

About the project and Problem Statement :

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.

The company wants to know:

- Which variables are significant in predicting the demand for shared electric cycles in the Indian market?
- How well those variables describe the electric cycle demands

Column Profiling:

- datetime: datetime
- season: season (1: spring, 2: summer, 3: fall, 4: winter)
- workingday: if day is neither weekend nor holiday is 1, otherwise is 0.
- weather:
 1. : Clear, Few clouds, partly cloudy, partly cloudy
 2. : Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist
 3. : Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds
 4. : Heavy Rain + Ice Pellets + Thunderstorm + Mist, Snow + Fog
- temp: temperature in Celsius
- atemp: feeling temperature in Celsius
- humidity: humidity
- windspeed: wind speed
- casual: count of casual users
- registered: count of registered users
- count: count of total rental bikes including both casual and registered

shape of the data :

```
(10886, 12)
```

10886 Records of bike Rented (each record shows howmany bikes were rented during that hour of the day.)

```
data.isna().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
```

no null values detected

unique values per columns:

```
data.nunique()
```

```
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
```

Data Pre-Processing Feature Extraction:

```
data["weather"].replace({1:"Clear",
                        2:"Cloudy",
                        3:"Little Rain",
                        4:"Heavy Rain"},inplace=True)
data["season"].replace({1:"Spring",
                        2:"Summer",
                        3:"Fall",
                        4:"Winter"},inplace=True)
data["workingday"].replace({1:"Yes",
                             0:"No"},inplace=True)
data["datetime"] = pd.to_datetime(data["datetime"])
data["holiday"].replace({1:"Yes",
                          0:"No"},inplace=True)
data["day"]=data["datetime"].dt.day_name()
data["date"] = data["datetime"].dt.date
data["hour"] = data["datetime"].dt.hour
data["Month"] = data["datetime"].dt.month
data["Month_name"] = data["datetime"].dt.month_name()
data["year"] = data["datetime"].dt.year
```

Describing Statistical summery of Independent Numerical Features :

Categorising Temperature And Humidity Levels and Windspeed column data :

	count	mean	std	min	25%	50%	75%	max
atemp	10886.0	23.655084	8.474601	0.76	16.665	24.24	31.06	45.455

	count	mean	std	min	25%	50%	75%	max
humidity	10886.0	61.88646	19.245033	0.0	47.0	62.0	77.0	100.0

	count	mean	std	min	25%	50%	75%	max
windspeed	10886.0	12.799395	8.164537	0.0	7.0015	12.998	16.9979	56.9969

Data information :

