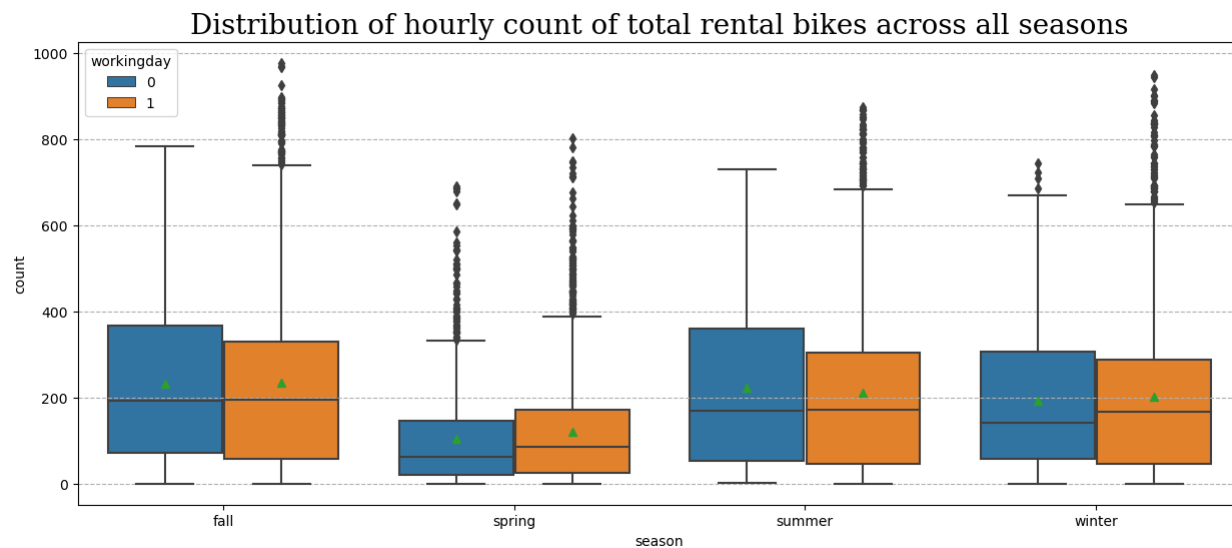
## Observation

- There is no outlier in the temp column.
- There are few outliers present in humidity column.
- There are many outliers present in each of the columns : windspeed, casual, registered, count.
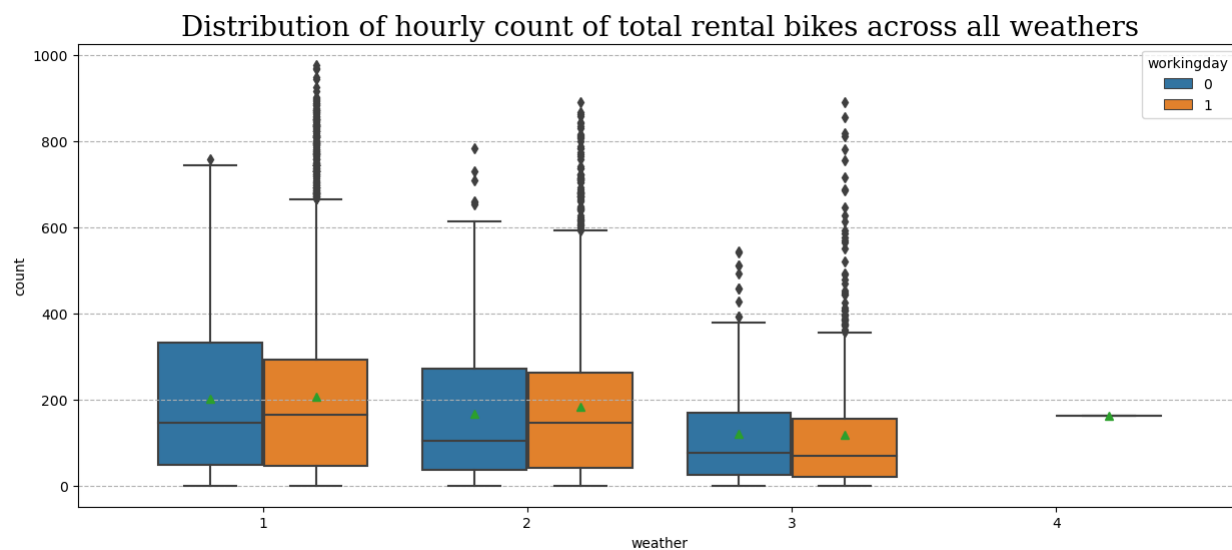
## Bivariate Analysis

```
In [51]:  plt.figure(figsize = (15, 6))
          plt.title('Distribution of hourly count of total rental bikes across all seasons',
                    fontdict = {'size' : 20,
                                'family' : 'serif'})
          sns.boxplot(data = df, x = 'season', y = 'count', hue = 'workingday', showmeans = True)
          plt.grid(axis = 'y', linestyle = '--')
          plt.show()
```

## Distribution of hourly count of total rental bikes across all seasons



```
In [52]:  plt.figure(figsize = (15, 6))
          plt.title('Distribution of hourly count of total rental bikes across all weathers',
                    fontdict = {'size' : 20,
                                'family' : 'serif'})
          sns.boxplot(data = df, x = 'weather', y = 'count', hue = 'workingday', showmeans = True)
          plt.grid(axis = 'y', linestyle = '--')
          plt.plot()

Out[52]:  []
```

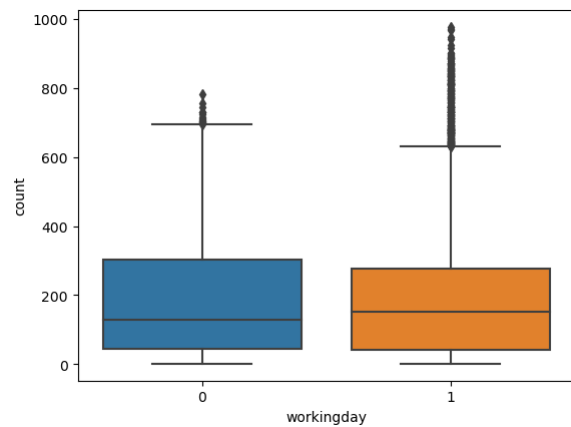## Distribution of hourly count of total rental bikes across all weathers



Is there any effect of Working Day on the number of electric cycles rented ?

```
In [53]:  df.groupby(by = 'workingday')['count'].describe()
```

Out[53]:

| workingday | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| 0 | 3474.0 | 188.506621 | 173.724015 | 1.0 | 44.0 | 128.0 | 304.0 | 783.0 |
| 1 | 7412.0 | 193.011873 | 184.513659 | 1.0 | 41.0 | 151.0 | 277.0 | 977.0 |

```
In [54]:  sns.boxplot(data = df, x = 'workingday', y = 'count')
          plt.show()
```

## Hypothesis Test

*STEP-1* : Set up Null Hypothesis

---

- **Null Hypothesis ( H0 )** - Working Day does not have any effect on the number of electric cycles rented.

- **Alternate Hypothesis ( HA )** - Working Day has some effect on the number of electric cycles rented

*STEP-2* : Checking for basic assumpitons for the hypothesis

---

- Distribution check using **QQ Plot**
- Homogeneity of Variances using **Levene's test**

*STEP-3*: Define Test statistics; Distribution of T under H0.

---

- If the assumptions of T Test are met then we can proceed performing T Test for independent samples else we will perform the non parametric test equivalent to T Test for independent sample i.e., Mann-Whitney U rank test for two independent samples.

