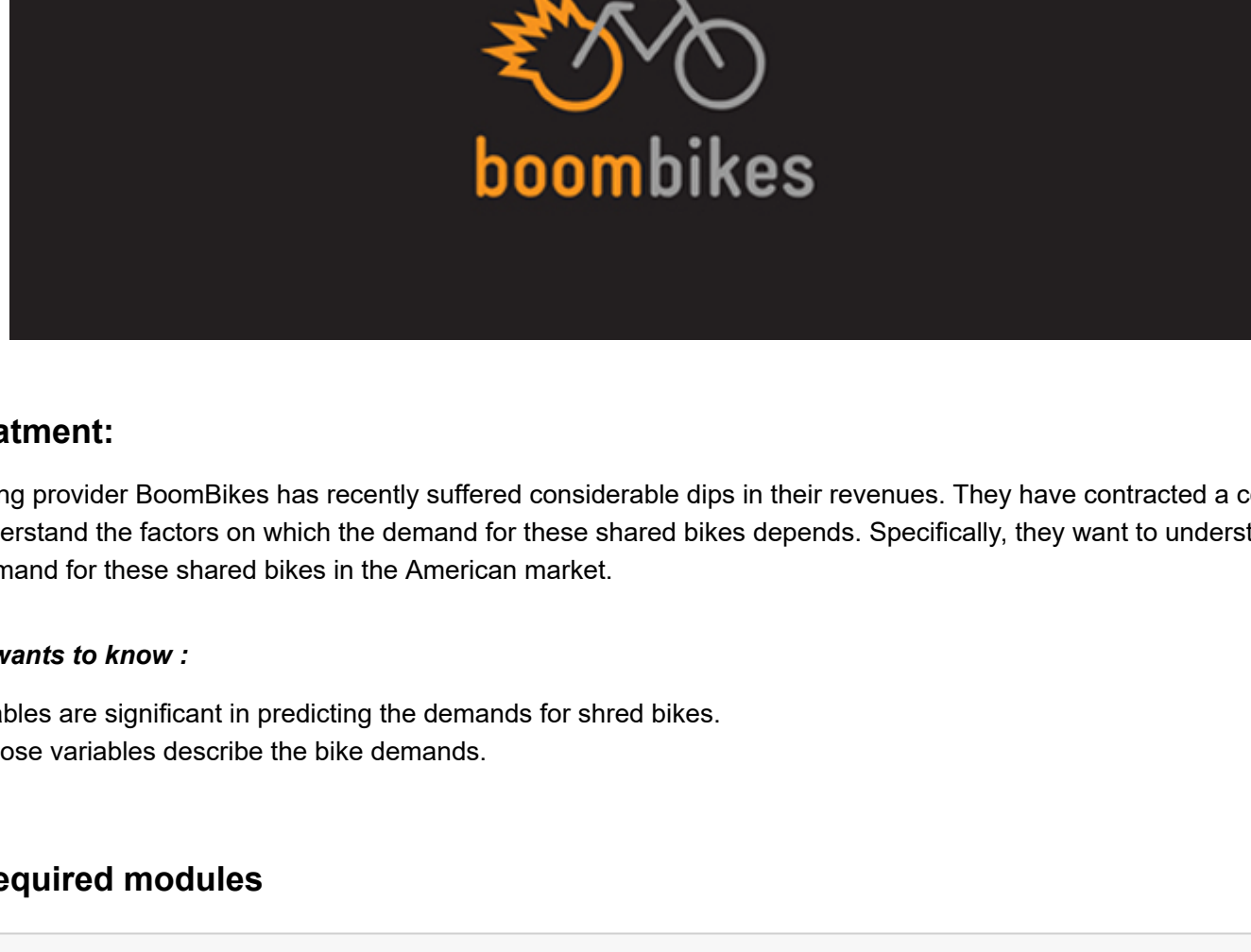


Boom Bikes



Problem Statement:

A US bike-sharing provider BoomBikes has recently suffered considerable dips in their revenues. They have contracted a consulting company to understand the factors on which the demand for these shared bikes depends. Specifically, they want to understand the factors affecting the demand for these shared bikes in the American market.

The company wants to know :

- Which variables are significant in predicting the demands for shared bikes.
- How well those variables describe the bike demands.

Importing required modules

```
In [41]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

Importing statsmodel to get the detailed statistics summary of the trained model

```
In [42]: import statsmodels.api as sm
from statsmodels.stats.outliers_influence import variance_inflation_factor
```

Importing sklearn required modules for performing Linear Regression

```
In [43]: import sklearn
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import Normalizer
from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
```

```
In [44]: # Reading Dataset
df = pd.read_csv('day.csv')
df.head()
```

```
Out [44]:
```

	instant	day	season	yr	mnth	holiday	weekday	workingday	weathersit	temp	atemp	hum	windspeed	casual	regis
0	1	01-01-2018	1	0	1	0	6	0	2	14.110847	18.18125	80.5833	10.749882	331	
1	2	02-01-2018	1	0	1	0	0	0	2	14.902598	17.68695	69.6087	16.652113	131	
2	3	03-01-2018	1	0	1	0	1	1	1	8.050924	9.47025	43.7273	16.636703	120	
3	4	04-01-2018	1	0	1	0	2	1	1	8.200000	10.60610	59.0435	10.739832	108	
4	5	05-01-2018	1	0	1	0	3	1	1	9.305237	11.46350	43.6957	12.522300	82	

```
In [45]: df.shape
Out [45]: (730, 16)
```

```
In [46]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 730 entries, 0 to 729
Data columns (total 16 columns):
#   Column      Non-Null Count  Dtype
---  --
0   instant      730 non-null    int64
1   day          730 non-null    object
2   season       730 non-null    int64
3   yr           730 non-null    int64
4   mnth         730 non-null    int64
5   holiday      730 non-null    int64
6   weekday      730 non-null    int64
7   workingday   730 non-null    int64
8   weathersit    730 non-null    int64
9   temp         730 non-null    float64
10  atemp        730 non-null    float64
11  hum          730 non-null    float64
12  windspeed    730 non-null    float64
13  casual       730 non-null    int64
14  registered   730 non-null    int64
15  cnt          730 non-null    int64
dtypes: float64(4), int64(11), object(1)
memory usage: 91.4+ KB
```

Data Cleaning

- All the columns contains 0 null values
- dayday columns needs to get converted into dateime format as well as the name of the column should also be standardised
- instant column contains the pre-populated index, we can drop this column as it is of no use
- Season and year need to be converted into its actual categorical value at first

```
In [47]: df.drop('instant', axis=1, inplace=True)
```

```
In [48]: df['dteday'] = pd.to_datetime(df['day'])
```

Mapping the numerical values to its respective description

- season : season (1:spring, 2:summer, 3:fall, 4:winter)
- yr : year (0:18, 1:2019)
- mnth : (1:Jan,2:Feb,3:Mar,...,12:Dec)
- weekdays : days ('Monday',2:'Tuesday',3:'Wednesday',.....,0:'Sunday')
- weathersit : climate (1:'Clear', 2:'Cloudy', 3:'LightSnow Rain', 4:'HeavySnow Rain')

```
In [49]: # Year
df.yr = df.yr.map([1: 2019, 0 : 2018])

# Season
df.season = df.season.map([1:'spring', 2:'summer', 3:'fall', 4:'winter'])

# Month
df.mnth = df.mnth.map([1:'Jan',2:'Feb',3:'March',4:'April',
                    5:'May',6:'June',7:'July',8:'Aug',9:'Sep',
                    10:'Oct',11:'Nov',12:'Dec'])

# weekday
df.weekday = df.weekday.map([1:'Monday',2:'Tuesday',3:'Wednesday',4:'Thursday',
                            5:'Friday',6:'Saturday',0:'Sunday',])
```

weathersit :

- 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

```
In [50]: df.weathersit = df.weathersit.map([1:'Clear', 2:'Cloudy', 3:'LightSnow Rain', 4:'HeavySnow Rain'])
```

```
In [51]: print("Season :- {}".format(df.season.unique()))
print("Years :- {}".format(df.yr.unique()))
print("Months :- {}".format(df.mnth.unique()))
print("Weekday :- {}".format(df.weekday.unique()))
print("Weather Situation :- {}".format(df.weathersit.unique()))

Season :- ['spring' 'summer' 'fall' 'winter']
Years :- ['2018' '2019']
Months :- ['Jan' 'Feb' 'March' 'April' 'May' 'June' 'July' 'Aug' 'Sep' 'Oct' 'Nov' 'Dec']
Weekday :- ['Saturday' 'Sunday' 'Monday' 'Tuesday' 'Wednesday' 'Thursday' 'Friday']
Weather Situation :- ['Cloudy' 'Clear' 'LightSnow Rain']
```

Looking to the weather situation we can say that the rider never rented a boom bike during a heavy rain or during storm which was expected

```
In [52]: df.head()
```

```
Out [52]:
```

	dteday	season	yr	mnth	holiday	weekday	workingday	weathersit	temp	atemp	hum	windspeed	casual	registe
0	2018-01-01	spring	18	Jan	0	Saturday	0	Cloudy	14.110847	18.18125	80.5833	10.749882	331	65
1	2018-02-01	spring	18	Jan	0	Sunday	0	Cloudy	14.902598	17.68695	69.6087	16.652113	131	67
2	2018-03-01	spring	18	Jan	0	Monday	1	Clear	8.050924	9.47025	43.7273	16.636703	120	122
3	2018-04-01	spring	18	Jan	0	Tuesday	1	Clear	8.200000	10.60610	59.0435	10.739832	108	145
4	2018-05-01	spring	18	Jan	0	Wednesday	1	Clear	9.305237	11.46350	43.6957	12.522300	82	151

Changing name of the columns as desired

```
In [53]: df.rename(columns = {'dteday':'date','yr':'year','mnth':'month','hum':'humidity','cnt':'count', inplace = True)
```

Let see the dataset information