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10886 entries, 0 to 10885
Data columns (total 21 columns):
#   Column                Non-Null Count  Dtype
---  -
0   datetime              10886 non-null  datetime64[ns]
1   season                10886 non-null  object
2   holiday                10886 non-null  object
3   workingday            10886 non-null  object
4   weather                10886 non-null  object
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
12  day                   10886 non-null  object
13  date                  10886 non-null  object
14  hour                  10886 non-null  int64
15  Month                 10886 non-null  int64
16  Month_name            10886 non-null  object
17  year                  10886 non-null  int64
18  temperature           10886 non-null  object
19  gethumidity           10886 non-null  object
20  windspeed_category    10886 non-null  object
dtypes: datetime64[ns](1), float64(3), int64(7), object(10)
memory usage: 1.7+ MB
```

statistical summery about categorical data :

	season	holiday	workingday	weather	day	date	Month_name	temperature	gethumidity	windspeed_category
count	10886	10886	10886	10886	10886	10886	10886	10886	10886	10886
unique	4	2	2	4	7	456	12	4	10	8
top	Winter	No	Yes	Clear	Saturday	2011-01-01	May	moderate	70%	(-0.001, 6.003]
freq	2734	10575	7412	7192	1584	24	912	4767	1845	2185

Moderate level Temperature frequency is highest in given data

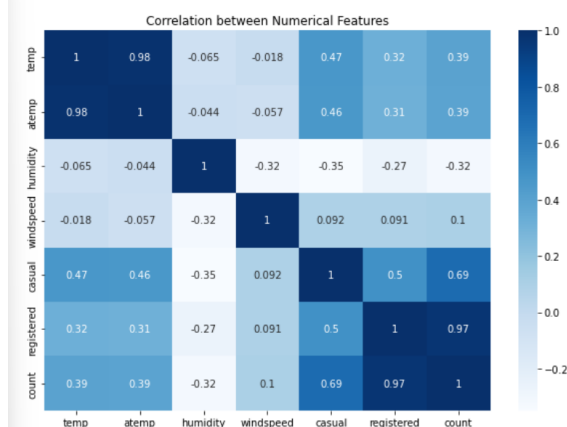
70% humdity

and most preferable windspeed 8-12

Correlation Matrix :

	temp	atemp	humidity	windspeed	casual	registered	count
temp	1.000000	0.984948	-0.064949	-0.017852	0.467097	0.318571	0.394454
atemp	0.984948	1.000000	-0.043536	-0.057473	0.462067	0.314635	0.389784
humidity	-0.064949	-0.043536	1.000000	-0.318607	-0.348187	-0.265458	-0.317371
windspeed	-0.017852	-0.057473	-0.318607	1.000000	0.092276	0.091052	0.101369
casual	0.467097	0.462067	-0.348187	0.092276	1.000000	0.497250	0.690414
registered	0.318571	0.314635	-0.265458	0.091052	0.497250	1.000000	0.970948
count	0.394454	0.389784	-0.317371	0.101369	0.690414	0.970948	1.000000

Heatmap (correlation between features)



Correlation between Temperature and Number of Cycles Rented for all customers : 0.39

Correlation between Temperature and Number of Cycles Rented for casual subscribers : 0.46