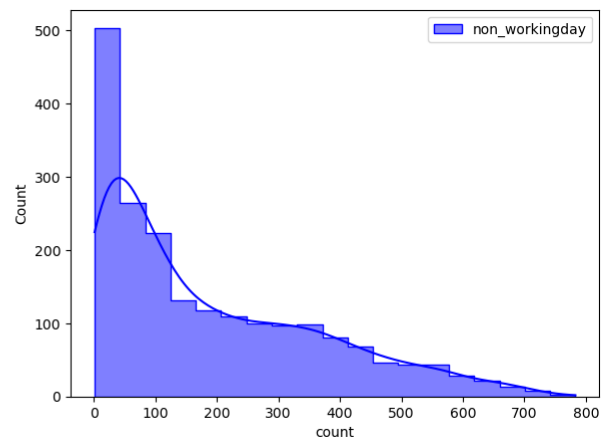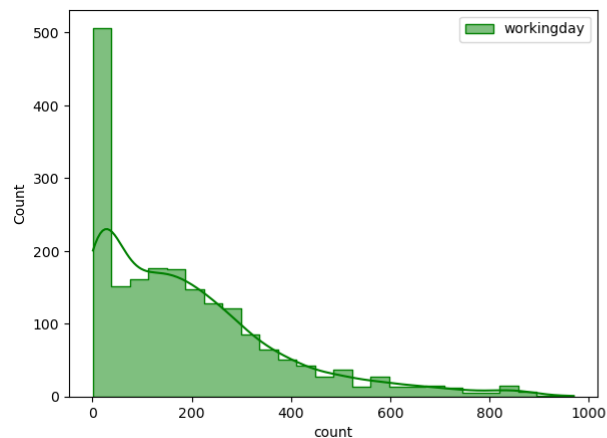*STEP-4*: Compute the p-value and fix value of alpha.

---

- We set our *alpha to be 0.05*

*STEP-5*: Compare p-value and alpha.

---

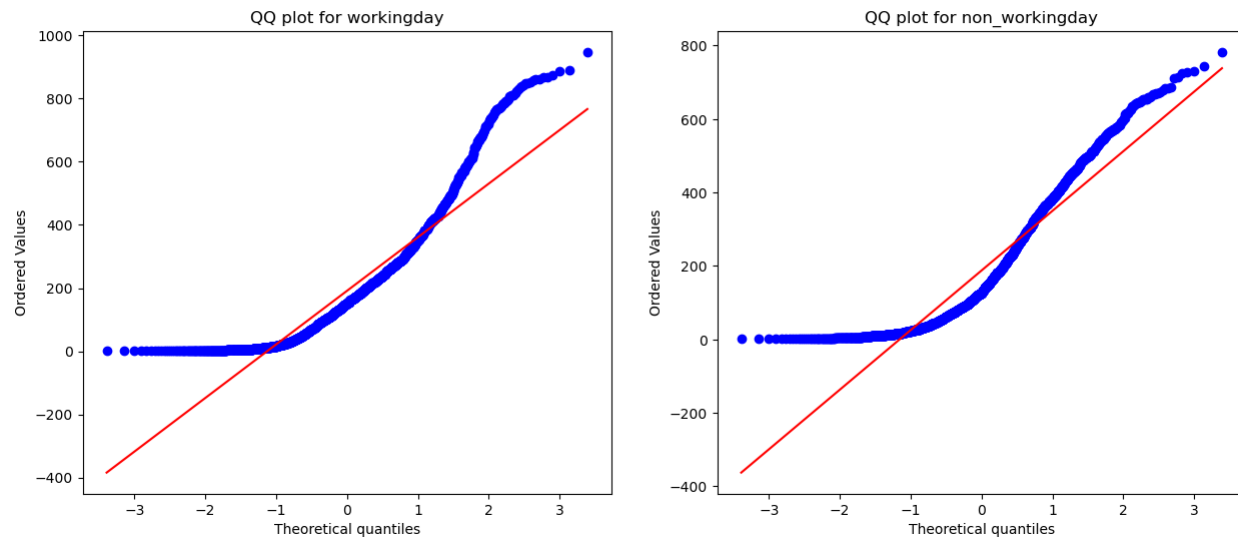- Based on p-value, we will accept or reject H0.

  1. **p-val > alpha** : Accept H0
  2. **p-val < alpha** : Reject H0

```
In [55]: plt.figure(figsize = (15, 5))
         plt.subplot(1, 2, 1)
         sns.histplot(df.loc[df['workingday'] == 1, 'count'].sample(2000),
                     element = 'step', color = 'green', kde = True, label = 'workingday')
         plt.legend()
         plt.subplot(1, 2, 2)
         sns.histplot(df.loc[df['workingday'] == 0, 'count'].sample(2000),
                     element = 'step', color = 'blue', kde = True, label = 'non_workingday')
         plt.legend()
         plt.show()
```



```
In [56]: plt.figure(figsize = (15, 6))
         plt.subplot(1, 2, 1)
         plt.suptitle('QQ plots for the count of electric vehicles rented in workingday and non_workingday')
         stats.probplot(df.loc[df['workingday'] == 1, 'count'].sample(2000), plot = plt, dist = 'norm')
         plt.title('QQ plot for workingday')
         plt.subplot(1, 2, 2)
         stats.probplot(df.loc[df['workingday'] == 0, 'count'].sample(2000), plot = plt, dist = 'norm')
         plt.title('QQ plot for non_workingday')
         plt.show()
```

## QQ plots for the count of electric vehicles rented in workingday and non_workingday



- It can be inferred from the above plot that the distributions do not follow normal distribution.

It can be seen from the above plots that the samples do not come from normal distribution.

- Applying Shapiro-Wilk test for normality

$H_0$ : The sample **follows normal distribution**

$H_1$ : The sample **does not follow normal distribution**

alpha = 0.05

Test Statistics : **Shapiro-Wilk test for normality**

```python
In [57]:  test_stat, p_value = stats.shapiro(df.loc[df['workingday'] == 1, 'count'].sample(2000))
          print('p-value', p_value)
          if p_value < 0.05:
              print('The sample does not follow normal distribution')
          else:
              print('The sample follows normal distribution')
```

```
p-value 9.798081546369001e-37
The sample does not follow normal distribution
```

```python
In [58]:  test_stat, p_value = stats.shapiro(df.loc[df['workingday'] == 0, 'count'].sample(2000))
          print('p-value', p_value)
          if p_value < 0.05:
              print('The sample does not follow normal distribution')
          else:
              print('The sample follows normal distribution')
```

```
p-value 1.4333352412862916e-36
The sample does not follow normal distribution
```

***Transforming the data using boxcox transformation and checking if the transformed data follows normal distribution.***

```python
In [59]:  transformed_workingday = stats.boxcox(df.loc[df['workingday'] == 1, 'count'])[0]
          test_stat, p_value = stats.shapiro(transformed_workingday)
          print('p-value', p_value)
          if p_value < 0.05:
              print('The sample does not follow normal distribution')
          else:
              print('The sample follows normal distribution')
```

```
p-value 1.6132153862898905e-33
The sample does not follow normal distribution
```
```
/Users/debnsuma/anaconda3/lib/python3.10/site-packages/scipy/stats/_morestats.py:1816: UserWarning: p-value may not be accurate for N > 5000.
  warnings.warn("p-value may not be accurate for N > 5000.")
```

- Even after applying the boxcox transformation on each of the "workingday" and "non_workingday" data, the samples do not follow normal distribution.

- Homogeneity of Variances using **Lavene's test**

```python
In [60]:  # Null Hypothesis(H0) - Homogenous Variance

          # Alternate Hypothesis(HA) - Non Homogenous Variance

          test_stat, p_value = stats.levene(df.loc[df['workingday'] == 1, 'count'].sample(2000),
                                            df.loc[df['workingday'] == 0, 'count'].sample(2000))
          print('p-value', p_value)
          if p_value < 0.05:
              print('The samples do not have  Homogenous Variance')
          else:
              print('The samples have Homogenous Variance ')
```

```
p-value 0.5233340046767797
The samples have Homogenous Variance
```

Since the samples are not normally distributed, T-Test cannot be applied here, we can perform its non parametric equivalent test i.e., Mann-Whitney U rank test for two independent samples.

```python
In [61]:  # Ho : Mean no.of electric cycles rented is same for working and non-working days
          # Ha : Mean no.of electric cycles rented is not same for working and non-working days
          # Assuming significance Level to be 0.05
          # Test statistics : Mann-Whitney U rank test for two independent samples

          test_stat, p_value = stats.mannwhitneyu(df.loc[df['workingday'] == 1, 'count'],
                                                  df.loc[df['workingday'] == 0, 'count'])
          print('P-value :',p_value)
```

```
if p_value < 0.05:
    print('Mean no.of electric cycles rented is not same for working and non-working days')
else:
    print('Mean no.of electric cycles rented is same for working and non-working days')
```

```
P-value : 0.9679139953914079
Mean no.of electric cycles rented is same for working and non-working days
```