```
In [54]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 730 entries, 0 to 729
Data columns (total 15 columns):
#   Column      Non-Null Count  Dtype
---  --
0   date         730 non-null    datetime64[ns]
1   season       730 non-null    object
2   year         730 non-null    int64
3   month        730 non-null    object
4   holiday      730 non-null    int64
5   weekday      730 non-null    object
6   workingday   730 non-null    int64
7   weathersit    730 non-null    object
8   temp         730 non-null    float64
9   atemp        730 non-null    float64
10  humidity     730 non-null    float64
11  windspeed    730 non-null    float64
12  casual       730 non-null    int64
13  registered   730 non-null    int64
14  count        730 non-null    int64
dtypes: datetime64[ns](1), float64(4), int64(6), object(4)
memory usage: 85.7+ KB
```

As we can see that the categorical columns have been changed from integer type to object type as well as the date data type has been changed to Datetime data type

```
In [55]: # Statistical Summary
df.describe()
```

```
Out [55]:
```

	year	holiday	workingday	temp	atemp	humidity	windspeed	casual	registered	count
count	730.000000	0.000000	730.000000	730.000000	730.000000	730.000000	730.000000	730.000000	730.000000	730.000000
mean	2018.500000	0.626767	0.683562	20.319259	23.728322	62.761575	12.763820	849.249315	3658.757534	4508.006849
std	0.500343	0.167266	0.465405	2.424346	3.953498	8.152038	1.500244	686.479875	1559.758728	1836.011647
min	2018.000000	0.000000	0.000000	14.138185	16.869713	52.000000	9.041650	316.250000	2502.250000	3169.750000
25%	2018.000000	0.000000	0.000000	20.486828	24.368225	62.825000	12.125325	717.000000	3684.500000	4548.500000
50%	2018.500000	0.000000	1.000000	20.468828	24.368225	62.825000	12.125325	717.000000	3684.500000	4548.500000
75%	2018.500000	0.000000	1.000000	20.880615	30.445775	72.989575	15.625589	1096.500000	4783.250000	5968.000000
max	2019.000000	1.000000	1.000000	35.328347	42.044800	97.250000	34.000021	3410.000000	6948.000000	8714.000000

Data Exploration

Do we have any repetitive date in the Dataset?

```
In [56]: df.groupby('date').agg(['date':'count']).shape
```

```
Out [56]: (730, 1)
```

So, the total no's of rows comes out to be 730 after grouping on the date column which is same as the nos of rows in the actual dataset so we don't have any repetitive dates in our dataset

Summarizing the date column in the dataset

```
In [57]: print("dataset starting date = {}".format(df.date.min().date()))
print("dataset ending date = {}".format(df.date.max().date()))
s = df.date.min().date()
e = df.date.max().date()
print("No's of days between max and min date = {}".format((e-s).days))

dataset starting date = 2018-01-01
dataset ending date = 2019-12-31
No's of days between max and min date = 729
```

```
In [58]: total_days = df.groupby(['weathersit','workingday']).agg(['workingday':'count'])
total_days.rename(columns=['workingday':'count of days', inplace = True)
```

```
In [59]: total_days.pivot_table(index='weathersit',columns='workingday',
                             margins=True,
                             margins_name='total', # defaults to 'All'
                             aggfunc=sum)

# total_days
```

```
Out [59]:
```

count of days		
weathersit	workingday	
Clear	0	158
	1	307
Cloudy	0	70
	1	178
LightSnow Rain	0	5
	1	16
total		730

Earlier while looking into the unique values of the weather situation we stated that the riders don't use the bike rental service during storm or heavy rain which seems to be right as expected but on looking to the total no's of days on which the bike has been rented we can say that the firm never faced the Heavy Rain or storm situation

So, we don't have enough data for the heavy rain situation to conclude any statement but ideally the providers should avoid renting out their bikes during this situation also take an additional steps to face this situation in the future

lets see the spread of the renters data across different weather situation using bee swarm plot

```
In [60]: plt.figure(figsize=(16,6))
sns.swarmplot(x='weathersit', y='count', hue='workingday',
              data=df, palette="Set1", dodge=True)
plt.title('Climatic Condition Vs Count of the Renters')
plt.show()
```



From above plot we can state that during clear weather the numbers of riders renting the Bikes are more as compare to the Mist/Cloudy weather

- we can see that the riders use to access the services more during Mist/Cloudy or Clear weather
- during light snow/rain and the user won't prefer to avail the services if it's non-working day but during working days they do have rented the bike a very few nos of times

Now let's analysed on which week day the services has been used more with in combination of weathar

```
In [61]: # Though we have seen the spread of the data lets also explore its density
df.weathersit.value_counts(normalize=True)
```

```
Out [61]:
```

Clear	0.634247
Cloudy	0.336986
LightSnow Rain	0.028767

Name: weathersit, dtype: float64

Though the spread across the weekdays was seemed to be the same but the density and user preference lies towards Clear weather

```
In [62]: df.weekday.value_counts(normalize=True).apply(lambda x:str(round(x*100,2)+'%'))

Out [62]:
```

Monday	14.38%
Sunday	14.38%
Saturday	14.38%
Friday	14.25%
Thursday	14.25%
Tuesday	14.25%
Wednesday	14.11%

Name: weekday, dtype: object

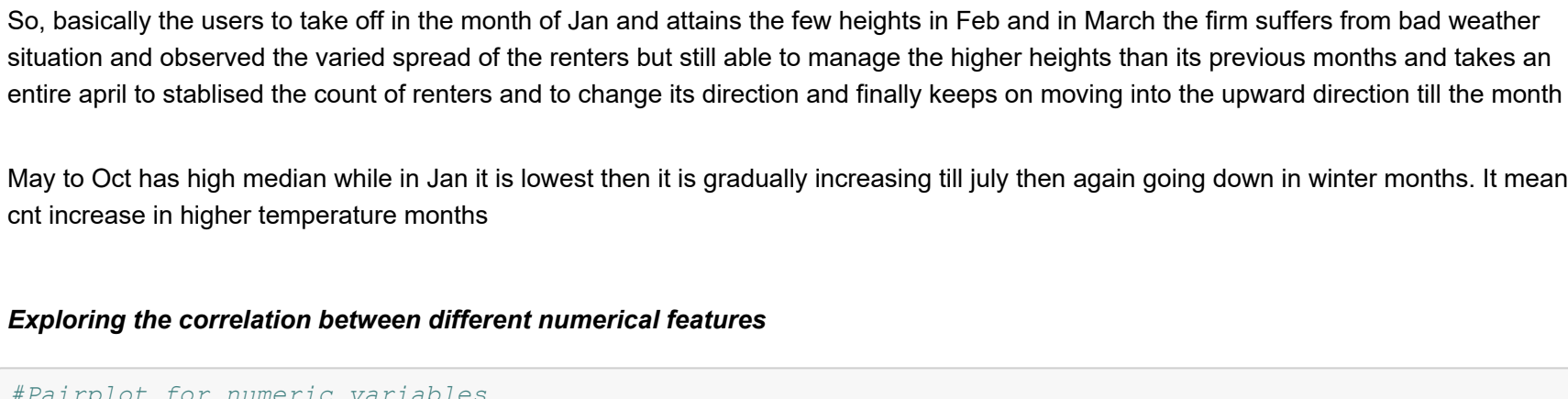
```
In [63]: %matplotlib inline

fig, axes = plt.subplots(1, 2, sharey=True, figsize=(15,8))
plt.suptitle('WEEKDAYS VS SPREAD OF RENTERS')

sns.swarmplot(x='weekday', y='count', hue='weathersit',
              data=df, palette=['LightSnow Rain', 'Clear', 'Cloudy'],
              order=[df['weathersit'].value_counts().index[2], df['weathersit'].value_counts().index[1], df['weathersit'].value_counts().index[0]],
              palette=['Green','Blue'],dodge=True,ax=axes[0])

sns.swarmplot(x='weekday', y='count', hue='weathersit',
              data=df, palette=['LightSnow Rain', 'Clear', 'Cloudy'],
              order=[df['weathersit'].value_counts().index[2], df['weathersit'].value_counts().index[1], df['weathersit'].value_counts().index[0]],
              palette=['Blue','Green'],dodge=True,ax=axes[1])

plt.xlabel('Weekend')
plt.show()
```



There is not significant increase in the numbers but if we compare the weekends as compare to any days of the weekdays excluding Friday(day=5) the numbers of the user renting the bikes have been increase and among weekdays the Friday is the most popular day irrespective of the weather

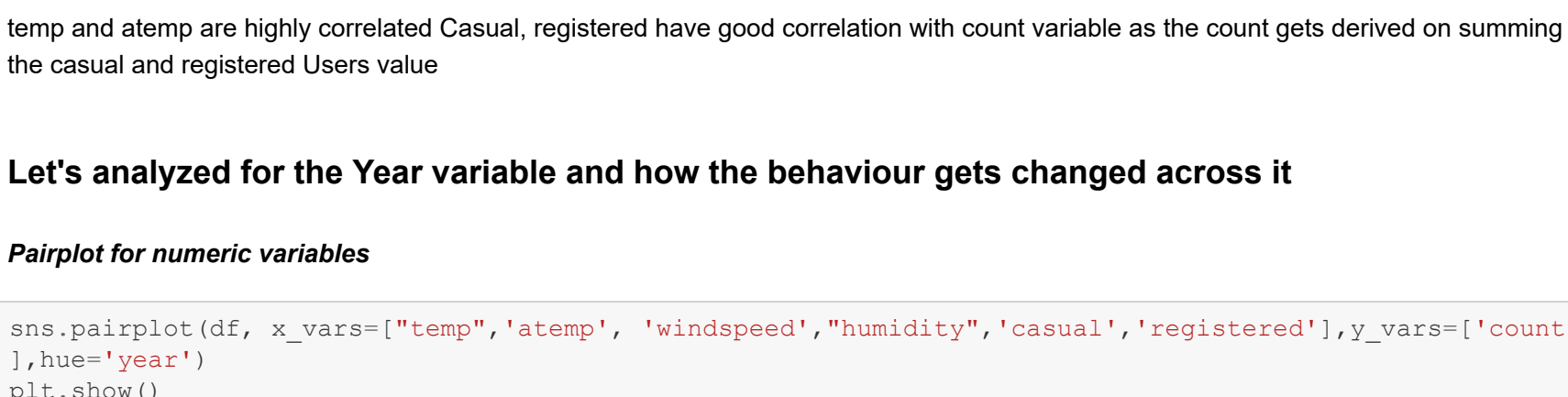
The count variable has been derived by summing up the registered users and the casual users. So, lets see the contribution of each class of users

