

▼ Business Problem

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.

Objective:

- To find which variables are significant in predicting the demand for shared electric cycles in the Indian market?
- To understand how well those variables describe the electric cycle demands

Concept Used:

- Bi-Variate Analysis
- 2-sample t-test: testing for difference across populations
- ANNOVA
- Chi-square

▼ import modules and load data

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5
6 from scipy.stats import norm, chi2, f # distributions
7
8 from scipy.stats import ttest_ind, ttest_rel, f_oneway, kruskal # numerical vs categorical
9 from scipy.stats import chisquare, chi2_contingency # categorical features
10 from scipy.stats import pearsonr, spearmanr # numeric vs numeric
11
12 from scipy.stats import kstest # cdf
13
14 from statsmodels.distributions.empirical_distribution import ECDF
15 # Empirical CDF
```

```
1 df=pd.read_csv('/content/bike_sharing.txt')
```

```
1 df
```

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	cc
0	2011-01-01 00:00:00	1	0	0	1	9.84	14.395	81	0.0000	3	13	
1	2011-01-01 01:00:00	1	0	0	1	9.02	13.635	80	0.0000	8	32	
2	2011-01-01 02:00:00	1	0	0	1	9.02	13.635	80	0.0000	5	27	
3	2011-01-01 03:00:00	1	0	0	1	9.84	14.395	75	0.0000	3	10	
4	2011-01-01 04:00:00	1	0	0	1	9.84	14.395	75	0.0000	0	1	
...
	2012-12-											

```
1 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
```

▼ Basic data exploration

```
1 df.shape
2 # 10886 rows, 12 columns
```

```
(10886, 12)
```

```
1 df.isna().sum()
2 # No missing value
```

```
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
```

```
1 df.dtypes
2 # checking datatype of columns
```

```
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
```

```
1 df.nunique()
2 # few int columns have very less unique values, we will convert them to object
```

```
datetime    10886
season       4
holiday      2
workingday   2
weather      4
temp         49
atemp        60
humidity     89
windspeed    28
casual       309
registered   731
count        822
dtype: int64
```

```
1 # Drop unnecessary columns
2 df.drop('datetime',axis=1, inplace=True )
```

```
1 columns=['season','holiday','workingday','weather']
2 df[columns] = df[columns].astype('object')
```

```
1 df.dtypes
```

```
season      object
holiday      object
workingday   object
weather      object
temp        float64
atemp        float64
humidity     int64
```

```
windspeed    float64
casual        int64
registered    int64
count         int64
dtype: object
```

```
1 ## separeate categorical and numeric features.
2 # catgeorical and numerical columns
3 cat_cols = df.dtypes == 'object'
4 cat_cols = list(cat_cols[cat_cols].index)
5 cat_cols
```

```
['season', 'holiday', 'workingday', 'weather']
```

```
1 num_cols = df.dtypes != 'object'
2 num_cols = list(num_cols[num_cols].index)
3 num_cols
```

```
['temp', 'atemp', 'humidity', 'windspeed', 'casual', 'registered', 'count']
```

▼ Univariate Data Analysis

```
1 #Checking how the data is spread on basis of distinct users (customer analysis)
2 df2=df.copy()
3 cat_count = df2[cat_cols].melt().groupby(['variable', 'value'])['value'].size().reset_index(name='counts')
4 s = df2[cat_cols].melt().variable.value_counts()
5 cat_count['Percent'] = cat_count['counts'].div(cat_count['variable'].map(s)).mul(100).round().astype('int')
6 cat_count.groupby(['variable', 'value']).first()
7
8 # only 1 electric cycles rented in extreame weather 4. (outlier)
9 # only 3% days are holidays.
```

		counts	Percent
variable	value		
holiday	0	10575	97
	1	311	3
season	1	2686	25
	2	2733	25
	3	2733	25
	4	2734	25
weather	1	7192	66
	2	2834	26
	3	859	8
	4	1	0
workingday	0	3474	32
	1	7412	68