Correlation between Temperature and Number of Cycles Rented for registered subscribers : 0.31

Correlation between Temperature and Number of Cycles Rented for registered subscribers : 0.31

Humidity has a negative correlation with the number of cycles rented which is -0.32

◆ About the features :

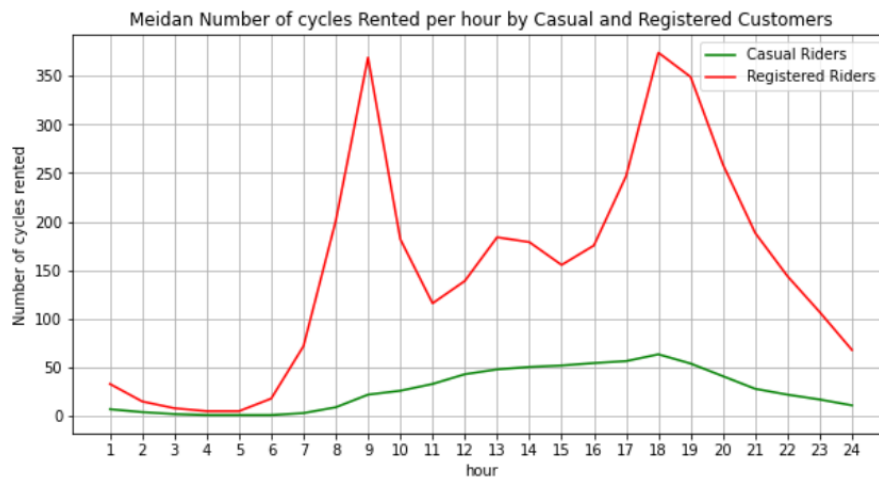
- dependent variables : count / registerd / casual
- independent variables : workingday / holiday / weather / seasons /temperature /humidity /windspeed.

Outlier detection in Dataset :

'0.0278 % Outliers data from input data found'

NUMBER OF CYCLES RENTED BY : CASUAL USERS AND REGISTERED USERS

Average Number of Cycles rented by Casual vs Registered Subscribes :



From above linplot :

- registered customers seems to be using rental cycles mostly for work-commute purposes.
- registered cycle counts seems to be much higher than the casual customers.

Casual Users (in %) :

18.8031413451893

Registered Users (in %) :

81.1968586548107

81% cycles had been rented by registered customers.

19% cycles had been rented by casual customers.

USING BOOTSTRAPPING : CONFIDENCE INTERVAL OF MEAN NUMBER OF CYCLES RENTED BY CASUAL AND REGISTERED CUSTOMERS :

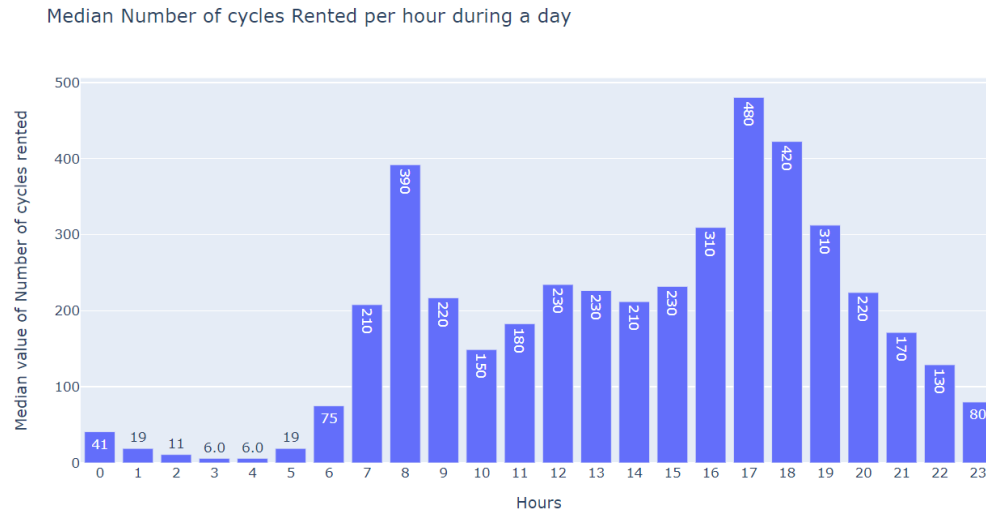
➤ Confidence Interval of Average Number of Cycles Rented by Registered Customers

Confidence Interval : (153.6716933943279, 157.0899466056721)

➤ Confidence Interval of Average Number of Cycles Rented by Casual Customers

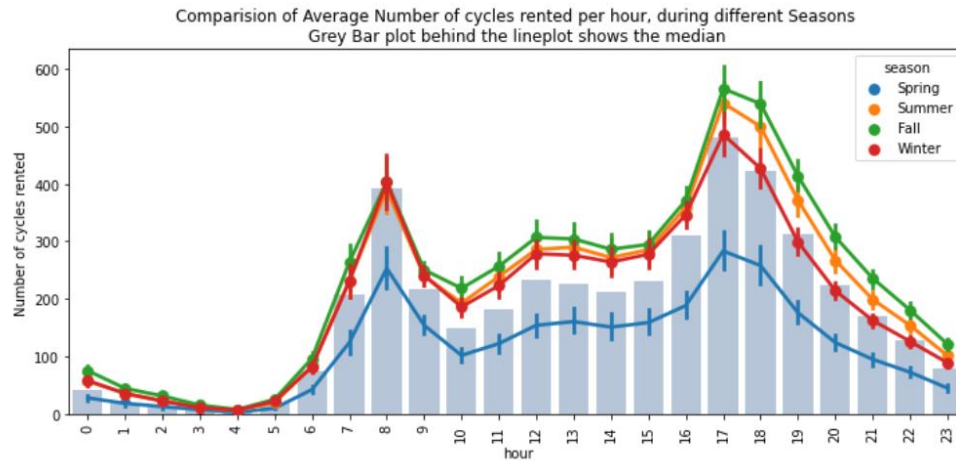
Confidence Interval : (35.44287786523068, 36.57356280143599)

HOURLY MEDIAN NUMBER OF CYCLES RENTED DURING THE DAY :



- from above bar chart :
- shows the median value of number of cycles were rented during perticular hour of the day.
- Median of number of cycles rented are higher during morning 7 to 9 am to evening 4 to 8pm .

EFFECT OF SEASONS ON NUMBER OF CYCLES RENTED DURING HOURS :

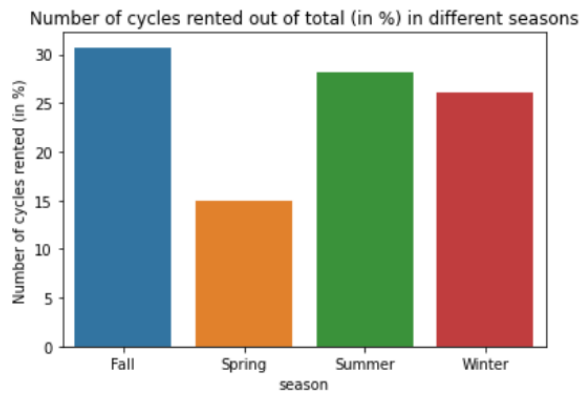


during the morning 7-9am and afternoon 4pm to 7pm , the cycles rent counts is increasing.

during the spring season , looks like people prefer less likely to rent the cycle.

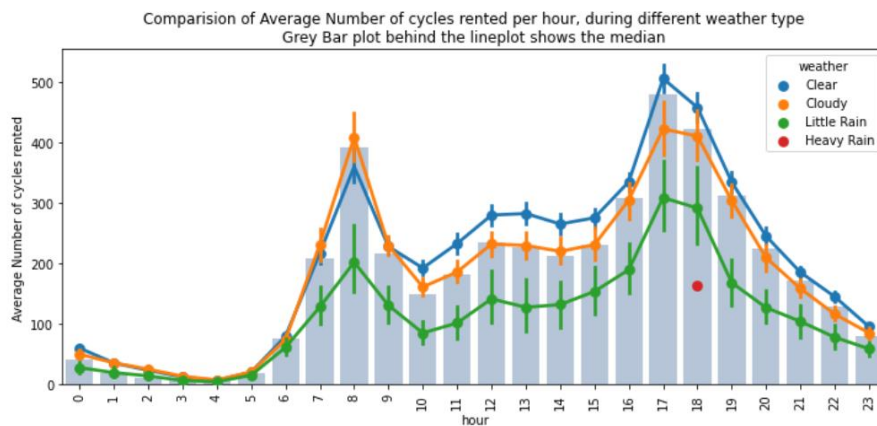
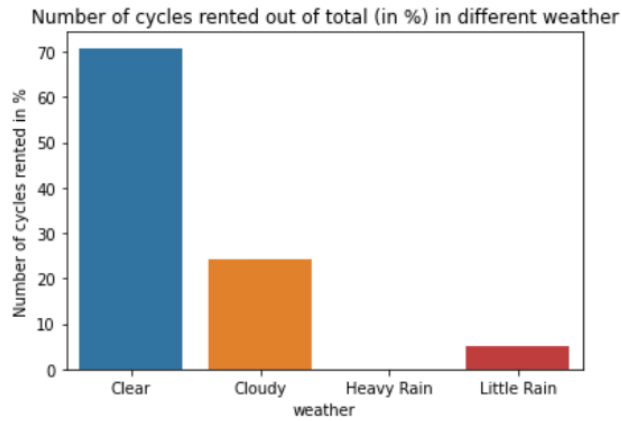
Number of cycles rented during different seasons (in %) :

- season
- Fall 30.720181
- Spring 14.984493
- Summer 28.208524
- Winter 26.086802



WEATHER EFFECT ON CYCLE RENTAL MEDIAN COUNTS HOURLY :