```
In [64]: rental_user = df.groupby('month', as_index=False, sort=False).agg(['count':'sum','registered':'sum','casual':'sum'])
rental_user
```

```
Out [64]:
```

	month	count	registered	casual
0	Jan	134933	122891	12042
1	Feb	143518	134620	14898
2	March	228620	184476	44444
3	April	269094	208292	60802
4	May	331686	256401	75285
5	June	346342	272438	73906
6	July	344948	266791	78157
7	Aug	351194	279155	72039
8	Sept	345991	278668	70323
9	Oct	322352	262592	59760
10	Nov	254831	218228	36603
11	Dec	211036	189343	21693

```
In [65]: plt.figure(figsize=(10,5))
sns.lineplot(x='month',month,rented_user.casual,color='b',label='Casual User',sort=False)
sns.lineplot(x='month',month,rented_user.registered,color='g',dashes=True,label='Registered User',sort=False)
plt.xlabel('Types of Users and their respective interactions')
plt.ylabel('Month')
plt.title('Casual User vs Registered User Monthly')
plt.show()
```



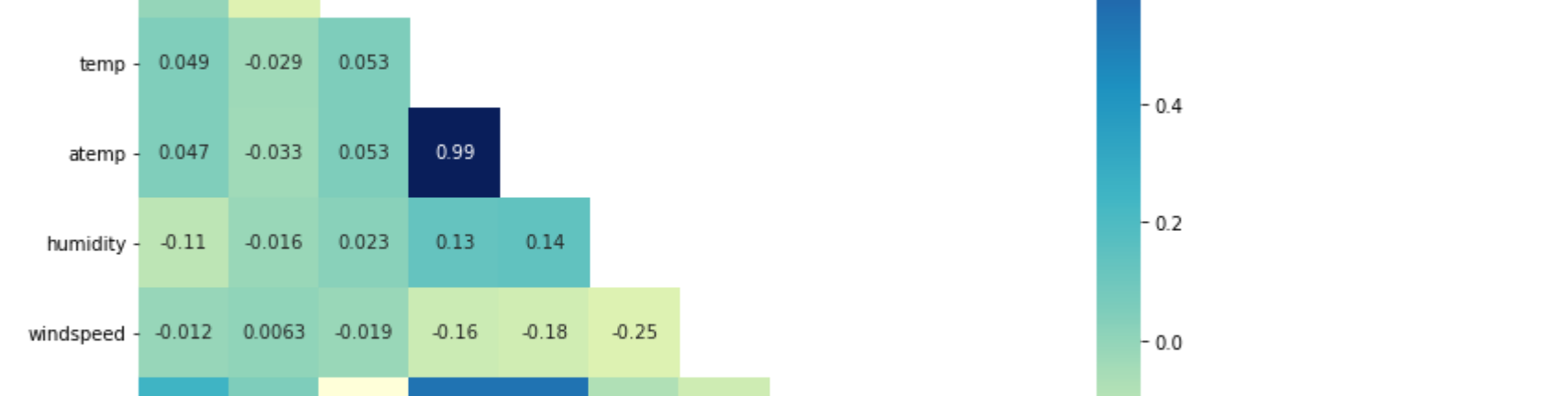
Clearly the Booms Bike get rented by the registered user more frequently as compared to the their casual users on the same note, the cnt of casual user is also good in numbers and can contribute or affect the firms revenue directly

Also, the firm can increase the margin of price for the non-registered user to earn the extra penny

The Box plot unlike the Boxes

On which month the wind is blowing and in what direction my bees are flying?

```
In [66]: plt.figure(figsize=(16,6))
# sns.lineplot(x='month', y='count', data=df, color='red')
sns.violinplot(x='month', y='count', data=df)
sns.boxplot(x='month', y='count', data=df)
plt.title('Renters Density Across Months')
plt.show()
```



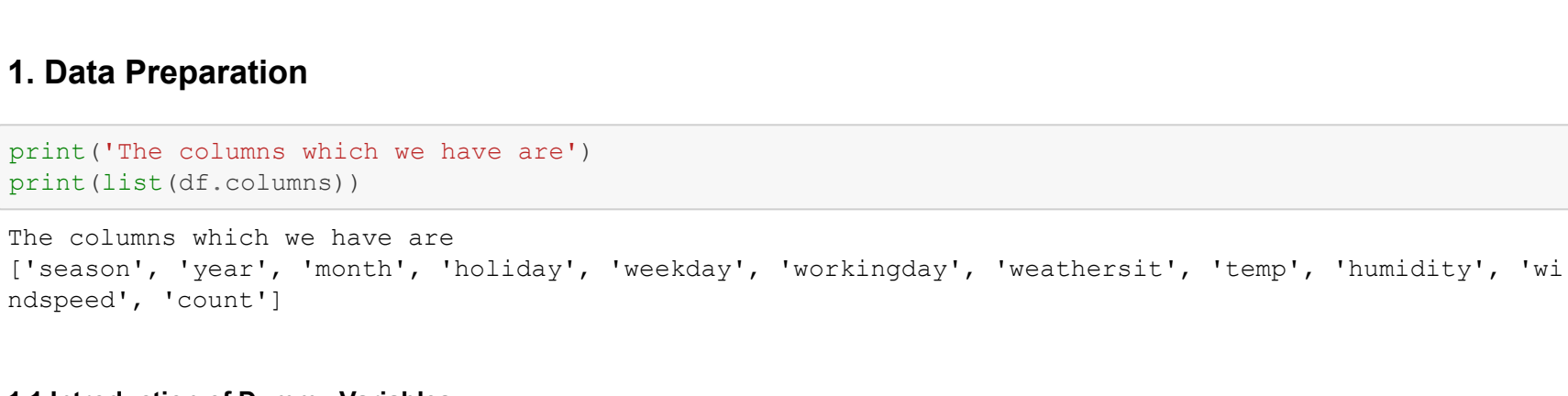
Consider the the renter as bees and tail representing the lower density of the renters and the head showing the higher density of the renters and more stable on the flowers bed(June, July)....

So, basically the users take off in the month of Jan and attains the few heights in Feb and in March the firm suffers from bad weather season and observed that the rented cases of the renters but still able to manage the higher heights than its previous months and takes an entire april to the stabilised count of renters and to change its direction finally keeps on moving into the upward direction till the month

May to Oct has high median while in Jan it is lowest then it is gradually increasing till July then again going down in winter months. It means cnt increase in higher temperature months

Exploring the correlation between different numerical features

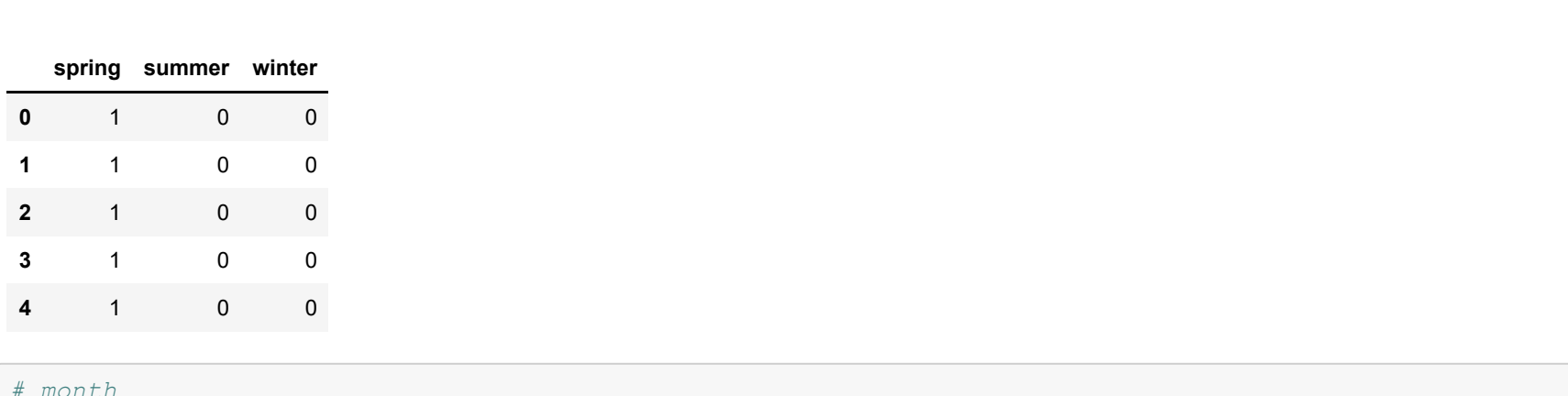
```
In [67]: #sns.pairplot for numeric variables
sns.pairplot(df, vars=['temp','atemp','windspeed','humidity','casual','registered','count'])
plt.show()
```



temp and atemp are highly correlated Casual, registered have good correlation with count variable as the count gets derived from summing the casual and registered Users value

Let's analyzed for the Year variable and how the behaviour gets changed across it

