

Importing Libraries

```
In [1]: import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
import seaborn as sns
from scipy import stats
from statsmodels.graphics.gcfplots import qqplot_2samples
import warnings
warnings.filterwarnings('ignore')
```

Loading the data

```
In [2]: df = pd.read_csv('C:/Users/aksh/Desktop/Scaler/Case_Study/Tulu/yulu.csv')
df.head()
```

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

No missing values

```
In [3]: 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   float64
 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
```

```
In [4]: print("season = %s" %df['season'].unique())
print("holiday = %s" %df['holiday'].unique())
print("workingday = %s" %df['workingday'].unique())
print("Weather = %s" %df['weather'].unique())

season = [1 2 3 4]
holiday = [0 1]
workingday = [0 1]
Weather = [1 2 3 4]
```

Converting Categorical values to object

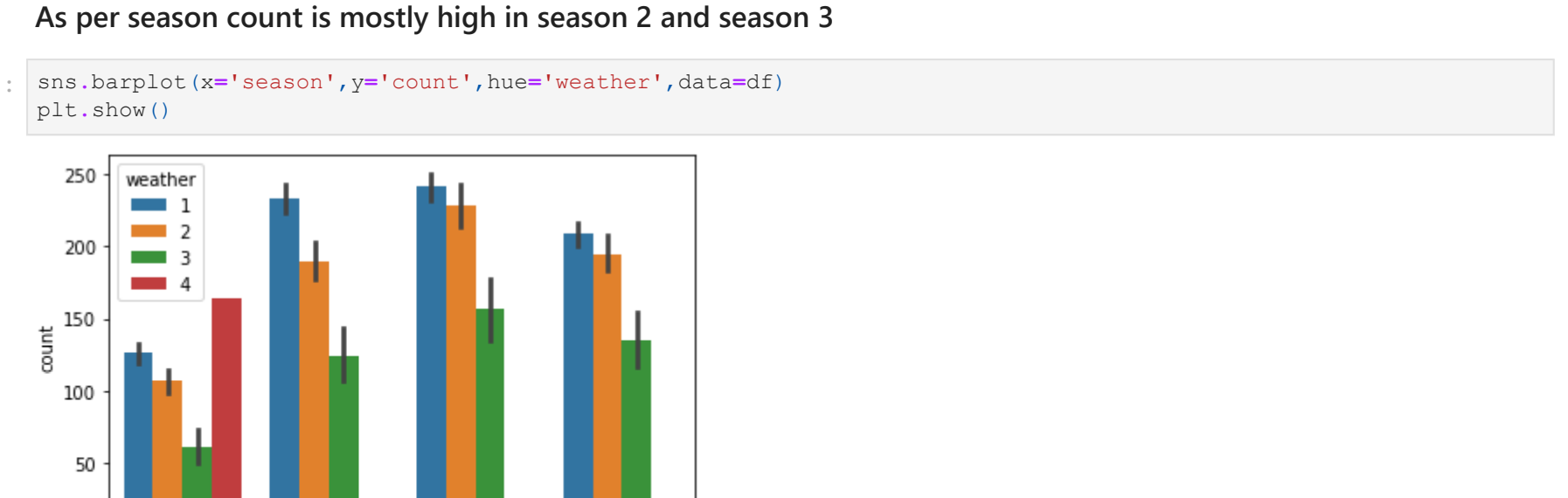
```
In [5]: df['season'] = df['season'].astype('object')
df['holiday'] = df['holiday'].astype('object')
df['workingday'] = df['workingday'].astype('object')
df['Weather'] = df['season'].astype('object')
```

```
In [6]: df.describe(include = ['int64', 'float64', 'object'])
```