```
weather
Clear      70.778230
Cloudy     24.318669
Heavy Rain  0.007864
Little Rain 4.895237
```



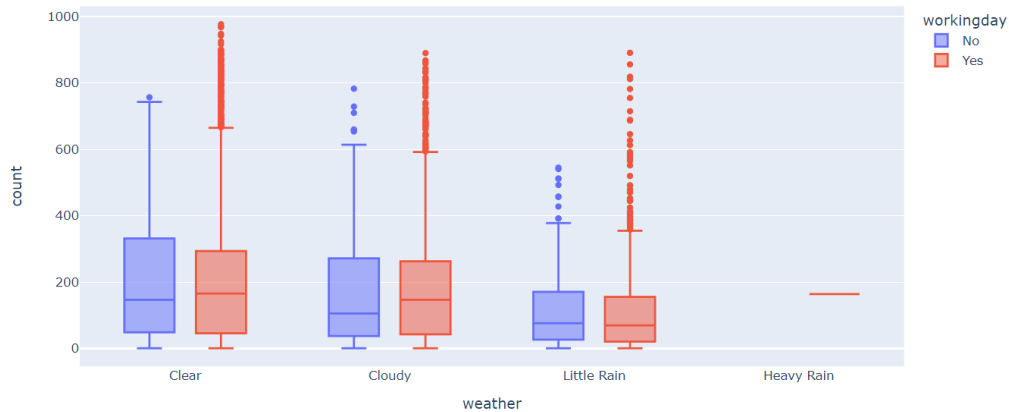
70% of the cycles were rented when it was clear weather.

24% when it was cloudy weather .

during rainy weather , only around 5% of the cycles were rented.

DISTRIBUTIONS and Comparision of number of cycles rented during working days and off day , across different seasons.

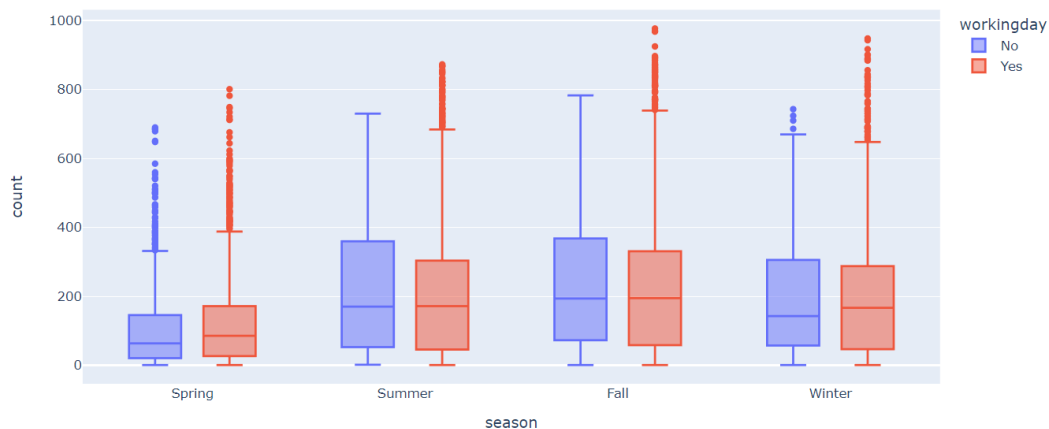
Number of cycles rented Boxplot during Workday and Offday as per different weather conditions



from above boxplot, we can say , there's no significant activity during heavy rain weather.
High activity during clear and cloudy weather.

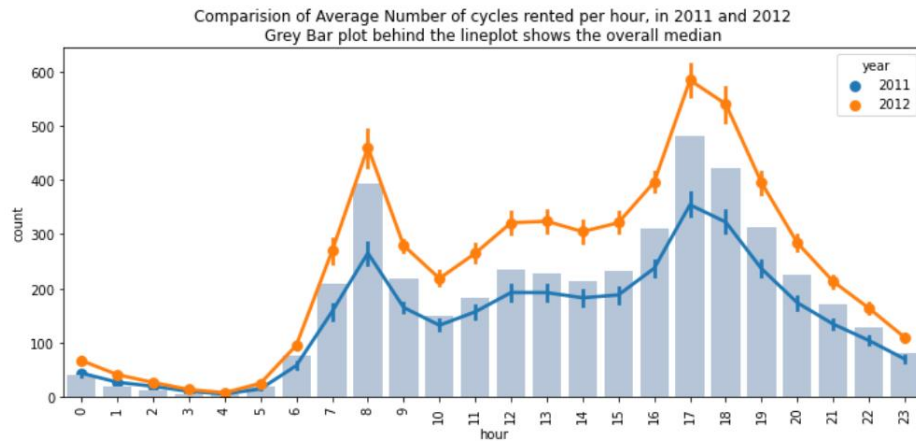
- **Boxplot - distribution of number of bike rented , during different seasons as per workingday or not!**

Number of cycles rented Boxplot during Workday and Offday as per different seasons



during spring season , number of bike rented were lower than summer and fall.

YEARLY DIFFERENCE IN NUMBER OF BIKE RENTAL :



hourly average bike rented in year 2011 and 2012

