

```
In [2]: import numpy as np
import pandas as pd
from scipy import stats
import matplotlib.pyplot as plt
import seaborn as sns

In [3]: df = pd.read_csv("yulu.csv")

In [4]: df.head()

Out [4]:
      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 [5]: # no of rows and columns in dataset
print(f"Rows: {df.shape[0]} \nW columns: {df.shape[1]}")
# Rows: 10886
# Columns: 12

In [6]: 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

Datatype of following attributes needs to be changed to proper data type

• datetime -> datetime
• season -> to categorical
• holiday -> to categorical
• workingday -> to categorical
• weather -> to categorical
```

```
In [7]: df['datetime'] = pd.to_datetime(df['datetime'])

cat_cols = ['season', 'holiday', 'workingday', 'weather']
for col in cat_cols:
    df[col] = df[col].astype('object')

In [8]: df.iloc[:, 1:].describe(include='all')

      season  holiday  workingday  weather  temp  atemp  humidity  windspeed  casual  registered  count
count 10886.0 10886.0 10886.0 10886.0 10886.000000 10886.000000 10886.000000 10886.000000 10886.000000 10886.000000 10886.000000
unique    4.0    2.0          2.0    4.0      NaN      NaN      NaN      NaN      NaN      NaN      NaN
top      4.0    0.0          1.0    1.0      NaN      NaN      NaN      NaN      NaN      NaN      NaN
freq    2734.0 10575.0      7412.0  7192.0      NaN      NaN      NaN      NaN      NaN      NaN      NaN
mean      NaN      NaN      NaN      NaN      20.23086  23.65094    61.89460    12.79395    36.02195    155.95217    191.57432
std      NaN      NaN      NaN      NaN      7.77159   8.474601   19.24503    8.164537   49.96047    151.03903    181.144454
min      NaN      NaN      NaN      NaN      0.82000   0.780000   0.000000   0.000000   0.000000   0.000000   1.000000
25%      NaN      NaN      NaN      NaN      13.94000   16.665000   47.000000   7.001500   4.000000   36.000000   42.000000
50%      NaN      NaN      NaN      NaN      20.50000   24.240000   62.000000  12.998000  17.000000  118.000000  145.000000
75%      NaN      NaN      NaN      NaN      26.24000   31.060000   77.000000  16.997900  49.000000  222.000000  284.000000
max      NaN      NaN      NaN      NaN      41.00000   45.455000   100.00000   56.996900  367.000000  886.000000  977.000000

• There are no missing values in the dataset.
• casual and registered attributes might have outliers because their mean and median are very far away to one another and the value of standard deviation is also high which tells us that there is high variance in the data of these attributes.
```

```
In [9]: # detecting missing values in the dataset
df.isnull().sum()

Out [9]:
datetime      0
season        0
holiday       0
workingday    0
weather       0
temp          0
atemp         0
humidity      0
windspeed     0
casual        0
registered    0
count         0
dtype: int64

There are no missing values present in the dataset.

In [10]: # minimum datetime and maximum datetime
df['datetime'].min(), df['datetime'].max()

Out [10]:
(Timestamp('2011-01-01 00:00:00'), Timestamp('2012-12-19 23:00:00'))

In [11]: # number of unique values in each categorical columns
df[cat_cols].melt().groupby(['variable', 'value'])['value'].count()

Out [11]:
variable value
holiday      0  10575
             1   311
season       1  2686
             2  2733
             3  2733
             4  2734
weather      1  7192
             2  2834
             3   859
             4    1
workingday   0  3474
             1  7412
```

