

```
'''
```

Business Case: Yulu - Hypothesis Testing

About Yulu: 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\*.

```
'''
```

```
\nBusiness Case: Yulu - Hypothesis Testing\n\nAbout Yulu: Yulu is India's leading *micro-mobility service provider, which offers un
in India, Yulu provides the safest commute solution through a\n user-friendly mobile app to enable shared, solo and sustainable comm
bus stands, office spaces,\n residential areas, corporate offices, etc*) to make those first and last miles smooth, affordable, and
d a consulting company to understand\n the factors on which the demand for these shared electric cycles depends. Specifically, they
ndian market*.\n\n'
```

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
import datetime as dt
```

```
from scipy import stats
from scipy.stats import f_oneway, ttest_ind, shapiro
from statsmodels.graphics.gofplots import qqplot
```

```
df_yulu = pd.read_csv("/content/Business Case -Yulu - Hypothesis Testing.csv_1642089089.txt")
```

```
df_yulu.shape
```

```
(10886, 12)
```

```
df_yulu.columns
```

```
Index(['datetime', 'season', 'holiday', 'workingday', 'weather', 'temp',
       'atemp', 'humidity', 'windspeed', 'casual', 'registered', 'count'],
      dtype='object')
```

```
df_yulu.head()
```

	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

```
df_yulu.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 - to datetime
season - to categorical
holiday - to categorical
workingday - to categorical
weather - to categorical
'''

'\nDatatype of following attributes needs to be changed to proper data type\n\ndatetime - to datetime\nseason - to categorical\nholiday

df_yulu['datetime'] = pd.to_datetime(df_yulu['datetime'])

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

#Converting the datatype of datetime column from object to datetime

df_yulu['datetime'] = pd.to_datetime(df_yulu['datetime'])

df_yulu['datetime'].min()

Timestamp('2011-01-01 00:00:00')

df_yulu['datetime'].max()

Timestamp('2012-12-19 23:00:00')

df_yulu['datetime'].max() - df_yulu['datetime'].min()

Timedelta('718 days 23:00:00')

df_yulu.iloc[:, 1:].describe(include='all')

```

	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	
<b>count</b>	10886.0	10886.0	10886.0	10886.0	10886.00000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.
<b>unique</b>	4.0	2.0	2.0	4.0	NaN	NaN	NaN	NaN	NaN	NaN	
<b>top</b>	4.0	0.0	1.0	1.0	NaN	NaN	NaN	NaN	NaN	NaN	
<b>freq</b>	2734.0	10575.0	7412.0	7192.0	NaN	NaN	NaN	NaN	NaN	NaN	
<b>mean</b>	NaN	NaN	NaN	NaN	20.23086	23.655084	61.886460	12.799395	36.021955	155.552177	191.
<b>std</b>	NaN	NaN	NaN	NaN	7.79159	8.474601	19.245033	8.164537	49.960477	151.039033	181.
<b>min</b>	NaN	NaN	NaN	NaN	0.82000	0.760000	0.000000	0.000000	0.000000	0.000000	1.
<b>25%</b>	NaN	NaN	NaN	NaN	13.94000	16.665000	47.000000	7.001500	4.000000	36.000000	42.
<b>50%</b>	NaN	NaN	NaN	NaN	20.50000	24.240000	62.000000	12.998000	17.000000	118.000000	145.
<b>75%</b>	NaN	NaN	NaN	NaN	26.24000	31.060000	77.000000	16.997900	49.000000	222.000000	284.
<b>max</b>	NaN	NaN	NaN	NaN	41.00000	45.455000	100.000000	56.996900	367.000000	886.000000	977.

```

# *There are no missing values in the dataset
# *casual and registered attributes might have outliers because their mean and median are
# very far away from 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.