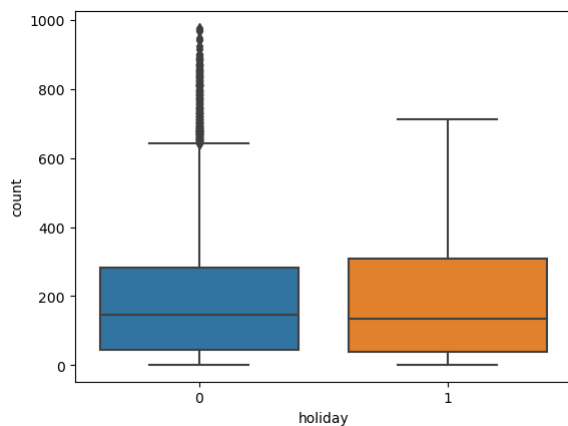## Is there any effect of holidays on the number of electric cycles rented ?¶

In [62]: `df.groupby(by = 'holiday')['count'].describe()`

Out[62]:

| holiday | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| 0 | 10575.0 | 191.741655 | 181.513131 | 1.0 | 43.0 | 145.0 | 283.0 | 977.0 |
| 1 | 311.0 | 185.877814 | 168.300531 | 1.0 | 38.5 | 133.0 | 308.0 | 712.0 |

In [63]: 
```
sns.boxplot(data = df, x = 'holiday', y = 'count')
plt.plot()
```

Out[63]: `[]`



***STEP-1*** : Set up Null Hypothesis

- **Null Hypothesis ( H0 )** - Holidays have no effect on the number of electric vehicles rented

- **Alternate Hypothesis ( HA )** - Holidays has some effect on the number of electric vehicles rented

***STEP-2*** : Checking for basic assumpitons for the hypothesis

- Distribution check using **QQ Plot**
- Homogeneity of Variances using **Levene's test**

***STEP-3***: Define Test statistics; Distribution of T under H0.

- If the assumptions of T Test are met then we can proceed performing T Test for independent samples else we will perform the non parametric test equivalent to T Test for independent sample i.e., Mann-Whitney U rank test for two independent samples.

***STEP-4***: Compute the p-value and fix value of alpha.

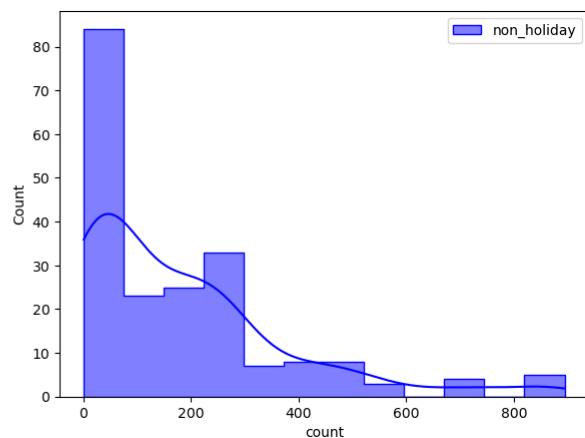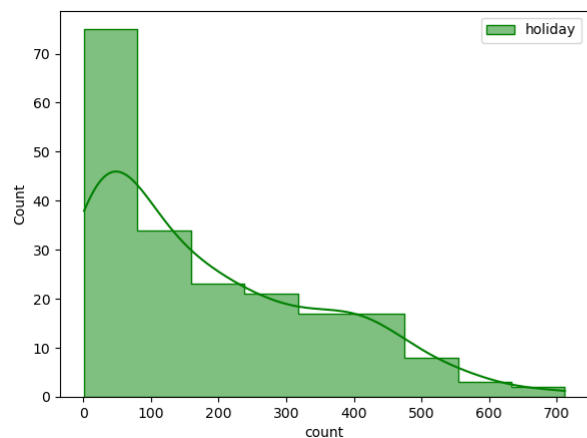- We set our ***alpha to be 0.05***

***STEP-5***: Compare p-value and alpha.

- Based on p-value, we will accept or reject H0.

1. **p-val > alpha** : Accept H0
2. **p-val < alpha** : Reject H0

***Visual Tests to know if the samples follow normal distribution***

In [64]: 
```
plt.figure(figsize = (15, 5))
plt.subplot(1, 2, 1)
sns.histplot(df.loc[df['holiday'] == 1, 'count'].sample(200),
             element = 'step', color = 'green', kde = True, label = 'holiday')
plt.legend()
plt.subplot(1, 2, 2)
sns.histplot(df.loc[df['holiday'] == 0, 'count'].sample(200),
             element = 'step', color = 'blue', kde = True, label = 'non_holiday')
plt.legend()
plt.show()
```
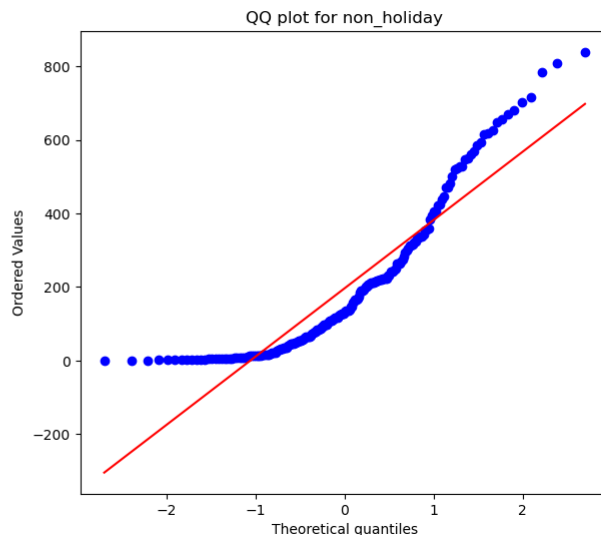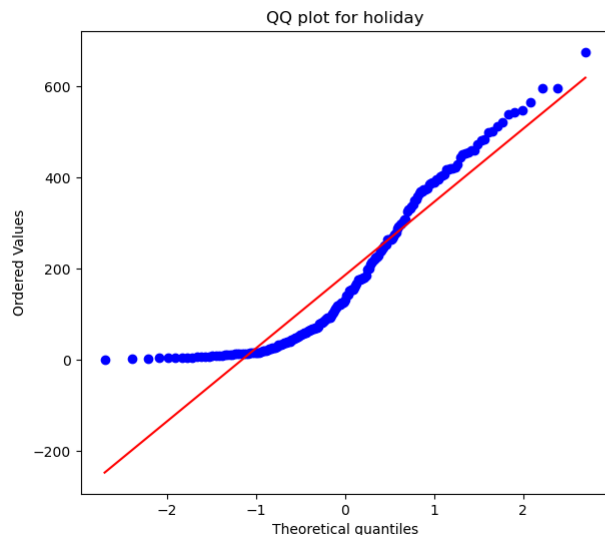
- It can be inferred from the above plot that the distributions do not follow normal distribution.

***Distribution check using QQ Plot***

```
In [65]: plt.figure(figsize = (15, 6))
         plt.subplot(1, 2, 1)
         plt.suptitle('QQ plots for the count of electric vehicles rented in holiday and non_holiday')
         stats.probplot(df.loc[df['holiday'] == 1, 'count'].sample(200), plot = plt, dist = 'norm')
         plt.title('QQ plot for holiday')
         plt.subplot(1, 2, 2)
         stats.probplot(df.loc[df['holiday'] == 0, 'count'].sample(200), plot = plt, dist = 'norm')
         plt.title('QQ plot for non_holiday')
         plt.show()
```

QQ plots for the count of electric vehicles rented in holiday and non_holiday



- It can be inferred from the above plot that the distributions do not follow normal distribution.

It can be seen from the above plots that the samples do not come from normal distribution.

- Applying Shapiro-Wilk test for normality

$H_0$ : The sample **follows normal distribution** $H_1$ : The sample **does not follow normal distribution**

alpha = 0.05

Test Statistics : **Shapiro-Wilk test for normality**

```
In [66]: test_stat, p_value = stats.shapiro(df.loc[df['holiday'] == 1, 'count'].sample(200))
         print('p-value', p_value)
         if p_value < 0.05:
             print('The sample does not follow normal distribution')
         else:
             print('The sample follows normal distribution')
```

```
p-value 2.962478595769369e-10
The sample does not follow normal distribution
```

```
In [67]: test_stat, p_value = stats.shapiro(df.loc[df['holiday'] == 0, 'count'].sample(200))
         print('p-value', p_value)
         if p_value < 0.05:
             print('The sample does not follow normal distribution')
         else:
             print('The sample follows normal distribution')
```

```
p-value 1.6528133028881342e-12
The sample does not follow normal distribution
```

***Transforming the data using boxcox transformation and checking if the transformed data follows normal distribution.***

```
In [68]: transformed_holiday = stats.boxcox(df.loc[df['holiday'] == 1, 'count'])[0]
         test_stat, p_value = stats.shapiro(transformed_holiday)
```

```
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 2.1349286782879062e-07
The sample does not follow normal distribution
```

In [69]:
```
transformed_non_holiday = stats.boxcox(df.loc[df['holiday'] == 0, 'count'].sample(5000))[0]
test_stat, p_value = stats.shapiro(transformed_non_holiday)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 9.483587965325829e-27
The sample does not follow normal distribution
```

- Even after applying the boxcox transformation on each of the "holiday" and "non_holiday" data, the samples do not follow normal distribution.