```
In [68]: sns.pairplot(df, x_vars=['temp','atemp','windspeed','humidity','casual','registered'], y_vars=['count',
                                                'year'])
plt.show()
```



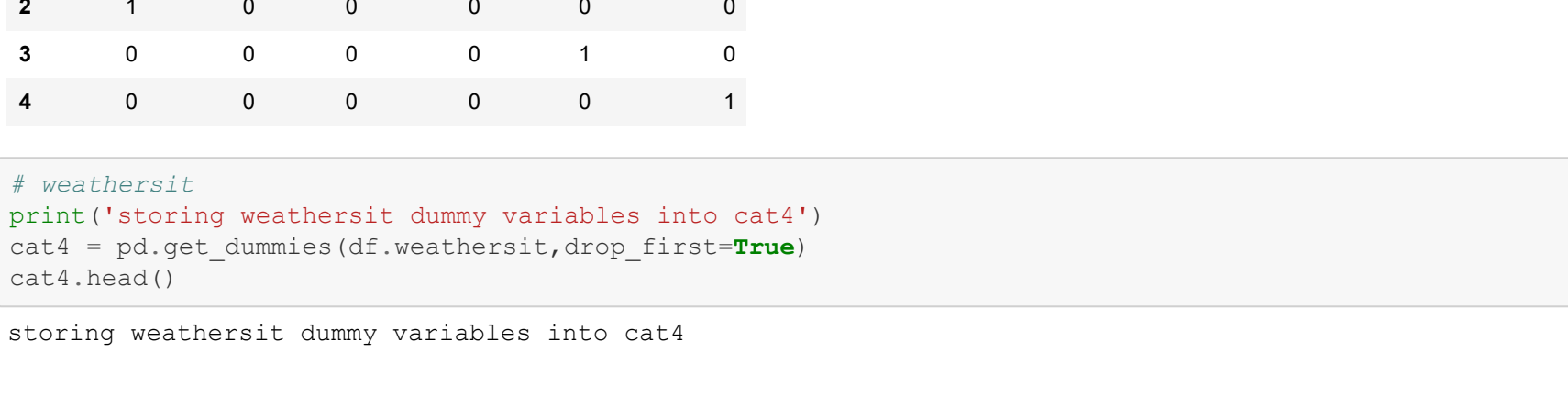
From above pairplot we can infer that the no's of casual users are getting increased as compared to the previous year while the registered user show slightly downward trend as compared to the previous year

Correlation Heat Map

```
In [69]: corr = df.corr()
# Generate a mask for the upper triangle
mask = np.triu(np.ones_like(corr, dtype=bool))

# Set up the matplotlib figure
fig, ax = plt.subplots(figsize=(11, 9))

# Draw the heatmap with the mask and correct aspect ratio
sns.heatmap(corr, mask=mask, cmap='YlGnBu', annot = True)
plt.show()
```



temp and atemp both have the same correlation coefficient as well as they are highly correlated to each other. So, it need to analyzed while building model and need to drop one variable while building the model to avoid multicollinearity using VIF and p-value

- On the other hand, year and casual users, and registered users are positively correlated to the count variable
- Also, the target variable gets derived from the casual and registered feature so moving forward we will not consider these columns

Drop the unnecessary variables from the dataset

We can see the dataset has some variables that are not required. We can drop date, casual, registered, atemp.

```
In [70]: #drop unnecessary columns
df.drop(['date','casual', 'registered','atemp'], axis=1, inplace=True)
df.head()
```

```
Out [70]:
```

	season	year	month	holiday	weekday	workingday	weathersit	temp	humidity	windspeed	count
0	spring	2018	Jan	0	Saturday	0	Cloudy	14.110847	80.5833	10.749882	985
1	spring	2018	Jan	0	Sunday	0	Cloudy	14.902598	69.6087	16.652113	801
2	spring	2018	Jan	0	Monday	1	Clear	8.050924	43.7273	16.636703	1349
3	spring	2018	Jan	0	Tuesday	1	Clear	8.200000	59.0435	10.739832	1562
4	spring	2018	Jan	0	Wednesday	1	Clear	9.305237	43.6957	12.522300	1600


```
In [79]: df_model = pd.concat([df_cat1,cat2,cat3,cat4], axis=1)

Out [81]: df_model.drop(['year', 'month', 'weekday', 'weatherisit'], axis =1, inplace = True)

In [82]: print('Concatenating the cat1,cat2,cat3,cat4 & dropping table as well as dropping unwanted repeated column')
df_model.drop(['year', 'month', 'weekday', 'weatherisit'], axis =1, inplace = True)

Out [83]: df_model.head()
```

Concatenating the cat1,cat2,cat3,cat4 & dropping table as well as dropping unwanted repeated columns

```
Out [81]: year                2
           holiday          2
           Tuesday         2
           Thursday        2
           Sunday          2
           Saturday        2
           Monday          2
           Sep              2
           Feb              2
           Nov              2
           May              2
           March            2
           Cloudy           2
           July             2
           Feb              2
           Dec              2
           Aug              2
           spring           2
           summer           2
           winter           2
           LightSnow Rain   2
           temp             496
           humidity         594
           windspeed        649
           count            695
           dtype: int64
```

```
In [82]: print('The Final Model Dataset')
The Final Model Dataset
```

```
In [83]: df_model.head()
```

```
Out [83]: year holiday workingday temp humidity windspeed count spring summer winter ... Oct Sep Monday Saturday Sunday
0 0 0 0 14.1108673 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
1 0 0 0 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
2 0 0 0 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
3 0 0 0 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
4 0 0 0 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
```

5 rows × 29 columns

2. Splitting data into Train and Test

```
In [84]: df_model_train, df_model_test = train_test_split(df_model, train_size= 0.7, test_size = 0.3, random_state = 100)
```

```
In [85]: print('Shape of Train Dataset (i)').format(df_model_train.shape)
df_model_train.head()
```

```
Out [85]: Shape of Train Dataset (510, 29)
df_model_train.head()
```

```
year holiday workingday temp humidity windspeed count spring summer winter ... Oct Sep Monday Saturday Sunday
653 1 0 1 14.1108673 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
576 1 0 1 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
426 1 0 0 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
728 1 0 0 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
482 1 0 0 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
```

5 rows × 29 columns

```
In [88]: print('Shape of Test Dataset (i)').format(df_model_test.shape)
df_model_test.head()
```

```
Out [88]: Shape of Test Dataset (219, 29)
df_model_test.head()
```

```
year holiday workingday temp humidity windspeed count spring summer winter ... Oct Sep Monday Saturday Sunday
184 0 1 1 14.1108673 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
576 1 0 1 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
299 0 0 0 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
221 0 0 0 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
192 0 0 0 14.902598 15.03333 10.780287 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
```

5 rows × 29 columns

3. Scaling the required variables to bring all of the dimension on the same scale

If a feature in the dataset is big in scale compared to others then in algorithms where Euclidean distance is measured this big scaled feature becomes dominating and needs to be normalized.

```
In [89]: scaler = MinMaxScaler()
variables = ['temp', 'humidity', 'windspeed']
df_model_train[variables] = scaler.fit_transform(df_model_train[variables])
```

```
In [90]: df_model_train.head()
```

```
Out [90]: year holiday workingday temp humidity windspeed count spring summer winter ... Oct Sep Monday Saturday Sunday
653 1 0 1 0.509887 0.575354 0.307094 7534 0 0 0 1 ... 1 0 0 0 0 0 0 0
576 1 0 1 0.815169 0.725633 0.265869 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
426 1 0 0 0.442393 0.640178 0.255342 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
728 1 0 0 0.245101 0.498067 0.063106 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
482 1 0 0 0.395666 0.504508 0.188475 7534 0 0 0 0 ... 0 0 0 0 0 0 0 0
```

5 rows × 29 columns

```
In [91]: print('Statistical Summary')
df_model_train.describe()
```

```
Out [91]: Statistical Summary
year holiday workingday temp humidity windspeed count spring summer winter ...
count 510.000000 510.000000 510.000000 510.000000 510.000000 510.000000 510.000000 510.000000 510.000000 ... 510
mean 0.507843 0.025490 0.676471 0.537262 0.650369 0.320768 0.48638233 0.243137 0.245098 0.24902 ... 0
std 0.500429 0.157763 0.468282 0.225844 0.145884 0.169797 0.952158739 0.429398 0.430668 0.43287 ... 0
min 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 22.000000 0.000000 0.000000 0.00000 ... 0
50% 0.000000 0.000000 0.000000 0.000000 0.540519 0.653714 0.296763 0.320.000000 0.000000 0.000000 ... 0
95% 1.000000 0.000000 1.000000 0.735215 0.754830 0.414447 5973.500000 0.000000 0.000000 0.00000 ... 0
max 1.000000 0.000000 1.000000 1.000000 1.000000 1.000000 8714.000000 1.000000 1.000000 1.00000 ... 1
```

8 rows × 29 columns

As we can see that the maximum value for all the variables is 1 while min value is 0 so the whole features are now in between the range of 0 to 1. So all the features are now on the same scale

4. Let's check the correlation coefficients to see which variables are highly correlated

```
In [273]: # Let's check the correlation coefficients to see which variables are highly correlated
plt.figure(figsize=(25,20))
ax = sns.heatmap(df_model_train.corr(), annot = True, cmap='Greens', fmt='.2g', linewidth=1)
plt.yticks(size=15)
plt.xticks(size=15)
# plt.show()
```

```
for lab in ax.get_xticklabels():
    text = lab.get_text()
    if text == 'count': # lets highlight row 2
        # set the properties of the ticklabel
        lab.set_weight('bold')
        lab.set_size(20)
        lab.set_color('purple')

for lab in ax.get_yticklabels():
    if text == 'count': # lets highlight row 2
        # set the properties of the ticklabel
        lab.set_weight('bold')
        lab.set_size(20)
        lab.set_color('purple')
```


inferences

- workingday variable has high negative correlation with Sat & Sun (where workingday=0)
- Spring is negatively correlated with cnt
- emp, atemp and yr has strong correlation with cnt
- misty weather and humidity has correlation
- various months and corresponding weather has correlation

We can see that temperature, Summer season, June to October months are in good correlation with the 'count' variable. And seem to have good influence on the number of bike rentals.

5. Splitting the Data into X and y train