# detecting missing values in the dataset
df_yulu.isnull().sum()

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.
```

### Univariate Analysis

```
# number of unique values in each categorical columns
df_yulu[cat_cols].melt().groupby(['variable', 'value'])['value'].count()
```

		value	
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	

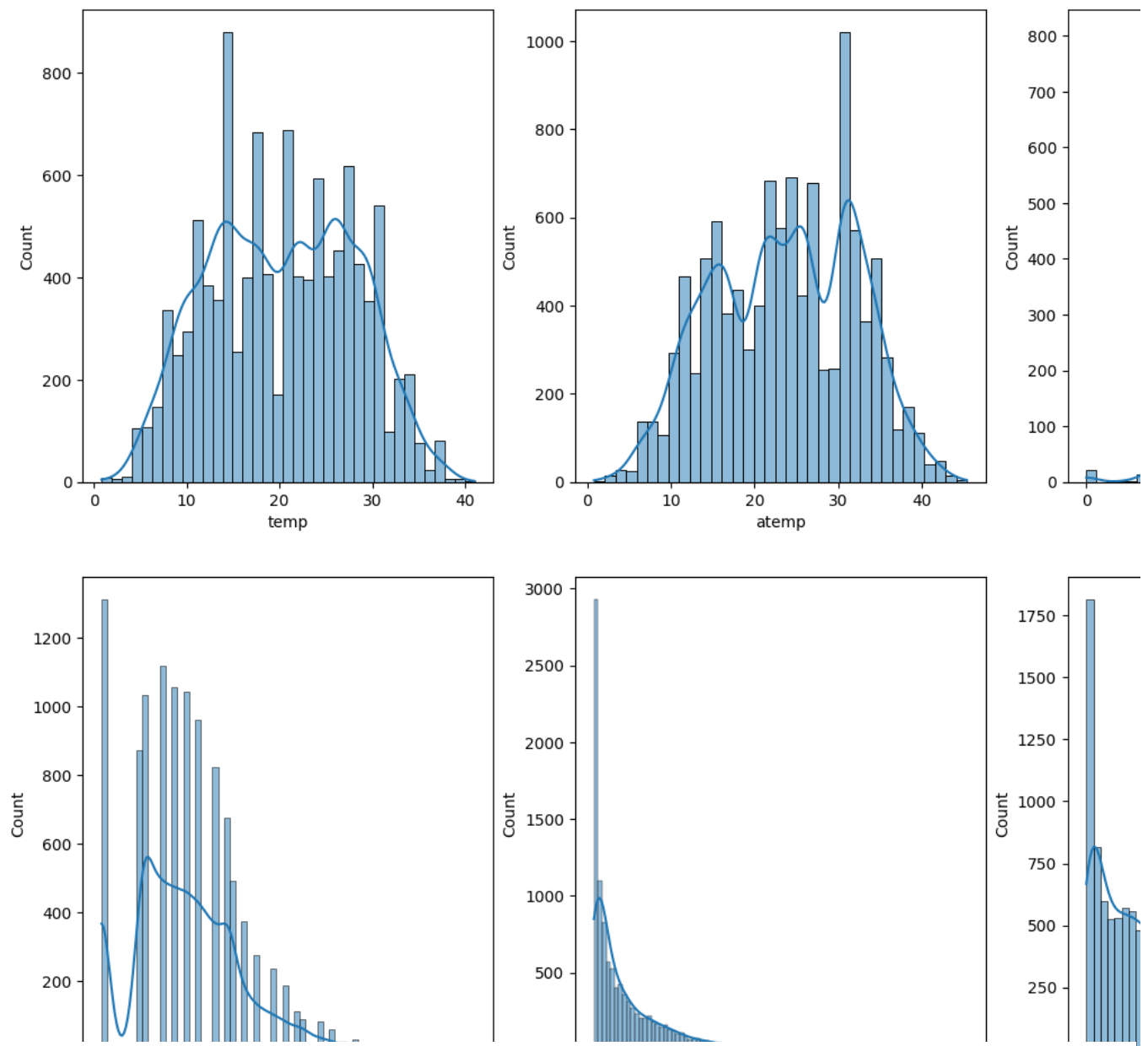
```
# understanding the distribution for numerical variables
```