```
data.groupby("year")["count"].median()
```

```
year
2011    111.0
2012    199.0
```

```
data.groupby("year")["casual"].median()
```

```
year
2011    13.0
2012    20.0
Name: casual, dtype: float64
```

```
data.groupby("year")["registered"].median()
```

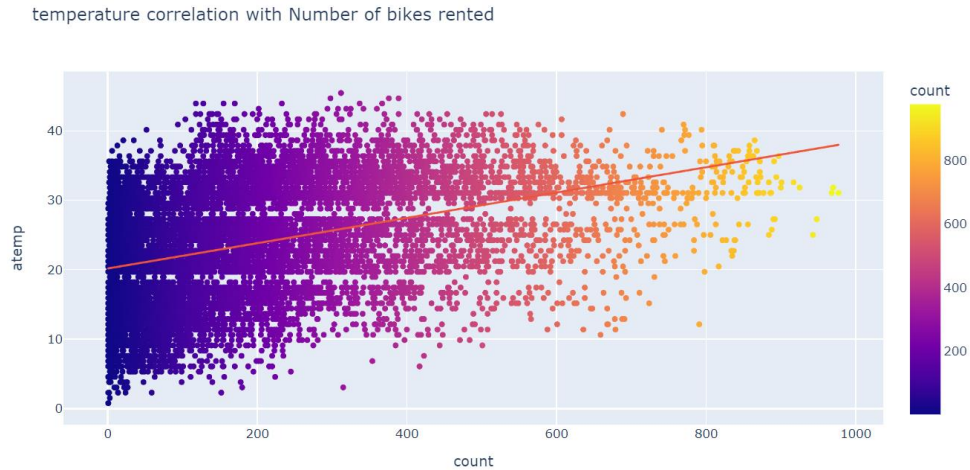
```
year
2011    91.0
2012   161.0
Name: registered, dtype: float64
```

from 2011 , there's 79.27% hike in overall hourly median number of bike rental in 2012.

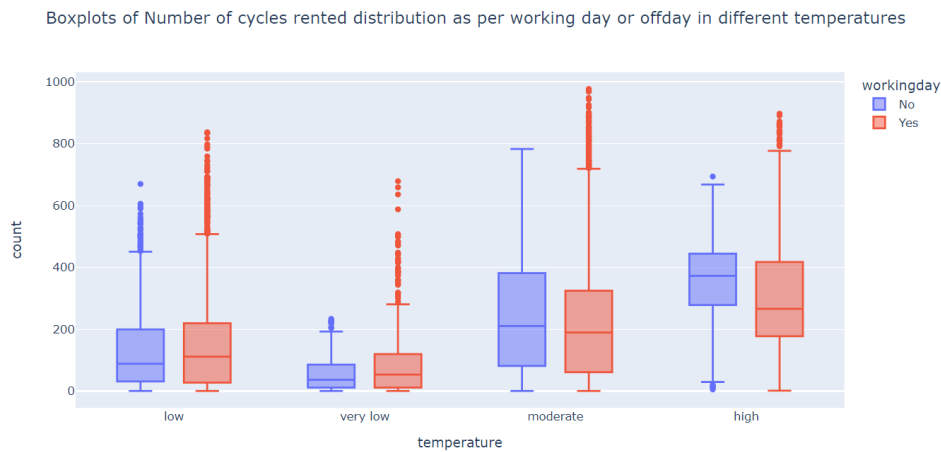
in registered customers , 76% hike in hourly median cycle rental from 2011 to 2012.

in 2011 , median number of hourly rental were 13 , and in 2012 , its 20. -

NUMBER AND CYCLES RENTED AND TEMPERATURE CORRELATION :



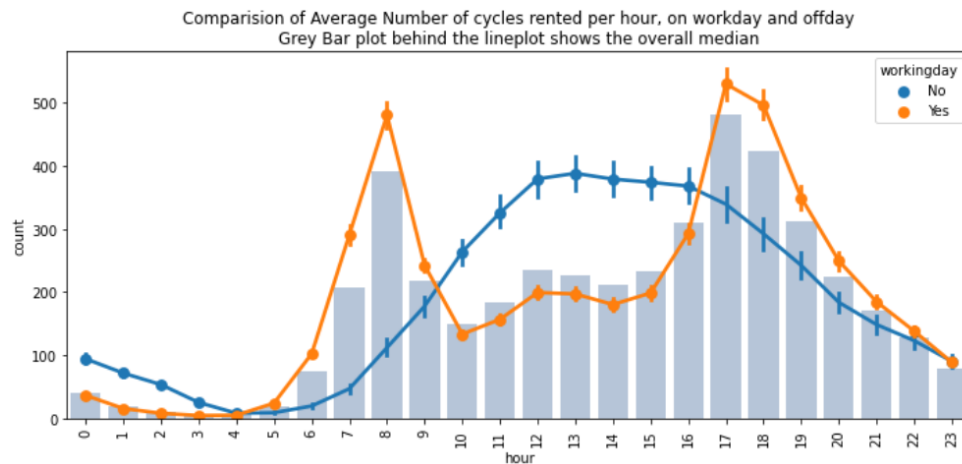
- from scatter plot , there's a positive correlation across temperature and number of bikes rented.
- After categorising the temperature as low, verylow, moderate, high :



from above boxplot :

number of bike rented during moderate to high temperature is significantly higher than lower temperature.

OFFDAY VS WORKING DAY NUMBER OF CYCLES RENTED TREND DURING A DAY :



number of cycles rented changed as per working day and off-day . trend is opposit.

on off days , number of cycles rented increases during the day time ! which is opposite of during working days.

from above plot it looks like, working day count of cycle rented seems to be higher than offday! lets do a AB test : weather mean of rented cycled on working day and offdays are same or not !

hourly median number of cycles rented during

```
data.groupby("workingday")["count"].median()
```

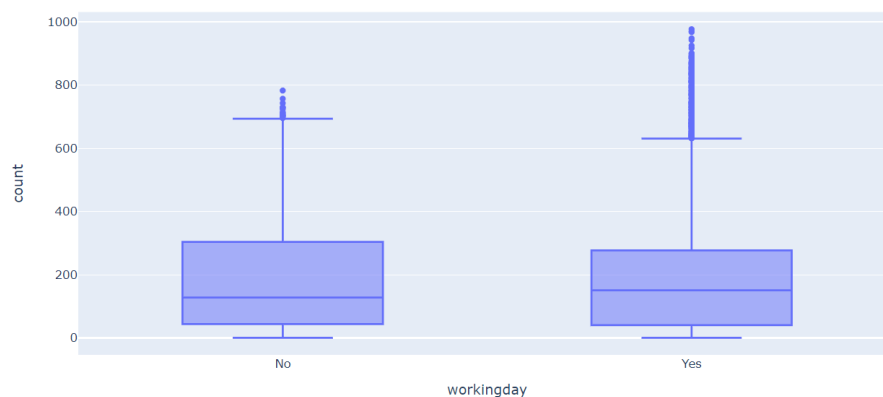
```
workingday
No      128.0
Yes     151.0
Name: count, dtype: float64
```

hourly average number of cycles rented during

```
data.groupby("workingday")["count"].mean()
```

```
workingday
No      188.506621
Yes     193.011873
Name: count, dtype: float64
```

Boxplot shows the distribution of number of bikes rented on offdays and workingdays



- from above boxplot ,
- distributions of hourly number of bike rented during working day and off day seems similar .
- though there are more outliers in workinday category.

testing if mean number of electric cycles rented on workday is equal to on offday !

t-test :

If working day and offday has an effect on the number of electric cycles rented.

distribution of number of bikes rented as per working day or offday (in percentages)

- Establishing Hypothesis :

H0: average # of cycles rented on workingdays = average # of cycles rented on offday

Ha: average # of cycles rented on workingdays != average # of cycles rented on offday

calculating Test Statistic :

```
T_observed =(m1-m2)/(np.sqrt(((s1**2)/n1)+((s2**2)/n2)))  
T_observed
```

1.236258041822322

p-Value :

```
p_value = 2*(1-stats.t.cdf(T_observed,n1+n2-2))  
p_value
```

0.2163893399034813

Extream Critical Value