```
In [94]: y_train = df_model_train.pop('count')
```

```
In [95]: X_train = df_model_train
```

```
In [97]: print('Count as Target variable (ytrain)')
pd.info(X_train.head())
```

```
Out [97]: Count as Target variable (ytrain)
count
653 7534
576 7216
426 4066
728 1796
482 4220
```

```
In [98]: print('Xtrain as')
X_train.head()
```

```
Out [98]: year holiday workingday temp humidity windspeed spring summer winter Aug ... Oct Sep Monday Saturday Sunday
653 1 0 1 0.509887 0.575354 0.307094 0 0 0 1 0 ... 1 0 0 0 0 0 0 0
576 1 0 1 0.815169 0.725633 0.265869 0 0 0 0 0 ... 0 0 0 0 0 0 0 0
426 1 0 0 0.442393 0.640178 0.255342 1 0 0 0 0 ... 0 0 0 0 0 0 0 0
728 1 0 0 0.245101 0.498067 0.063106 1 0 0 0 0 ... 0 0 0 0 0 0 0 0
482 1 0 0 0.395666 0.504508 0.188475 0 0 0 0 0 ... 0 0 0 0 0 0 0 0
```

5 rows × 28 columns

6. Training Model

Model 1

Considering all the variables

```
In [99]: # Add Constant
X_train_sm = sm.add_constant(X_train)
# Create a fitted model in one line
lr_1 = sm.OLS(y_train, X_train_sm).fit()
print('Model 1 summary')
print(lr_1.summary())
```

```
Model 1 summary
OLS Regression Results
=====
Dep. Variable: count R-squared: 0.853
Model: OLS Adj. R-squared: 0.845
Method: Least Squares F-statistic: 103.8
Date: Thu, 01 Oct 2020 Prob (F-statistic): 8.74e-182
Time: 12:25:11 Log-Likelihood: -4097.8
No. Observations: 510 AIC: 8252.
Df Residuals: 482 BIC: 8370.
Df Model: 27
Covariance Type: nonrobust
=====
```

```
coef std err t >|t| [0.025 0.975]
const 2147.8745 306.272 7.013 0.000 1546.081 2749.668
year 0.000 0.000 0.000 1.000 -0.000 0.000
holiday 2017.5498 70.005 28.820 0.000 1879.997 2155.103
workingday 90.9621 208.622 0.436 0.663 -318.958 500.883
temp -1620.6522 223.329 -7.257 0.000 -2059.472 -1181.833
humidity -418.2822 260.478 -1.570 0.115 -939.472 152.908
windspeed 336.5941 227.730 1.478 0.140 -110.872 784.061
spring 920.0214 242.512 3.794 0.000 -1972.595 -658.801
summer 125.1994 292.671 0.428 0.669 -449.870 700.268
winter -258.2326 291.836 -0.883 0.375 -869.668 352.195
Jan -280.5935 285.798 -0.982 0.327 -842.157 280.970
Feb -545.6541 291.390 -1.873 0.062 -1118.206 26.998
Mar -350.7234 324.698 -1.081 0.280 -999.424 287.978
Apr -205.8360 217.957 -0.945 0.342 -644.101 232.419
May -9.0845 213.174 0.043 0.966 -409.781 427.950
Jun 328.0850 182.530 1.800 0.072 59.217 596.953
Jul -363.7716 315.893 -1.152 0.250 -984.668 256.125
Aug 65.4819 310.994 0.211 0.833 -545.589 676.553
Sep 704.5726 278.125 2.533 0.012 158.086 1251.059
Oct -181.0213 333.066 -0.543 0.589 -849.889 367.799
Nov 842.2102 119.136 7.066 0.000 608.002 1076.418
Dec 368.1868 121.011 3.043 0.002 130.413 605.960
LightSnow Rain 134.832 324.457 0.416 0.675 -345.217 614.645
temp -164.1905 314.665 -0.523 0.599 -982.794 100.413
humidity -46.8282 129.428 -0.362 0.718 -301.140 207.484
windspeed 30.7891 90.677 -0.340 0.735 -109.261 352.617
count -237.6527 228.832 -1.038 0.300 -685.364 178.022
Omnibus: 84.475 Durbin-Watson: 2.040
Prob(Omnibus): 0.000 Jarque-Bera (JB): 158.454
Skew: -0.804 Prob(JB): 7.72e-16
Kurtosis: 5.914 Cond. No. 1.34e+16
```

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The smallest eigenvalue is 8.83e-30. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

From Above P-values:

year, workingday, temp, humidity, windspeed, spring, summer, winter, Sep, Sunday, Saturday, Cloudy, LightSnow Rain are statistically significant

Recursive Feature Elimination

RFE is an efficient approach for eliminating features from a training dataset for feature selection.

Using RFE we will select top 15 features which mostly describe our dependent/target variables

```
In [100]: from sklearn.feature_selection import RFE
lr = LinearRegression()
rfe = RFE(lr, 15)
rfe = rfe.fit(X_train, y_train)
```

```
In [101]: #list of variables selected
list(zip(X_train.columns, rfe.support_, rfe.ranking_))
```

```
Out [101]: (('year', True, 1),
('holiday', True, 1),
('workingday', True, 1),
('temp', True, 1),
('humidity', True, 1),
('windspeed', True, 1),
('spring', True, 1),
('summer', True, 1),
('winter', True, 1),
('Aug', False, 6),
('Dec', False, 3),
('Feb', False, 4),
('Jan', True, 1),
('July', True, 1),
('May', False, 13),
('March', False, 1),
('Sep', True, 1),
('Nov', False, 2),
('Oct', False, 10),
('Saturday', True, 1),
('Sunday', False, 8),
('Thursday', False, 11),
('Tuesday', False, 9),
('Wednesday', False, 12),
('Cloudy', True, 1),
('LightSnow Rain', True, 1))
```

```
In [281]: #Columns where RFE support is True
feature = X_train.columns[rfe.support_]
feature = list(feature)
```

```
Out [281]: ['year',
'holiday',
'workingday',
'temp',
'humidity',
'windspeed',
'spring',
'summer',
'winter',
'Jan',
'July',
'Nov',
'Oct',
'Monday',
'Sunday',
'Tuesday',
'Thursday',
'Wednesday',
'Cloudy',
'LightSnow Rain']
```

```
In [282]: #Columns where RFE support is False
col = X_train.columns[~rfe.support_]
col = list(col)
```

```
Out [283]: ['Aug',
'Dec',
'Feb',
'June',
'March',
'May',
'Nov',
'Oct',
'Monday',
'Sunday',
'Tuesday',
'Thursday',
'Wednesday',
'Cloudy',
'LightSnow Rain']
```

Model 2

Build using 15 Featured Variables using Statsmodel for detail analysis and further optimization of the model

```
In [284]: # Adding constant as statsmodel wont provide constant by default and assumes line to be passed through the origin
X_train_2 = sm.add_constant(X_train[feature])
```

```
In [286]: lr = sm.OLS(y_train, X_train_2).fit()
print(lr.summary())
```

```
OLS Regression Results
=====
Dep. Variable: count R-squared: 0.844
Model: OLS Adj. R-squared: 0.844
Method: Least Squares F-statistic: 184.6
Date: Tue, 29 Sep 2020 Prob (F-statistic): 2.31e-191
Time: 23:15:14 Log-Likelihood: -4105.9
No. Observations: 510 AIC: 8246.
Df Residuals: 496 BIC: 8322.
Df Model: 15
Covariance Type: nonrobust
=====
```

```
coef std err t >|t| [0.025 0.975]
const 2204.0556 314.032 7.019 0.000 1587.052 2821.059
year 2003.5240 69.550 28.807 0.000 1866.873 2140.175
holiday -506.7381 233.568 -2.178 0.030 -967.646 -49.830
workingday 78.2822 100.557 0.778 0.438 -181.693 547.872
temp -1431.4379 238.389 -6.007 0.000 -1918.228 -884.507
humidity -1344.7234 324.457 -4.145 0.000 -1982.280 -707.167
windspeed -466.3289 221.379 -2.107 0.034 -905.289 -18.369
spring -441.6037 179.462 -2.461 0.014 -794.207 -89.001
summer 820.0214 242.512 3.382 0.000 335.217 1308.829
winter -258.2326 291.836 -0.883 0.375 -869.668 352.195
Jan -280.5935 285.798 -0.982 0.327 -842.157 280.970
Feb -545.6541 291.390 -1.873 0.062 -1118.206 26.998
Mar -350.7234 324.698 -1.081 0.280 -999.424 287.978
Apr -205.8360 217.957 -0.945 0.342 -644.101 232.419
May -9.0845 213.174 0.043 0.966 -409.781 427.950
Jun 328.0850 182.530 1.800 0.072 59.217 596.953
Jul -363.7716 315.893 -1.152 0.250 -984.668 256.125
Aug 65.4819 310.994 0.211 0.833 -545.589 676.553
Sep 704.5726 278.125 2.533 0.012 158.086 1251.059
Oct -181.0213 333.066 -0.543 0.589 -849.889 367.799
Nov 842.2102 119.136 7.066 0.000 608.002 1076.418
Dec 368.1868 121.011 3.043 0.002 130.413 605.960
LightSnow Rain 134.832 324.457 0.416 0.675 -345.217 614.645
temp -164.1905 314.665 -0.523 0.599 -982.794 100.413
humidity -46.8282 129.428 -0.362 0.718 -301.140 207.484
windspeed 30.7891 90.677 -0.340 0.735 -109.261 352.617
count -237.6527 228.832 -1.038 0.300 -685.364 178.022
Omnibus: 66.344 Durbin-Watson: 2.071
Prob(Omnibus): 0.000 Jarque-Bera (JB): 158.454
Skew: -0.661 Prob(JB): 7.72e-16
Kurtosis: 5.403 Cond. No. 22.1
```

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

The difference b/w Aug, Resquare and Resquare score is 0.005 which is very less

Jan followed by holiday has high p-value among selected variables

Let's check VIF method to further justify this model

VIF

Variance Inflation Factor or VIF, gives a basic quantitative idea about how much the feature variables are correlated with each other. It is an extremely important parameter to test our linear model. The formula for calculating VIF is:

$$VIF_i = \frac{1}{1 - R_i^2}$$