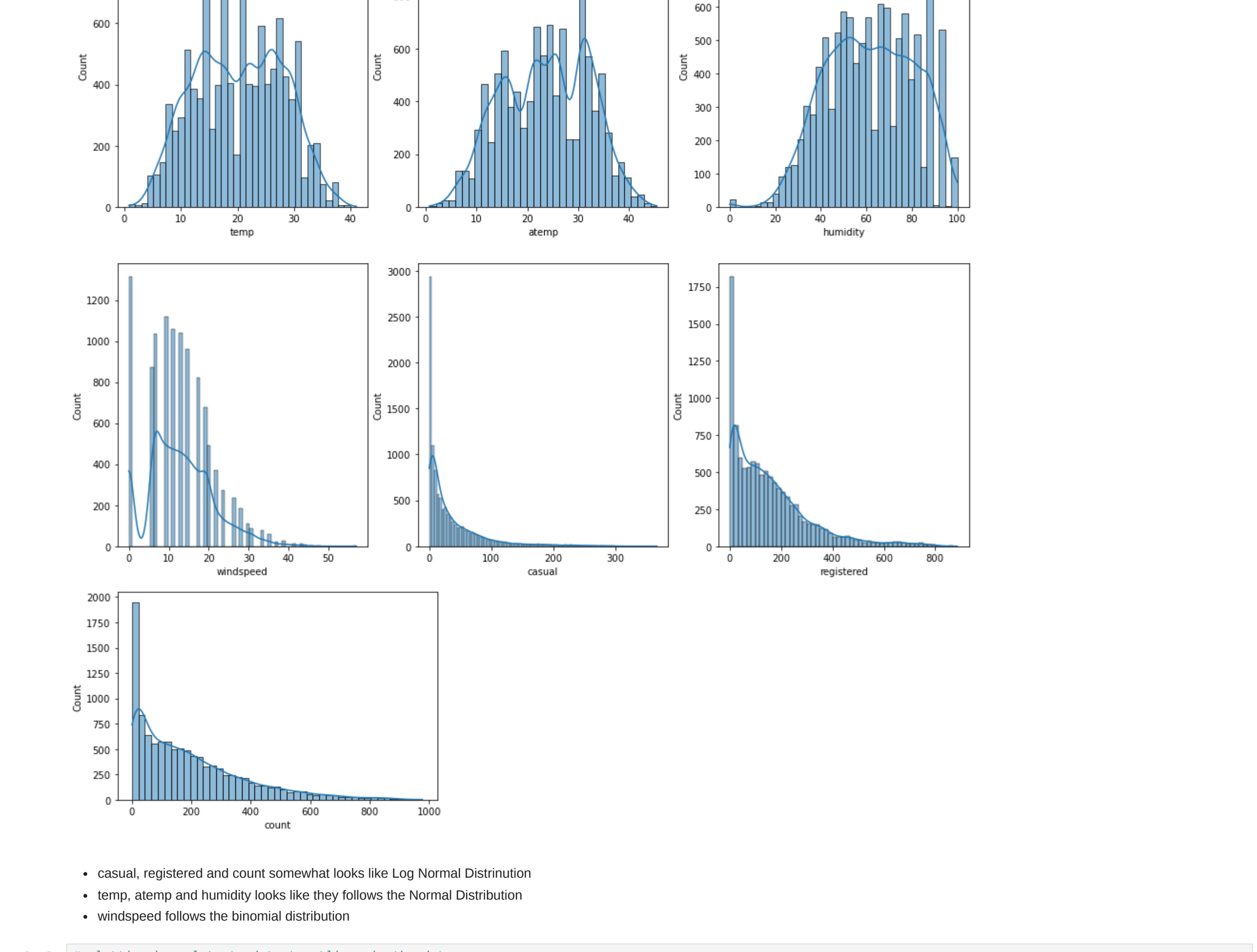
Univariate Analysis

```
In [12]: # understanding the distribution for numerical variables
num_cols = ['temp', 'atemp', 'humidity', 'windspeed', 'casual', 'registered', 'count']

fig, axis = plt.subplots(nrows=2, ncols=3, figsize=(16, 12))

index = 0
for row in range(2):
    for col in range(3):
        sns.histplot(df[num_cols[index]], ax=axis[row, col], kde=True)
        index += 1

plt.show()
sns.histplot(df[num_cols[-1]], kde=True)
plt.show()
```