```
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_yulu[num_cols[index]], ax=axis[row, col], kde=True)
        index += 1
```

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



#1: casual, registered and count somewhat looks like Log Normal Distribution  
 #2: temp, atemp and humidity looks like they follows the Normal Distribution  
 #3: windspeed follows the binomial distribution

```

40000  _
|

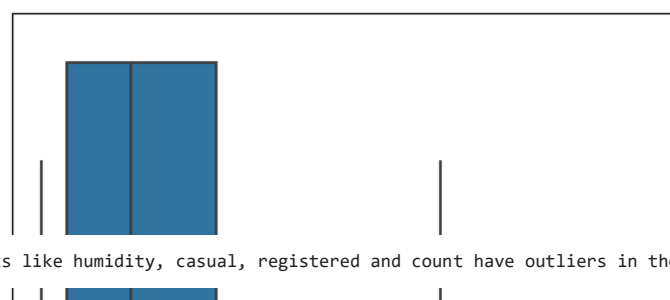
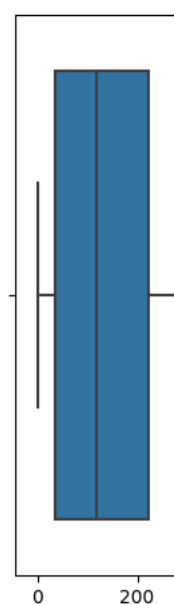
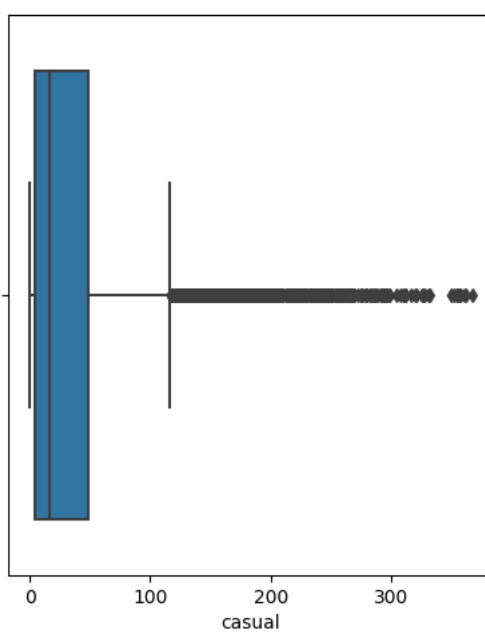
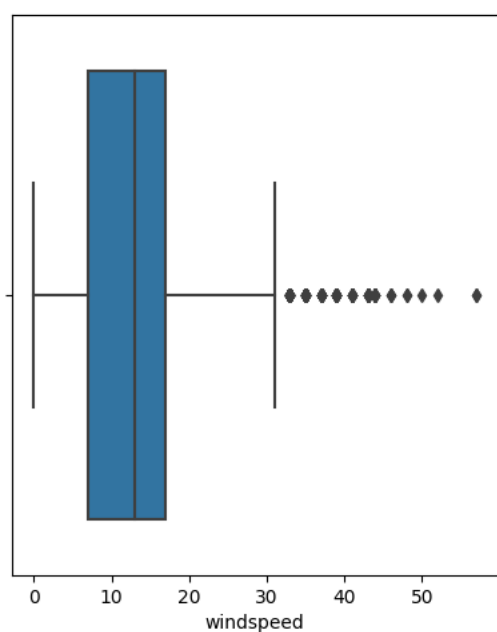
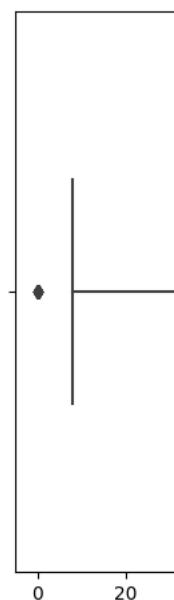
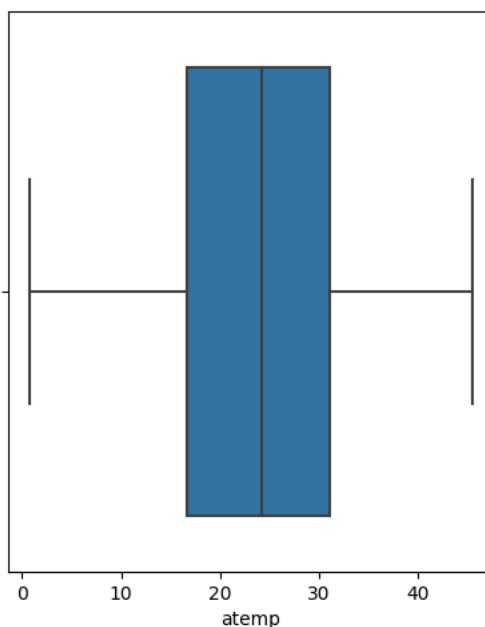
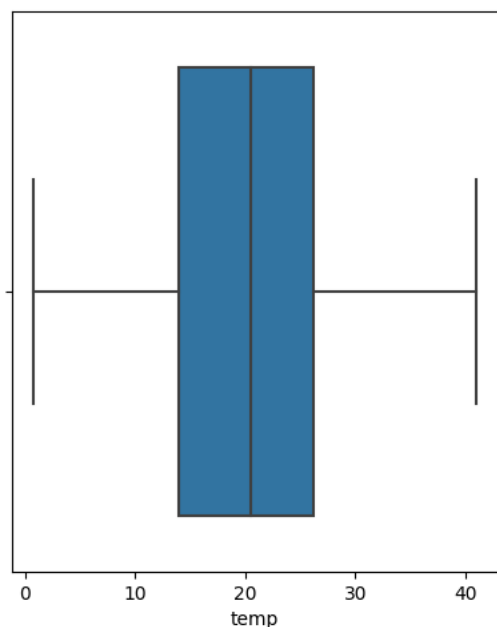
#plotting boxplots to detect outliers in the data

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_yulu[num_cols[index]], ax=axis[row,col])
        index +=1
plt.show()
sns.boxplot(x=df_yulu[num_cols[-1]])
plt.show()

```



# looks like humidity, casual, registered and count have outliers in the data

#countplot of each categorical column

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

index = 0

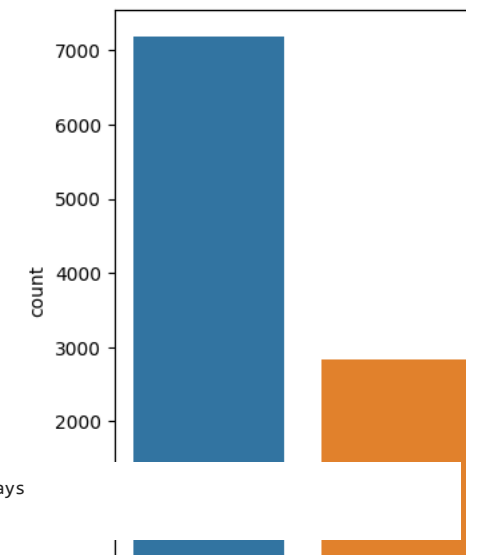
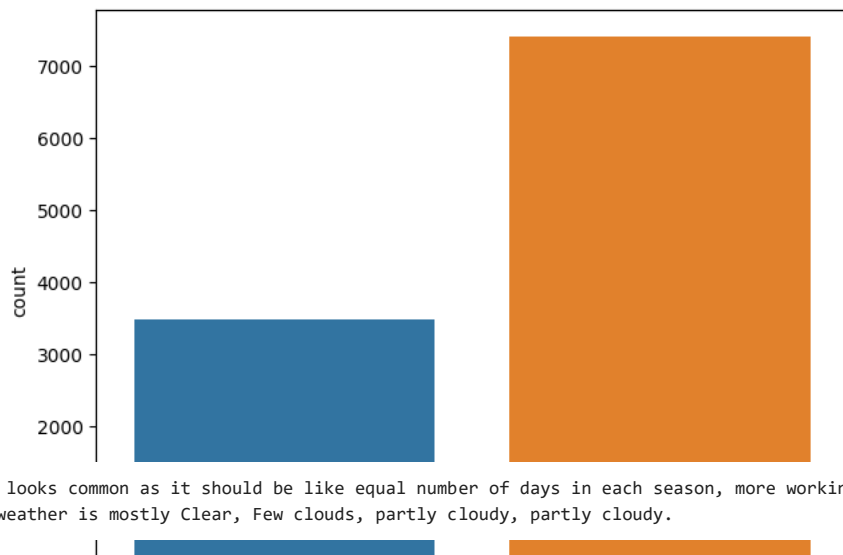
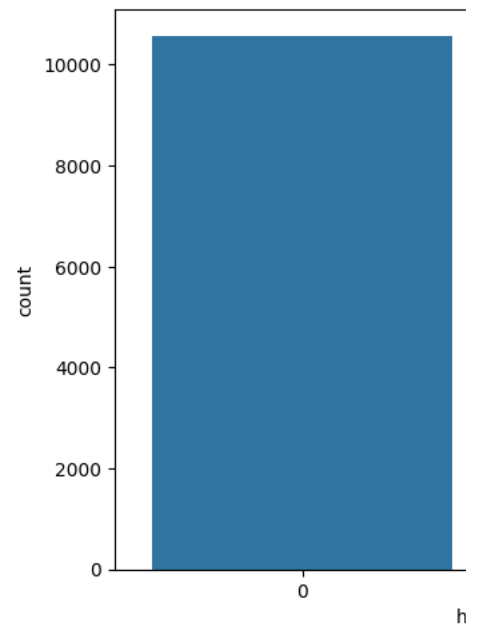
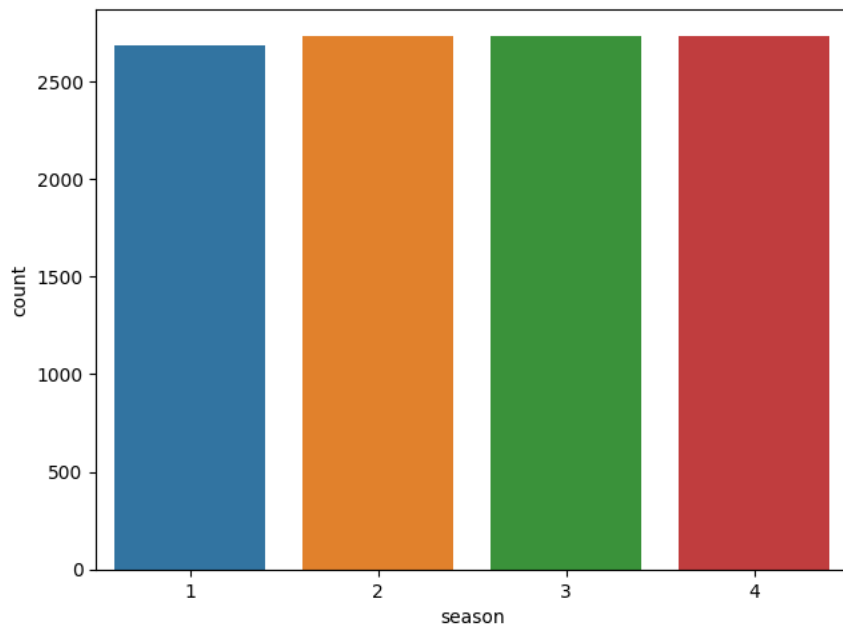
for row in range(2):

for col in range(2):

sns.countplot(data=df\_yulu , x= cat\_cols[index], ax=axis[row,col])

index += 1

plt.show()

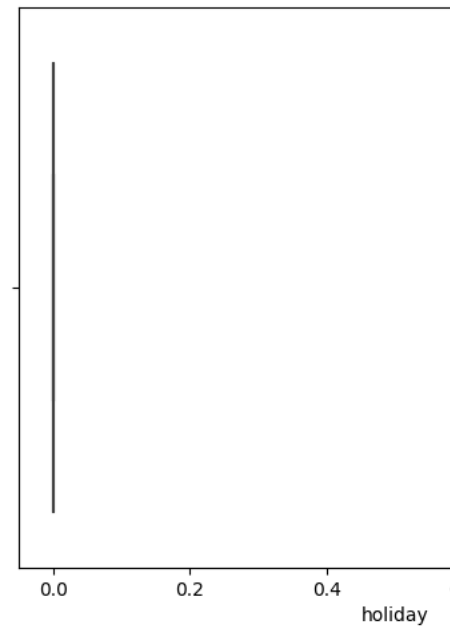
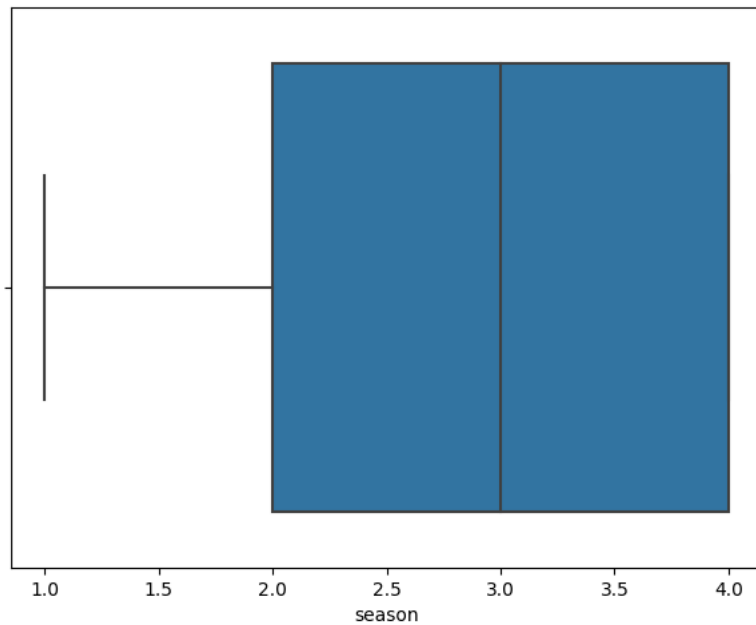


#Data looks common as it should be like equal number of days in each season, more working days  
#and weather is mostly Clear, Few clouds, partly cloudy, partly cloudy.

## ▼ Bi- variate analysis

```
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_yulu , x= cat_cols[index], ax=axis[row,col])
        index += 1
plt.show()
```



```
...
```

- In summer and fall seasons more bikes are rented as compared to other seasons.
- Whenever its 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

```
...
```

```
'\n• In summer and fall seasons more bikes are rented as compared to other seasons.\n• Whenever its a holiday more bikes are rented.\nmore bikes were rented.\n• Whenever there is rain, thunderstorm, snow or fog, there were less bikes were rented\n'
```



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

```
index = 0
```

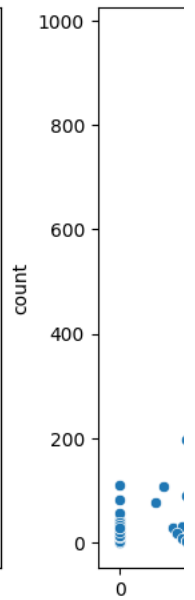
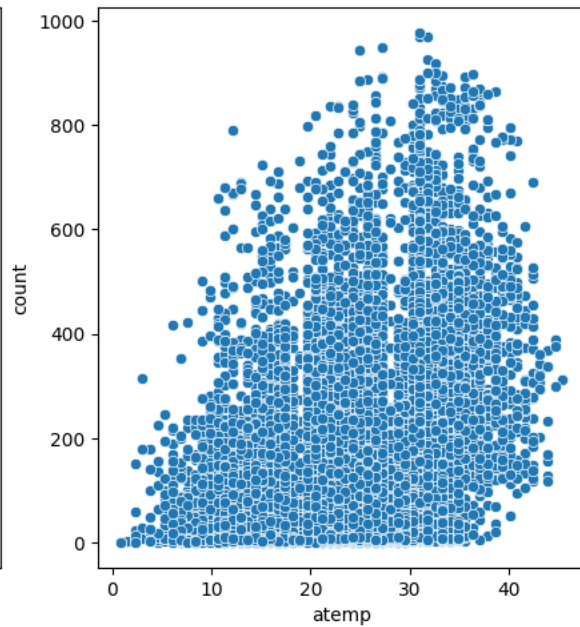
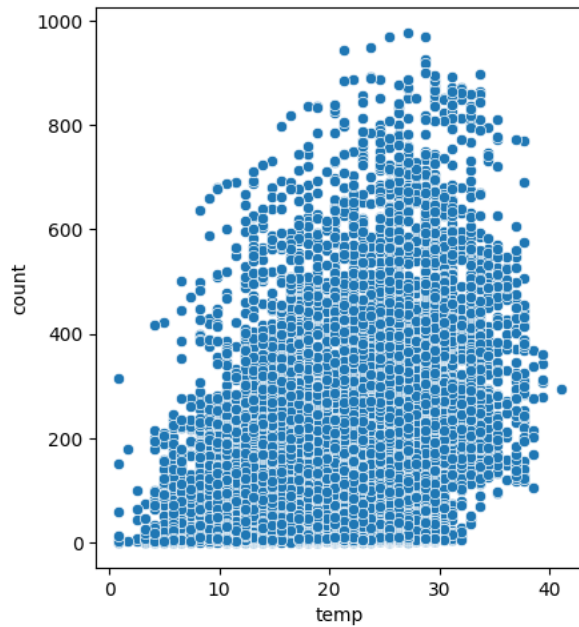
```
for row in range(2):
```

```
    for col in range(3):
```

```
        sns.scatterplot(data=df_yulu , 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.

```
'''
```

```
'\n• Whenever the humidity is less than 20, number of bikes rented is very very low.\n• Whenever the temperature is less than 10, number of bikes rented is less.\n'
```

```
600 | 600 | 600 |
```

```
df_yulu.corr()['count']
```

```
<ipython-input-37-ba9ced1f72e6>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, this will raise an error.
df_yulu.corr()['count']
```

```
temp      0.394454
atemp     0.389784
humidity  -0.317371
windspeed  0.101369
casual     0.690414
registered 0.970948
count     1.000000
Name: count, dtype: float64
```



```
sns.heatmap(df_yulu.corr(), annot=True)
plt.show()
```



```
<ipython-input-38-9ebb3d142672>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future ver
enc heatmap(df_yulu.corr(), annot=True)
```

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'



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

Observed values:



data\_table

weather	1	2	3	4
season				
1	1759	715	211	1
2	1801	708	224	0
3	1930	604	199	0
4	1702	807	225	0

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

```
array([[1.77454639e+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],
       [1.80625831e+03, 7.11754180e+02, 2.15736359e+02, 2.51148264e-01]])
```

```
nrows, ncols = 4, 4
dof = (nrows-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)])
chi_sqr_statistic = chi_sqr[0] + chi_sqr[1]
print("chi-square test statistic: ", chi_sqr_statistic)
```

```
print('f critical value: {critical_val}')
p_val = 1-stats.chi2.cdf(x=chi_sqr_statistic, df_yulu=dof)
```

```
critical_val = stats.chi2.ppf(q=1-alpha, df_yulu=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
f critical value: {critical_val}
```

```
-----
TypeError                                Traceback (most recent call last)
<ipython-input-57-0bd012b42824> in <cell line: 10>()
8
9 print('f critical value: {critical_val}')
```

```
----> 10 p_val = 1-stats.chi2.cdf(x=chi_sqr_statistic, df_yulu=dof)
11
12 critical_val = stats.chi2.ppf(q=1-alpha, df_yulu=dof)
```

```
/usr/local/lib/python3.10/dist-packages/scipy/stats/_distn_infrastructure.py in cdf(self, x, *args, **kws)
2059
2060 """
-> 2061 args, loc, scale = self._parse_args(*args, **kws)
2062 x, loc, scale = map(asarray, (x, loc, scale))
2063 args = tuple(map(asarray, args))
```

```
TypeError: _parse_args() got an unexpected keyword argument 'df_yulu'
```

SEARCH STACK OVERFLOW

```

nrows, ncols = 4, 4
dof = (nrows-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)])
chi_sqr_statistic = chi_sqr[0] + chi_sqr[1]
print("chi-square test statistic: ", chi_sqr_statistic)

critical_val = stats.chi2.ppf(q=1-alpha, df_yulu=dof)
print(f"critical value: {critical_val}")

p_val = 1-stats.chi2.cdf(x=chi_sqr_statistic, df_yulu=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
-----
TypeError                                 Traceback (most recent call last)
<ipython-input-52-e91238ae576c> in <cell line: 11>()
      9 print("chi-square test statistic: ", chi_sqr_statistic)
     10
----> 11 critical_val = stats.chi2.ppf(q=1-alpha, df_yulu=dof)
     12 print(f"critical value: {critical_val}")
     13

/usr/local/lib/python3.10/dist-packages/scipy/stats/_distn_infrastructure.py in ppf(self, q, *args, **kwargs)
    2228
    2229     """
-> 2230     args, loc, scale = self._parse_args(*args, **kwargs)
    2231     q, loc, scale = map(asarray, (q, loc, scale))
    2232     args = tuple(map(asarray, args))

TypeError: _parse_args() got an unexpected keyword argument 'df_yulu'

```

SEARCH STACK OVERFLOW

## 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

...

'\nNull Hypothesis: Working day has no effect on the number of cycles being rented.\n\nAlternate Hypothesis: Working day has effect
ample T-Test to test the hypothesis defined above\n\n'

data_group1 = df_yulu[df_yulu['workingday']==0]['count'].values
data_group2 = df_yulu[df_yulu['workingday']==1]['count'].values

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

(30171.346098942427, 34040.69710674686)

...
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

'\nBefore conducting the two-sample T-Test we need to find if the given data groups have the same variance. If the ratio of the larg
groups have equal variance.'

...

Here, the ratio is 34040.70 / 30171.35 which is less than 4:1
...

```

```
'\n\nHere, the ratio is 34040.70 / 30171.35 which is less than 4:1\n'
```

```
stats.ttest_ind(a=data_group1, b=data_group2, equal_var=True)
```

```
TtestResult(statistic=-1.2096277376026694, pvalue=0.22644804226361348, df=10884.0)
```

Since pvalue 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 hypothess defined above
```

```
'''
```

```
'\n\nNull Hypothesis: Number of cycles rented is similar in different weather and season.\n\nAlternate Hypothesis: Number of cycles r
ented is not similar in different weather and season.\n\nSignificance level (alpha): 0.05\n\nHere, we will use the ANOVA to test th
e hypothess defined above\n\n'
```

```
#defining the data groups for the ANOVA
```

```
gp1 = df_yulu[df_yulu['weather']==1]['count'].values
gp2 = df_yulu[df_yulu['weather']==2]['count'].values
gp3 = df_yulu[df_yulu['weather']==3]['count'].values
gp4 = df_yulu[df_yulu['weather']==4]['count'].values
```

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

```
# conduct the one-way anova
```

```
stats.f_oneway(gp1, gp2, gp3, gp4, gp5, gp6, gp7, gp8)
```

```
F_onewayResult(statistic=127.96661249562491, pvalue=2.8074771742434642e-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 wea
```

```
'''
```

```
'\n)\nSince p-value is less than 0.05, we reject the null hypothesis. This implies that Number of cycles rented is not similar in di
```

```
'''
```

### Insights

In summer and fall seasons more bikes are rented as compared to other seasons.

Whenever its 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 com

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 temprature 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.

```
'''
```

```
'\nInsights\n\nIn summer and fall seasons more bikes are rented as compared to other seasons.\nWhenever its a holiday more bikes are ghtly more bikes were rented.\nWhenever there is rain, thunderstorm, snow or fog, there were less bikes were rented.\nWhenever the h  
ess than 10 number of bikes rented is less \nWhenever the windeneed is greater than 35 number of bikes rented is less \n\n\nPeromng
```