```
T_critical = stats.t.ppf(0.975,n1+n2-2)  
T_critical
```

1.9601819678713073

```
p_value > 0.05
```

True

```
-T_critical < T_observed < T_critical
```

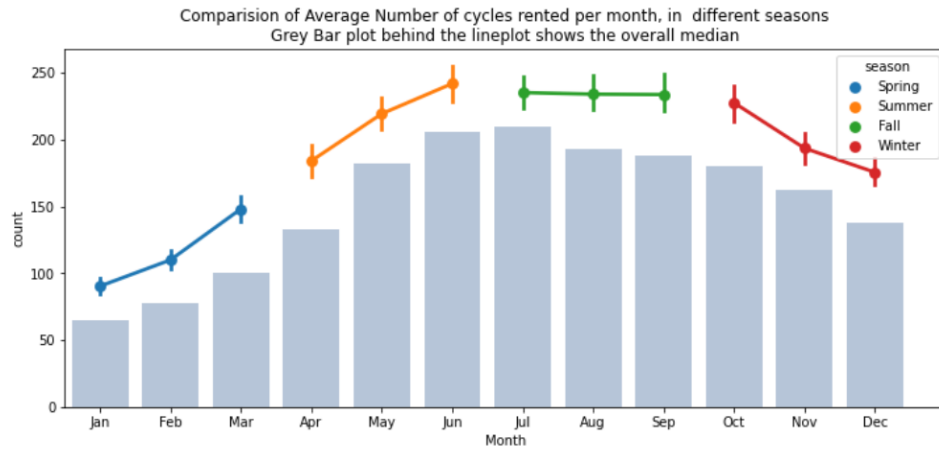
True

we failed to reject null Hypothesis

mean of number of cycles rented on

working days are equal as the cycles rented on offdays.

MONTH AND SEASON WISE , EFFECT ON MEDIAN AND AVERAGE NUMBER OF CYCLES RENTED .

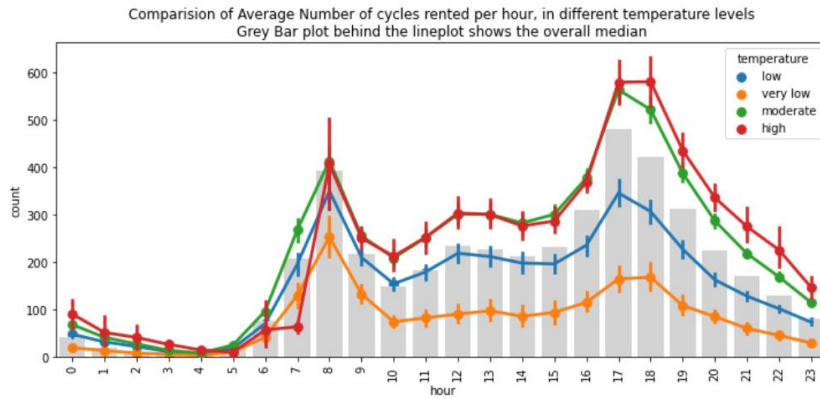


cycle rental counts decreased during winter season and opening spring season .

During Summer season , count increase and stays a constant till pre-winter season .

From May to November the number of cycles rented are increasing

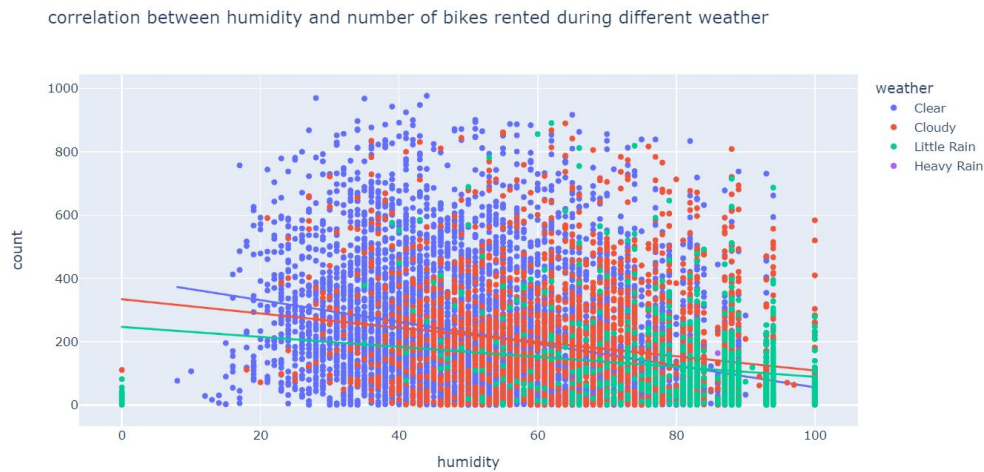
TEMPERATURE EFFECT ON CYCLE RENTAL



Average Number of Bikes rented are higher in moderate to high temperature.

which decreases when temperature is low to very low!

HUMIDITY VS COUNT

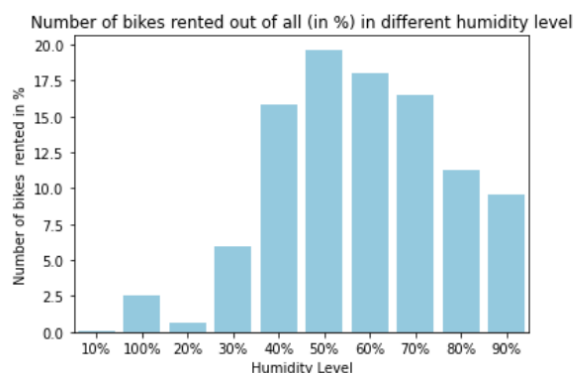


Scatter plot above , shows kind of a negative correlation , between humidity and number of bikes rented. After Categorising Humidity level , we can see

10%	0.038696
100%	2.565314
20%	0.635970
