

1) PROBLEM STATEMENT

The objective of this case study is to help Yulu find which variables are significant in predicting the demand for shared electric cycles in the Indian market. How well these variables describe the electric cycle demands.

```
In [126]: import pandas as pd
import sys
from ydata.profitiling import ProfileReport
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import ttest_ind,ttest_samp,kstest,norm
from scipy.stats import chi2_contingency,chi2
from scipy.stats import pearsonr,f_oneway

In [13]: import warnings
warnings.filterwarnings("ignore")

In [2]: df=pd.read_csv("yulu_dataset.txt")

In [5]: df.head(5)

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]: df.shape
Out[6]: (18886, 12)

In [11]: df.describe()

Out[11]:
   season    holiday  workingday  weather  temp  atemp  humidity  windspeed  casual  registered  winds
count 10886.000000 10886.000000 10886.000000 10886.000000 10886.000000 10886.000000 10886.000000 10886.000000 10886.000000 10886.000000 10886.000000
mean 2.506614 0.028569 0.068875 0.680875 1.418427 20.23006 23.655084 61.886460 61.886460 12.79932
std 1.116174 0.166599 0.466159 0.63839 0.63839 7.79159 8.474601 19.245033 8.164531
min 1.000000 0.000000 0.000000 1.000000 0.82000 13.94000 16.66500 0.00000 0.00000 0.00000
25% 2.000000 0.000000 0.000000 1.000000 1.00000 18.00000 20.50000 47.00000 7.00500
50% 3.000000 0.000000 1.000000 1.000000 20.50000 24.24000 62.00000 12.99000
75% 4.000000 0.000000 1.000000 2.000000 26.24000 31.96000 77.00000 16.9979
max 4.000000 1.000000 1.000000 4.00000 41.0000 45.45000 100.00000 56.9969

In [12]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18886 entries, 0 to 18885
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 datetime 18886 non-null object
1 season 18886 non-null int64
2 holiday 18886 non-null int64
3 workingday 18886 non-null int64
4 weather 18886 non-null float64
5 temp 18886 non-null float64
6 atemp 18886 non-null float64
7 humidity 18886 non-null float64
8 windspeed 18886 non-null float64
9 casual 18886 non-null int64
10 registered 18886 non-null int64
11 count 18886 non-null int64
dtypes: float64(8), int64(8), object(1)
memory usage: 3020.7 KB

In [14]: df.duplicated().sum()
Out[14]: 0

In [15]: df.isnull().sum().sum()
Out[15]: 0

Observation
The dataset has neither duplicate rows nor has it any duplicate value.
```

```
In [3]: dfi=df.copy()

In [10]: dfi.head()

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

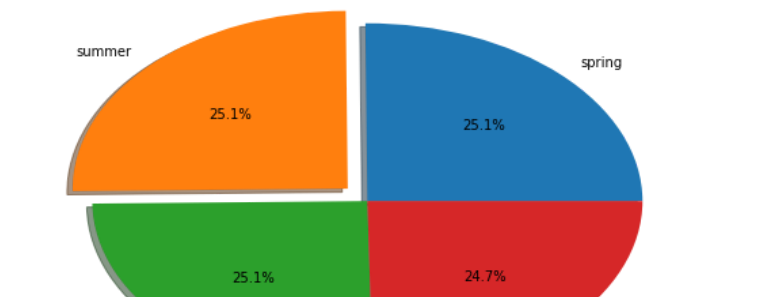
In [4]: dfi['datetime']=pd.to_datetime(dfi['datetime'])

In [50]: dfi[['datetime']].info()

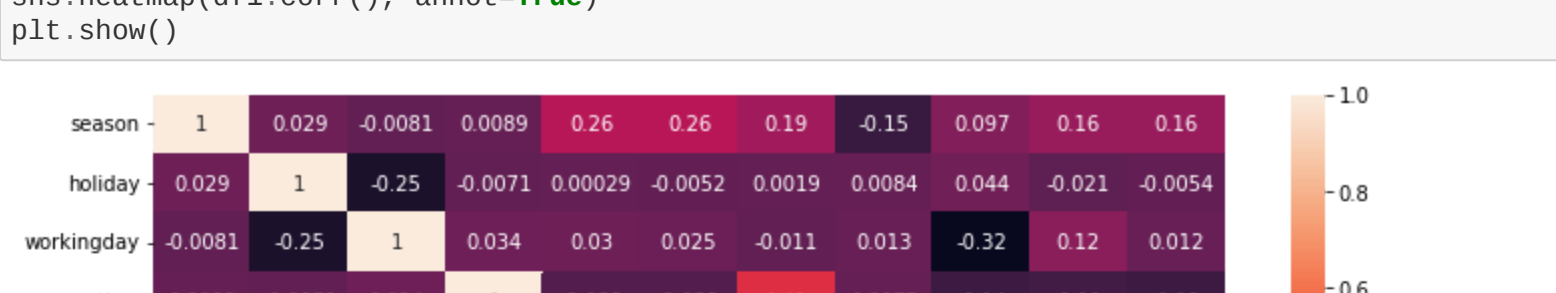
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18886 entries, 0 to 18885
Data columns (total 1 columns):
# Column Non-Null Count Dtype
---
0 datetime 18886 non-null datetime64[ns]
dtypes: datetime64[ns](1)
memory usage: 85.2 KB
```