***Homogeneity of Variances using Levene's test***

In [70]:
```
# Null Hypothesis(H0) - Homogenous Variance

# Alternate Hypothesis(HA) - Non Homogenous Variance

test_stat, p_value = stats.levene(df.loc[df['holiday'] == 0, 'count'].sample(200),
                                  df.loc[df['holiday'] == 1, 'count'].sample(200))
print('p-value', p_value)
if p_value < 0.05:
    print('The samples do not have  Homogenous Variance')
else:
    print('The samples have Homogenous Variance ')
```

```
p-value 0.9696611694288528
The samples have Homogenous Variance
```

Since the samples are not normally distributed, T-Test cannot be applied here, we can perform its non parametric equivalent test i.e., Mann-Whitney U rank test for two independent samples.

In [71]:
```
# Ho : No.of electric cycles rented is similar for holidays and non-holidays
# Ha : No.of electric cycles rented is not similar for holidays and non-holidays days
# Assuming significance Level to be 0.05
# Test statistics : Mann-Whitney U rank test for two independent samples

test_stat, p_value = stats.mannwhitneyu(df.loc[df['holiday'] == 0, 'count'].sample(200),
                                        df.loc[df['holiday'] == 1, 'count'].sample(200))
print('P-value :',p_value)
if p_value < 0.05:
    print('No.of electric cycles rented is not similar for holidays and non-holidays days')
else:
    print('No.of electric cycles rented is similar for holidays and non-holidays')
```

```
P-value : 0.7710050333031868
No.of electric cycles rented is similar for holidays and non-holidays
```

# Is weather dependent on the season ?¶

In [72]:
```
df['season'].describe()
```

Out[72]:
```
count      10886
unique         4
top       winter
freq        2734
Name: season, dtype: object
```

In [73]:
```
df['weather'].describe()
```

Out[73]:
```
count   10886.000000
mean        1.418427
std         0.633839
min         1.000000
25%         1.000000
50%         1.000000
75%         2.000000
max         4.000000
Name: weather, dtype: float64
```

***STEP-1*** : Set up Null Hypothesis

1. **Null Hypothesis ( H0 )** - weather is independent of season

2. **Alternate Hypothesis ( HA )** - weather is dependent of seasons.

***STEP-2***: Define Test statistics

Since we have two categorical features, the Chi- square test is applicable here. Under H0, the test statistic should follow **Chi-Square Distribution**.

***STEP-3***: Checking for basic assumptons for the hypothesis (Non-Parametric Test)

1. The data in the cells should be **frequencies**, or **counts** of cases.
2. The levels (or categories) of the variables are **mutually exclusive**. That is, a particular subject fits into one and only one level of each of the variables.
3. There are 2 variables, and both are measured as **categories**.
4. The **value of the cell expecteds should be 5 or more** in at least 80% of the cells, and no cell should have an expected of less than one (3).

***STEP-4***: Compute the p-value and fix value of alpha.

we will be computing the chi square-test p-value using the chi2_contingency function using scipy.stats. We set our **alpha to be 0.05**

***STEP-5***: Compare p-value and alpha.

Based on p-value, we will accept or reject H0.

1. **p-val > alpha** : Accept H0
2. **p-val < alpha** : Reject H0

The **Chi-square statistic is a non-parametric** (distribution free) tool designed to analyze group differences when the dependent variable is measured at a nominal level. Like all non-parametric statistics, the Chi-square is robust with respect to the distribution of the data. Specifically, it does not require equality of variances among the study groups or homoscedasticity in the data.

```python
In [74]:  # First, finding the contingency table such that each value is the total number of total bikes rented
          # for a particular season and weather
          cross_table = pd.crosstab(index = df['season'],
                                    columns = df['weather'],
                                    values = df['count'],
                                    aggfunc = np.sum).replace(np.nan, 0)
          cross_table
```

Out[74]:

| weather | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| **season** | | | | |
| **fall** | 470116 | 139386 | 31160 | 0 |
| **spring** | 223009 | 76406 | 12919 | 164 |
| **summer** | 426350 | 134177 | 27755 | 0 |
| **winter** | 356588 | 157191 | 30255 | 0 |

Since the above contingency table has one column in which the count of the rented electric vehicle is less than 5 in most of the cells, we can remove the weather 4 and then proceed further.

```python
In [75]:  cross_table = pd.crosstab(index = df['season'],
                                    columns = df.loc[df['weather'] != 4, 'weather'],
                                    values = df['count'],
                                    aggfunc = np.sum).to_numpy()[:, :3]
          cross_table
```

Out[75]:
```
array([[470116, 139386,  31160],
       [223009,  76406,  12919],
       [426350, 134177,  27755],
       [356588, 157191,  30255]])
```

```python
In [76]:  chi_test_stat, p_value, dof, expected = stats.chi2_contingency(observed = cross_table)
          print('Test Statistic =', chi_test_stat)
          print('p value =', p_value)
          print('-' * 65)
          print("Expected : '\n'", expected)
```

```
Test Statistic = 10838.372332480214
p value = 0.0
-----------------------------------------------------------------
Expected : '
' [[453484.88557396 155812.72247031  31364.39195574]
 [221081.86259035  75961.44434981  15290.69305984]
 [416408.3330293  143073.60199337  28800.06497733]
 [385087.91880639 132312.23118651  26633.8500071 ]]
```

```python
In [77]:  alpha = 0.05
          if p_value < alpha:
              print('Reject Null Hypothesis')
          else:
              print('Failed to reject Null Hypothesis')
```
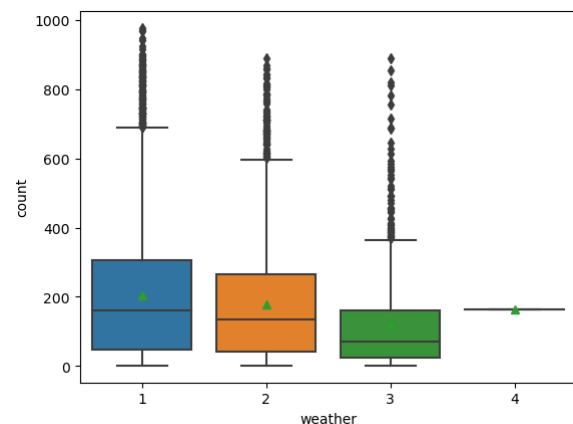
```
Reject Null Hypothesis
```

Therefore, there is statistically significant dependency of weather and season based on the number of number of bikes rented.

## Is the number of cycles rented is similar or different in different weather ?

```python
In [78]:  df.groupby(by = 'weather')['count'].describe()
```

Out[78]:

| | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| **weather** | | | | | | | | |
| **1** | 7192.0 | 205.236791 | 187.959566 | 1.0 | 48.0 | 161.0 | 305.0 | 977.0 |
| **2** | 2834.0 | 178.955540 | 168.366413 | 1.0 | 41.0 | 134.0 | 264.0 | 890.0 |
| **3** | 859.0 | 118.846333 | 138.581297 | 1.0 | 23.0 | 71.0 | 161.0 | 891.0 |
| **4** | 1.0 | 164.000000 | NaN | 164.0 | 164.0 | 164.0 | 164.0 | 164.0 |

```python
In [79]:  sns.boxplot(data = df, x = 'weather', y = 'count', showmeans = True)
          plt.show()
```



```python
In [80]:  df_weather1 = df.loc[df['weather'] == 1]
          df_weather2 = df.loc[df['weather'] == 2]
          df_weather3 = df.loc[df['weather'] == 3]
          df_weather4 = df.loc[df['weather'] == 4]
          len(df_weather1), len(df_weather2), len(df_weather3), len(df_weather4)
```

Out[80]:  (7192, 2834, 859, 1)

**STEP-1** : Set up Null Hypothesis

- **Null Hypothesis ( H0 )** - Mean of cycle rented per hour is same for weather 1, 2 and 3. (We wont be considering weather 4 as there in only 1 data point for weather 4 and we cannot perform a ANOVA test with a single data point for a group)

- **Alternate Hypothesis ( HA )** -Mean of cycle rented per hour is not same for season 1,2,3 and 4 are different.