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
count	10886	10886.0	10886.0	10886.0	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000
unique	10886	4	2.0	2.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
top	10-01-2012 09:00	4.0	0.0	1.0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
freq	1	2734.0	1057.50	7412.0	NaN	1.418427	20.23086	23.655084	61.886460	12.799995	36.021955	155.552
mean	NaN	NaN	NaN	NaN	NaN	0.633839	7.79159	8.474601	19.245033	8.164537	49.960477	151.039
std	NaN	NaN	NaN	NaN	NaN	1.000000	0.82000	0.766000	0.000000	0.000000	0.000000	0.000
min	NaN	NaN	NaN	NaN	NaN	1.000000	13.94000	16.665000	47.000000	7.001500	4.000000	36.000
50%	NaN	NaN	NaN	NaN	NaN	1.000000	20.50000	24.240000	62.000000	12.998000	17.000000	118.000
75%	NaN	NaN	NaN	NaN	NaN	2.000000	26.34000	31.060000	77.000000	16.997900	49.000000	222.000
max	NaN	NaN	NaN	NaN	NaN	4.000000	41.00000	45.455000	100.000000	56.996900	367.000000	886.000

Checking for outliers in casual and registered as their mean and mean have significant difference.

```
In [7]: i = ['casual','registered','count']
fig,axs = plt.subplots(ncols = 3, rows = 1, figsize=(15,5))
index=0
for i in i:
    sns.boxplot(y = i, data = df, ax = axs[index])
    index=index+1
```



```
In [8]: Upper_bound_C = (1.5*stats.iqr(df['casual']))+np.quantile(df['casual'],0.75)
Upper_bound_R = (1.5*stats.iqr(df['registered']))+np.quantile(df['registered'],0.75)
print("No of Outliers for casual = %s" %len(df.loc(df['casual']>Upper_bound_C)))
print("No of Outliers for registered = %s" %len(df.loc(df['registered']>Upper_bound_R)))
print("No of Outliers for count = %s" %len(df.loc(df['count']>Upper_bound_Co)))
```

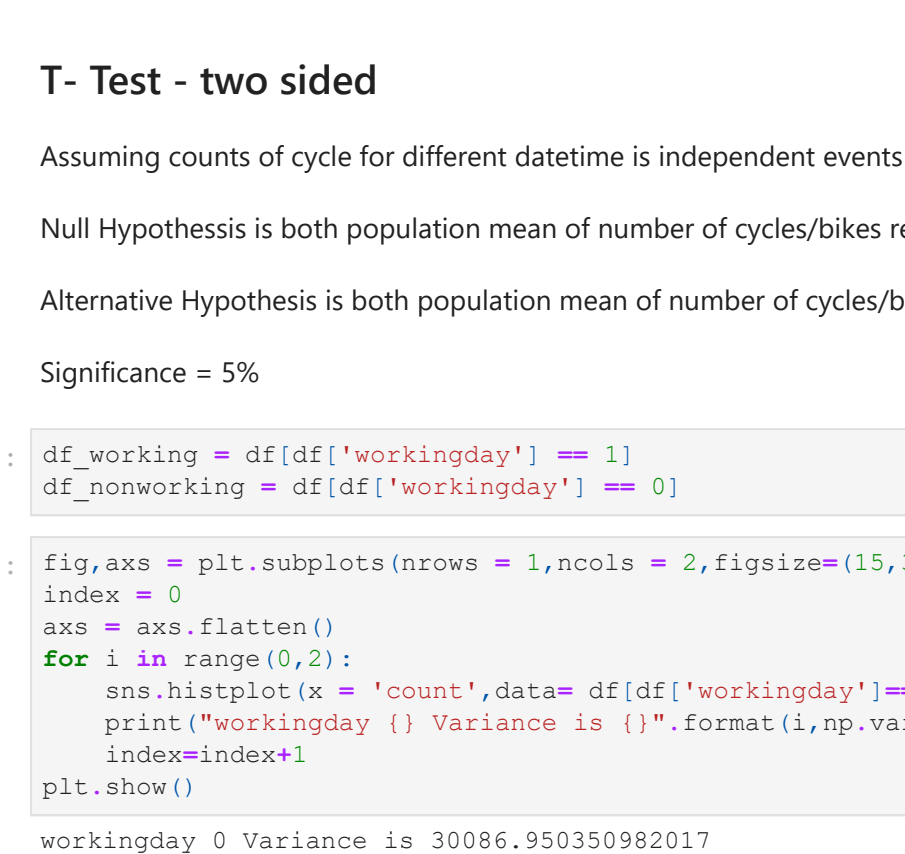
No of Outliers for casual = 749
No of Outliers for registered = 423
No of Outliers for count = 300