30%	5.942528
40%	15.798887
50%	19.659541
60%	18.030512
70%	16.507215
80%	11.268459
90%	9.552879

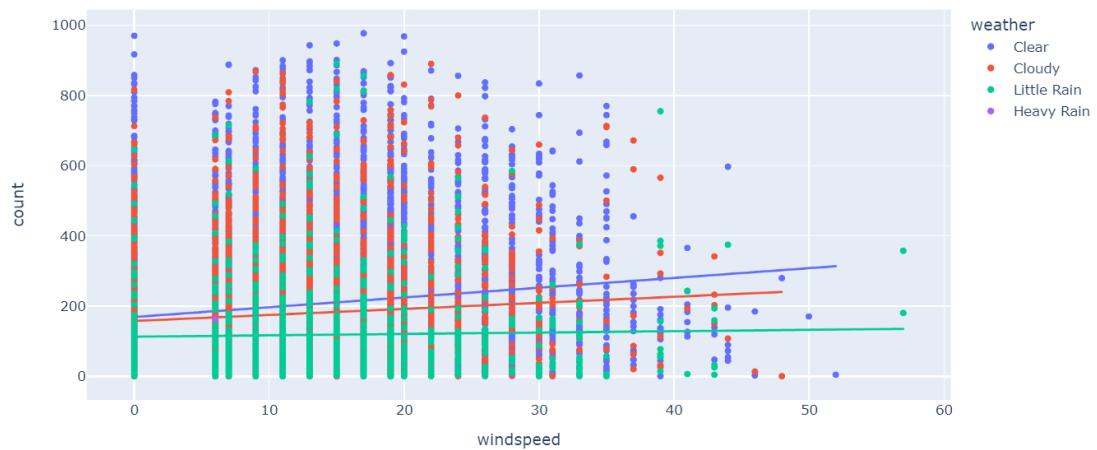
Counts are increasing from humidity level of 40% to 70% .

40 to 70% humidity level seems to be most comfortable for cycling.

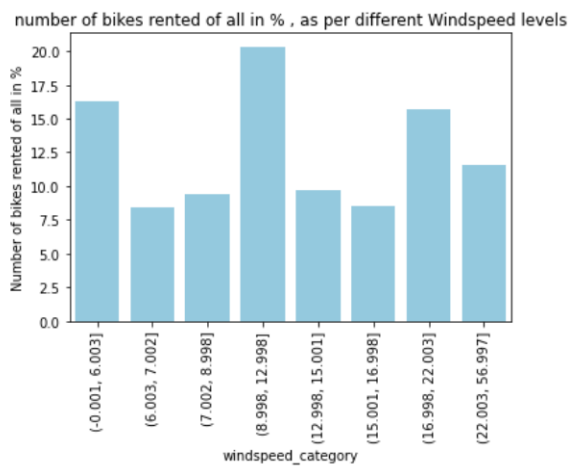


WINDSPEED VS COUNT:

Correlation of Windspeed with Count of bikes rented during different weather



(-0.001, 6.003]	16.325482
(6.003, 7.002]	8.421435
(7.002, 8.998]	9.433002
(8.998, 12.998]	20.356743
(12.998, 15.001]	9.715336
(15.001, 16.998]	8.488901
(16.998, 22.003]	15.682703
(22.003, 56.997]	11.576398



from above, plot:

windspeed are categorised in different groups .

Windspeed increases , the number of bike rented are decreases.

Most often windspeed is 8 to 24.

TEST FOR INDEPENDENCE BETWEEN FEW CATEGORICAL FEATURES. :

IF WEATHER IS DEPENDENT ON THE SEASON

chi-square test : for independence :

H0: weather and seasons are independent

Ha: weather and seasons are dependent

weather	Clear	Cloudy	Little Rain
season			
Fall	470116	139386	31160
Spring	223009	76406	12919
Summer	426350	134177	27755
Winter	356588	157191	30255