*STEP-2* : Checking for basic assumpitons for the hypothesis

Normality check using **QQ Plot**. If the distribution is not normal, use **BOX-COX transform** to transform it to normal distribution.

Homogeneity of Variances using **Levene's test**

Each observations are **independent**.

*STEP-3*: Define **Test statistics**

The test statistic for a One-Way ANOVA is denoted as F. For an independent variable with k groups, the F statistic evaluates whether the group means are significantly different.

**F=MSB / MSW**

Under H0, the test statistic should follow **F-Distribution**.

*STEP-4*: Decide the kind of test.

We will be performing **right tailed f-test**

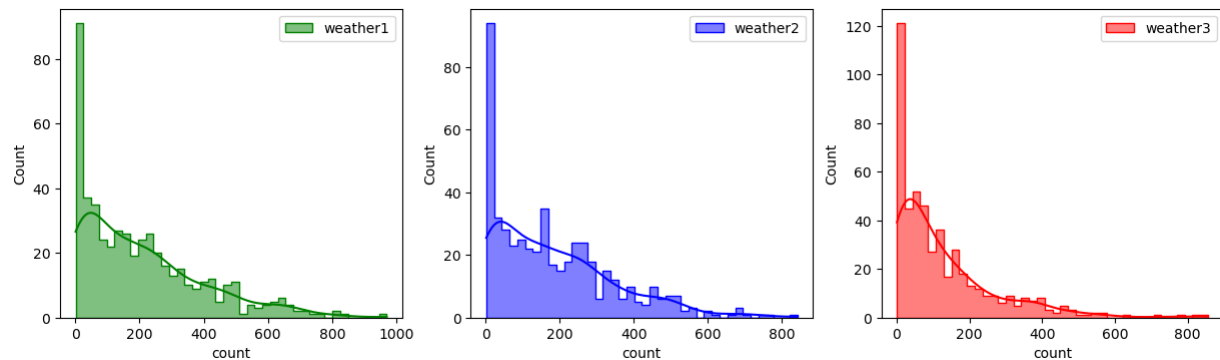*STEP-5*: Compute the **p-value** and fix value of alpha.

we will be computing the anova-test p-value using the f_oneway function using scipy.stats. We set our **alpha to be 0.05**

*STEP-6*: Compare p-value and alpha.

Based on p-value, we will accept or reject H0.

- **p-val > alpha** : Accept H0
- **p-val < alpha** : Reject H0

```
In [81]: plt.figure(figsize = (15, 4))
plt.subplot(1, 3, 1)
sns.histplot(df_weather1.loc[:, 'count'].sample(500), bins = 40,
             element = 'step', color = 'green', kde = True, label = 'weather1')
plt.legend()
plt.subplot(1, 3, 2)
sns.histplot(df_weather2.loc[:, 'count'].sample(500), bins = 40,
             element = 'step', color = 'blue', kde = True, label = 'weather2')
plt.legend()
plt.subplot(1, 3, 3)
sns.histplot(df_weather3.loc[:, 'count'].sample(500), bins = 40,
             element = 'step', color = 'red', kde = True, label = 'weather3')
plt.legend()
plt.show()
```
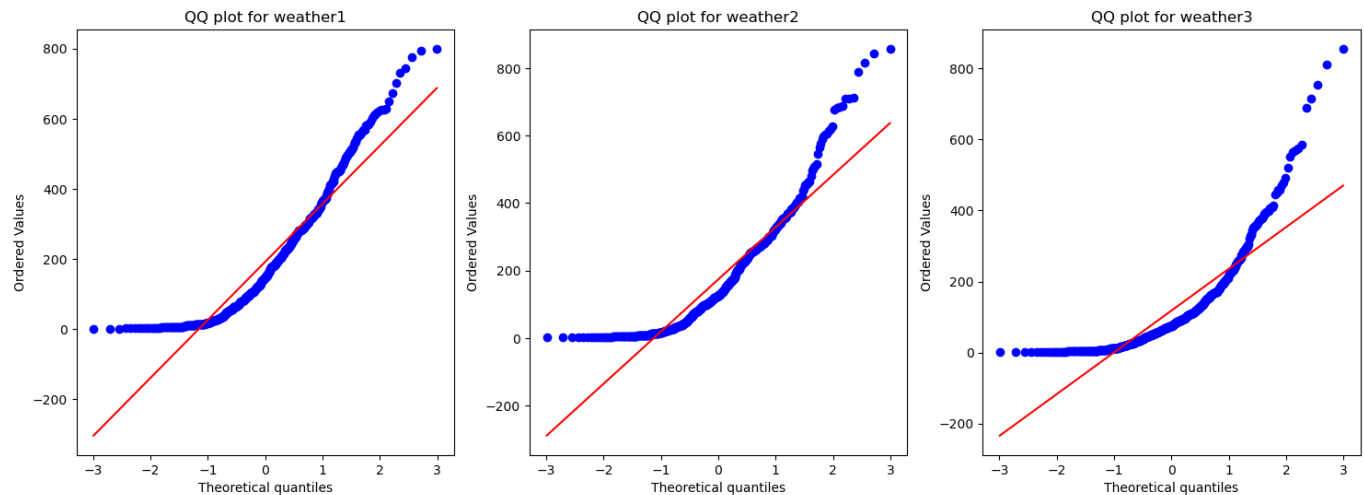


- It can be inferred from the above plot that the distributions do not follow normal distribution.

*Distribution check using QQ Plot*

```
In [82]: plt.figure(figsize = (18, 6))
plt.subplot(1, 3, 1)
plt.suptitle('QQ plots for the count of electric vehicles rented in different weathers')
stats.probplot(df_weather1.loc[:, 'count'].sample(500), plot = plt, dist = 'norm')
plt.title('QQ plot for weather1')
plt.subplot(1, 3, 2)
stats.probplot(df_weather2.loc[:, 'count'].sample(500), plot = plt, dist = 'norm')
plt.title('QQ plot for weather2')
plt.subplot(1, 3, 3)
stats.probplot(df_weather3.loc[:, 'count'].sample(500), plot = plt, dist = 'norm')
plt.title('QQ plot for weather3')
plt.show()
```

QQ plots for the count of electric vehicles rented in different weathers



- It can be inferred from the above plot that the distributions do not follow normal distribution.

  ###### It can be seen from the above plots that the samples do not come from normal distribution.

  - Applying Shapiro-Wilk test for normality

$H_0$ : The sample **follows normal distribution** $H_1$ : The sample **does not follow normal distribution**

alpha = 0.05

Test Statistics : **Shapiro-Wilk test for normality**

```
In [83]:  test_stat, p_value = stats.shapiro(df_weather1.loc[:, 'count'].sample(500))
          print('p-value', p_value)
          if p_value < 0.05:
              print('The sample does not follow normal distribution')
          else:
              print('The sample follows normal distribution')
```

```
p-value 2.923493189289835e-17
The sample does not follow normal distribution
```

```
In [84]:  test_stat, p_value = stats.shapiro(df_weather2.loc[:, 'count'].sample(500))
          print('p-value', p_value)
          if p_value < 0.05:
              print('The sample does not follow normal distribution')
          else:
              print('The sample follows normal distribution')
```

```
p-value 1.6449148225106998e-19
The sample does not follow normal distribution
```

```
In [85]:  test_stat, p_value = stats.shapiro(df_weather3.loc[:, 'count'].sample(500))
          print('p-value', p_value)
          if p_value < 0.05:
              print('The sample does not follow normal distribution')
          else:
              print('The sample follows normal distribution')
```

```
p-value 9.19182421581874e-27
The sample does not follow normal distribution
```

***Transforming the data using boxcox transformation and checking if the transformed data follows normal distribution.***

```
In [86]:  transformed_weather1 = stats.boxcox(df_weather1.loc[:, 'count'].sample(5000))[0]
          test_stat, p_value = stats.shapiro(transformed_weather1)
          print('p-value', p_value)
          if p_value < 0.05:
              print('The sample does not follow normal distribution')
          else:
              print('The sample follows normal distribution')
```

```
p-value 3.4772520701165714e-28
The sample does not follow normal distribution
```

```
In [87]:  transformed_weather2 = stats.boxcox(df_weather2.loc[:, 'count'])[0]
          test_stat, p_value = stats.shapiro(transformed_weather2)
          print('p-value', p_value)
          if p_value < 0.05:
              print('The sample does not follow normal distribution')
          else:
              print('The sample follows normal distribution')
```

```
p-value 1.9212615187509174e-19
The sample does not follow normal distribution
```

```
In [88]:  transformed_weather3 = stats.boxcox(df_weather3.loc[:, 'count'])[0]
          test_stat, p_value = stats.shapiro(transformed_weather3)
          print('p-value', p_value)
          if p_value < 0.05:
              print('The sample does not follow normal distribution')
          else:
              print('The sample follows normal distribution')
```

```
p-value 1.4131142052065115e-06
The sample does not follow normal distribution
```

- Even after applying the boxcox transformation on each of the weather data, the samples do not follow normal distribution.