```
1 df[df['weather']==4]
2
3 # Only one row where weather is 4, so will replace it's value with 3.
```

	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
5631	1	0	1	4	8.2	11.365	86	6.0032	6	158	164

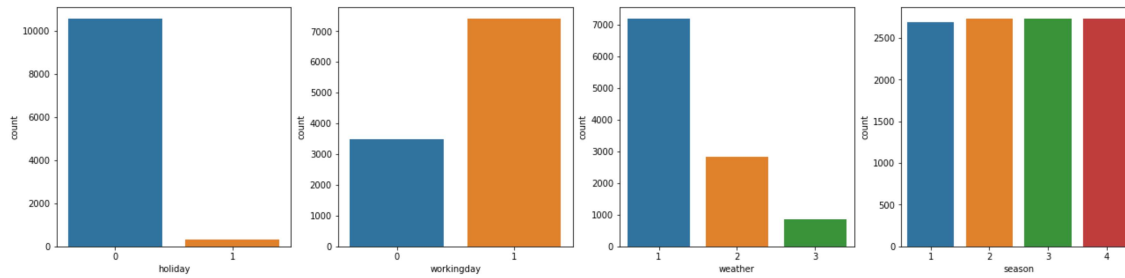
```
1 df['weather'].replace(4,3,inplace=True)
2 df['weather'].value_counts()
```

```
1    7192
2    2834
3     860
Name: weather, dtype: int64
```

```
1 plt.figure(figsize = [22,5])
2 cat_cols = ['holiday', 'workingday', 'weather', 'season']
3 for i in range (len(cat_cols)):
4     plt.subplot(1, 4, i+1)
5     sns.countplot(data=df, x=cat_cols[i])
```

6

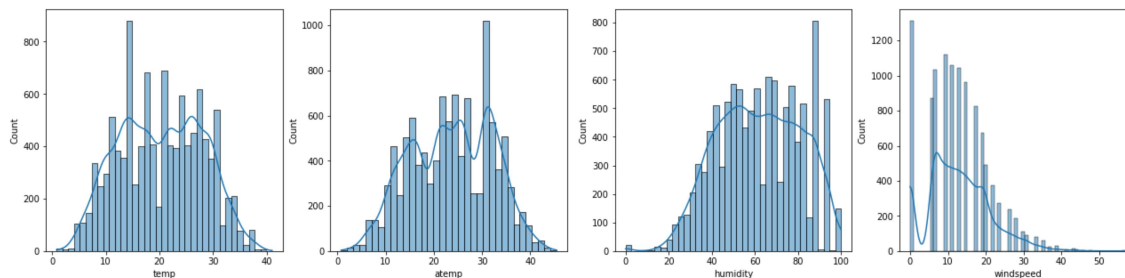
7 # Weather 1 is most liklihood weather.



```
1 df.describe(include='object')
```

	season	holiday	workingday
count	10886	10886	10886
unique	4	2	2
top	4	0	1
freq	2734	10575	7412

```
1 plt.figure(figsize = [22,5])
2 num_cols = ['temp', 'atemp', 'humidity', 'windspeed']
3 for i in range (len(num_cols)):
4     plt.subplot(1, 4, i+1)
5     sns.histplot(data=df, x=num_cols[i], kde=True)
```



```
1 for i in (num_cols):
2     print(i, round(df[i].skew(),1))
3
4 # windspeed distribution is right skewed, means it has some outliers in right.
```

```
temp 0.0
atemp -0.1
humidity -0.1
windspeed 0.6
```

```
1 from scipy.stats import shapiro
2 num_cols = ['temp', 'atemp', 'humidity', 'windspeed']
3 for i in (num_cols):
4     print(shapiro(df[i]).pvalue)
5
6 # Since the p-value is less than .05, we reject the null hypothesis.
7 # By shapiro test we can say that the sample data does not come from a normal distribution.
8 # We could normalize it but We will consider then normally distributed in further analysis.
9 # Because in visual plots we can see that variables are looking kind of normally distributed.
```