Temp have positive correlation with count

humidity have negative correlation with count

```
In [9]: sns.heatmap(df.corr(method='spearman'),annot=True)

<AxesSubplot>
```



We cant remove all of these as it is 10 percent of our data.

```
In [10]: len(df[(df['casual']>Upper_bound_C) | (df['registered']>Upper_bound_R) | (df['count']>Upper_bound_Co)])
10.490538306081206
```

Lets check for common rows which have outlier in all 3 columns -- As it is only .275 percent ,this we can remove it from dataframe

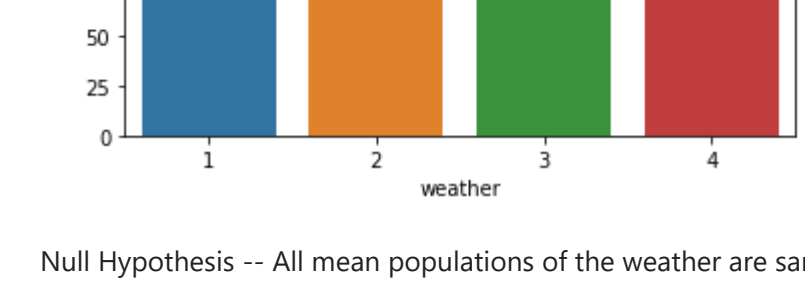
```
In [11]: len(df[(df['casual']>Upper_bound_C) & (df['registered']>Upper_bound_R) & (df['count']>Upper_bound_Co)])
0.27553318023149
```

```
In [12]: drop = df[(df['casual']>Upper_bound_C) & (df['registered']>Upper_bound_R) & (df['count']>Upper_bound_Co)].index
df.drop(drop,inplace=True)
```

As we see here count is mostly high in weather 1 and weather 2

As per season count is mostly high in season 2 and season 3

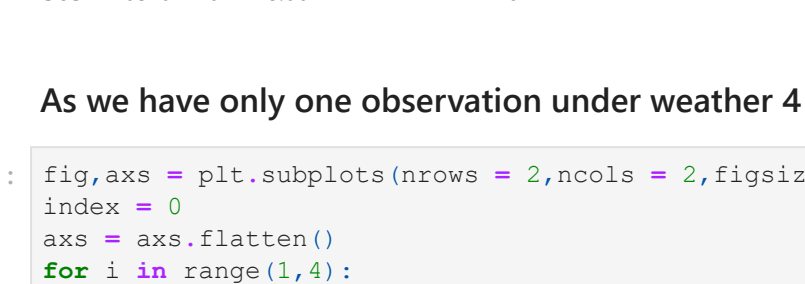
```
In [13]: sns.barplot(x='season',y='count',hue='weather',data=df)
plt.show()
```



Working Day has effect on number of electric cycles rented

The mean looks more are less same as per the graph lets do Hypothesis test to find out

```
In [14]: sns.barplot(y='count',x='workingday',data=df)
plt.show()
```



T- Test - two sided

Assuming counts of cycle for different datetime is independent events

Null Hypothesis is both population mean of number of cycles/bikes rented are same for Working and non working.