***Homogeneity of Variances using Levene's test***

```
In [89]:  # Null Hypothesis(H0) - Homogenous Variance
```

```
# Alternate Hypothesis(HA) - Non Homogenous Variance

test_stat, p_value = stats.levene(df_weather1.loc[:, 'count'].sample(500),
                                  df_weather2.loc[:, 'count'].sample(500),
                                  df_weather3.loc[:, 'count'].sample(500))
print('p-value', p_value)
if p_value < 0.05:
    print('The samples do not have  Homogenous Variance')
else:
    print('The samples have Homogenous Variance ')
```

```
p-value 1.1664938638862327e-16
The samples do not have  Homogenous Variance
```

Since the samples are not normally distributed and do not have the same variance, f_oneway test cannot be performed here, we can perform its non parametric equivalent test i.e., Kruskal-Wallis H-test for independent samples.

### Is the number of cycles rented is similar or different in different season ?
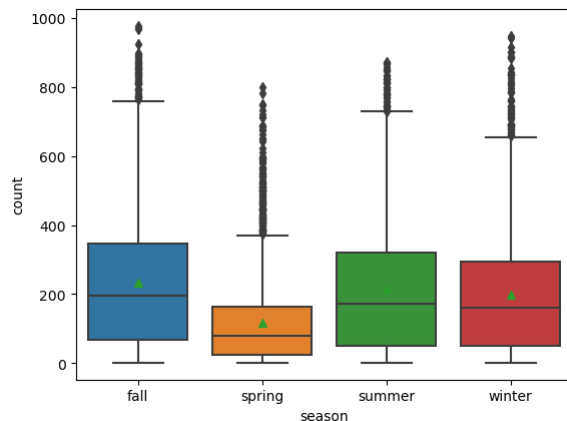
In [90]: `df.groupby(by = 'season')['count'].describe()`

Out[90]:

| season | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| fall | 2733.0 | 234.417124 | 197.151001 | 1.0 | 68.0 | 195.0 | 347.0 | 977.0 |
| spring | 2686.0 | 116.343261 | 125.273974 | 1.0 | 24.0 | 78.0 | 164.0 | 801.0 |
| summer | 2733.0 | 215.251372 | 192.007843 | 1.0 | 49.0 | 172.0 | 321.0 | 873.0 |
| winter | 2734.0 | 198.988296 | 177.622409 | 1.0 | 51.0 | 161.0 | 294.0 | 948.0 |

In [91]:
```
df_season_spring = df.loc[df['season'] == 'spring', 'count']
df_season_summer = df.loc[df['season'] == 'summer', 'count']
df_season_fall = df.loc[df['season'] == 'fall', 'count']
df_season_winter = df.loc[df['season'] == 'winter', 'count']
len(df_season_spring), len(df_season_summer), len(df_season_fall), len(df_season_winter)
```

Out[91]: `(2686, 2733, 2733, 2734)`

In [92]:
```
sns.boxplot(data = df, x = 'season', y = 'count', showmeans = True)
plt.show()
```



*STEP-1* : Set up Null Hypothesis

- **Null Hypothesis ( H0 )** - Mean of cycle rented per hour is same for season 1,2,3 and 4.

- **Alternate Hypothesis ( HA )** -Mean of cycle rented per hour is different for season 1,2,3 and 4.

*STEP-2* : Checking for basic assumpitons for the hypothesis

1. **Normality check** using QQ Plot. If the distribution is not normal, use **BOX-COX transform** to transform it to normal distribution.

2. Homogeneity of Variances using **Levene's test**

3. Each observations are **independent**.

*STEP-3*: Define Test statistics

The test statistic for a One-Way ANOVA is denoted as F. For an independent variable with k groups, the F statistic evaluates whether the group means are significantly different.

**F=MSB/MSW**

Under H0, the test statistic should follow **F-Distribution**.

*STEP-4*: Decide the kind of test.

We will be performing **right tailed f-test**

*STEP-5*: Compute the p-value and fix value of alpha.

we will be computing the anova-test p-value using the **f_oneway** function using scipy.stats. We set our alpha to be **0.05**
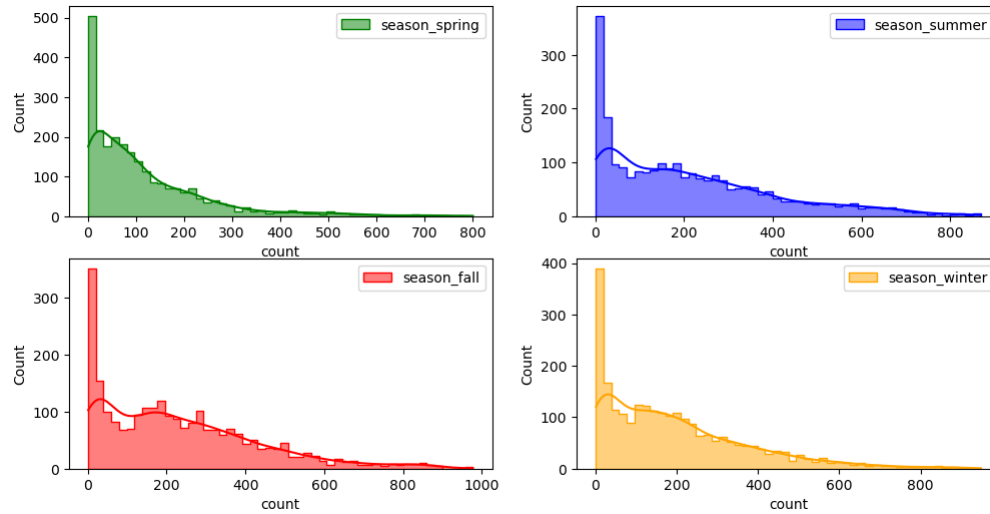
*STEP-6*: Compare p-value and alpha.

Based on p-value, we will accept or reject H0. p-val > alpha : Accept H0 p-val < alpha : Reject H0

*Visual Tests to know if the samples follow normal distribution*

```
In [93]: plt.figure(figsize = (12, 6))
         plt.subplot(2, 2, 1)
         sns.histplot(df_season_spring.sample(2500), bins = 50,
                      element = 'step', color = 'green', kde = True, label = 'season_spring')
         plt.legend()
         plt.subplot(2, 2, 2)
         sns.histplot(df_season_summer.sample(2500), bins = 50,
                      element = 'step', color = 'blue', kde = True, label = 'season_summer')
         plt.legend()
         plt.subplot(2, 2, 3)
         sns.histplot(df_season_fall.sample(2500), bins = 50,
                      element = 'step', color = 'red', kde = True, label = 'season_fall')
         plt.legend()
         plt.subplot(2, 2, 4)
         sns.histplot(df_season_winter.sample(2500), bins = 50,
                      element = 'step', color = 'orange', kde = True, label = 'season_winter')
         plt.legend()
         plt.plot()
```

Out[93]: []

- It can be inferred from the above plot that the distributions do not follow normal distribution.