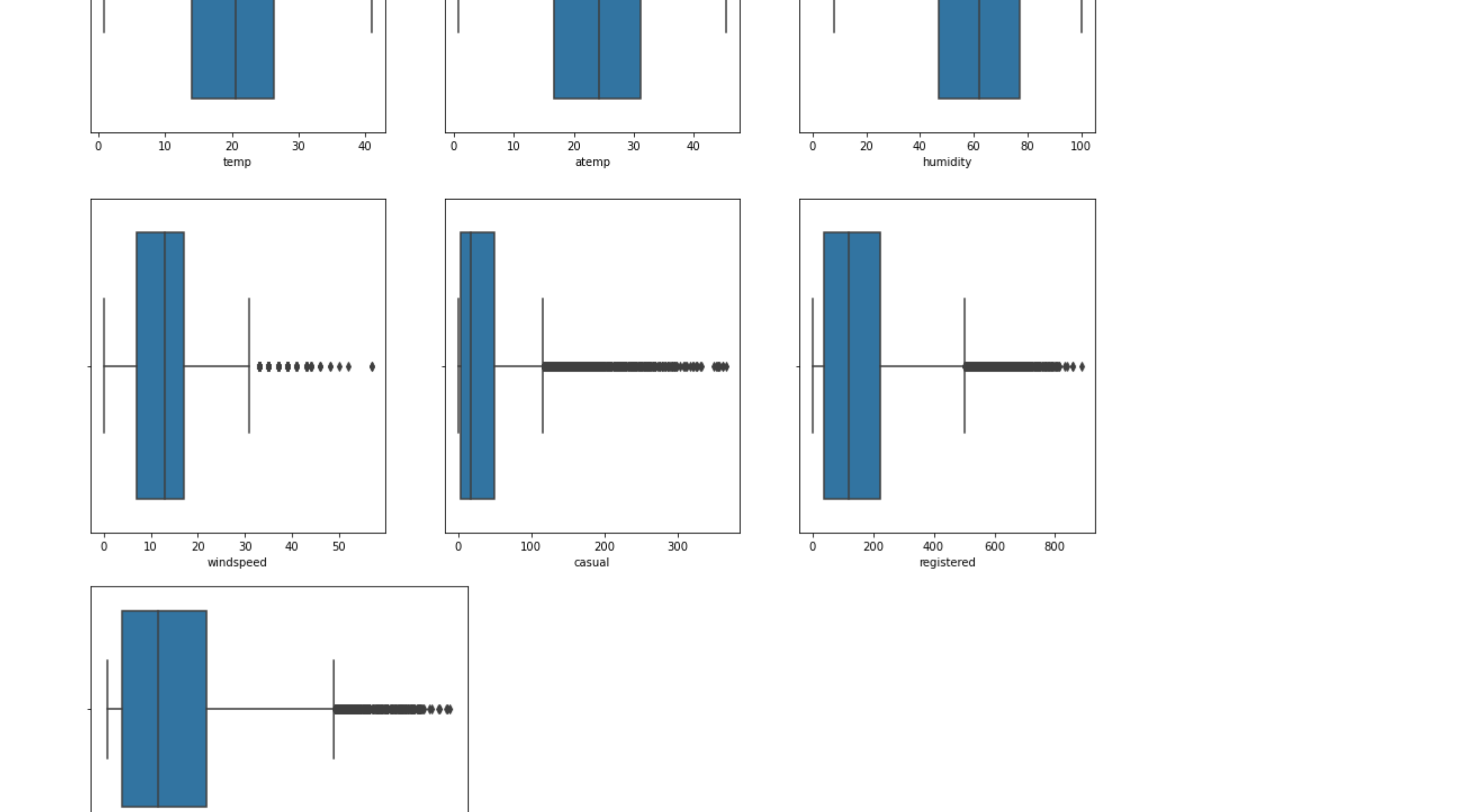


- casual, registered and count somewhat looks like Log Normal Distribution
- temp, atemp and humidity looks like they follow the Normal Distribution
- windspeed follows the binomial distribution

```
In [13]: # plotting box plots to detect outliers in the data
fig, axis = plt.subplots(nrows=2, ncols=3, figsize=(16, 12))

index = 0
for row in range(2):
    for col in range(3):
        sns.boxplot(x=df[num_cols[index]], ax=axis[row, col])
        index += 1

plt.show()
sns.boxplot(x=df[num_cols[-1]])
plt.show()
```

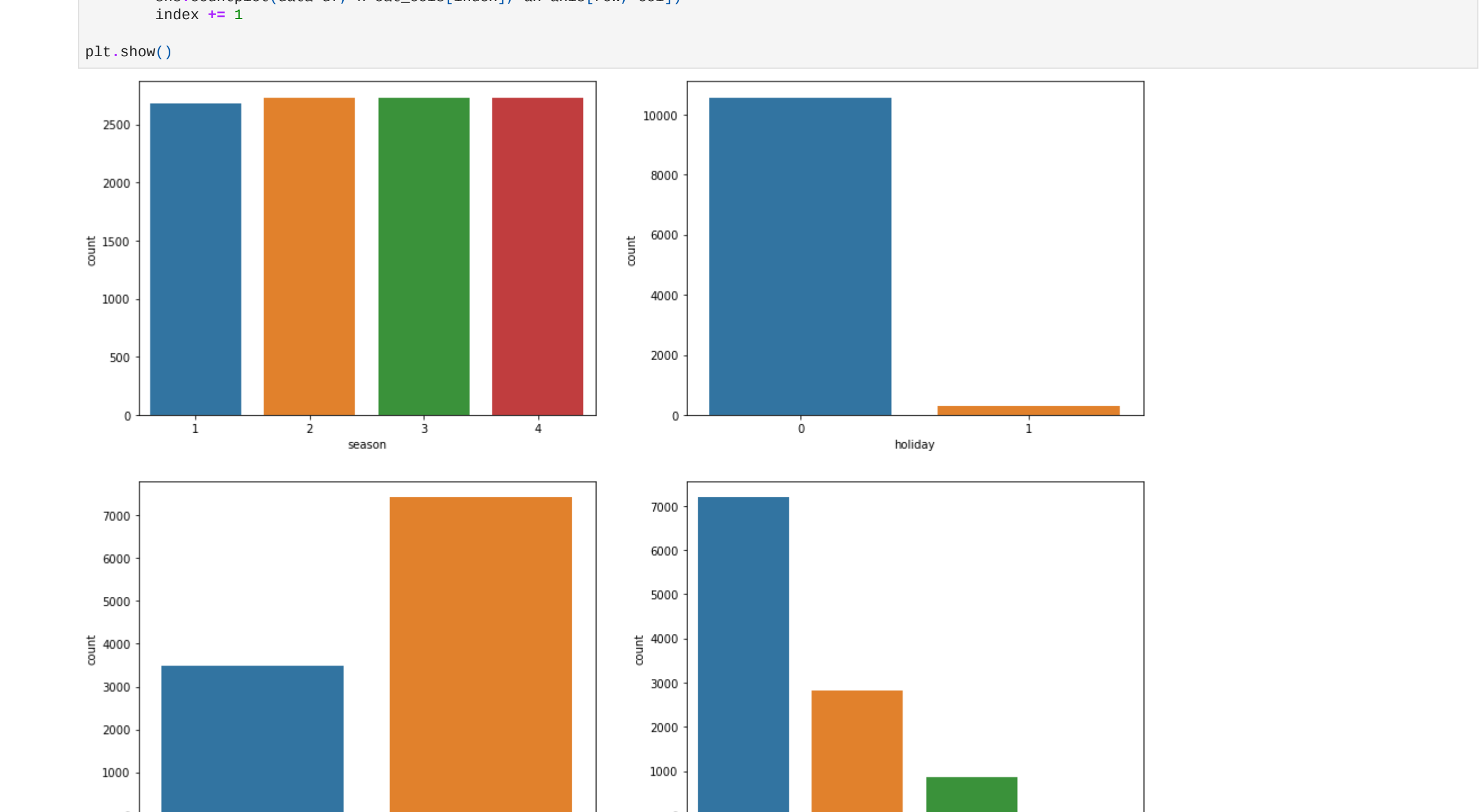


Looks like humidity, casual, registered and count have outliers in the data.

```
In [14]: # countplot of each categorical column
fig, axis = plt.subplots(nrows=2, ncols=2, figsize=(16, 12))

index = 0
for row in range(2):
    for col in range(2):
        sns.countplot(data=df, x=cat_cols[index], ax=axis[row, col])