```
In [287]: # Creating dataframe to store a VIF information
X_train_vif = X_train[feature]
vif_df.DataFrame()
vif['features'] = X_train_vif.columns
vif['VIF'] = [variance_inflation_factor(X_train_vif.values, i) for i in range(15)]
vif['VIF'] = round(vif['VIF'], 2)
vif = vif.sort_values(by = "VIF", ascending = False).reset_index(drop=True)
vif
```

```
Out [287]: features VIF
0 humidity 29.40
1 temp 17.77
2 workingday 5.31
3 windspeed 4.73
4 spring 4.53
5 winter 3.48
6 summer 2.84
7 holiday 2.29
8 year 2.09
9 Saturday 1.98
10 Jan 1.67
11 July 1.58
12 Sep 1.36
13 LightSnow Rain 1.25
14 holiday 1.18
```

The Rules to eliminate the Features

- High p-value and High VIF
- Low p-value and Low VIF
- Low p-value and High VIF

Since the Jan has high p-value and low VIF, so dropping the Jan column as it is insignificant in presence of other variables

```
In [288]: feature.remove('Jan')
```

Since, to get the better model we need to repeat the above step so, lets create a function to perform the above operation

```
In [289]: def get_model(cnt, feature, X_train, final=False):
global X_train_sm
global vif
global vif_df
# add constant
X_train_sm = sm.add_constant(X_train[feature])
# fit straight line
lr = sm.OLS(y_train, X_train_sm).fit()
# print summary
print(lr.summary())
vif = vif.sort_values(by = "VIF", ascending = False).reset_index(drop=True)
vif['VIF'] = [variance_inflation_factor(X_train_vif.values, i) for i in range(X_train_vif.shape[1])]
vif['VIF'] = round(vif['VIF'], 2)
print('-----')
print('VIF: center(75)')
print('-----')
print(vif)
```

Model3

```
In [290]: get_model(3, feature, X_train)
```

```
OLS Regression Results
=====
Dep. Variable: count R-squared: 0.842
Model: OLS Adj. R-squared: 0.843
Method: Least Squares F-statistic: 196.3
Date: Tue, 29 Sep 2020 Prob (F-statistic): 1.13e-191
Time: 23:15:14 Log-Likelihood: -4117.1
No. Observations: 510 AIC: 8246.
Df Residuals: 496 BIC: 8322.
Df Model: 14
Covariance Type: nonrobust
=====
```

```
coef std err t >|t| [0.025 0.975]
const 2192.6839 310.741 7.067 0.000 1587.052 2821.059
year 2003.5240 69.550 28.807 0.000 1866.873 2140.175
holiday -506.7381 233.568 -2.178 0.030 -967.646 -49.830
workingday 78.2822 100.557 0.778 0.438 -181.693 547.872
temp -1431.4379 238.389 -6.007 0.000 -1918.228 -884.507
humidity -1344.7234 324.457 -4.145 0.000 -1982.280 -707.167
windspeed -466.3289 221.379 -2.107 0.034 -905.289 -18.369
spring -441.6037 179.462 -2.461 0.014 -794.207 -89.001
summer 820.0214 242.512 3.382 0.000 335.217 1308.829
winter -258.2326 291.836 -0.883 0.375 -869.668 352.195
Jan -280.5935 285.798 -0.982 0.327 -842.157 280.970
Feb -545.6541 291.390 -1.873 0.062 -1118.206 26.998
Mar -350.7234 324.698 -1.081 0.280 -999.424 287.978
Apr -205.8360 217.957 -0.945 0.342 -644.101 232.419
May -9.0845 213.174 0.043 0.966 -409.781 427.950
Jun 328.0850 182.530 1.800 0.072 59.217 596.953
Jul -363.7716 315.893 -1.152 0.250 -984.668 256.125
Aug 65.4819 310.994 0.211 0.833 -545.589 676.553
Sep 704.5726 278.125 2.533 0.012 158.086 1251.059
Oct -181.0213 333.066 -0.543 0.589 -849.889 367.799
Nov 842.2102 119.136 7.066 0.000 608.002 1076.418
Dec 368.1868 121.011 3.043 0.002 130.413 605.960
LightSnow Rain 134.832 324.457 0.416 
```



```
In [293]: feature.append('temp')
feature.remove('July')
get_model(6,feature,X_train)

=====
OLS Regression Results
=====
Dep. Variable: count R-squared: 0.839
Model: OLS Adj. R-squared: 0.836
Method: Least Squares F-statistic: 216.5
Date: Tue, 29 Sep 2020 Prob (F-statistic): 1.48e-168
Time: 23:35:17 Log-Likelihood: -4120.9
No. Observations: 510 AIC: 8268.
DF Residuals: 497 BIC: 8323.
DF Model: 12
Covariance Type: nonrobust

=====
coef std err t >|t| (0.025 0.975)
-----
const 1332.2548 267.028 4.914 0.000 787.612 1836.897
year 2035.3614 70.646 28.729 0.000 1896.107 2174.496
holiday -478.8258 239.645 -1.998 0.046 -949.669 -71.983
workingday 413.2153 102.447 4.033 0.000 211.933 614.498
windspeed -1297.0616 216.282 -5.997 0.000 -1722.001 -872.122
spring -469.1358 178.254 -2.632 0.009 -819.360 -118.911
summer 534.2226 122.036 4.378 0.000 294.452 773.993
winter 853.3356 143.152 5.961 0.000 572.079 1134.592
Cloudy -469.1358 178.254 -2.632 0.009 -819.360 -118.911
LightSnow Rain 1.08
temp -1297.0616 216.282 -5.997 0.000 -1722.001 -872.122
Omnibus: 71.158 Durbin-Watson: 1.092
Prob(Omnibus): 0.000 Jarque-Bera (JB): 170.059
Skew: -0.729 Prob(JB): 1.18e-37
Kurtosis: 5.424 Cond. No. 19.5

=====

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

=====
VIF
=====
features VIF
0 temp 5.70
1 workingday 2.20
2 windspeed 4.63
3 spring 2.40
4 year 2.07
5 summer 2.00
6 Saturday 1.96
7 winter 1.83
8 Cloudy 1.56
9 LightSnow Rain 1.24
10 holiday 1.17
11 LightSnow Rain 1.08
```

Adjusted R2_score changes by 0.002 (0.2%), and won't affect our model

Now, we will eliminate the holiday column as all other variables except holiday have a value of 0

Model7

```
In [294]: feature.remove('holiday')
get_model(7,feature,X_train)

=====
OLS Regression Results
=====
Dep. Variable: count R-squared: 0.838
Model: OLS Adj. R-squared: 0.835
Method: Least Squares F-statistic: 237.5
Date: Tue, 29 Sep 2020 Prob (F-statistic): 6.86e-189
Time: 23:35:18 Log-Likelihood: -4140.8
No. Observations: 510 AIC: 8304.
DF Residuals: 498 BIC: 8321.
DF Model: 11
Covariance Type: nonrobust

=====
coef std err t >|t| (0.025 0.975)
-----
const 1243.5563 265.599 4.682 0.000 721.723 1765.390
year 2036.7262 71.055 28.664 0.000 1897.122 2176.330
workingday 413.2153 102.447 4.033 0.000 211.933 614.498
windspeed -1297.0616 216.282 -5.997 0.000 -1722.001 -872.122
spring -469.1358 178.254 -2.632 0.009 -819.360 -118.911
summer 534.2226 122.036 4.378 0.000 294.452 773.993
winter 853.3356 143.152 5.961 0.000 572.079 1134.592
Cloudy -469.1358 178.254 -2.632 0.009 -819.360 -118.911
LightSnow Rain 1.08
temp -1297.0616 216.282 -5.997 0.000 -1722.001 -872.122
Omnibus: 76.386 Durbin-Watson: 2.073
Prob(Omnibus): 0.000 Jarque-Bera (JB): 190.041
Skew: -0.765 Prob(JB): 5.41e-42
Kurtosis: 5.569 Cond. No. 19.4

=====

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

=====
VIF
=====
features VIF
0 temp 5.48
1 workingday 4.63
2 windspeed 4.63
3 spring 2.30
4 year 2.07
5 summer 1.99
6 Saturday 1.83
7 winter 1.77
8 Cloudy 1.56
9 Sep 1.23
10 LightSnow Rain 1.08
```

The Adjusted R-Square gets reduced by 0.001 or 0.1% which again won't affect our model also, all other variable have high statistical significance

Now, if we look to correlation coefficient we find a negative coefficient among all the features windspeed and LightSnow Rain have high negative coefficient and also VIF value for the windspeed is high.

So, let's try dropping windspeed and evaluate the model

Model8

```
In [295]: feature.remove('windspeed')
get_model(8,feature,X_train)

=====
OLS Regression Results
=====
Dep. Variable: count R-squared: 0.826
Model: OLS Adj. R-squared: 0.823
Method: Least Squares F-statistic: 237.5
Date: Tue, 29 Sep 2020 Prob (F-statistic): 1.55e-182
Time: 23:35:18 Log-Likelihood: -4140.8
No. Observations: 499 AIC: 8304.
DF Residuals: 499 BIC: 8350.
DF Model: 10
Covariance Type: nonrobust