```
4.47221826500091e-36
3.4538982852050647e-35
1.245496990918048e-34
0.0
/usr/local/lib/python3.8/dist-packages/scipy/stats/morestats.py:1760: UserWarning: p-value may not be accurate for N > 5000.
  warnings.warn("p-value may not be accurate for N > 5000.")
```

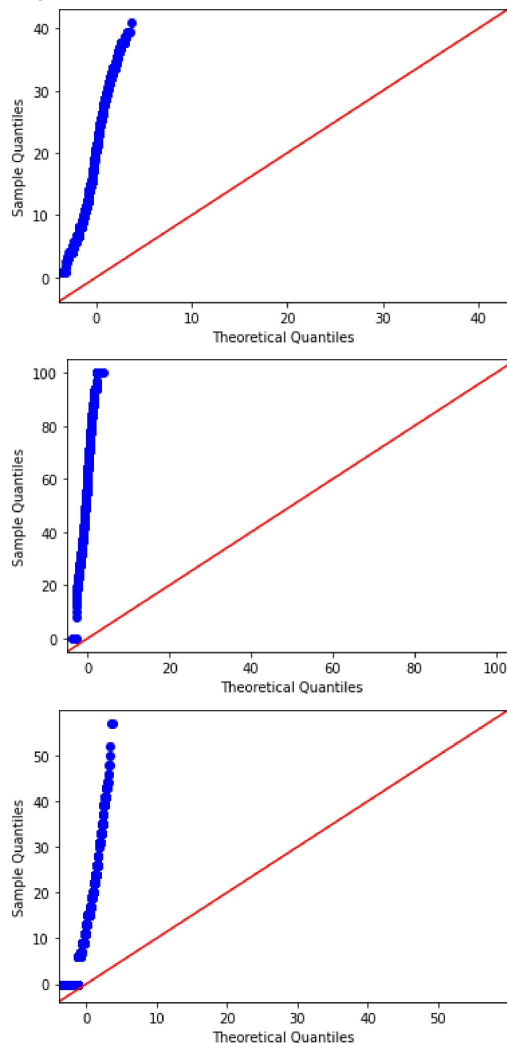
```
1 import statsmodels.api as sm
2 import matplotlib.pyplot as plt
3 plt.figure(figsize = [22,5])
4 num_cols = ['temp', 'humidity', 'windspeed']
5 for i in (num_cols):
```

```

6 sm.qqplot(df[i], line='45')
7
8 #The data values clearly do not follow the red 45-degree line,
9 # which is an indication that they do not follow a normal distribution.

```

<Figure size 1584x360 with 0 Axes>



```

1 from scipy.stats import levene
2 stat, p = levene(df['temp'],df['atemp'], center='median')
3 p

```

1.3036286748857844e-16

```

1 from scipy.stats import levene
2 stat, p = levene(df['temp'],df['atemp'], center='mean')
3 p

```

3.23529399853922e-17

```

1 alpha =0.05
2 # now we pass the groups and center value from the following
3 # ('trimmed mean', 'mean', 'median')
4 w_stats, p_value =levene(df['temp'],df['atemp'], center='median')
5
6 if p_value > alpha :
7     print("We do not reject the null hypothesis")
8 else:
9     print("Reject the Null Hypothesis")
10
11 # This means we do have sufficient evidence to say that the variance between temp and atemp is significantly different.

```

Reject the Null Hypothesis

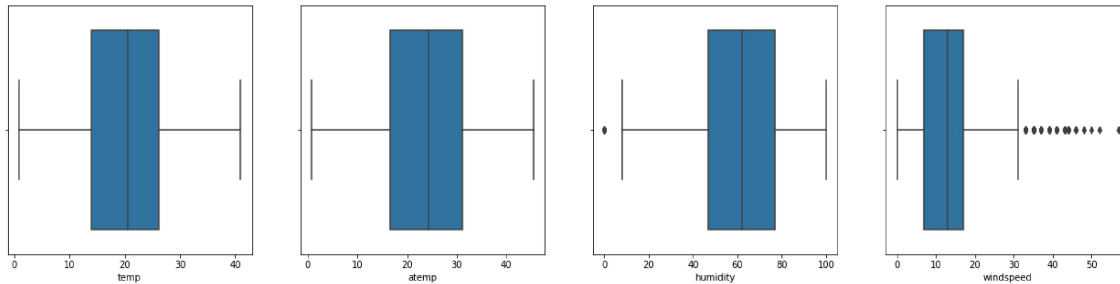
```
1 np.percentile(df['temp'],97.5)-np.percentile(df['temp'],2.5)
```

27.7775

```
1 np.percentile(df['atemp'],97.5)-np.percentile(df['atemp'],2.5)
```

31.06

```
1 plt.figure(figsize = [22,5])
2 num_cols = ['temp', 'atemp', 'humidity', 'windspeed']
3 for i in range (len(num_cols)):
4     plt.subplot(1, 4, i+1)
5     sns.boxplot(data=df, x=num_cols[i])
```