        expected_values = df[cat_cols[index]].value_counts()
        index += 1

plt.show()
```



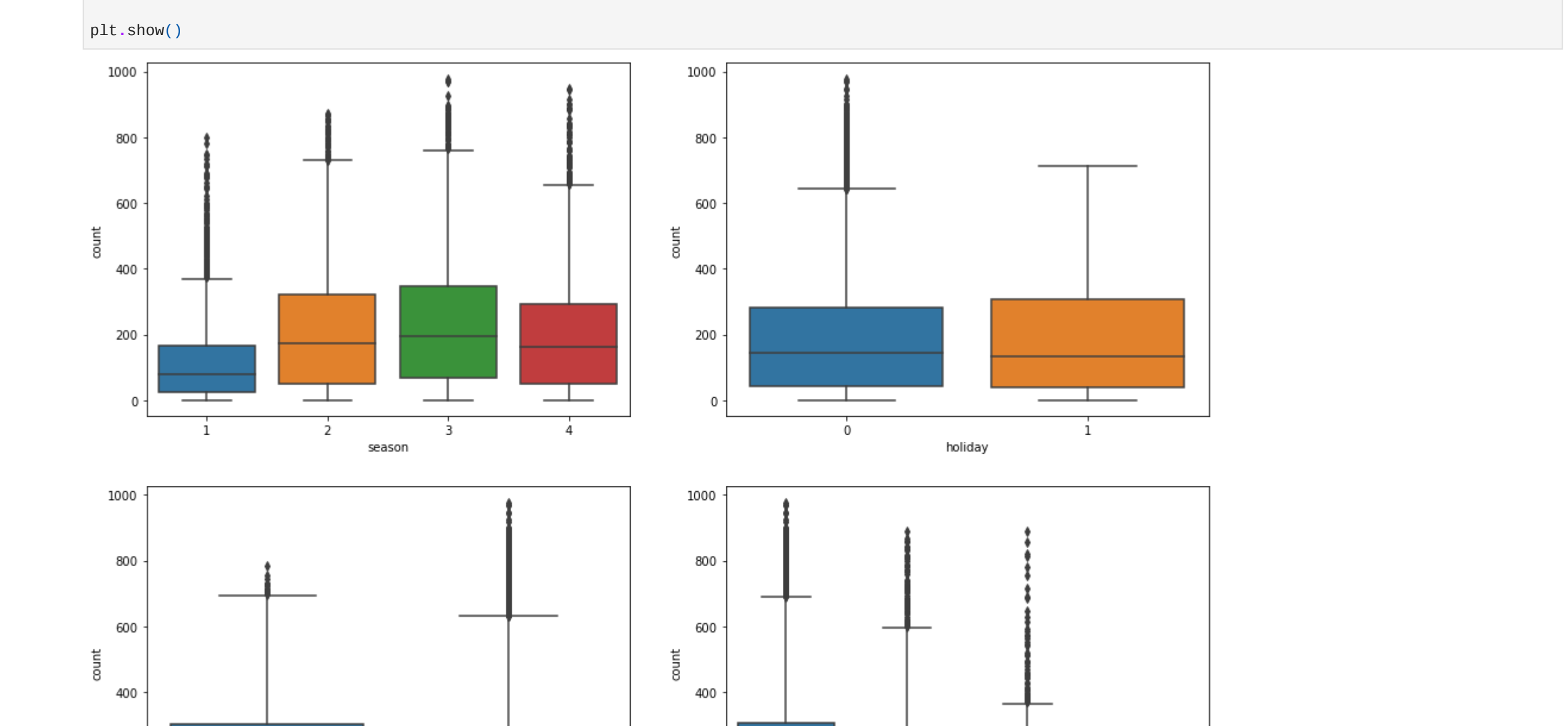
Data looks common as it should be like equal number of days in each season, more working days and weather is mostly Clear, Few clouds, partly cloudy, partly cloudy.