```
T_observed
```

```
10838.372332480216
```

```
df = (len(observed)-1)*(len(observed.columns)-1)
```

```
T_critical = stats.chi2.ppf(0.95,df)
```

```
T_critical
```

```
12.591587243743977
```

```
p_value = 1-stats.chi2.cdf(T_observed,df)
```

```
p_value
```

```
0.0
```

```
if T_observed > T_critical:
    print("Reject Null Hypothesis : \nWeather and Season are dependent variables")
else:
    print("Failed to Reject Null Hypothesis : \nWeather and Season are independent Variables")
```

```
Reject Null Hypothesis :
Weather and Season are dependent variables
```

From ChiSquare test of independence :

We reject Null hypothesis as independence:

Conclude that weather and seasons are Dependent Features.

IF WEATHER AND TEMPERATURE ARE DEPENDENT :

for dependency : chi square test :

H0: weather and temperature are independent

Ha: weather and temperature are dependent

```
: chi2Test_of_independence(observed_temp_weather)
```

temperature	high	low	moderate	very low
Clear	52538	56379	177592	3391
Cloudy	11496	23163	51780	807
Little Rain	1726	3249	9869	139

temperature	high	low	moderate	very low
Clear	48616.205381	61207.181565	176870.279678	3206.333375
Cloudy	14631.146791	18420.426916	53229.473683	964.952610
Little Rain	2512.647828	3163.391519	9141.246638	165.714015

T_statistic : 2979.804

p_value : 0.0

Reject Null Hypothesis

"Weather and Ttemperature are dependent variables"

IF WEATHER AND HUMIDITY LEVEL ARE DEPENDENT :

for dependency : chi square test :

H0: weather and Humidity are independent

Ha: weather and Humidity are dependent

T_statistic : 75755.823

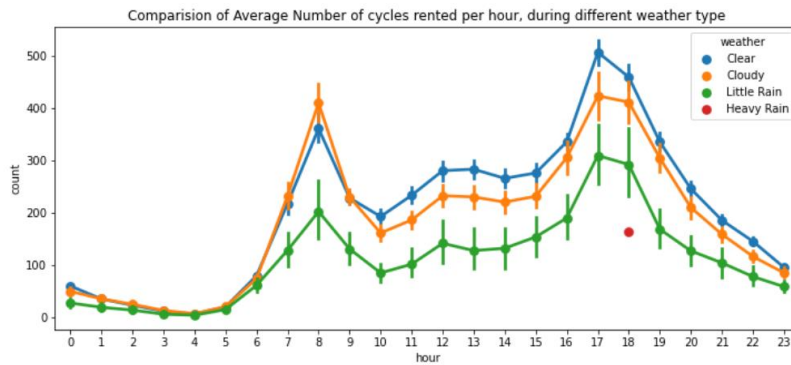
p_value : 0.0

Reject Null Hypothesis

From the dependency test :

we can conclude that weather and humidity are dependent features.

IF THE DISTRIBUTION OF NUMBER OF CYCLES RENTED ARE SIMILAR IN DIFFERENT WEATHER.



- we have 4 different weather here, to check if there's significant difference between 4 weathers, we can perform anova test :

H0: population mean of number of cycles rented in different seasons are same

Ha: population mean of number of cycles rented in different seasons are different

Since the datasets for tests, are not normally distributed, and having significance variance between weathers ,

we cannot perform anova test .

non parametric test : Kruskal Wallis test :

```
H = ((12/(N*(N+1))) * (np.sum(((rank_sum**2)/(kr.groupby("weather")["rank"].count())))) - (3*(N+1)))
H
```

```
: 204.95101790400076
```

```
p_value = 1-stats.chi2.cdf(205.073,degree_of_freedom)
p_value
```

```
: 0.0
```

```
H_critical = stats.chi2.ppf(0.95,2)
H_critical
```

```
: 5.991464547107979
```

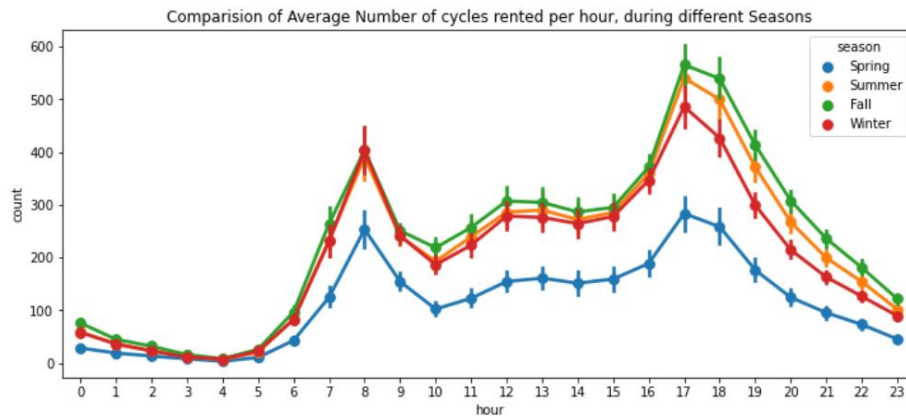
H statistic from Kruskal Wallis test , is higher than the Critical Value ,

p_value is smaller than significant value 0.05 ,

we reject Null Hypothesis.

Hence we conclude that the Population mean number of cycles rented across different weather are not same.

IF NO. OF CYCLES RENTED IS SIMILAR OR DIFFERENT IN DIFFERENT SEASONS



Since the datasets for tests, are not normally distributed, and having significance varinace between all seasons ,

we can use non parametric test : Kruskal Wallis test :

```
H = (((12/(N*(N+1))))*(np.sum(((rank_sum**2)/(kr.groupby("season")["rank"].count())))))-(3*(N+1))
H
699.6499424783542

p_value = 1-stats.chi2.cdf(205.073,degree_of_freedom)
p_value
0.0

H_critical = stats.chi2.ppf(0.95,degree_of_freedom)
H_critical
7.814727903251179

H > H_critical
True
```

H statistic from Kruskal Wallis test , is higher than the Critical Value , p_value is smaller than significant value 0.05 ,

we reject Null Hypothesis.

Hence we conclude that the Population mean number of cycles rented across different Seasons are not same.

INFERENCES AND RECOMMENDATIONS :

There is a positive Correlation between Temperature and Number of cycles rented.
Demand increases with the rise in the temperature from moderate to not very high.

As per shows in the charts in the file , till certain level of humidity level , demand increases , when humidity is too low or very high , there are very few observations.
Humidity level , 40% to 70% highest records have been observed.

As per hourly average number of cycles rented by registered and casual customer plots ,
Registered Customers seems to be using rental cycles mostly for work commute purposes.

registered customers are much higher than the casual customers. 81% customers are Registered and 19% only are casual riders. Which is good thing for a consistent business. Though it is recommended to introduce more go-to offers and strategical execution to attract more casual riders, that further increase chances of converting to consistent users.
Confidence interval of average number of cycles rented by registered customers is (153,157) and casual customers is (35,37).

Demand for cycles increases during the rush hours specifically during working days , from morning 7 to 9 am and in evening 4 to 8pm.
on off days demands are higher from 10 am to evening 7pm.
Though it is concluded from statistical tests, that demand on weekdays and off-days are similar. We can say demand is equal with 95% confidence.

During spring season , customers prefer less likely to rent cycle. demand increases in summer and fall season.
From May to October, demand is increasing .

During clear and cloudy weather demand is higher than in rainy weather.

in 2012 , there's 180% hike in demand , from 2011.
in registered customers , its been 176% hike , where casual customers in 2013 were average 13 to in 2012 are 20.

statistical test results shows,

average number of cycles rented during working days and off days are significantly similar.
weather and seasons are dependent.
Weather and temperature , Weather and humidity level are also dependent .

There's significance difference in demand during different weather and seasons .