```
1 df[num_cols].describe()
```

	temp	atemp	humidity	windspeed
count	10886.00000	10886.000000	10886.000000	10886.000000
mean	20.23086	23.655084	61.886460	12.799395
std	7.79159	8.474601	19.245033	8.164537
min	0.82000	0.760000	0.000000	0.000000
25%	13.94000	16.665000	47.000000	7.001500
50%	20.50000	24.240000	62.000000	12.998000
75%	26.24000	31.060000	77.000000	16.997900
max	41.00000	45.455000	100.000000	56.996900



```
1 R_whisker = np.percentile(df['windspeed'],75)+(np.percentile(df['windspeed'],75)-np.percentile(df['windspeed'],25))*1.5
2 R_whisker
```

31.992500000000003

```
1 df[df['windspeed']>R_whisker].index.size
2 # 227 outliers showing in windspeed column (using 1.5*IQR method).
```

227

```
1 L_whisker = np.percentile(df['humidity'],25)-(np.percentile(df['humidity'],75)-np.percentile(df['humidity'],25))*1.5
2 L_whisker
```

2.0

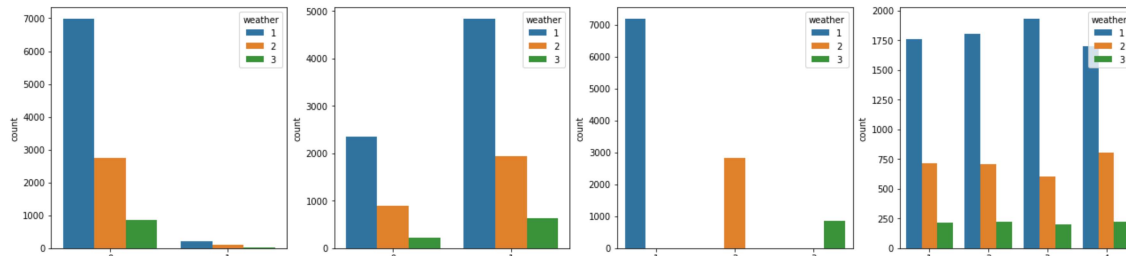
```
1 df[df['humidity']<L_whisker].index.size
2 # 22 outliers showing in humidity column (using 1.5*IQR method).
```

22

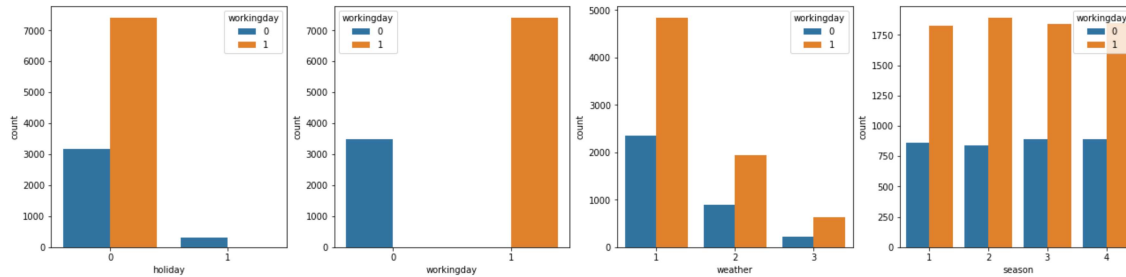
```
1 # We have found outliers and we could remove,
2 # Since we have limited data, we will keep outliers in further exploration and test.
3 # earlier we detected single outlier data point in weather column and we fixed it by imputation.
```

▼ Bivariate Data Analysis

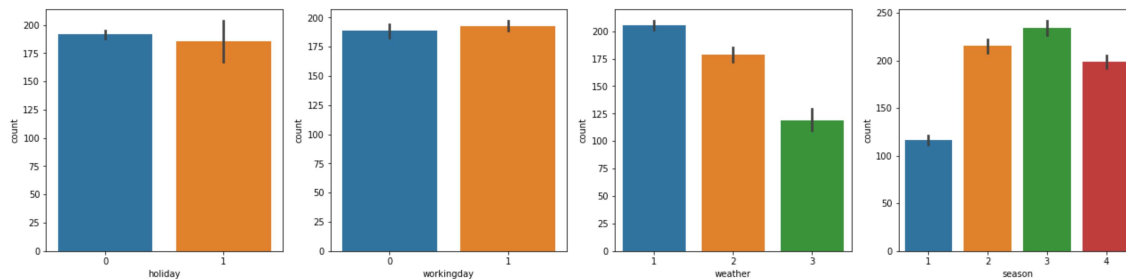
```
1 plt.figure(figsize = [22,5])
2 cat_cols = ['holiday', 'workingday', 'weather', 'season']
3 for i in range (len(cat_cols)):
4     plt.subplot(1, 4, i+1)
5     sns.countplot(data=df, x=cat_cols[i], hue=cat_cols[2])
```