Bi-variate Analysis

```
In [15]: # plotting categorical variables against count using boxplots
fig, axis = plt.subplots(nrows=2, ncols=2, figsize=(16, 12))

index = 0
for row in range(2):
    for col in range(2):
        sns.boxplot(data=df, x=cat_cols[index], y='count', ax=axis[row, col])
        index += 1

plt.show()
```

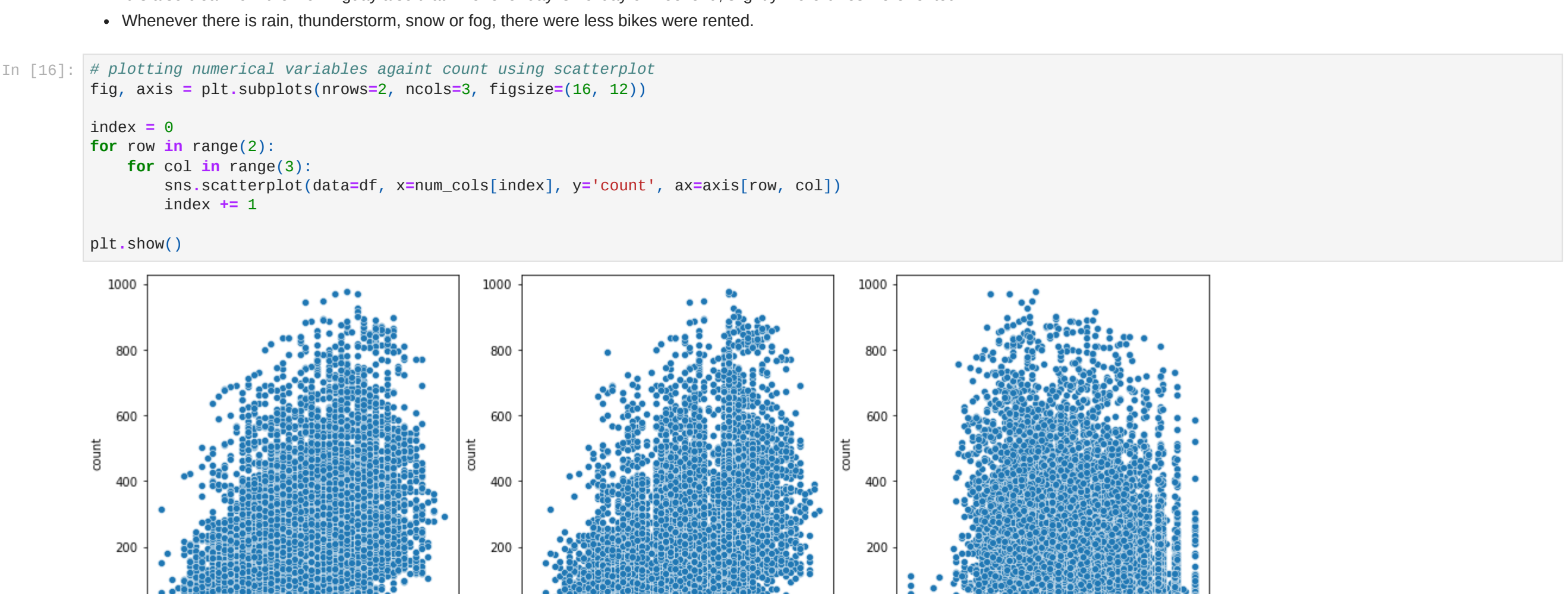


- In summer and fall seasons more bikes are rented as compared to other seasons.
- Whenever it is a holiday more bikes are rented.
- It is also clear from the workingday also that whenever day is holiday or weekend, slightly more bikes were rented.
- Whenever there is rain, thunderstorm, snow or fog, there were less bikes were rented.

```
In [16]: # plotting numerical variables against count using scatterplot
fig, axis = plt.subplots(nrows=2, ncols=3, figsize=(16, 12))

index = 0
for row in range(2):
    for col in range(3):
        sns.scatterplot(data=df, x=num_cols[index], y='count', ax=axis[row, col])
        index += 1

plt.show()
```