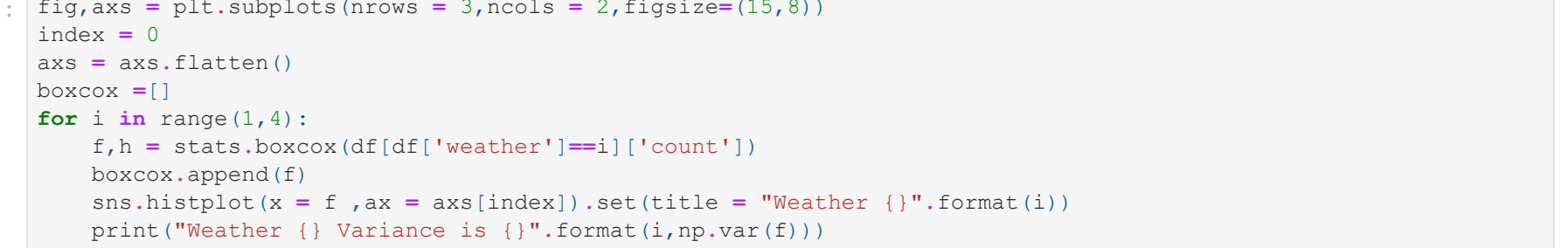
Alternative Hypothesis is both population mean of number of cycles/bikes rented are not equal for Working and non working.

Significance = 5%

```
In [15]: df_working = df[df['workingday'] == 1]
df_nonworking = df[df['workingday'] == 0]
```

```
In [16]: fig,axs = plt.subplots(nrows = 3,ncols = 2,figsize=(15,3))
index = 0
axs = axs.flatten()
for i in range(0,2):
    sns.histplot(x = 'count',data= df[df['workingday']==i],ax = axs[index]).set(title = "Workingday {}".format(i))
    print("Workingday {} Variance is {}".format(i,np.var(df[df['workingday']==i]['count'])))
    index=index+1
plt.show()
```

workingday 0 Variance is 30086.95035098207
workingday 1 Variance is 32668.88306255997



```
In [17]: stats.ttest_ind(df_working['count'],df_nonworking['count'])
Ttest_indResult(statistic=0.616762184303081, pvalue=0.5377911440927007)
```

pvalue is significantly greater than .05, so we can conclude with it that both population mean are same.

Anova

Weather -- No. of cycles rented is similar or different in different weathers

```
In [18]: df['weather'].unique()
```

```
Out[18]: array([1, 2, 3, 4], dtype=int64)
```

We see noticeable difference in sample mean

```
In [19]: sns.barplot(y='count',x='weather',data=df)
plt.show()
```



Null Hypothesis -- All mean populations of the weather are same

Alternative Hypothesis -- Not all are equal

Significance level = 5% percent

```
In [20]: df_1=df[df['weather']==1]
df_2=df[df['weather']==2]
df_3=df[df['weather']==3]
df_4=df[df['weather']==4]
```

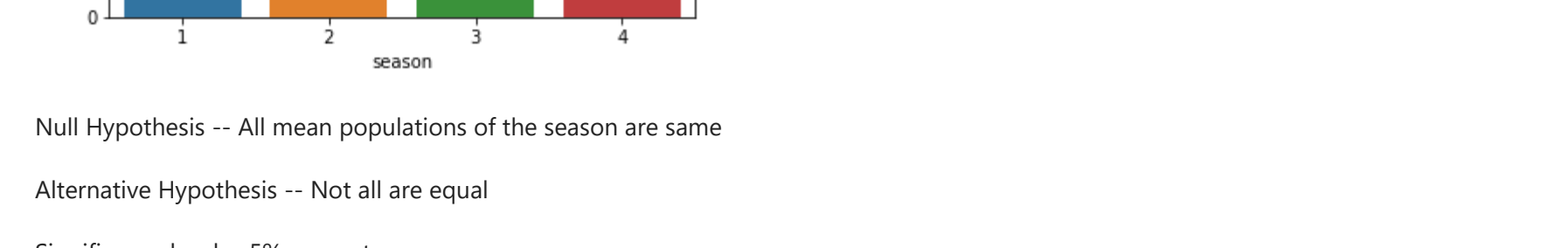
```
In [21]: df_4
```

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	Weather
5631	09-01-2012 18:00	1	0	1	4	8.2	11.365	86	6.0032	6	158	164	1

As we have only one observation under weather 4 will check for other weathers

```
In [22]: fig,axs = plt.subplots(nrows = 3,ncols = 2,figsize=(15,8))
index = 0
axs = axs.flatten()
for i in range(1,4):
    sns.histplot(x = 'count',data= df[df['weather']==i],ax = axs[index]).set(title = "Weather {}".format(i))
    print("Weather {} Variance is {}".format(i,np.var(df[df['weather']==i]['count'])))
    index=index+1
plt.show()
```