```
1 plt.figure(figsize = [22,5])
2 cat_cols = ['holiday','workingday','weather','season']
3 for i in range (len(cat_cols)):
4     plt.subplot(1, 4, i+1)
5     sns.countplot(data=df, x=cat_cols[i],hue=cat_cols[1])
```



```
1 plt.figure(figsize = [22,5])
2 cat_cols = ['holiday','workingday','weather','season']
3 for i in range (len(cat_cols)):
4     plt.subplot(1, 4, i+1)
5     sns.barplot(data=df, x=cat_cols[i], y='count')
```



▼ 2 Sample T-Test

```
1 # Checking: Does number of electric cycles rented on weekday greater than weekend?
2
3 null_hypothesis = 'number of electric cycles rented are similar on weekday and weekend'
4 alternative_hypothesis = 'number of electric cycles rented higher on weekday than weekend'
5
6 sample1 = df[df['workingday']==0]['count']
7 sample2 = df[df['workingday']==1]['count']
8 t_stat, p_value = ttest_ind(sample1, sample2, equal_var=False, alternative='greater')
9 print(t_stat, p_value)
10
11 if(p_value < 0.05):
12     print('Since, p-value < 0.05, the null hypothesis is rejected')
13     print(alternative_hypothesis)
14 else:
15     print('Since p-value > 0.05, we fail to reject null hypothesis')
16     print(null_hypothesis)
17
18 # conclusion: number of electric cycles rented are similar on weekday and weekend.
```

```
-1.2362580418223226 0.8917984385965245
Since p-value > 0.05, we fail to reject null hypothesis
number of electric cycles rented are similar on weekday and weekend
```

```
1 # Checking: Does the Working Day has an effect on number of electric cycles rented?
2
3 null_hypothesis = 'Working Day has no effect on number of electric cycles rented'
```

```

4 alternative_hypothesis = 'Working Day has effect on number of electric cycles rented'
5
6 sample1 = df[df['workingday']==0]['count']
7 sample2 = df[df['workingday']==1]['count']
8 t_stat, p_value = ttest_ind(sample1, sample2)
9 print(t_stat, p_value)
10
11 if(p_value < 0.05):
12     print('Since, p-value < 0.05, the null hypothesis is rejected')
13     print(alternative_hypothesis)
14 else:
15     print('Since p-value > 0.05, we fail to reject null hypothesis')
16     print(null_hypothesis)
17
18 # conclusion: workingday has no effect on number of electric cycles rented.

```

```

-1.2096277376026694 0.22644804226361348
Since p-value > 0.05, we fail to reject null hypothesis
Working Day has no effect on number of electric cycles rented

```

```

1 print(np.std(sample1),np.std(sample2))
2 # There is not huge difference in standerd deviation, so no need to equalize size of samples.

```

```

173.69901006897658 184.501211667422

```

```

1 df['workingday'].value_counts()
2 # We could make sample size equal for both variable (by considering lesser variable size).
3 # But we are not doing it because it is giving difference p_value and different result when doing same test multiple times.
4 # So we are considering unequal sample size to get fixed p_value (it might increase type 1 error, but it will give sure result).

```

```

1    7412
0    3474
Name: workingday, dtype: int64

```

▼ ANNOVA test

```

1 # Checking : Are the number of cycles rented similar or different in the different seasons?
2
3 null_hypothesis = 'There is no significant difference in the No. of cycles rented in the different seasons'
4 alternative_hypothesis = 'There is significant difference in No. of cycles rented in the different seasons'
5
6 sample1 = df[df['season']==1]['count']
7 sample2 = df[df['season']==2]['count']
8 sample3 = df[df['season']==3]['count']
9 sample4 = df[df['season']==4]['count']
10 t_stat, p_value = f_oneway(sample1,sample2,sample3,sample4)
11 print(t_stat, p_value)
12
13 if(p_value < 0.05):
14     print('Since, p-value < 0.05, the null hypothesis is rejected')
15     print(alternative_hypothesis)
16 else:
17     print('Since p-value > 0.05, we fail to reject null hypothesis')
18     print(null_hypothesis)
19
20 # conclusion: Number of cycles rented are similar in different seasons.

```

```

236.94671081032106 6.164843386499654e-149
Since, p-value < 0.05, the null hypothesis is rejected
There is significant difference in No. of cycles rented in the different seasons

```

```

1 # Checking : Are the number of cycles rented similar or different in the different weather?
2
3 null_hypothesis = 'There is no significant difference in the No. of cycles rented in the different weather'
4 alternative_hypothesis = 'There is significant difference in No. of cycles rented in the different weather'
5
6 sample1 = df[df['weather']==1]['count']
7 sample2 = df[df['weather']==2]['count']
8 sample3 = df[df['weather']==3]['count']
9 t_stat, p_value = f_oneway(sample1,sample2,sample3)
10 print(t_stat, p_value)
11
12 if(p_value < 0.05):
13     print('Since, p-value < 0.05, the null hypothesis is rejected')
14     print(alternative_hypothesis)
15 else:
16     print('Since p-value > 0.05, we fail to reject null hypothesis')
17     print(null_hypothesis)

```



```

18
19 # Since, p-value < 0.05, the null hypothesis is rejected
    Since, p-value < 0.05, the null hypothesis is rejected
    There is significant difference in No. of cycles rented in the different weather

```

▼ Chi-square test

```

1 # Checking : Is Weather is dependent on season?
2
3 null_hypothesis = 'There is no relationship between Weather and Season'
4 alternative_hypothesis = 'There is relationship between Weather and Season'
5
6 contingency = pd.crosstab(df['weather'], df['season'])
7 print(contingency)
8
9 # p-value calculation
10 p_value = chi2_contingency(contingency)[1]
11 print('p-value:', round(p_value, 4))
12
13 if(p_value < 0.05):
14     print('Since, p-value < 0.05, the null hypothesis is rejected')
15     print(alternative_hypothesis)
16 else:
17     print('Since p-value > 0.05, we fail to reject null hypothesis')
18     print(null_hypothesis)
19
20 # conclusion: Weather and Season are dependent.

```

```

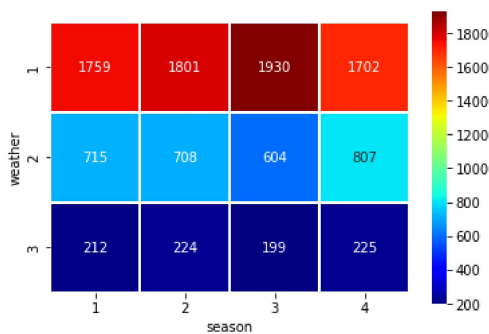
season      1      2      3      4
weather
1         1759   1801   1930   1702
2          715    708    604    807
3          212    224    199    225
p-value: 0.0
Since, p-value < 0.05, the null hypothesis is rejected
There is relationship between Weather and Season

```

```

1 sns.heatmap(pd.crosstab(df['weather'], df['season']), cmap= "jet", annot = True, fmt = 'd', square=1, linewidth=1.)
2 plt.show()
3
4 # Proportion of season differ across different weather.

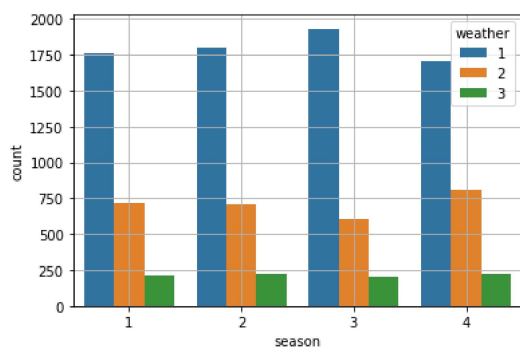
```



```

1 sns.countplot(data=df, x='season', hue='weather')
2 plt.grid()
3 plt.show()
4
5 # In season 3 most of the days have weather 1 and less days have with weather 3, Compared to other season.
6 # likelihood of weather 1 is different in different seasons.

```



All appropriate insights are mentioned with relevent code as comments. All the results are printed for relevent test.

✓ 0s completed at 11:34 PM