=====
coef std err t >|t| (0.025 0.975)
-----
const 805.3959 264.256 3.048 0.002 266.203 1324.588
year 2027.6600 73.499 27.587 0.000 1883.254 2172.066
workingday 413.2153 102.447 4.033 0.000 211.933 614.498
spring -562.0888 184.357 -3.049 0.002 -924.300 -199.877
summer 448.6179 125.768 3.567 0.000 201.517 695.719
winter 854.9460 148.475 5.758 0.000 563.234 1146.658
Cloudy -469.1358 178.254 -2.632 0.009 -819.360 -118.911
LightSnow Rain 1.08
temp -1297.0616 216.282 -5.997 0.000 -1722.001 -872.122
Omnibus: 77.818 Durbin-Watson: 2.057
Prob(Omnibus): 0.000 Jarque-Bera (JB): 209.021
Skew: -0.753 Prob(JB): 4.09e-46
Kurtosis: 5.751 Cond. No. 18.9

=====

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

=====
VIF
=====
features VIF
0 temp 4.72
1 workingday 4.63
2 year 2.07
3 Saturday 1.81
4 summer 1.74
5 spring 1.70
6 winter 1.59
7 Cloudy 1.56
8 Sep 1.23
9 LightSnow Rain 1.07
```

On removing Windspeed

R-square get change from 0.838 to 0.826 with the difference of 0.012 (1.2%)

Adj. R-square get change from 0.835 to 0.823 with the difference of 0.012(1.2%)

So, now the p-value and the VIF of the features are within the acceptable range

- With VIF values (< 0.05) indicates that we can reject the null hypothesis.
- All other VIF values less than 5, we are in good shape, and can proceed with our regression without having multicollinearity
- We have the R-squared value of 0.826 or 82.6%
- The adjusted R-squared came to be 0.823 or 82.3%. The low differences between R-square and the adj. R-square signifies that all the variables in the model are describing our data points

Also, let's check by removing temp and keeping windspeed into the model and

Model9

```
In [296]: feature.append('windspeed')
feature.remove('temp')
get_model(9,feature,X_train)

=====
OLS Regression Results
=====
Dep. Variable: count R-squared: 0.768
Model: OLS Adj. R-squared: 0.763
Method: Least Squares F-statistic: 237.5
Date: Tue, 29 Sep 2020 Prob (F-statistic): 3.40e-151
Time: 23:35:19 Log-Likelihood: -4214.9
No. Observations: 510 AIC: 8452.
DF Residuals: 499 BIC: 8490.
DF Model: 10
Covariance Type: nonrobust

=====
coef std err t >|t| (0.025 0.975)
-----
const 4630.4923 158.146 29.280 0.000 4319.779 4941.206
year 2150.8308 74.495 28.875 0.000 1984.234 2317.428
workingday 491.7071 115.328 4.264 0.000 265.120 718.295
spring -2580.2866 128.370 -20.100 0.000 -2832.494 -2328.068
summer -469.1358 178.254 -2.632 0.009 -819.360 -118.911
winter -634.8625 122.413 -5.186 0.000 -875.370 -394.355
Sep 627.5042 166.868 3.760 0.000 299.453 955.355
Saturday 534.2226 122.036 4.378 0.000 294.452 773.993
Cloudy -469.1358 178.254 -2.632 0.009 -819.360 -118.911
LightSnow Rain -2617.7535 255.318 -10.253 0.000 -3119.385 -2216.122
windspeed -1496.4872 259.051 -5.778 0.000 -2005.373 -987.602
Omnibus: 32.745 Durbin-Watson: 2.016
Prob(Omnibus): 0.000 Jarque-Bera (JB): 70.965
Skew: -0.356 Prob(JB): 3.89e-16
Kurtosis: 4.683 Cond. No. 10.0

=====

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

=====
VIF
=====
features VIF
0 temp 4.72
1 workingday 4.63
2 year 2.07
3 Saturday 1.81
4 summer 1.74
5 spring 1.70
6 winter 1.59
7 Cloudy 1.56
8 Sep 1.23
9 LightSnow Rain 1.08
```

As we can see, for the above model on removing temp variable and keeping windspeed leads to huge drop in R-Square and adj. R-square which is about (0.5.8%) and (10%) respectively.

```
In [297]: feature.append('temp')
feature.remove('windspeed')
get_model(9,feature,X_train)

=====
OLS Regression Results
=====
Dep. Variable: count R-squared: 0.826
Model: OLS Adj. R-squared: 0.823
Method: Least Squares F-statistic: 1.55e-182
Date: Tue, 29 Sep 2020 Prob (F-statistic): 1.55e-182
Time: 23:35:19 Log-Likelihood: -4140.8
No. Observations: 510 AIC: 8304.
DF Residuals: 499 BIC: 8350.
DF Model: 10
Covariance Type: nonrobust

=====
coef std err t >|t| (0.025 0.975)
-----
const 805.3959 264.256 3.048 0.002 266.203 1324.588
year 2027.6600 73.499 27.587 0.000 1883.254 2172.066
workingday 413.2153 102.447 4.033 0.000 211.933 614.498
spring -562.0888 184.357 -3.049 0.002 -924.300 -199.877
summer 448.6179 125.768 3.567 0.000 201.517 695.719
winter 854.9460 148.475 5.758 0.000 563.234 1146.658
Sep 795.5240 144.460 5.507 0.000 511.700 1079.348
Saturday 562.5839 128.432 4.380 0.000 310.250 814.918
Cloudy -469.1358 178.254 -2.632 0.009 -819.360 -118.911
LightSnow Rain -2651.4976 219.793 -12.064 0.000 -3083.331 -2219.664
temp -4274.7202 293.051 14.587 0.000 3698.955 4850.486
Omnibus: 77.818 Durbin-Watson: 2.057
Prob(Omnibus): 0.000 Jarque-Bera (JB): 209.021
Skew: -0.753 Prob(JB): 4.09e-46
Kurtosis: 5.751 Cond. No. 18.9

=====

Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

=====
VIF
=====
features VIF
0 temp 4.72
1 workingday 4.63
2 year 2.07
3 Saturday 1.81
4 summer 1.74
5 spring 1.70
6 winter 1.59
7 Cloudy 1.56
8 Sep 1.23
9 LightSnow Rain 1.08
```

We will move with the model8 as our final model

y_pred = 0.0901 + year (0.2333) + workingday (0.0568) + spring (-0.0647) + summer (0.0516) + winter (0.0984) + Sep (0.0915) + Saturday (0.0647) + Cloudy (-0.0800) + LightSnow Rain (-0.3051) + temp (0.4918)

R-Square: 0.826

Adj. R-Square: 0.823

Residual Analysis of the train data

So, now to check if the error terms are also normally distributed (which is infact, one of the major assumptions of linear regression), let's plot the histogram of the error terms and see what it looks like.

```
In [298]: # residual = y_pred - y_train
# y_train = y_train.reshape(y_train.shape[0])

In [299]: y_train.shape

Out[299]: (510,)

In [300]: X_train_sml.head(3)

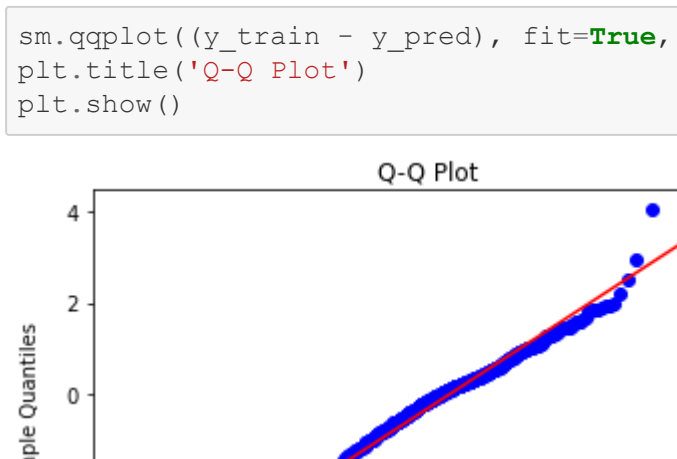
Out[300]:
```

	const	year	workingday	spring	summer	winter	Sep	Saturday	Cloudy	LightSnow Rain	temp
576	1.0	1	1	1	0	0	1	0	0	0	0.500887
653	1.0	1	1	0	0	0	0	0	0	0	0.018169
426	1.0	1	0	1	0	0	0	1	1	0	0.044283

Normal distribution of the Error term around mean 0

```
In [301]: y_pred = lr.predict(X_train_sml)

fig = plt.figure()
sns.distplot((y_train - y_pred), bins = 20)
fig.suptitle('Error Terms', fontsize = 20)
plt.xlabel('Errors', fontsize = 18)
plt.show()
```

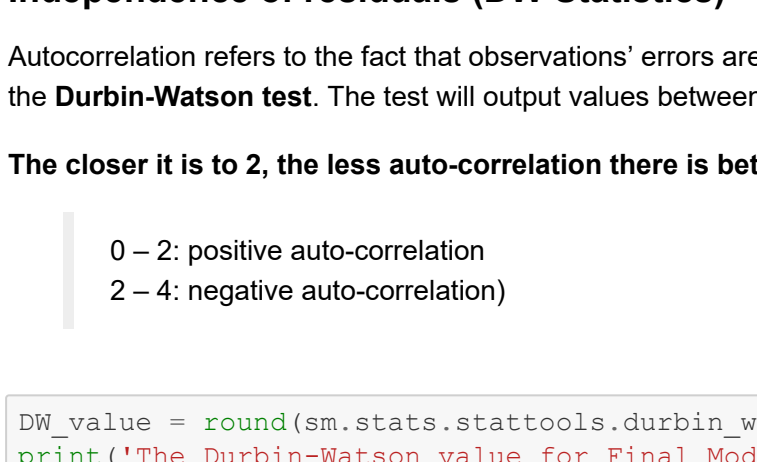


The error terms are normally Distributed with mean around 0

Homoscedasticity

Homoscedasticity means that the residuals have constant variance no matter the level of the dependent variable.

```
In [302]: residual = y_train - y_pred
sns.scatterplot(y_train,residual)
plt.plot(y_train,(y_train - y_train), '-')
plt.xlabel('Count')
plt.ylabel('Residual')
plt.show()
```



Insight: There is no visible pattern in residual values, thus homoscedasticity is well preserved

Normality of Error terms

```
In [303]: sm.qqplot(y_train - y_pred, fit=True, line='45')
plt.title('Q-Q Plot')
plt.show()
```



If the errors are normally distributed, the plot would show fairly straight line. Absence of normality in the errors can be seen with deviation in the straight line.

Insight:

The error terms are normally distributed

Independence of residuals (DW Statistics)

Autocorrelation refers to the fact that observations' errors are correlated. To verify that the observations are not auto-correlated, we can use the Durbin-Watson test. The test will output values between 0 and 4.

The closer it is to 2, the less auto-correlation there is between the various variables.

0 = 2 positive auto-correlation

2 = 4 negative auto-correlation

```
In [304]: DW_value = round(sm.stats.stattools.durbin_watson(y_train - y_pred),4)
print("The Durbin-Watson value for Final Model8 is ",DW_value)

The Durbin-Watson value for Final Model8 is 3.0567

Insight: There is almost no autocorrelation.
```

Multicollinearity Check

VIF value = 4 suggests no multicollinearity

VIF value of >= 10 implies serious multicollinearity

```
In [305]: # Calling VIF from the function for the final model define within the function
vif

Out[305]:
```

features	VIF
0 temp	4.72
1 workingday	4.63
2 year	2.07
3 Saturday	1.81
4 summer	1.74
5 spring	1.70
6 winter	1.59
7 Cloudy	1.56
8 Sep	1.23
9 LightSnow Rain	1.07

All the predictor variables have VIF value less than 5. So we can say that there is insignificant multicollinearity among the predictor variables.

Thus, the model stands on all the assumption of the Linear regression

Making Predictions

```
In [306]: df_model_test.head()

Out[306]:
```

	year	holiday	workingday	temp	humidity	windspeed	count	spring	summer	winter	...	Oct	Sep	Monday	Saturday	Sunday
184	0	1	0	29.793347	63.7917	5.459106	6043	0	0	0	...	0	0	1	0	0
535	1	0	1	32.082500	59.2083	7.625404	8211	0	1	0	...	0	0	0	0	0
299	0	0	1	19.270000	81.2917	13.250121	2659	0	0	1	...	1	0	0	0	0
121	0	0	1	31.433347	42.4167	13.417286	4780	0	0	0	...	0	0	0	0	0
221	0	0	1	29.315000	30.5000	19.883229	4968	0	1	0	...	0	0	0	0	0

5 rows × 29 columns

Scaling

Apply scaler() to all numeric variables in test dataset.

Note: we will only use scaler.transform, as we want to use the metrics that the model learned from the training data to be applied on the test data.

In other words, we want to prevent the information leak from train to test dataset.

```
In [307]: df_model_test.columns

Out[307]: Index(['year', 'holiday', 'workingday', 'temp', 'humidity', 'windspeed', 'count', 'spring', 'summer', 'winter', '...', 'Oct', 'Sep', 'Monday', 'Saturday', 'Sunday', 'type'], dtype='object')

In [308]: df_model_test[variables] = scaler.transform(df_model_test[variables])

In [309]: df_model_test.head()

Out[309]:
```

	year	holiday	workingday	temp	humidity	windspeed	count	spring	summer	winter	...	Oct	Sep	Monday	Saturday	Sunday
184	0	1	0	0.831783	0.657384	0.084219	6043	0	0	0	...	0	0	1	0	0
535	1	0	1	0.901354	0.610133	0.153728	8211	0	1	0	...	0	0	0	0	0
299	0	0	1	0.511964	0.837699	0.334206	2659	0	0	1	...	1	0	0	0	0
121	0	0	1	0.881625	0.437098	0.336570	4780	0	0	0	...	0	0	0	0	0
152	0	0	1	0.817246	0.314296	0.537414	4968	0	1	0	...	0	0	0	0	0

5 rows × 29 columns

Creating X_test and y_test dataset

```
In [310]: X_test = df_model_test.pop('count')
y_test = df_model_test

In [311]: # Add Constant
X_test = sm.add_constant(X_test[feature])
# create a fitted model in one line
lr_test = sm.OLS(y_test,X_test).fit()
y_test_pred = lr_test.predict(X_test_sm)

In [312]: #MSE
print('Mean Square Error: {}'.format(np.sqrt(mean_squared_error(y_test, y_test_pred))))
print('R Square: {}'.format(r_squared))
print('Adj. R Square: {}'.format(Adj_r2))

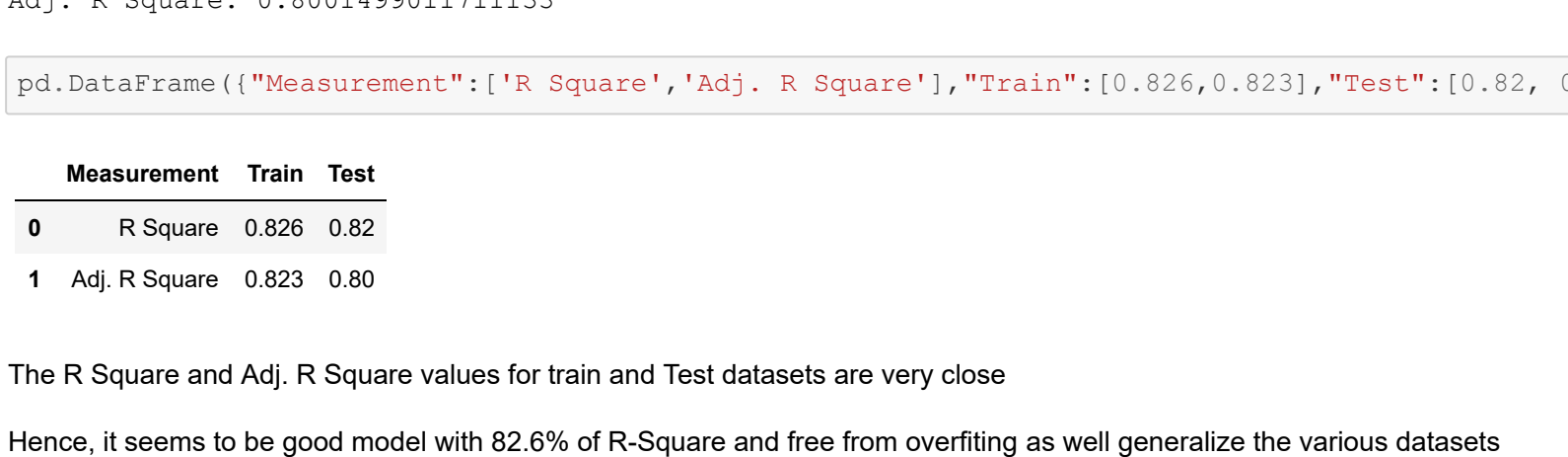
Adj_r2=1-(1-0.820134910540019)*(51-1)/(51-1-1)
print('Adj. R Square: {}'.format(Adj_r2))

Mean Square Error: 803.3795949561533
R Square: 0.820134910540019
Adj. R Square: 0.800149901171133
```

Since, the R2_score for train and test datasets are approximately same

Model Evaluation

```
In [313]: # Plotting y_test and y_pred to understand the spread.
plt.figure(figsize=(15,6))
plt.scatter(y_test,y_test_pred,color='blue')
plt.title('y_test vs y_pred', fontsize=15)
plt.xlabel('y_test', fontsize=15)
plt.ylabel('y_pred', fontsize=15)
plt.show()
```



```
In [314]: from math import sqrt
round(sqrt(mean_squared_error(y_test, y_test_pred)),4)

Out[314]: 803.3796
```

```
In [315]: r_squared = r2_score(y_test, y_test_pred)
print('R Square: {}'.format(r_squared))
Adj_r2=1-(1-0.820134910540019)*(51-1)/(51-1-1)
print('Adj. R Square: {}'.format(Adj_r2))

R Square: 0.820134910540019
Adj. R Square: 0.800149901171133
```

```
In [316]: pd.DataFrame({'Measurement': ['R Square', 'Adj. R Square'], 'Train': [0.826, 0.823], 'Test': [0.82, 0.80]})

Out[316]:
```

Measurement	Train	Test
0 R Square	0.826	0.82
1 Adj. R Square	0.823	0.80

The R Square and Adj. R Square values for train and Test datasets are very close

Hence, it seems to be good model with 82.6% of R-Square and free from overfitting as well generalize the various datasets

```
In [317]: #RMSE
print('Root Mean Square Error: {}'.format(np.sqrt(mean_squared_error(y_test, y_test_pred))))
#MAE
print('Mean Absolute Error: {}'.format(np.sqrt(mean_absolute_error(y_test, y_test_pred))))

Root Mean Square Error: 803.3795949561533
Mean Absolute Error: 24.74813594167846
```

The Root Mean Squared Error value for the test dataset based on final model is 0.0924 and Mean Absolute Error is 0.26544, which indicates that the model is really good.

```
In [318]: #Regression plot
plt.figure(figsize=(15,6))
sns.regplot(x=y_test, y=y_test_pred, ci=8, fit_reg=True, scatter_kws={"color": "blue", line_kws={"color": "red"}})
plt.title('y_test vs y_pred', fontsize=20)
plt.xlabel('y_test', fontsize=19)
plt.ylabel('y_pred', fontsize=16)
plt.show()
```



Hence, We reach to the decent model for the the demand for shared bikes with the significant variables

Best Fit Line (y_pred) = 0.0901 + year (0.2333) + workingday (0.0568) + spring (-0.0647) + summer (0.0516) + winter (0.0984) + Sep (0.0915) + Saturday (0.0647) + Cloudy (-0.0800) + LightSnow Rain (-0.3051) + temp (0.4918)

R-Square: 82.6%

Adj. R-Square: 82.3%

F-Statistics: 237.5

Important & Significant Features(predictors) which effect demand for shared bikes:

- Temperature has a high coefficient of 0.4918, which means that if the temperature increases by 1 unit the demand goes up by 0.4918 units

- During