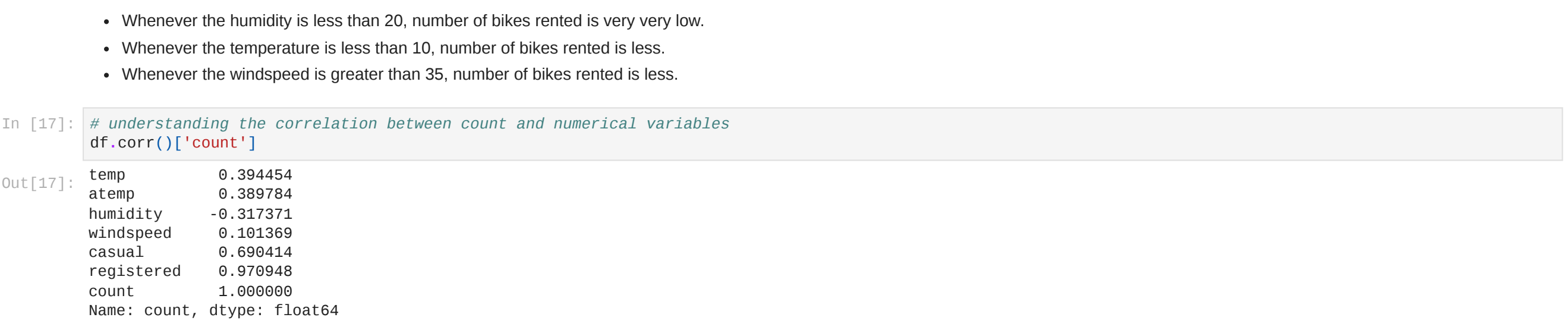


- Whenever the humidity is less than 20, number of bikes rented is very very low.
- Whenever the temperature is less than 10, number of bikes rented is less.
- Whenever the windspeed is greater than 35, number of bikes rented is less.

```
In [17]: # understanding the correlation between count and numerical variables
df.corr()['count']

Out [17]:
temp      0.394454
atemp     0.389784
humidity  -0.317374
windspeed  0.191369
casual     0.698414
registered  0.578948
count      1.000000
Name: count, dtype: float64
```

```
In [18]: sns.heatmap(df.corr(), annot=True)
plt.show()
```



Hypothesis Testing - 1

Null Hypothesis (H0): - Weather is independent of the season
Alternate Hypothesis (H1): Weather is not independent of the season
Significance level (alpha): 0.05

We will use chi-square test to test hypothesis defined above.

```
In [19]: data_table = pd.crosstab(df['season'], df['weather'])
print("Observed values:")
data_table

Observed values:
season  1  2  3  4
weather
1  1759  715  221  1
2  1801  708  224  0
3  1930  604  199  0
4  1702  807  225  0

In [20]: val = stats.chi2_contingency(data_table)
expected_values = val[3]
```

```
array([1.77454630e+03, 6.99258130e+02, 2.11948742e+02, 2.46738931e+01,
       1.80559765e+03, 7.11493845e+02, 2.15657450e+02, 2.51056403e+01,
       1.80559765e+03, 7.11493845e+02, 2.15657450e+02, 2.51056403e+01,
       1.80559765e+03, 7.11493845e+02, 2.15657450e+02, 2.51056403e+01])

In [22]: rows, ncols = 4, 4
dof = (rows-1)*(ncols-1)
print("degrees of freedom: ", dof)
alpha = 0.05

chi_sqr = sum([(o-e)**2/e for o, e in zip(data_table.values, expected_values)])
p_val = 1-stats.chi2.cdf(x=chi_sqr_statistic, df=dof)
print(f"P-value: {p_val}")

if p_val <= alpha:
    print("\nSince p-value is less than the alpha 0.05, We reject the Null Hypothesis. Meaning that\
Weather is dependent on the season.")
else:
    print("Since p-value is greater than the alpha 0.05, We do not reject the Null Hypothesis")
```

degrees of freedom: 9
chi-square test statistic: 44.09441248632364
critical value: 16.91897704629448
p-value: 1.356060157971317e-06

Since p-value is less than the alpha 0.05, We reject the Null Hypothesis. Meaning that Weather is dependent on the season.

Hypothesis Testing - 2

Null Hypothesis: Working day has no effect on the number of cycles being rented.
Alternate Hypothesis: Working day has effect on the number of cycles being rented.

Significance level (alpha): 0.05

We will use the 2-Sample T-Test to test the hypothesis defined above

```
In [23]: data_group1 = df[df['workingday']==0]['count'].values
data_group2 = df[df['workingday']==1]['count'].values

np.var(data_group1), np.var(data_group2)

Out [23]:
(38171.346098942427, 34840.69710674086)
```

Before conducting the two-sample T-test we need to find if the given data groups have the same variance. If the ratio of the larger data groups to the small data group is less than 4:1 then we can consider that the given data groups have equal variance.

Here, the ratio is 34040.70 / 37171.35 which is less than 4:1

```
In [24]: stats.ttest_ind(a=data_group1, b=data_group2, equal_var=True)

Out [24]:
Ttest_indResult(statistic=-1.209627376826694, pvalue=0.22644804226361348)
```

Since p-value is greater than 0.05 so we can not reject the Null hypothesis. We don't have the sufficient evidence to say that working day has effect on the number of cycles being rented.

Hypothesis Testing - 3

Null Hypothesis: Number of cycles rented is similar in different weather and season.
Alternate Hypothesis: Number of cycles rented is not similar in different weather and season.

Significance level (alpha): 0.05

Here, we will use the ANOVA to test the hypothesis defined above

```
In [25]: # defining the data groups for the ANOVA
gp1 = df[df['weather']==1]['count'].values
gp2 = df[df['weather']==2]['count'].values
gp3 = df[df['weather']==3]['count'].values
gp4 = df[df['weather']==4]['count'].values

gp5 = df[df['season']==1]['count'].values
gp6 = df[df['season']==2]['count'].values
gp7 = df[df['season']==3]['count'].values
gp8 = df[df['season']==4]['count'].values

# conduct the one-way anova
stats.f_oneway(gp1, gp2, gp3, gp4, gp5, gp6, gp7, gp8)

Out [25]:
F_onewayResult(statistic=127.96661249562491, pvalue=2.8674771742434642e-185)
```

Since p-value is less than 0.05, we reject the null hypothesis. This implies that Number of cycles rented is not similar in different weather and season conditions

Insights

- In summer and fall seasons more bikes are rented as compared to other seasons.
- Whenever it is a holiday more bikes are rented.
- It is also clear from the workingday also that whenever day is holiday or weekend, slightly more bikes were rented.
- Whenever there is rain, thunderstorm, snow or fog, there were less bikes were rented.
- Whenever the humidity is less than 20, number of bikes rented is very very low.
- Whenever the temperature is less than 10, number of bikes rented is less.
- Whenever the windspeed is greater than 35, number of bikes rented is less.

Recommendations

- In summer and fall seasons the company should have more bikes in stock to be rented. Because the demand in these seasons is higher as compared to other seasons.
- With a significance level of 0.05, workingday has no effect on the number of bikes being rented.
- In very low humid days, company should have less bikes in the stock to be rented.
- Whenever temperature is less than 10 or in very cold days, company should have less bikes.
- Whenever the windspeed is greater than 35 or in thunderstorms, company should have less bikes in stock to be rented.

```
In [ ]:
```