Weather 1 Variance is 34101.0262353682
Weather 2 Variance is 27882.22433892511
Weather 3 Variance is 19182.41876129777



```
In [17]: stats.ttest_ind(df_working['count'],df_nonworking['count'])
Ttest_indResult(statistic=0.616762184303081, pvalue=0.5377911440927007)
```

pvalue is significantly greater than .05, so we can conclude with it that both population mean are same.

Anova

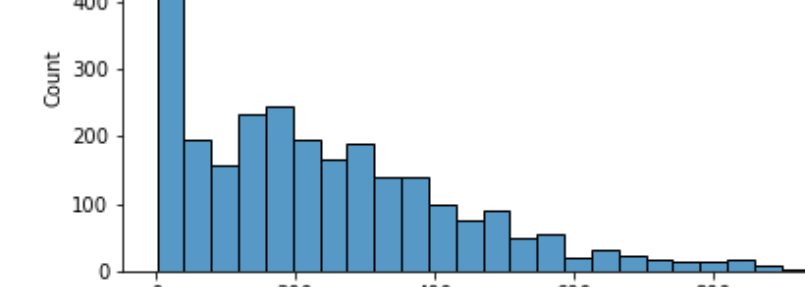
Weather -- No. of cycles rented is similar or different in different weathers

```
In [18]: df['weather'].unique()
```

```
Out[18]: array([1, 2, 3, 4], dtype=int64)
```

We see noticeable difference in sample mean

```
In [19]: sns.barplot(y='count',x='weather',data=df)
plt.show()
```



Null Hypothesis -- All mean populations of the weather are same

Alternative Hypothesis -- Not all are equal

Significance level = 5% percent

```
In [20]: df_1=df[df['weather']==1]
df_2=df[df['weather']==2]
df_3=df[df['weather']==3]
df_4=df[df['weather']==4]
```

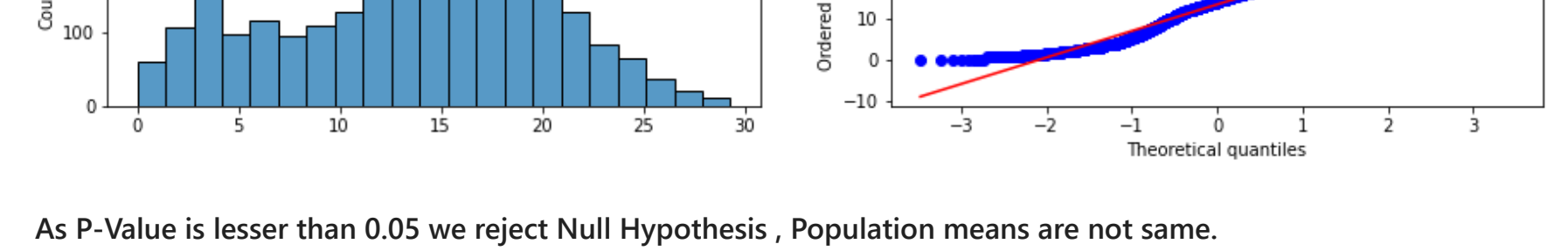
```
In [21]: df_4
```

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	Weather
5631	09-01-2012 18:00	1	0	1	4	8.2	11.365	86	6.0032	6	158	164	1

As we have only one observation under weather 4 will check for other weathers

```
In [22]: fig,axs = plt.subplots(nrows = 3,ncols = 2,figsize=(15,8))
index = 0
axs = axs.flatten()
for i in range(1,4):
    sns.histplot(x = 'count',data= df[df['weather']==i],ax = axs[index]).set(title = "Weather {}".format(i))
    print("Weather {} Variance is {}".format(i,np.var(df[df['weather']==i]['count'])))
    index=index+1
plt.show()
```

Weather 1 Variance is 34101.0262353682
Weather 2 Variance is 27882.22433892511
Weather 3 Variance is 19182.41876129777



```
In [17]: stats.ttest_ind(df_working['count'],df_nonworking['count'])
Ttest_indResult(statistic=0.616762184303081, pvalue=0.5377911440927007)
```

pvalue is significantly greater than .05, so we can conclude with it that both population mean are same.