Univariate analysis

```
In [63]: sns.histplot(dfi[['temp'],kde=True])
plt.show()
```



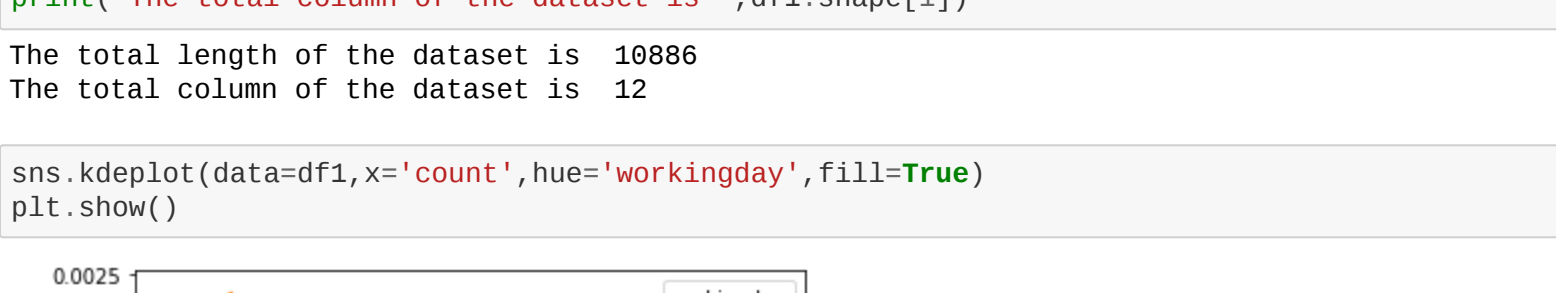
```
In [86]: fig, axis = plt.subplots(nrows=2, ncols=2, figsize=(20, 6))
total_seasons=dfi[season].value_counts()
sns.countplot(dfi[workingday],ax=axis[0])
plt.legend()
sns.countplot(dfi[holiday],ax=axis[1])
plt.show()
```



Observation

1. I.e working days/non weekend and non holidays) has more users of Yulu bikes .
2. I.e non holiday day has more yulu bikes users.
3. Users prefer temperature between 11 and 30 to use Yulu bikes.

```
In [41]: fig, axis = plt.subplots(nrows=1, ncols=2, figsize=(28, 6))
total_seasons=dfi[season].value_counts()
total_weather=dfi[weather].value_counts()
axis[0].pie(total_seasons,explode=explode,labels=label,autopct='%1.1f%%',shadow=True)
axis[1].pie(total_weather,explode=explode,labels=label2,autopct='%1.1f%%',shadow=True)
plt.show()
```

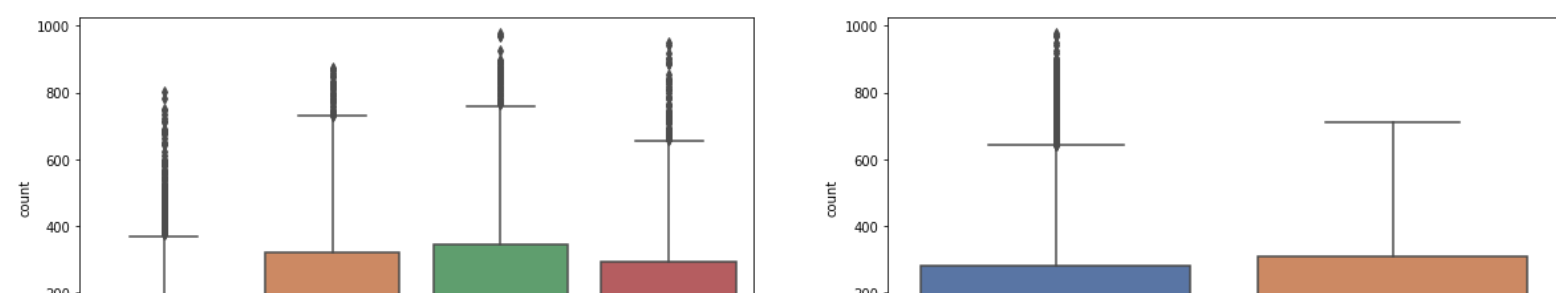


Observation

The dataset is almost distributed equally among seasons while most data has weather1. We can conclude people prefer Yulu irrespective of seasons but if the weather is weather_1 i.e Clear, Few clouds, partly cloudy, partly cloudy, then they prefer using YULU bikes.

Bivariate Analysis

```
In [44]: plt.figure(figsize=(12,6))
sns.heatmap(dfi.corr(),annot=True)
plt.show()
```



Observation

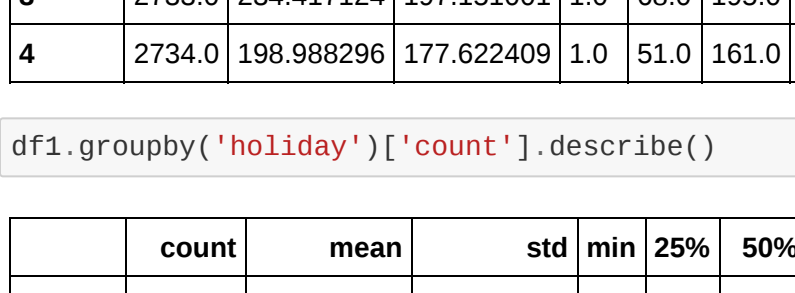
Here Correlation analysis can reveal meaningful relationships between different metrics or groups of metrics.

1. Temperature and feeling temperature are strongly positively correlated.
2. While season and weather or season and workingday has almost no correlation.

```
In [98]: print("The total length of the dataset is ",dfi.shape[0])
print("The total column of the dataset is ",dfi.shape[1])

The total length of the dataset is 18886
The total column of the dataset is 12

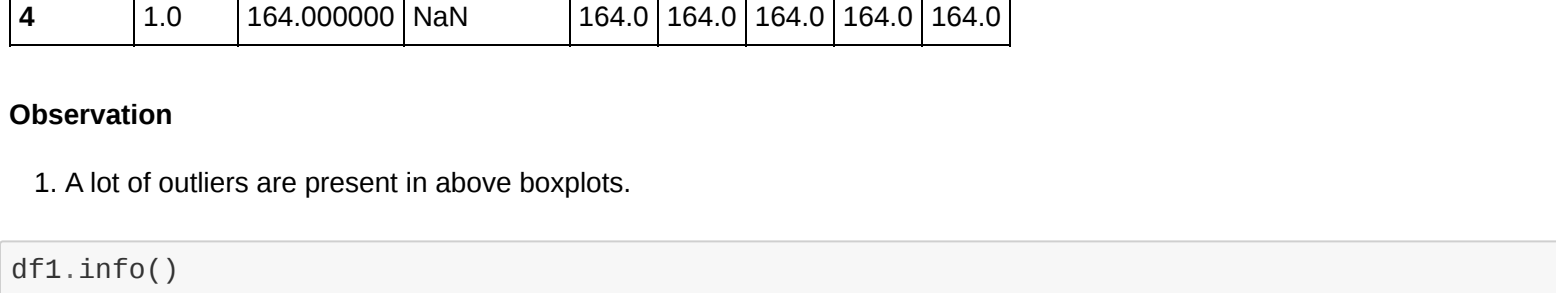
In [182]: sns.kdeplot(data=dfi,x='count',hue='workingday',fill=True)
plt.show()
```



Observation

1. Non holiday days i.e 0 have significantly higher users of Yulu bikes over days when its weekend or its a holiday.

```
In [111]: fig, axis = plt.subplots(nrows=2, ncols=2, figsize=(20, 6))
fig.subplots_adjust(top=1)
sns.boxplot(data=dfi,y='count',x='season',palette='deep',ax=axis[0,0])
sns.boxplot(data=dfi,y='count',x='holiday',palette='deep',ax=axis[0,1])
sns.boxplot(data=dfi,y='count',x='workingday',palette='deep',ax=axis[1,0])
sns.boxplot(data=dfi,y='count',x='weather',palette='deep',ax=axis[1,1])
plt.show()
```



```
In [113]: dfi.groupby("season")["count"].describe()

Out[113]:
season count mean std min 25% 50% 75% max
0 2686.0 116.343261 125.273974 1.0 24.0 78.0 164.0 801.0
1 2733.0 215.261372 192.007843 1.0 48.0 172.0 321.0 873.0
2 2734.0 234.417124 197.151001 1.0 68.0 195.0 347.0 977.0
4 2734.0 198.988296 177.622409 1.0 51.0 161.0 294.0 948.0
```

```
In [116]: dfi.groupby("holiday")["count"].describe()

Out[116]:
holiday count mean std min 25% 50% 75% max
0 10575.0 191.744655 181.513131 1.0 43.0 145.0 283.0 977.0
1 311.0 185.877814 169.306531 1.0 38.5 133.0 308.0 712.0
```

```
In [117]: dfi.groupby("workingday")["count"].describe()

Out[117]:
workingday count mean std min 25% 50% 75% max
0 3474.0 198.506621 173.724015 1.0 44.0 128.0 304.0 783.0
1 7432.0 193.011873 184.512659 1.0 41.0 151.0 277.0 977.0
```

```
In [118]: dfi.groupby("weather")["count"].describe()

Out[118]:
weather count mean std min 25% 50% 75% max
0 7192.0 205.236791 187.959566 1.0 48.0 161.0 305.0 977.0
1 2834.0 178.955540 168.366413 1.0 41.0 134.0 264.0 890.0
2 659.0 181.846533 136.581297 1.0 23.0 71.0 181.0 591.0
4 1.0 164.000000 NaN 164.0 164.0 164.0 164.0
```

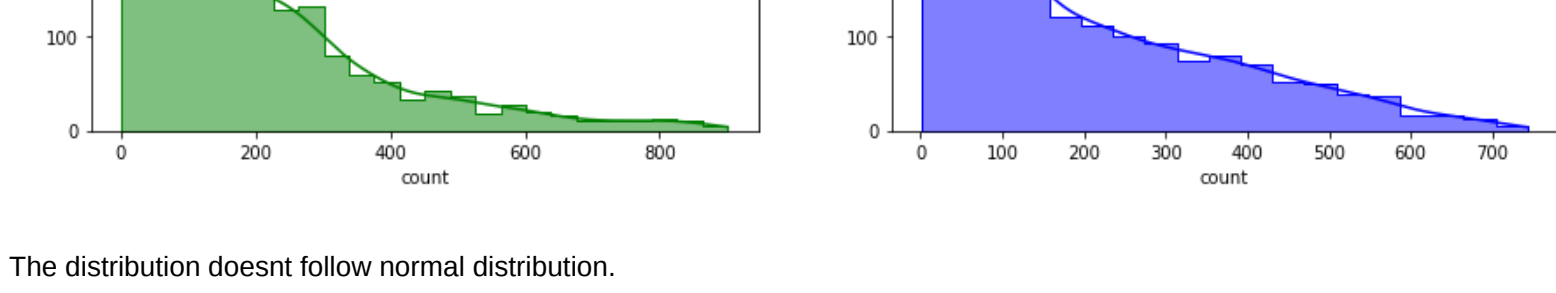
Observation

1. A lot of outliers are present in above boxplots.

```
In [65]: dfi.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18886 entries, 0 to 18885
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 datetime 18886 non-null datetime64[ns]
1 season 18886 non-null int64
2 holiday 18886 non-null int64
3 workingday 18886 non-null int64
4 weather 18886 non-null float64
5 temp 18886 non-null float64
6 atemp 18886 non-null float64
7 humidity 18886 non-null float64
8 windspeed 18886 non-null float64
9 casual 18886 non-null int64
10 registered 18886 non-null int64
11 count 18886 non-null int64
dtypes: datetime64[ns](1), float64(8), int64(8)
memory usage: 3020.7 KB

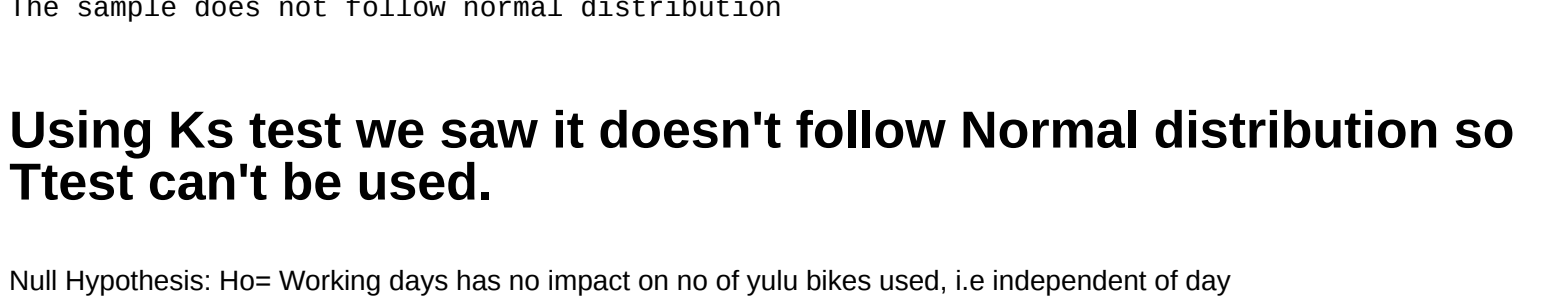
In [120]: sns.pairplot(dfi)
plt.show()
```



Hypothesis Testing

Q1. Whether working days has effect on no of yulu bikes users?

```
In [146]: plt.figure(figsize = (15, 5))
plt.subplot(1, 2, 1)
sns.histplot(dfi[['workingday'] == 1, 'count'].sample(2000),
             element = 'step', color = 'green', kde = True, label = 'workingday')
plt.legend()
plt.subplot(1, 2, 2)
sns.histplot(dfi[['workingday'] == 0, 'count'].sample(2000),
             element = 'step', color = 'blue', kde = True, label = 'non_workingday')
plt.legend()
plt.plot()
```



The distribution doesn't follow normal distribution.

Check for normality.

Ks test.

Ho: Gaussian distribution. Ha: Not a gaussian distribution.

```
In [28]: nonwork=dfi[dfi['workingday']==0]['count']
nonwork=nonwork.mean()/nonwork.std()
x_stat,t,p_value=stats.ttest_1samp(nonwork,norm.cdf)
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 4.7244827562341644e-69
The sample does not follow normal distribution
```

```
In [30]: work=dfi[dfi['workingday']==1]['count']
work=work.mean()/work.std()
x_stat,t,p_value=stats.ttest_1samp(work,norm.cdf)
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 3.7696531505378e-144
The sample does not follow normal distribution
```

Using Ks test we saw it doesn't follow Normal distribution so Test can't be used.

Null Hypothesis: Ho: Working days has no impact on no of yulu bikes used, i.e independent of day

Alternate Hypothesis: Ha: Working days has impact on no of yulu bikes users.