*Distribution check using QQ Plot*

```
In [94]: plt.figure(figsize = (12, 12))
         plt.subplot(2, 2, 1)
         plt.suptitle('QQ plots for the count of electric vehicles rented in different seasons')
         stats.probplot(df_season_spring.sample(2500), plot = plt, dist = 'norm')
         plt.title('QQ plot for spring season')

         plt.subplot(2, 2, 2)
         stats.probplot(df_season_summer.sample(2500), plot = plt, dist = 'norm')
         plt.title('QQ plot for summer season')

         plt.subplot(2, 2, 3)
         stats.probplot(df_season_fall.sample(2500), plot = plt, dist = 'norm')
         plt.title('QQ plot for fall season')

         plt.subplot(2, 2, 4)
         stats.probplot(df_season_winter.sample(2500), plot = plt, dist = 'norm')
         plt.title('QQ plot for winter season')
         plt.plot()
```
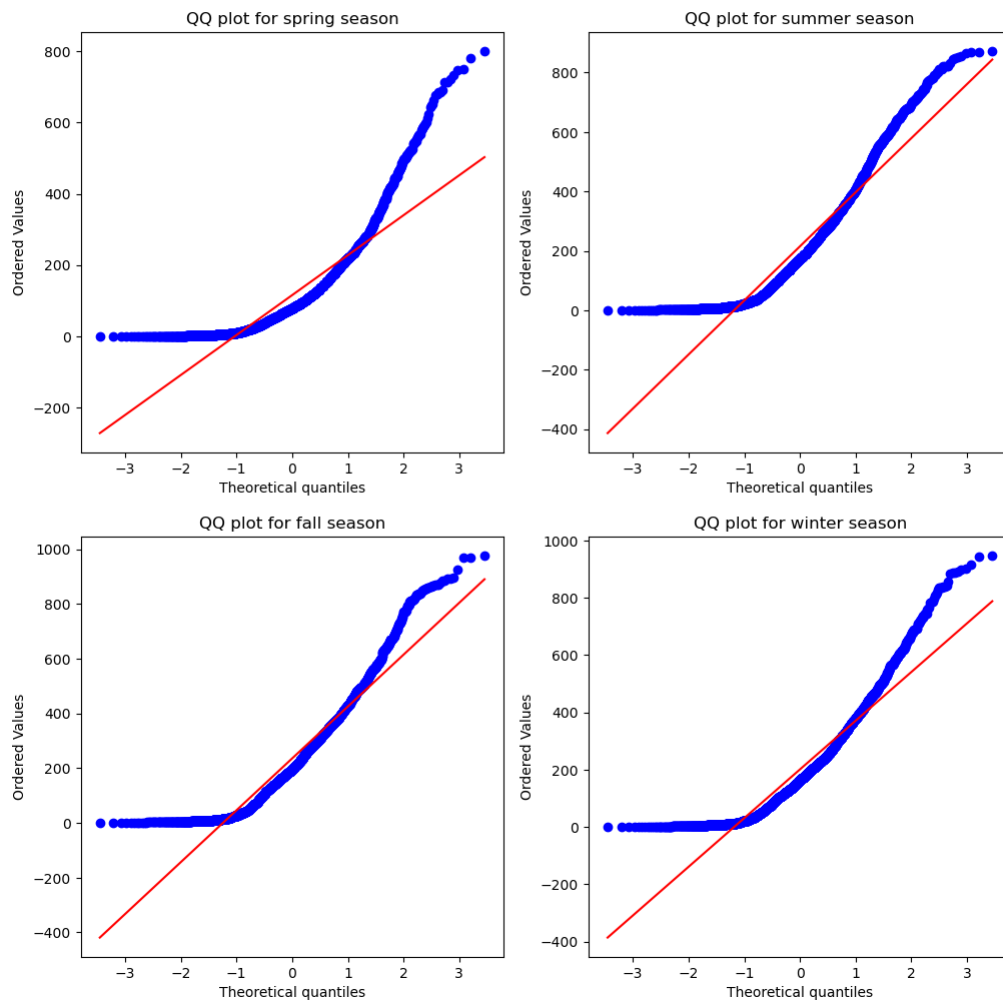
Out[94]: []

QQ plots for the count of electric vehicles rented in different seasons



- It can be inferred from the above plots that the distributions do not follow normal distribution.

It can be seen from the above plots that the samples do not come from normal distribution.

- Applying Shapiro-Wilk test for normality

$H_0$ : The sample **follows normal distribution** $H_1$ : The sample **does not follow normal distribution**

alpha = 0.05

Test Statistics : **Shapiro-Wilk test for normality**

```
In [95]: test_stat, p_value = stats.shapiro(df_season_spring.sample(2500))
         print('p-value', p_value)
         if p_value < 0.05:
             print('The sample does not follow normal distribution')
         else:
             print('The sample follows normal distribution')

         p-value 0.0
         The sample does not follow normal distribution
```

```
In [96]: test_stat, p_value = stats.shapiro(df_season_summer.sample(2500))
         print('p-value', p_value)
         if p_value < 0.05:
             print('The sample does not follow normal distribution')
         else:
             print('The sample follows normal distribution')

         p-value 1.414941169861852e-37
         The sample does not follow normal distribution
```

```
In [97]: test_stat, p_value = stats.shapiro(df_season_fall.sample(2500))
         print('p-value', p_value)
         if p_value < 0.05:
             print('The sample does not follow normal distribution')
         else:
             print('The sample follows normal distribution')

         p-value 1.4432535107445954e-35
         The sample does not follow normal distribution
```

```
In [98]: test_stat, p_value = stats.shapiro(df_season_winter.sample(2500))
         print('p-value', p_value)
         if p_value < 0.05:
             print('The sample does not follow normal distribution')
         else:
             print('The sample follows normal distribution')
```

```
p-value 1.7233151698685028e-38
The sample does not follow normal distribution
```

***Transforming the data using boxcox transformation and checking if the transformed data follows normal distribution.***

In [99]:
```python
transformed_df_season_spring = stats.boxcox(df_season_spring.sample(2500))[0]
test_stat, p_value = stats.shapiro(transformed_df_season_spring)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 8.352732844451035e-17
The sample does not follow normal distribution
```

In [100…
```python
transformed_df_season_summer = stats.boxcox(df_season_summer.sample(2500))[0]
test_stat, p_value = stats.shapiro(transformed_df_season_summer)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 1.0806075678529822e-21
The sample does not follow normal distribution
```

In [101…
```python
transformed_df_season_fall = stats.boxcox(df_season_fall.sample(2500))[0]
test_stat, p_value = stats.shapiro(transformed_df_season_fall)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 2.730513508514741e-21
The sample does not follow normal distribution
```

In [102…
```python
transformed_df_season_winter = stats.boxcox(df_season_winter.sample(2500))[0]
test_stat, p_value = stats.shapiro(transformed_df_season_winter)
print('p-value', p_value)
if p_value < 0.05:
    print('The sample does not follow normal distribution')
else:
    print('The sample follows normal distribution')
```

```
p-value 1.057489214708414e-19
The sample does not follow normal distribution
```

- Even after applying the boxcox transformation on each of the season data, the samples do not follow normal distribution.

***Homogeneity of Variances using Levene's test***

In [103…
```python
# Null Hypothesis(H0) - Homogenous Variance

# Alternate Hypothesis(HA) - Non Homogenous Variance

test_stat, p_value = stats.levene(df_season_spring.sample(2500),
                                  df_season_summer.sample(2500),
                                  df_season_fall.sample(2500),
                                  df_season_winter.sample(2500))
print('p-value', p_value)
if p_value < 0.05:
    print('The samples do not have  Homogenous Variance')
else:
    print('The samples have Homogenous Variance ')
```

```
p-value 3.110478071768311e-111
The samples do not have  Homogenous Variance
```

Since the samples are not normally distributed and do not have the same variance, f_oneway test cannot be performed here, we can perform its non parametric equivalent test i.e., Kruskal-Wallis H-test for independent samples.