Anova

Weather -- No. of cycles rented is similar or different in different weathers

```
In [18]: df['weather'].unique()
```

```
Out[18]: array([1, 2, 3, 4], dtype=int64)
```

We see noticeable difference in sample mean

```
In [19]: sns.barplot(y='count',x='weather',data=df)
plt.show()
```



Null Hypothesis -- All mean populations of the weather are same

Alternative Hypothesis -- Not all are equal

Significance level = 5% percent

```
In [20]: df_1=df[df['weather']==1]
df_2=df[df['weather']==2]
df_3=df[df['weather']==3]
df_4=df[df['weather']==4]
```

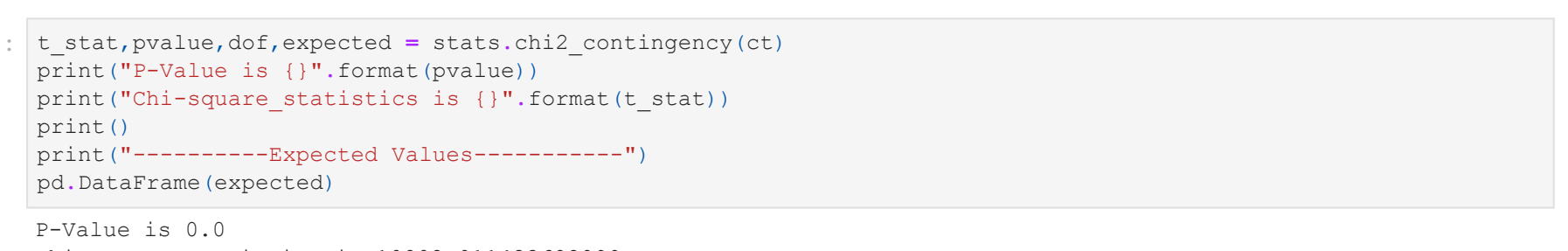
```
In [21]: df_4
```

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	Weather
5631	09-01-2012 18:00	1	0	1	4	8.2	11.365	86	6.0032	6	158	164	1

As we have only one observation under weather 4 will check for other weathers

```
In [22]: fig,axs = plt.subplots(nrows = 3,ncols = 2,figsize=(15,8))
index = 0
axs = axs.flatten()
for i in range(1,4):
    sns.histplot(x = 'count',data= df[df['weather']==i],ax = axs[index]).set(title = "Weather {}".format(i))
    print("Weather {} Variance is {}".format(i,np.var(df[df['weather']==i]['count'])))
    index=index+1
plt.show()
```

Weather 1 Variance is 34101.0262353682
Weather 2 Variance is 27882.22433892511
Weather 3 Variance is 19182.41876129777



```
In [17]: stats.ttest_ind(df_working['count'],df_nonworking['count'])
Ttest_indResult(statistic=0.616762184303081, pvalue=0.5377911440927007)
```

pvalue is significantly greater than .05, so we can conclude with it that both population mean are same.

Anova

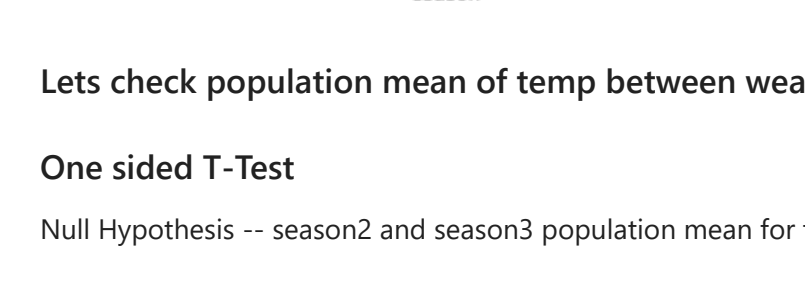
Weather -- No. of cycles rented is similar or different in different weathers

```
In [18]: df['weather'].unique()
```

```
Out[18]: array([1, 2, 3, 4], dtype=int64)
```

We see noticeable difference in sample mean

```
In [19]: sns.barplot(y='count',x='weather',data=df)
plt.show()
```



Null Hypothesis -- All mean populations of the weather are same

Alternative Hypothesis -- Not all are equal

Significance level = 5% percent

```
In [20]: df_1=df[df['weather']==1]
df_2=df[df['weather']==2]
df_3=df[df['weather']==3]
df_4=df[df['weather']==4]
```

```
In [21]: df_4
```

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	Weather
5631	09-01-2012 18:00	1	0	1	4	8.2	11.365	86	6.0032	6	158	164	1

As we have only one observation under weather 4 will check for other weathers

```
In [22]: fig,axs = plt.subplots(nrows = 3,ncols = 2,figsize=(15,8))
index = 0
axs = axs.flatten()
for i in range(1,4):
    sns.histplot(x = 'count',data= df[df['weather']==i],ax = axs[index]).set(title = "Weather {}".format(i))
    print("Weather {} Variance is {}".format(i,np.var(df[df['weather']==i]['count'])))
    index=index+1
plt.show()
```

Weather 1 Variance is 34101.0262353682
Weather 2 Variance is 27882.22433892511
Weather 3 Variance is 19182.41876129777



```
In [17]: stats.ttest_ind(df_working['count'],df_nonworking['count'])
Ttest_indResult(statistic=0.616762184303081, pvalue=0.5377911440927007)
```

pvalue is significantly greater than .05, so we can conclude with it that both population mean are same.

Anova

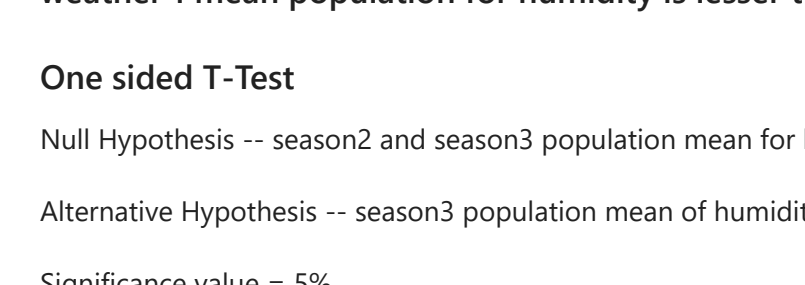
Weather -- No. of cycles rented is similar or different in different weathers

```
In [18]: df['weather'].unique()
```

```
Out[18]: array([1, 2, 3, 4], dtype=int64)
```

We see noticeable difference in sample mean

```
In [19]: sns.barplot(y='count',x='weather',data=df)
plt.show()
```



Null Hypothesis -- All mean populations of the weather are same

Alternative Hypothesis -- Not all are equal

Significance level = 5% percent

```
In [20]: df_1=df[df['weather']==1]
df_2=df[df['weather']==2]
df_3=df[df['weather']==3]
df_4=df[df['weather']==4]
```

```
In [21]: df_4
```

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	Weather
5631	09-01-2012 18:00	1	0	1	4	8.2	11.365	86	6.0032	6	158	164	1

As we have only one observation under weather 4 will check for other weathers

```
In [22]: fig,axs = plt.subplots(nrows = 3,ncols = 2,figsize=(15,8))
index = 0
axs = axs.flatten()
for i in range(1,4):
    sns.histplot(x = 'count',data= df[df['weather']==i],ax = axs[index]).set(title = "Weather {}".format(i))
    print("Weather {} Variance is {}".format(i,np.var(df[df['weather']==i]['count'])))
    index=index+1
plt.show()
```

Weather 1 Variance is 34101.0262353682
Weather 2 Variance is 27882.22433892511
Weather 3 Variance is 19182.41876129777