Mathematically: Ho: $\mu_1 = \mu_2$

Ha: $\mu_1 \neq \mu_2$

We will take 95% confidence, i.e. $\alpha = 0.05$

```
In [41]: temp=pd.crosstab([columns=dfi['count'],index=dfi[workingday]])
chi_stat,p_value,d_f,exp_freq=chi2_contingency(temp)
print('p_value =',p_value)
if p_value < 0.05:
    print('Mean no. of electric cycles rented is same for working and non-working days')
else:
    print('Mean no. of electric cycles rented is not same for working and non-working days')

p_value 1.2403697441548e-05
Mean no. of electric cycles rented is same for working and non-working days
```

Alternative method

```
In [80]: work and nonwork
print("The variance of working day is :",np.var(work))
print("The variance of non-working day is :",np.var(nonwork))

The variance of working day is : 34949.69710674693
The variance of non-working day is : 38173.34668984243
```

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 smaller data group is less than 4:1 then we can consider that the given data groups have equal variance.

Here, the ratio is 34949.70 / 38173.35 which is less than 4:1

```
In [132]: f_stat,p_value=ttest_ind(work, nonwork, equal_var=True)
print('p_value =',p_value)
if p_value < 0.05:
    print('Mean no. of electric cycles rented is 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 0.2264894226351348
Mean no. of electric cycles rented is same for working and non-working days
```

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.

ANOVA to check if No. of cycles rented is similar or different in

1. weather

2. season

```
In [43]: dfi.head()

Out[43]:
  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 [44]: dfi.groupby('season').value_counts()

Out[44]:
season count
fall 2733
spring 2686
summer 2734
winter 2734
Name: season, dtype: int64
```

```
In [51]: def season_rename(x):
    if x == 0:
        return 'spring'
    elif x == 1:
        return 'summer'
    elif x == 2:
        return 'fall'
    else:
        return 'winter'
dfi['season'] = dfi['season'].apply(season_rename)

In [52]: dfi.groupby(by = 'season')['count'].describe()

Out[52]:
season count mean std min 25% 50% 75% max
fall 2733 234.417124 197.151001 1.0 68.0 195.0 347.0 977.0
spring 2686 215.261372 192.007843 1.0 48.0 172.0 321.0 873.0
summer 2734 198.988296 177.622409 1.0 51.0 161.0 294.0 948.0
winter 2734 198.988296 177.622409 1.0 51.0 161.0 294.0 948.0
```

```
In [59]: spring_season = dfi[dfi['season'] == 'spring']['count']
summer_season = dfi[dfi['season'] == 'summer']['count']
fall_season = dfi[dfi['season'] == 'fall']['count']
winter_season = dfi[dfi['season'] == 'winter']['count']

In [82]: dfi.groupby('season').value_counts()

Out[82]:
season count
fall 2733
summer 2733
spring 2686
Name: season, dtype: int64
```

Null Hypothesis: Ho: Season has no impact on sales

Alternate Hypothesis: Ha: Season has impact on sales

Mathematically: Ho: $\mu_1 = \mu_2$ (i.e mean no of cycles is same for all seasons)

Ha: $\mu_1 \neq \mu_2$ (i.e mean no of cycles is different for all season)

We will take 95% confidence, i.e. $\alpha = 0.05$

So if $p\text{-value} > \alpha \Rightarrow$ Fail to reject Ho and if $p\text{-value} < \alpha \Rightarrow$ Reject Ho.

```
In [139]: f_stat,p_value,f_oneway(spring_season,summer_season,fall_season,winter_season)
print('p_value =',p_value)
if p_value < 0.05:
    print('Season has impact on sales of bicycle')
else:
    print('Season has no impact on sales of bicycle')

p_value 6.1648433864955e-149
Season has impact on sales of bicycle
```

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

```
In [64]: weather1 = dfi[dfi['weather'] == 1]['count']
weather2 = dfi[dfi['weather'] == 2]['count']
weather3 = dfi[dfi['weather'] == 3]['count']
weather4 = dfi[dfi['weather'] == 4]['count']

In [90]: dfi.groupby('weather').value_counts()

Out[90]:
weather count
fall 2733
spring 2686
Name: weather, dtype: int64
```

Null Hypothesis: Ho: Weather has no impact on sales

Alternate Hypothesis: Ha: Weather has impact on sales

Mathematically: Ho: $\mu_1 = \mu_2$ (i.e mean no of cycles is same for all weather)

Ha: $\mu_1 \neq \mu_2$ (i.e mean no of cycles is different for all weather)

We will take 95% confidence, i.e. $\alpha = 0.05$

So if $p\text{-value} > \alpha \Rightarrow$ Fail to reject Ho and if $p\text{-value} < \alpha \Rightarrow$ Reject Ho.

Since weather 4 has only 1 data point so we will ignore it for the test.

```
In [131]: f_stat,p_value,f_oneway(weather1,weather2,weather3)
print('p_value =',p_value)
if p_value < 0.05:
    print('Weather has impact on sales of bicycle')
else:
    print('Weather has no impact on sales of bicycle')

p_value 4.97644550994195e-43
Weather has impact on sales of bicycle
```

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

Chi-square test to check if Weather is dependent on the season

```
In [78]: dfi[['weather','season']].describe(include='all')

Out[78]:
weather season
count 10886.000000 10886
unique NaN 4
top NaN winter
freq NaN 2734
mean 1.418427 NaN
std 0.63839 NaN
min 1.000000 NaN
25% 2.000000 NaN
50% 3.000000 NaN
75% 4.000000 NaN
```

```
In [120]: cross_table = pd.crosstab(index = dfi['season'],
                                  columns = dfi['weather'],
                                  values = dfi['count'],
                                  aggfunc = np.sum).to_numpy()

print(cross_table)
real_cross=cross_table[0,0:3]
real_cross

[[4.7812e+05 5.3908e+05 3.1168e+04 nan]
 [2.2389e+05 7.4460e+04 1.2919e+04 1.6400e+02]
 [4.2630e+05 1.3417e+05 2.7758e+04 nan]
 [3.6508e+05 1.5713e+05 3.6058e+04 nan]]

Out[120]: array([[470116., 133988., 31168.],
       [223890., 74460., 12919.],
       [426300., 134170., 27755.],
       [365080., 157190., 36058.]])
```

Null Hypothesis: Ho: Weather is independent of season

Alternate Hypothesis: Ha: Weather is dependent on season

Mathematically: Ho: $\mu_1 = \mu_2$ (i.e mean no of cycles is same for all seasons)

Ha: $\mu_1 \neq \mu_2$ (i.e mean no of cycles is different for all season)

We will take 95% confidence, i.e. $\alpha = 0.05$

So if $p\text{-value} > \alpha \Rightarrow$ Fail to reject Ho and if $p\text{-value} < \alpha \Rightarrow$ Reject Ho.

```
In [133]: chi_test_stat, p_value, dof, expected = chi2_contingency(real_cross)
print('Test Statistic =', chi_test_stat)
print('p_value =', p_value)
if p_value < 0.05:
    print('Weather is dependent of season')
else:
    print('Weather is independent of season')

Test Statistic = 18838.37232489214
p_value = 0.0
Weather is dependent of season
```

Observation

1. The dataset is almost distributed equally among seasons while most data has weather1. We can conclude people prefer Yulu irrespective of seasons but if the weather is weather_1 i.e Clear, Few clouds, partly cloudy, partly cloudy, then they prefer using YULU bikes.

2. Temperature and feeling temperature are strongly positively correlated.

3. While season and weather or season and workingday has almost no correlation.

4. Working days(non weekend and non holidays) has more users of Yulu bikes.

5. Non holidays days have more yulu bikes users.

6. Users prefer temperature between 11 and 30 to use Yulu bikes.

7. 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.

-> Weather is dependent of season.

Recommendation

1. Since working days has more no of users, so Yulu can adjust its pricing according to its goals . and have more bikes availability during these days.
2. Temperature is ranging between 11 and 30 has highest number of users, so Yulu can have dynamic pricing during these times to earn extra profit.

3. Encourage customers to provide feedback and reviews on their biking experience. Collecting feedback can help the company assess its performance, understand customer preferences, and tailor the services to better meet customer expectations.
4. Recognize the impact of weather on bike rentals. Create weather-based promotions that target customers during clear and sunny weather, or those conditions show the highest rental counts. Yulu can offer weather-specific discounts to attract more customers during these favorable weather conditions.