In [104…
```python
# Ho : Mean no. of cycles rented is same for different weather
# Ha : Mean no. of cycles rented is different for different weather
# Assuming significance Level to be 0.05
alpha = 0.05
test_stat, p_value = stats.kruskal(df_season_spring, df_season_summer, df_season_fall,df_season_winter)
print('Test Statistic =', test_stat)
print('p value =', p_value)
```

```
Test Statistic = 699.6668548181988
p value = 2.479008372608633e-151
```
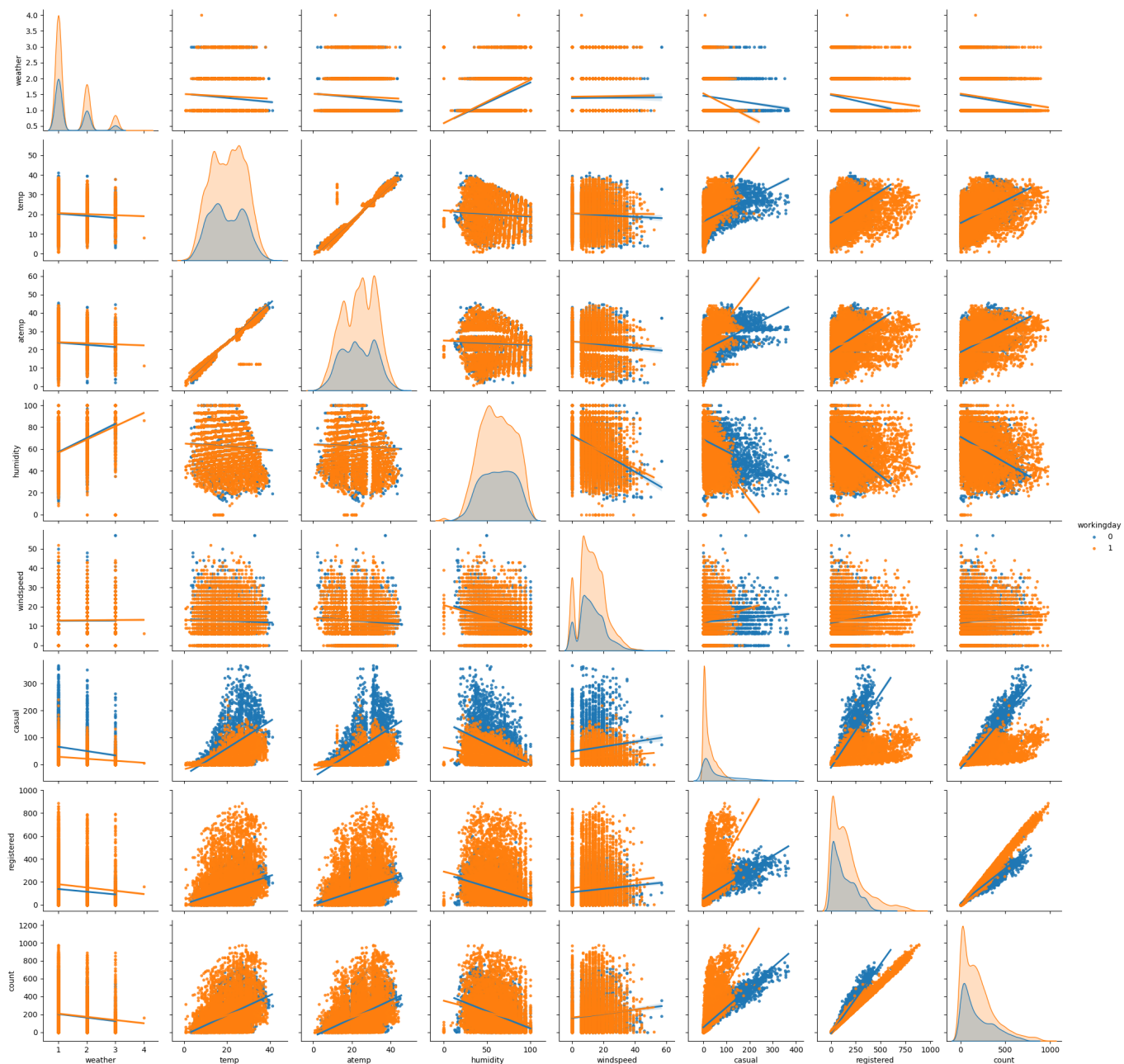
In [105…
```python
if p_value < alpha:
    print('Reject Null Hypothesis')
else:
    print('Failed to reject Null Hypothesis')
```

```
Reject Null Hypothesis
```

Therefore, the average number of rental bikes is statistically different for different seasons.

In [106…
```python
sns.pairplot(data = df,
             kind = 'reg',
             hue = 'workingday',
             markers = '.')
plt.show()
```

```
In [107…  corr_data = df.corr()
          corr_data
```
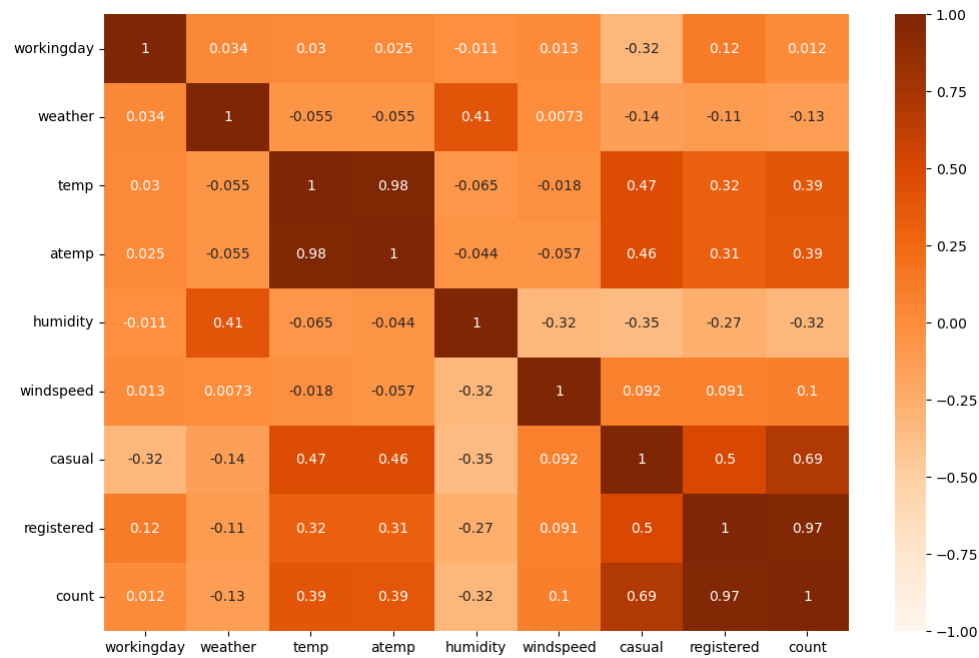
```
/var/folders/sg/qf1dw3cs4q5007gb2_9zd8600000gr/T/ipykernel_91140/919268980.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future
version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
    corr_data = df.corr()
```

Out[107]:

|  | workingday | weather | temp | atemp | humidity | windspeed | casual | registered | count |
|---|---|---|---|---|---|---|---|---|---|
| workingday | 1.000000 | 0.033772 | 0.029966 | 0.024660 | -0.010880 | 0.013373 | -0.319111 | 0.119460 | 0.011594 |
| weather | 0.033772 | 1.000000 | -0.055035 | -0.055376 | 0.406244 | 0.007261 | -0.135918 | -0.109340 | -0.128655 |
| temp | 0.029966 | -0.055035 | 1.000000 | 0.984948 | -0.064949 | -0.017852 | 0.467097 | 0.318571 | 0.394454 |
| atemp | 0.024660 | -0.055376 | 0.984948 | 1.000000 | -0.043536 | -0.057473 | 0.462067 | 0.314635 | 0.389784 |
| humidity | -0.010880 | 0.406244 | -0.064949 | -0.043536 | 1.000000 | -0.318607 | -0.348187 | -0.265458 | -0.317371 |
| windspeed | 0.013373 | 0.007261 | -0.017852 | -0.057473 | -0.318607 | 1.000000 | 0.092276 | 0.091052 | 0.101369 |
| casual | -0.319111 | -0.135918 | 0.467097 | 0.462067 | -0.348187 | 0.092276 | 1.000000 | 0.497250 | 0.690414 |
| registered | 0.119460 | -0.109340 | 0.318571 | 0.314635 | -0.265458 | 0.091052 | 0.497250 | 1.000000 | 0.970948 |
| count | 0.011594 | -0.128655 | 0.394454 | 0.389784 | -0.317371 | 0.101369 | 0.690414 | 0.970948 | 1.000000 |

```
In [108…  plt.figure(figsize = (12, 8))
          sns.heatmap(data = corr_data, cmap = 'Oranges', annot = True, vmin = -1, vmax = 1)
          plt.show()
```

- Very High Correlation (> 0.9) exists between columns [atemp, temp] and [count, registered]
- High positively / negatively correlation (0.7 - 0.9) does not exist between any columns.
- Moderate positive correlation (0.5 - 0.7) exists between columns [casual, count], [casual, registered].
- Low Positive correlation (0.3 - 0.5) exists between columns [count, temp], [count, atemp], [casual, atemp]
- Negligible correlation exists between all other combinations of columns.

In [ ]: