

In [163]: 1 ! pip install stats

Requirement already satisfied: stats in c:\users\dell\anaconda3\lib\site-packages (0.1.2a0)

In [165]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from scipy.stats import norm,t,binom,expon,chi2,chi-square,chi2_contingency,ttest_1samp,ttest_rel,ttest_:

import datetime as dt
7 from statsmodels.graphics.gofplots import qqplot

In [2]: 1 df=pd.read_csv('yulu_bike_sharing.csv')

In [100]: 1 df_original=pd.read_csv('yulu_bike_sharing.csv')

In [102]: 1 df_original['datetime']=pd.to_datetime(df['datetime'])

In [3]: 1 df

Out[3]:

localhost:8888/notebooks/Business_case_Yulu_v00.ipynb#

In [3]: 1 df

Out[3]:

		datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
0	2011-01-01 00:00:00		1	0	0	1	9.84	14.395	81	0.0000	3	13	16
1	2011-01-01 01:00:00		1	0	0	1	9.02	13.635	80	0.0000	8	32	40
2	2011-01-01 02:00:00		1	0	0	1	9.02	13.635	80	0.0000	5	27	32
3	2011-01-01 03:00:00		1	0	0	1	9.84	14.395	75	0.0000	3	10	13
4	2011-01-01 04:00:00		1	0	0	1	9.84	14.395	75	0.0000	0	1	1
...
10881	2012-12-19 19:00:00		4	0	1	1	15.58	19.695	50	26.0027	7	329	336
10882	2012-12-19 20:00:00		4	0	1	1	14.76	17.425	57	15.0013	10	231	241
10883	2012-12-19 21:00:00		4	0	1	1	13.94	15.910	61	15.0013	4	164	168
10884	2012-12-19 22:00:00		4	0	1	1	13.94	17.425	61	6.0032	12	117	129
10885	2012-12-19 23:00:00		4	0	1	1	13.12	16.665	66	8.9981	4	84	88

10886 rows × 12 columns

In [4]: 1 df['datetime']=pd.to_datetime(df['datetime'])

In [5]: 1 df

Out[5]:

datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
----------	--------	---------	------------	---------	------	-------	----------	-----------	--------	------------	-------

In [5]: 1 df

Out[5]:

		datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
0	2011-01-01 00:00:00	1	0	0	1	1	9.84	14.395	81	0.0000	3	13	16
1	2011-01-01 01:00:00	1	0	0	1	1	9.02	13.635	80	0.0000	8	32	40
2	2011-01-01 02:00:00	1	0	0	1	1	9.02	13.635	80	0.0000	5	27	32
3	2011-01-01 03:00:00	1	0	0	1	1	9.84	14.395	75	0.0000	3	10	13
4	2011-01-01 04:00:00	1	0	0	1	1	9.84	14.395	75	0.0000	0	1	1
...
10881	2012-12-19 19:00:00	4	0	1	1	1	15.58	19.695	50	26.0027	7	329	336
10882	2012-12-19 20:00:00	4	0	1	1	1	14.76	17.425	57	15.0013	10	231	241
10883	2012-12-19 21:00:00	4	0	1	1	1	13.94	15.910	61	15.0013	4	164	168
10884	2012-12-19 22:00:00	4	0	1	1	1	13.94	17.425	61	6.0032	12	117	129
10885	2012-12-19 23:00:00	4	0	1	1	1	13.12	16.665	66	8.9981	4	84	88

10886 rows × 12 columns

In [6]: 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   datetime64[ns]
 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 
```

In [7]: dtypes: datetime64[ns](1), float64(7), int64(8)
memory usage: 1020.7 KB

In [8]: 1 df

```
In [7]: █ 1 count          10886 non-null int64
      dtypes: datetime[1], float64(3), int64(8)
      memory usage: 1020.7 KB
```

```
In [8]: █ 1 df
```

Out[8]:

		datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	date
0	2011-01-01 00:00:00		1	0	0	1	9.84	14.395	81	0.0000	3	13	16	2011-01-01
1	2011-01-01 01:00:00		1	0	0	1	9.02	13.635	80	0.0000	8	32	40	2011-01-01
2	2011-01-01 02:00:00		1	0	0	1	9.02	13.635	80	0.0000	5	27	32	2011-01-01
3	2011-01-01 03:00:00		1	0	0	1	9.84	14.395	75	0.0000	3	10	13	2011-01-01
4	2011-01-01 04:00:00		1	0	0	1	9.84	14.395	75	0.0000	0	1	1	2011-01-01
...
10881	2012-12-19 19:00:00		4	0	1	1	15.58	19.695	50	26.0027	7	329	336	2012-12-19
10882	2012-12-19 20:00:00		4	0	1	1	14.76	17.425	57	15.0013	10	231	241	2012-12-19
10883	2012-12-19 21:00:00		4	0	1	1	13.94	15.910	61	15.0013	4	164	168	2012-12-19
10884	2012-12-19 22:00:00		4	0	1	1	13.94	17.425	61	6.0032	12	117	129	2012-12-19
10885	2012-12-19 23:00:00		4	0	1	1	13.12	16.665	66	8.9981	4	84	88	2012-12-19

10886 rows × 13 columns

```
In [9]: █ 1 df['time']=df['datetime'].dt.time
```

```
In [10]: █ 1 df
```

Out[10]:

datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	date	time
----------	--------	---------	------------	---------	------	-------	----------	-----------	--------	------------	-------	------	------

In [10]: 1 df

Out[10]:

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	date	time
0	2011-01-01 00:00:00	1	0	0	1	9.84	14.395	81	0.0000	3	13	16	2011-01-01	00:00:00
1	2011-01-01 01:00:00	1	0	0	1	9.02	13.635	80	0.0000	8	32	40	2011-01-01	01:00:00
2	2011-01-01 02:00:00	1	0	0	1	9.02	13.635	80	0.0000	5	27	32	2011-01-01	02:00:00
3	2011-01-01 03:00:00	1	0	0	1	9.84	14.395	75	0.0000	3	10	13	2011-01-01	03:00:00
4	2011-01-01 04:00:00	1	0	0	1	9.84	14.395	75	0.0000	0	1	1	2011-01-01	04:00:00
...
10881	2012-12-19 19:00:00	4	0	1	1	15.58	19.695	50	26.0027	7	329	336	2012-12-19	19:00:00
10882	2012-12-19 20:00:00	4	0	1	1	14.76	17.425	57	15.0013	10	231	241	2012-12-19	20:00:00
10883	2012-12-19 21:00:00	4	0	1	1	13.94	15.910	61	15.0013	4	164	168	2012-12-19	21:00:00
10884	2012-12-19 22:00:00	4	0	1	1	13.94	17.425	61	6.0032	12	117	129	2012-12-19	22:00:00
10885	2012-12-19 23:00:00	4	0	1	1	13.12	16.665	66	8.9981	4	84	88	2012-12-19	23:00:00

10886 rows × 14 columns



In [11]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10886 entries, 0 to 10885
Data columns (total 14 columns):
 #   Column  Non-Null Count  Dtype  
 0   datetime    10886 non-null  datetime64[ns]
 1   season      10886 non-null  int8    
 2   holiday     10886 non-null  int8    
 3   workingday  10886 non-null  int8    
 4   weather     10886 non-null  int8    
 5   temp        10886 non-null  float64
 6   atemp       10886 non-null  float64
 7   humidity    10886 non-null  float64
 8   windspeed   10886 non-null  float64
 9   casual      10886 non-null  float64
 10  registered  10886 non-null  float64
 11  count       10886 non-null  float64
 12  date        10886 non-null  datetime64[ns]
 13  time        10886 non-null  datetime64[ns]

```

In [11]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10886 entries, 0 to 10885
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype  
---  --  
 0   datetime    10886 non-null   datetime64[ns]
 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  
 12  date         10886 non-null   object  
 13  time         10886 non-null   object  
dtypes: datetime64[ns](1), float64(3), int64(8), object(2)
memory usage: 1.2+ MB
```

In [12]: 1 df.isna().sum()

```
Out[12]: datetime      0
          season       0
          holiday      0
          workingday   0
          weather      0
          temp         0
          atemp        0
          humidity     0
          windspeed    0
          casual       0
          registered   0
          count        0
          date         0
          time         0
          dtype: int64
```

1 **### No missing values in teh data**

In [13]: 1 df.describe()

1 **### No missing values in teh data**

In [13]: 1 df.describe()

Out[13]:

season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	reg
10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.000000	10886.
2.506614	0.028569	0.680875	1.418427	20.23086	23.655084	61.886460	12.799395	36.021955	155.
1.000000	0.000000	0.000000	1.000000	0.82000	0.760000	0.000000	0.000000	0.000000	0.
2.000000	0.000000	0.000000	1.000000	13.94000	16.665000	47.000000	7.001500	4.000000	36.
3.000000	0.000000	1.000000	1.000000	20.50000	24.240000	62.000000	12.998000	17.000000	118.
4.000000	0.000000	1.000000	2.000000	26.24000	31.060000	77.000000	16.997900	49.000000	222.
4.000000	1.000000	1.000000	4.000000	41.00000	45.455000	100.000000	56.996900	367.000000	886.
1.116174	0.166599	0.466159	0.633839	7.79159	8.474601	19.245033	8.164537	49.960477	151.

◀ ▶

In [14]: 1 df['date'].nunique()

Out[14]: 456

In [15]: 1 df.describe(include=object)

Out[15]:

	date	time
count	10886	10886
unique	456	24
top	2011-01-01	12:00:00
freq	24	456

In [16]: 1 df['date'].max(),df['date'].min()

In [17]: 1 df['year_month']=df['datetime'].dt.to_period('M')
Out[16]: (datetime.date(2012, 12, 19), datetime.date(2011, 1, 1))

In [18]: 1 df

In [17]: 1
Out[16]: (datetime.date(2012, 12, 19), datetime.date(2011, 1, 1))

In [18]: 1 df

Out[18]:

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	date	time
0	2011-01-01 00:00:00	1	0	0	1	9.84	14.395	81	0.0000	3	13	16	2011-01-01	00:00:00
1	2011-01-01 01:00:00	1	0	0	1	9.02	13.635	80	0.0000	8	32	40	2011-01-01	01:00:00
2	2011-01-01 02:00:00	1	0	0	1	9.02	13.635	80	0.0000	5	27	32	2011-01-01	02:00:00
3	2011-01-01 03:00:00	1	0	0	1	9.84	14.395	75	0.0000	3	10	13	2011-01-01	03:00:00
4	2011-01-01 04:00:00	1	0	0	1	9.84	14.395	75	0.0000	0	1	1	2011-01-01	04:00:00
...
10881	2012-12-19 19:00:00	4	0	1	1	15.58	19.695	50	26.0027	7	329	336	2012-12-19	19:00:00
10882	2012-12-19 20:00:00	4	0	1	1	14.76	17.425	57	15.0013	10	231	241	2012-12-19	20:00:00
10883	2012-12-19 21:00:00	4	0	1	1	13.94	15.910	61	15.0013	4	164	168	2012-12-19	21:00:00
10884	2012-12-19 22:00:00	4	0	1	1	13.94	17.425	61	6.0032	12	117	129	2012-12-19	22:00:00
10885	2012-12-19 23:00:00	4	0	1	1	13.12	16.665	66	8.9981	4	84	88	2012-12-19	23:00:00

10886 rows × 15 columns

In [68]: 1 df1=df[df['datetime']>='2012-01-01']

In [69]: 1 df1

```
In [68]: 1 df1=df[df['datetime']>='2012-01-01']
```

```
In [69]: 1 df1
```

Out[69]:

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	date	time
5422	2012-01-01 00:00:00	1	0	0	1	14.76	18.940	66	0.0000	5	43	48	2012-01-01	00:00:00
5423	2012-01-01 01:00:00	1	0	0	1	14.76	17.425	66	8.9981	15	78	93	2012-01-01	01:00:00
5424	2012-01-01 02:00:00	1	0	0	1	13.12	17.425	76	0.0000	16	59	75	2012-01-01	02:00:00
5425	2012-01-01 03:00:00	1	0	0	1	12.30	16.665	81	0.0000	11	41	52	2012-01-01	03:00:00
5426	2012-01-01 04:00:00	1	0	0	1	11.48	15.150	81	6.0032	0	8	8	2012-01-01	04:00:00
...
10881	2012-12-19 19:00:00	4	0	1	1	15.58	19.695	50	26.0027	7	329	336	2012-12-19	19:00:00
10882	2012-12-19 20:00:00	4	0	1	1	14.76	17.425	57	15.0013	10	231	241	2012-12-19	20:00:00
10883	2012-12-19 21:00:00	4	0	1	1	13.94	15.910	61	15.0013	4	164	168	2012-12-19	21:00:00
10884	2012-12-19 22:00:00	4	0	1	1	13.94	17.425	61	6.0032	12	117	129	2012-12-19	22:00:00
10885	2012-12-19 23:00:00	4	0	1	1	13.12	16.665	66	8.9981	4	84	88	2012-12-19	23:00:00

5464 rows × 15 columns

```
In [78]: 1 df11=df[(df['datetime']>='2012-01-01') & (df['datetime']<='2012-01-31')]
```

```
In [79]: 1 df11
```

In [78]: 1 df11=df[(df['datetime'] >= '2012-01-01') & (df['datetime'] <= '2012-01-31')]

In [79]: 1 df11

Out[79]:

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	date	time
5422	2012-01-01 00:00:00	1	0	0	1	14.76	18.940	66	0.0000	5	43	48	2012-01-01	00:00:00
5423	2012-01-01 01:00:00	1	0	0	1	14.76	17.425	66	8.9981	15	78	93	2012-01-01	01:00:00
5424	2012-01-01 02:00:00	1	0	0	1	13.12	17.425	76	0.0000	16	59	75	2012-01-01	02:00:00
5425	2012-01-01 03:00:00	1	0	0	1	12.30	16.665	81	0.0000	11	41	52	2012-01-01	03:00:00
5426	2012-01-01 04:00:00	1	0	0	1	11.48	15.150	81	6.0032	0	8	8	2012-01-01	04:00:00
...
5870	2012-01-19 19:00:00	1	0	1	2	10.66	11.365	48	19.9995	5	229	234	2012-01-19	19:00:00
5871	2012-01-19 20:00:00	1	0	1	1	10.66	10.605	52	30.0026	7	184	191	2012-01-19	20:00:00
5872	2012-01-19 21:00:00	1	0	1	1	10.66	11.365	44	23.9994	4	117	121	2012-01-19	21:00:00
5873	2012-01-19 22:00:00	1	0	1	2	10.66	11.365	44	27.9993	0	90	90	2012-01-19	22:00:00
5874	2012-01-19 23:00:00	1	0	1	1	10.66	11.365	48	22.0028	1	55	56	2012-01-19	23:00:00

453 rows × 15 columns

In [83]: 1 df111=df[(df['datetime'] >= '2012-01-01') & (df['datetime'] <= '2012-01-07')]

In [84]: 1 df111

In [83]: 1 df111=df[(df['datetime']>='2012-01-01') & (df['datetime']<='2012-01-07')]

In [84]: 1 df111

Out[84]:

	datetime	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count	date	time
5422	2012-01-01 00:00:00	1	0	0	1	14.76	18.940	66	0.0000	5	43	48	2012-01-01	00:00:00
5423	2012-01-01 01:00:00	1	0	0	1	14.76	17.425	66	8.9981	15	78	93	2012-01-01	01:00:00
5424	2012-01-01 02:00:00	1	0	0	1	13.12	17.425	76	0.0000	16	59	75	2012-01-01	02:00:00
5425	2012-01-01 03:00:00	1	0	0	1	12.30	16.665	81	0.0000	11	41	52	2012-01-01	03:00:00
5426	2012-01-01 04:00:00	1	0	0	1	11.48	15.150	81	6.0032	0	8	8	2012-01-01	04:00:00
...
5561	2012-01-06 20:00:00	1	0	1	1	17.22	21.210	41	12.9980	14	163	177	2012-01-06	20:00:00
5562	2012-01-06 21:00:00	1	0	1	1	16.40	20.455	43	15.0013	17	137	154	2012-01-06	21:00:00
5563	2012-01-06 22:00:00	1	0	1	1	14.76	17.425	50	12.9980	12	123	135	2012-01-06	22:00:00
5564	2012-01-06 23:00:00	1	0	1	1	14.76	18.940	50	0.0000	13	88	101	2012-01-06	23:00:00
5565	2012-01-07 00:00:00	1	0	0	1	14.76	17.425	50	11.0014	2	77	79	2012-01-07	00:00:00

144 rows × 15 columns

In [19]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
   count    mean    std
datetime  144.0  14.76  18.940
   count    mean    std
season       1.0    1.0  0.0000
   count    mean    std
holiday      0.0    0.0  0.0000
   count    mean    std
workingday   1.0    1.0  0.0000
   count    mean    std
weather       1.0    1.0  17.425
   count    mean    std
temp        144.0  14.76  18.940
   count    mean    std
atemp       144.0  13.12  16.665
   count    mean    std
humidity     144.0  76.00  8.9981
   count    mean    std
windspeed    144.0  6.0032  12.9980
   count    mean    std
casual      144.0  11.00  11.0014
   count    mean    std
registered   144.0  48.00  15.0013
   count    mean    std
count       144.0  101.0  101.0000
   count    mean    std
date        144.0  2012-01-06  2012-01-06
   count    mean    std
time        144.0  00:00:00  00:00:00
```

In [19]: 1 df.info()

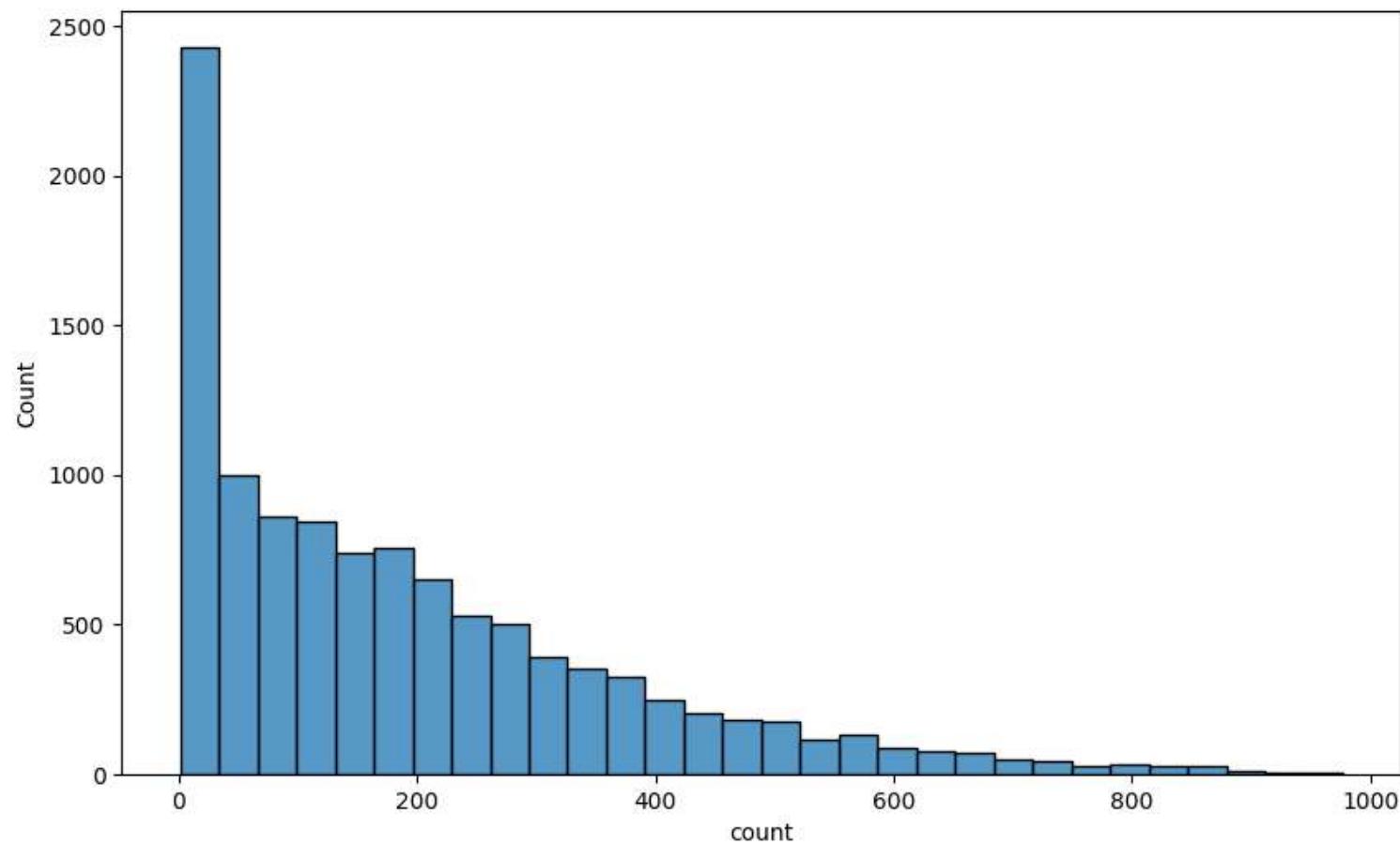
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10886 entries, 0 to 10885
Data columns (total 15 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   datetime    10886 non-null   datetime64[ns]
 1   season      10886 non-null   int64  
 2   holiday     10886 non-null   int64  
 3   workingday  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  
 12  date        10886 non-null   object  
 13  time        10886 non-null   object  
 14  year_month  10886 non-null   period[M] 
dtypes: datetime64[ns](1), float64(3), int64(8), object(2), period[M](1)
memory usage: 1.2+ MB
```

In [32]: 1 plt.figure(figsize=(10,6))

```
2 sns.histplot(data=df,x='count',bins=30)
3 plt.show()
```

In [32]:

```
1 plt.figure(figsize=(10,6))
2 sns.histplot(data=df,x='count',bins=30)
3 plt.show()
```

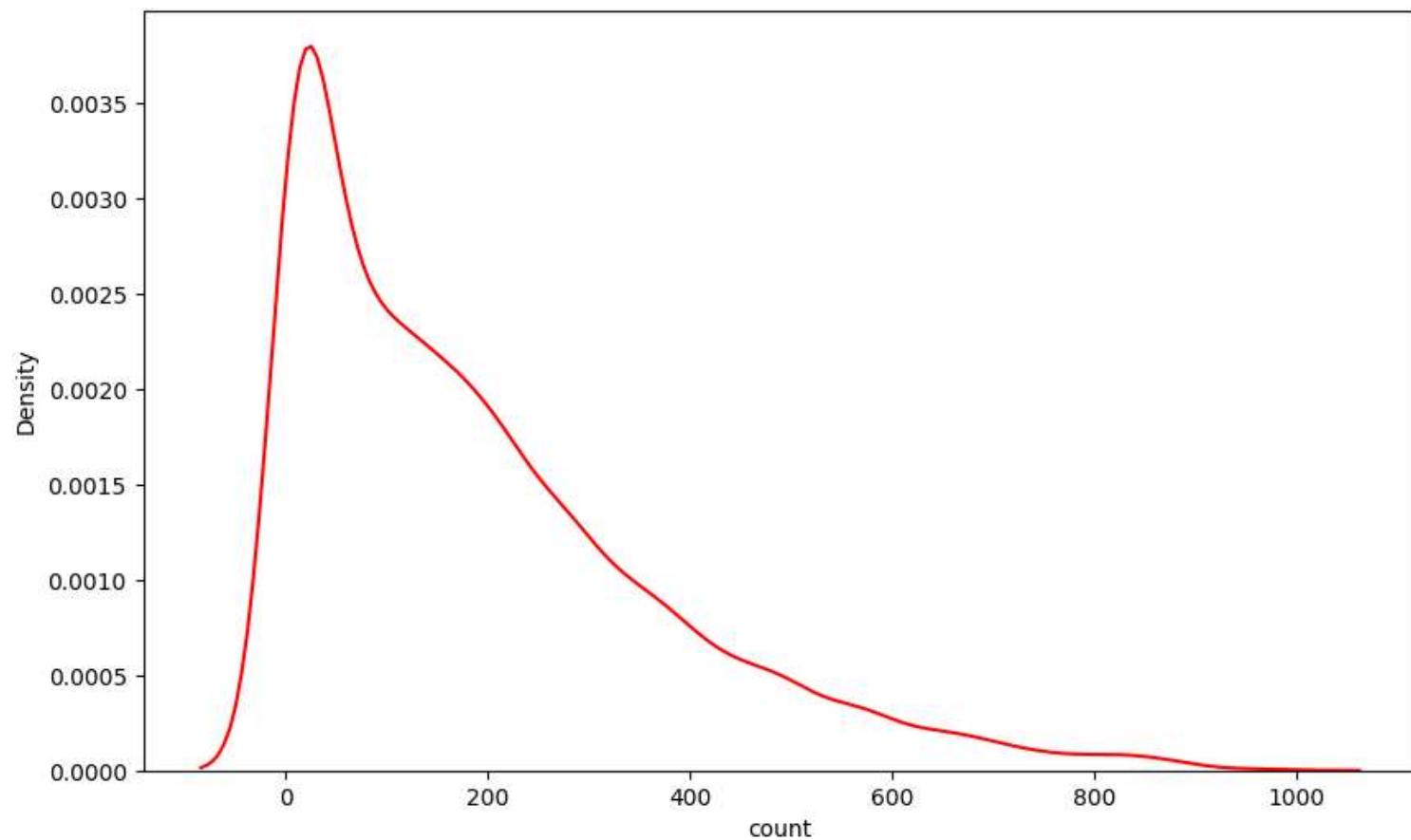


In [179]:

```
1 plt.figure(figsize=(10,6))
2 sns.kdeplot(x=df['count'],c='red')
3 plt.show()
```

In [179]:

```
1 plt.figure(figsize=(10,6))
2 sns.kdeplot(x=df['count'],c='red')
3 plt.show()
```

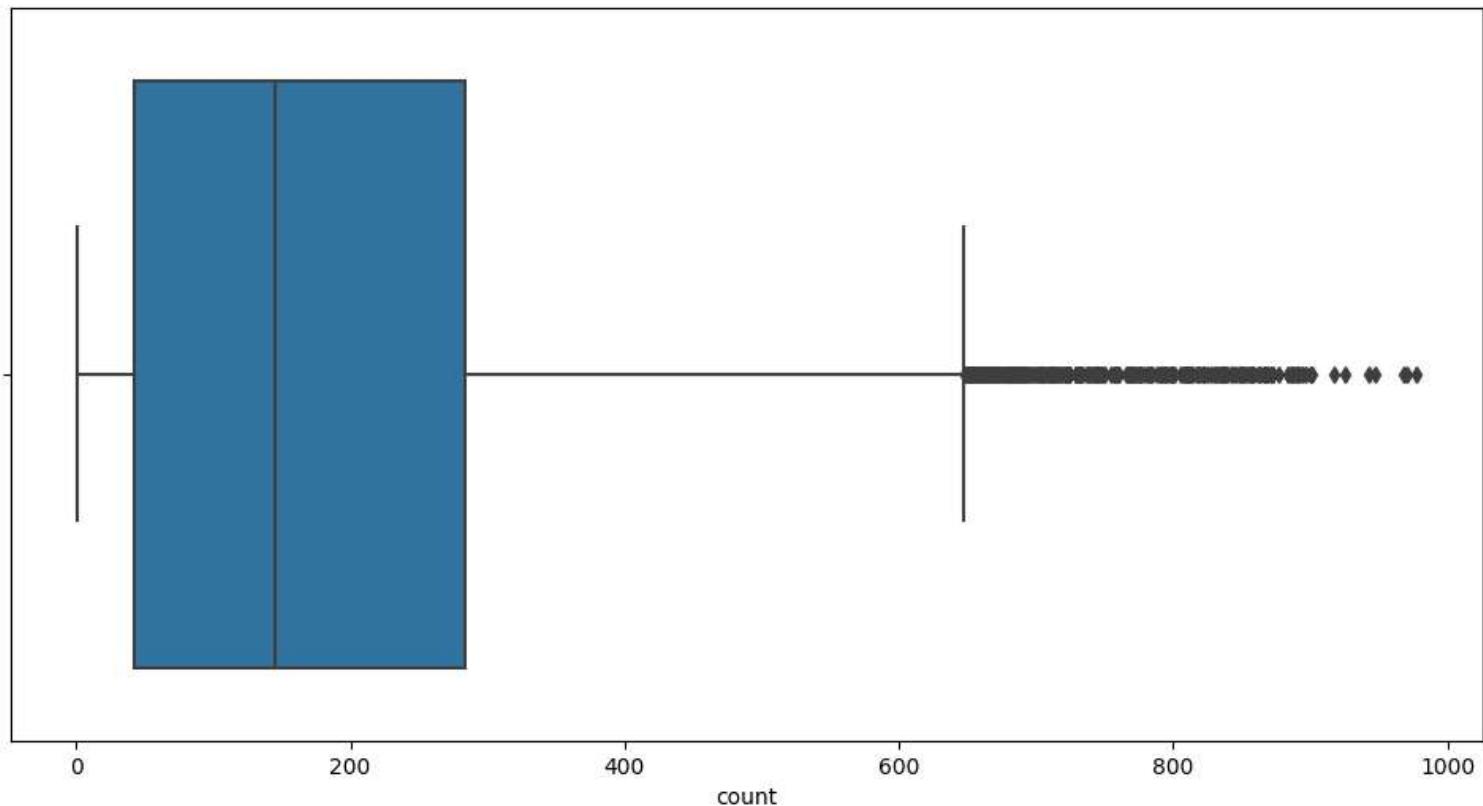


In [180]:

```
1 plt.figure(figsize=(12,6))
2 sns.boxplot(x=df['count'])
3 plt.show()
```

In [180]:

```
1 plt.figure(figsize=(12,6))
2 sns.boxplot(x=df['count'])
3 plt.show()
```



```
1 ## many outliers observed in various features; but no treatment of outliers
done...
```

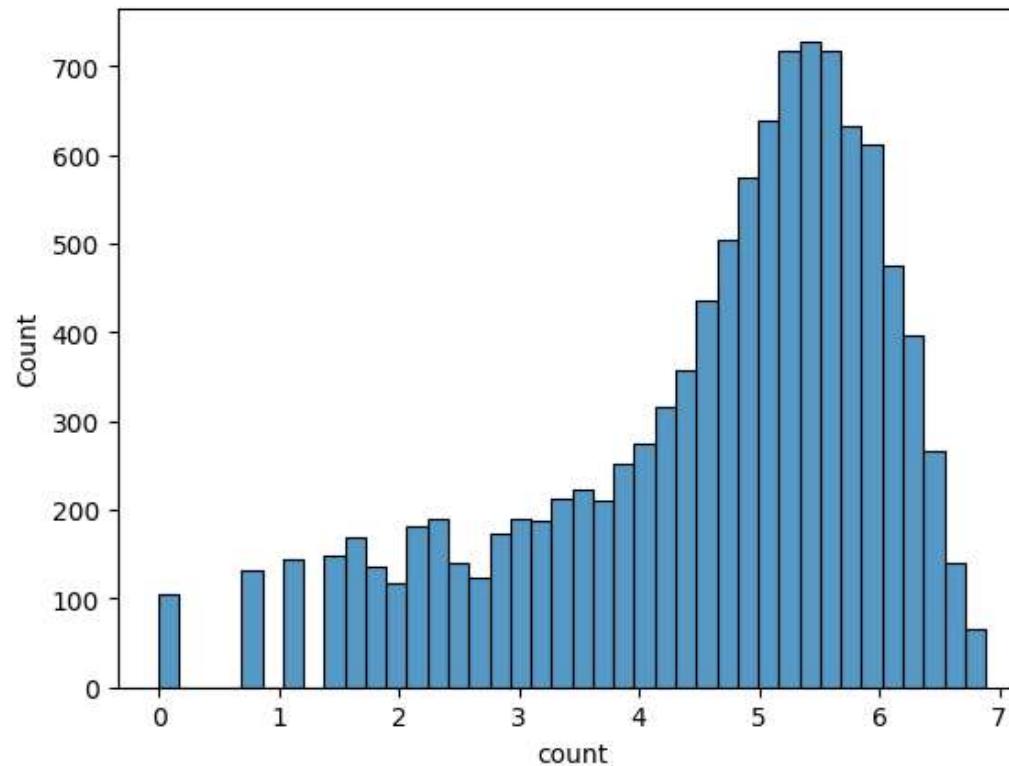
In [178]:

```
1 sns.histplot(np.log(df['count']))
```

Out[178]: <Axes: xlabel='count', ylabel='Count'>

In [178]: 1 sns.histplot(np.log(df['count']))

Out[178]: <Axes: xlabel='count', ylabel='Count'>



1 ## bike usage is not normally distributed. looks like it is slightly lognormally distributed

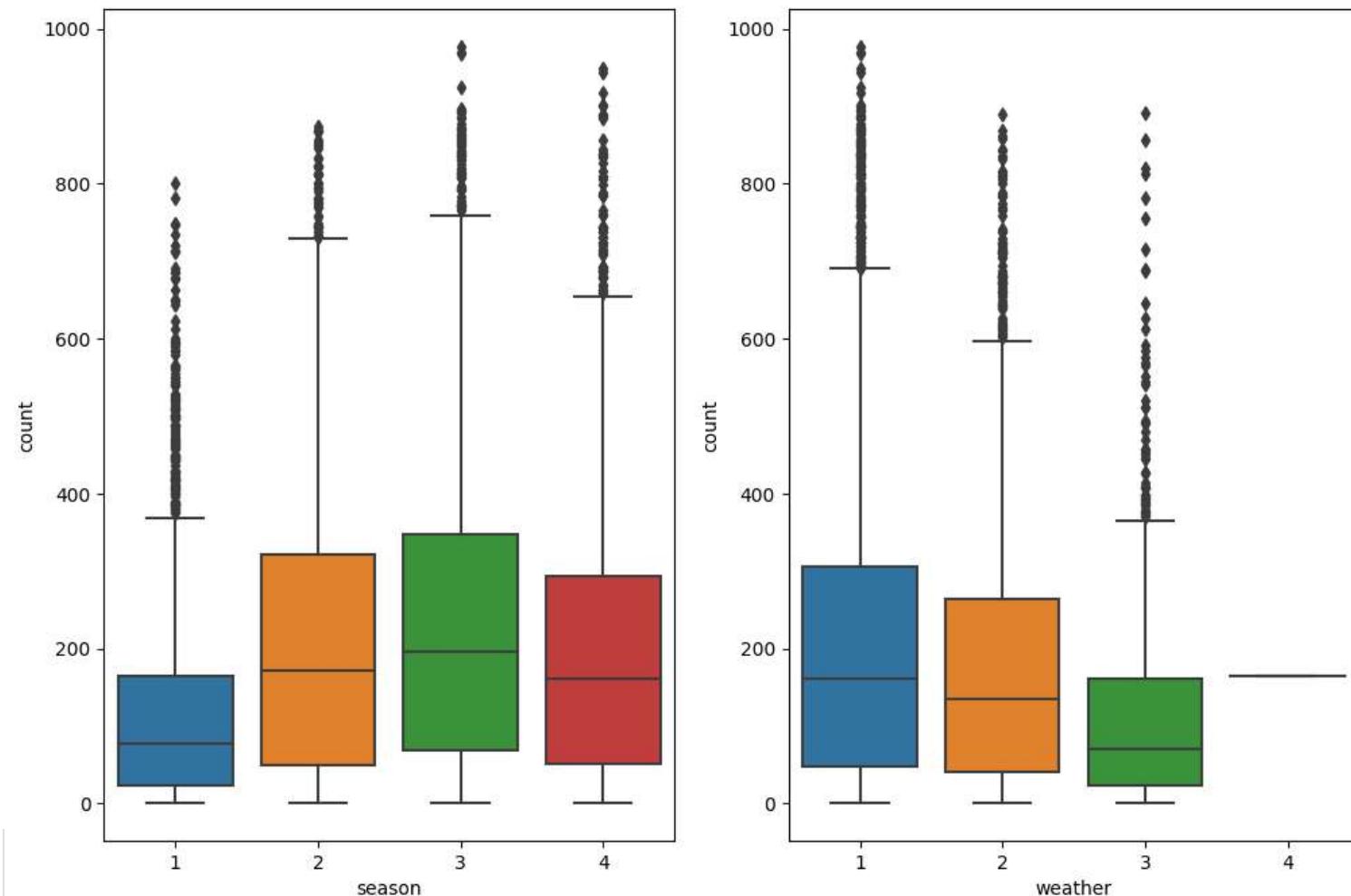
In [63]: 1 plt.figure(figsize=(12,8))

2 plt.subplot(1,2,1)

3 sns.boxplot(data=df,x='season', y='count')

In [63]:

```
1 plt.figure(figsize=(12,8))
2 plt.subplot(1,2,1)
3 sns.boxplot(data=df,x='season', y='count')
4
5 plt.subplot(1,2,2)
6 sns.boxplot(data=df,x='weather', y='count')
7 plt.show()
8
```



3 **## season 3 shows more bike usage compared to other seasons; season 1 shows the least**
4 **## weather 1 shows more bike usage than other weathers;**
5 **## heavy rain weather shows much less bike usage than other weathers;**
6 **## weather 1: clear weather 2: mist weather 3: light rain weather 4: heavy rain**

```
3 ### season 3 shows more bike usage compared to other seasons; season 1 shows the least
4 ### clear weather shows more bike usage than other weathers;
5 ### season rains spring weather is not as much as summer season bike usage season 4: winter
6 ### weather 1: clear weather 2: mist weather 3: light rain weather 4: heavy rain
7
```

In []:

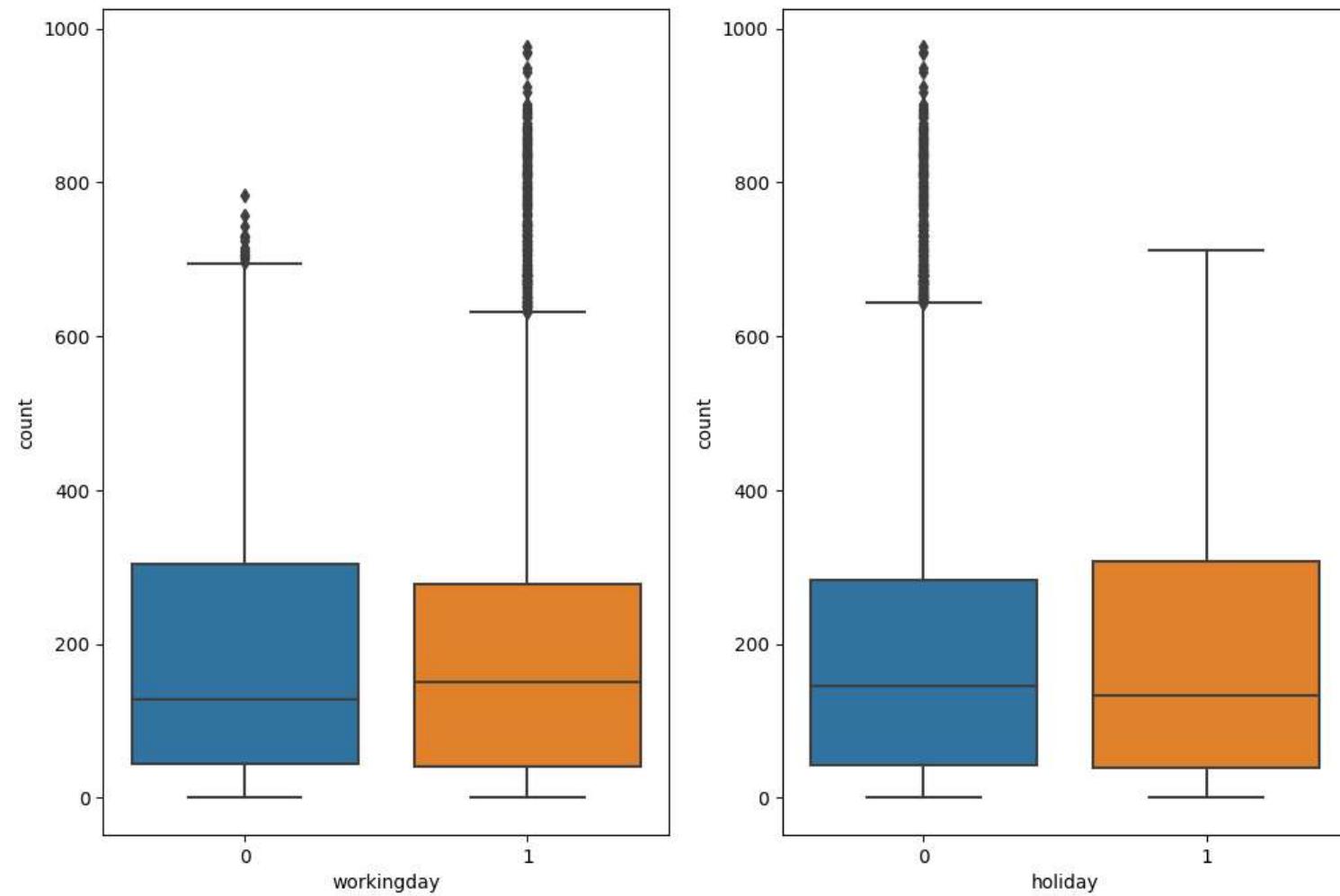
```
1 plt.figure(figsize=(12,8))
2 plt.subplot(1,2,1)
3 sns.boxplot(data=df,x='season', y='count')
4
5 plt.subplot(1,2,2)
6 sns.boxplot(data=df,x='weather', y='count')
7 plt.show()
8
```

In [62]:

```
1 plt.figure(figsize=(12,8))
2 plt.subplot(1,2,1)
3 sns.boxplot(data=df,x='workingday', y='count')
```

In [62]:

```
1 plt.figure(figsize=(12,8))
2 plt.subplot(1,2,1)
3 sns.boxplot(data=df,x='workingday', y='count')
4
5 plt.subplot(1,2,2)
6 sns.boxplot(data=df,x='holiday', y='count')
7 plt.show()
```



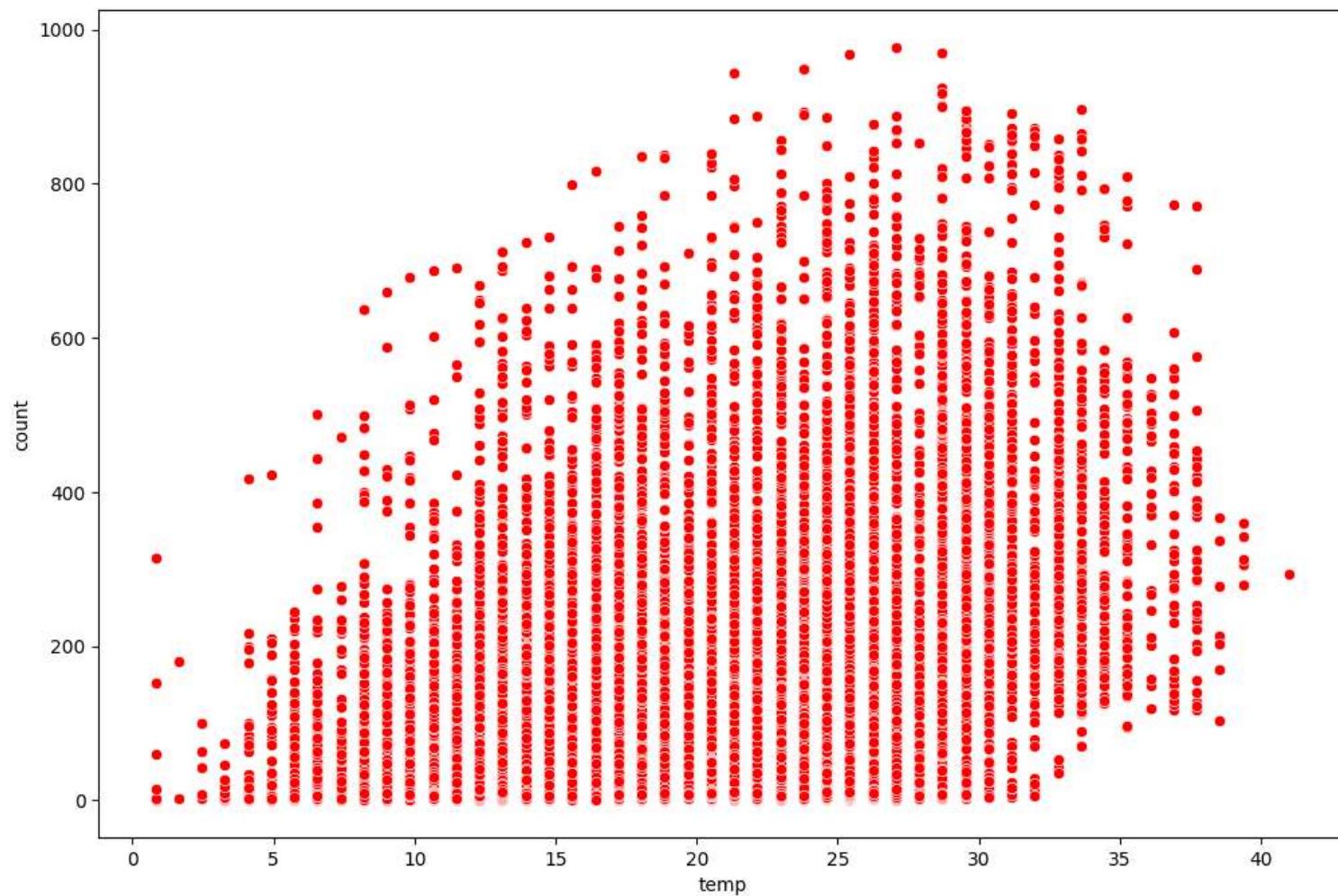
In [57]:

```
1 plt.figure(figsize=(12,8))
2 sns.scatterplot(data=df,x='temp',y='count',c='red')
3 plt.plot()
```

In [57]:

```
1 plt.figure(figsize=(12,8))
2 sns.scatterplot(data=df,x='temp',y='count',c='red')
3 plt.plot()
4
```

Out[57]: []



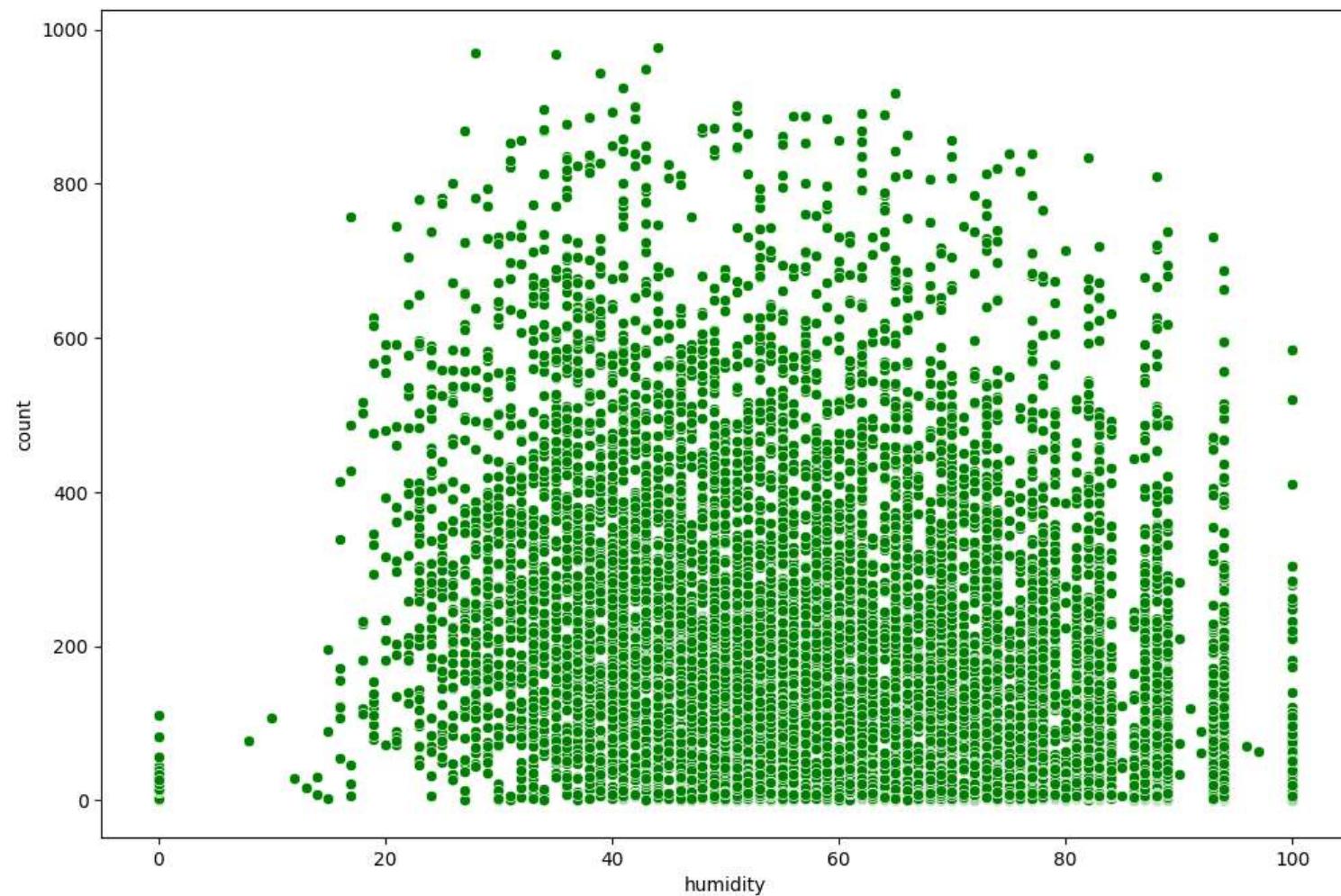
1 **### low temperature (less than 5 degree celcius) shows less bike usage based on the scatter plot**

In [56]:

```
1 plt.figure(figsize=(12,8))
2 sns.scatterplot(data=df,x='humidity',y='count',c='green')
3 plt.plot()
```

```
In [56]: 1 plt.figure(figsize=(12,8))
2 sns.scatterplot(data=df,x='humidity',y='count',c='green')
3 plt.plot()
4
```

Out[56]: []

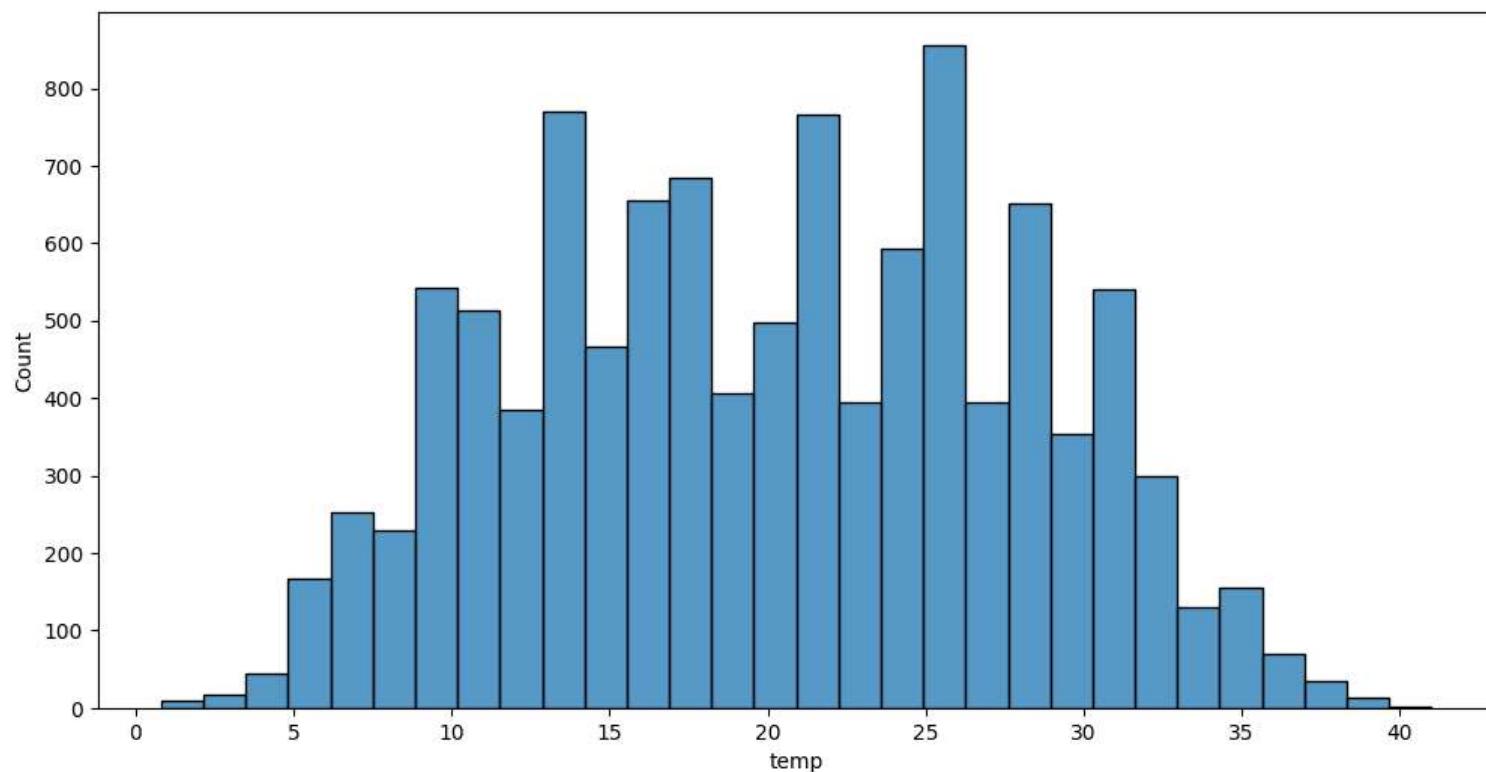


```
1 ### low humidity ( less than 20%) shows less bike usage;
2 ### very high (>90% humidity also shows lesser bike usage
```

```
In [58]: 1 plt.figure(figsize=(12,6))
2 sns.histplot(data=df,x='temp',bins=30)
3 plt.show()
```

In [58]:

```
1 plt.figure(figsize=(12,6))
2 sns.histplot(data=df,x='temp',bins=30)
3 plt.show()
```

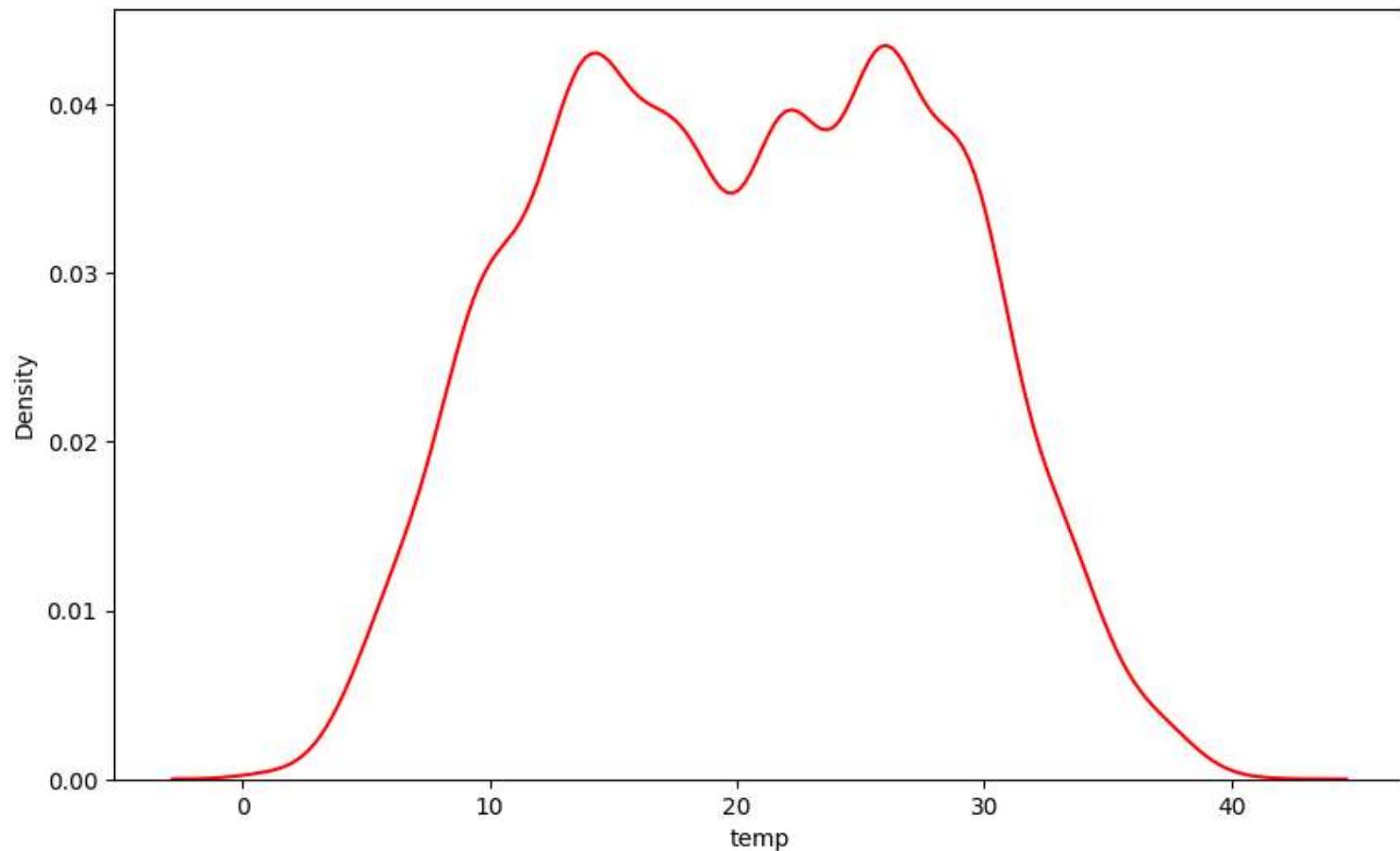


In [95]:

```
1 plt.figure(figsize=(10,6))
2 sns.kdeplot(x=df['temp'],c='red')
3 plt.show()
```

In [95]:

```
1 plt.figure(figsize=(10,6))
2 sns.kdeplot(x=df['temp'],c='red')
3 plt.show()
```



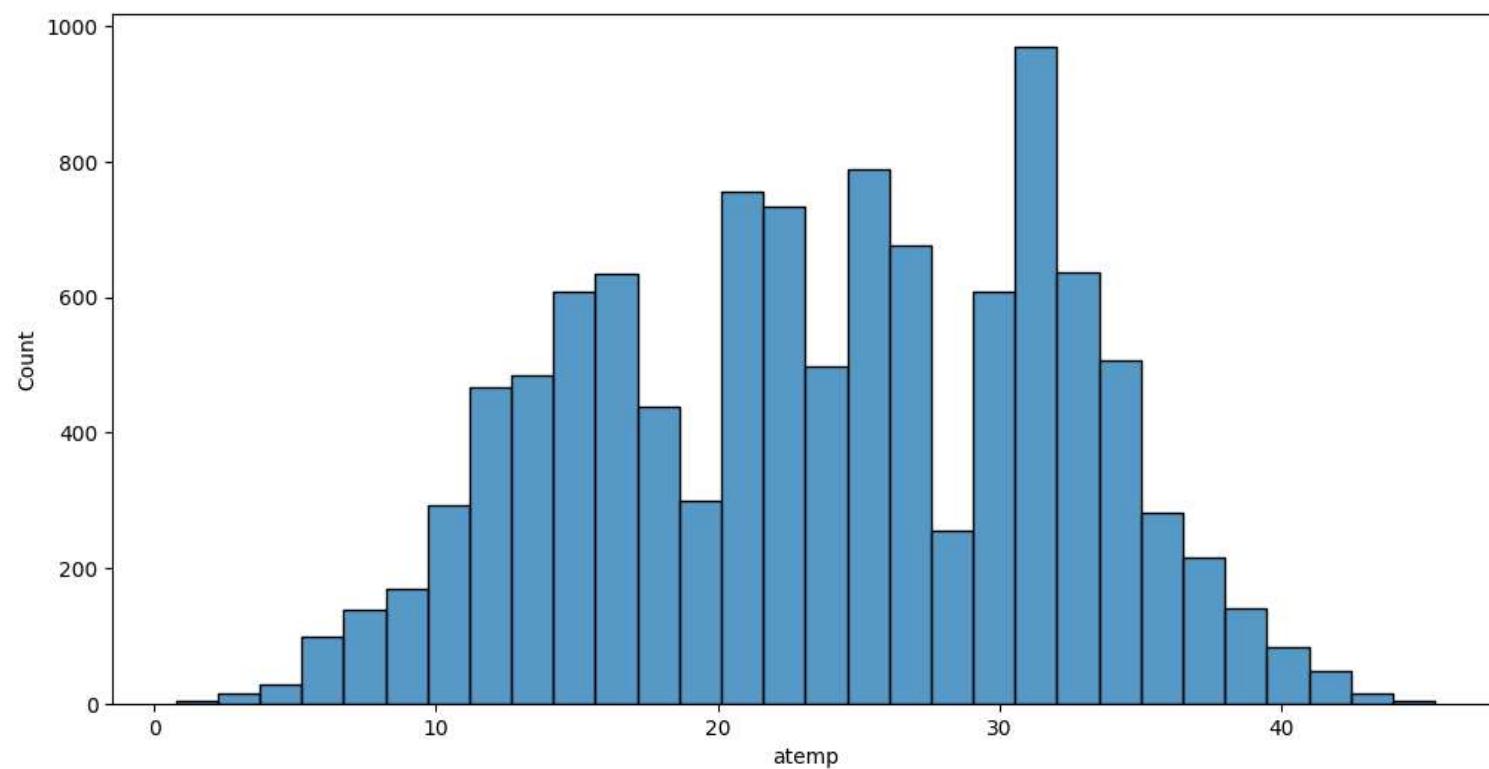
```
1 ### temp distribution can be broadly considered as normally distributed
```

In [94]:

```
1 plt.figure(figsize=(12,6))
2 sns.histplot(data=df,x='atemp',bins=30)
3 plt.show()
```

In [94]:

```
1 plt.figure(figsize=(12,6))
2 sns.histplot(data=df,x='atemp',bins=30)
3 plt.show()
```

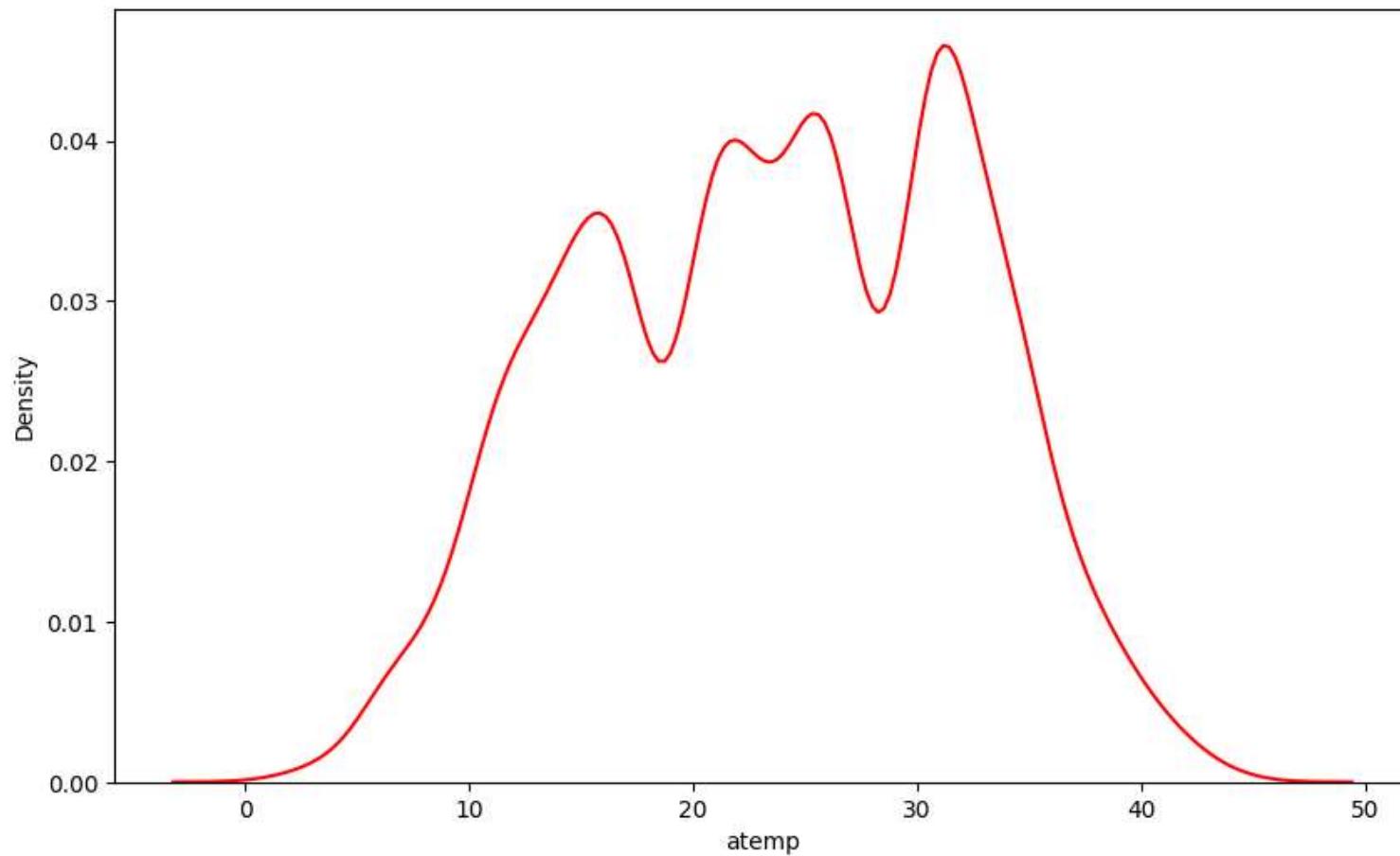


In [96]:

```
1 plt.figure(figsize=(10,6))
2 sns.kdeplot(x=df['atemp'],c='red')
3 plt.show()
```

In [96]:

```
1 plt.figure(figsize=(10,6))
2 sns.kdeplot(x=df['atemp'],c='red')
3 plt.show()
```



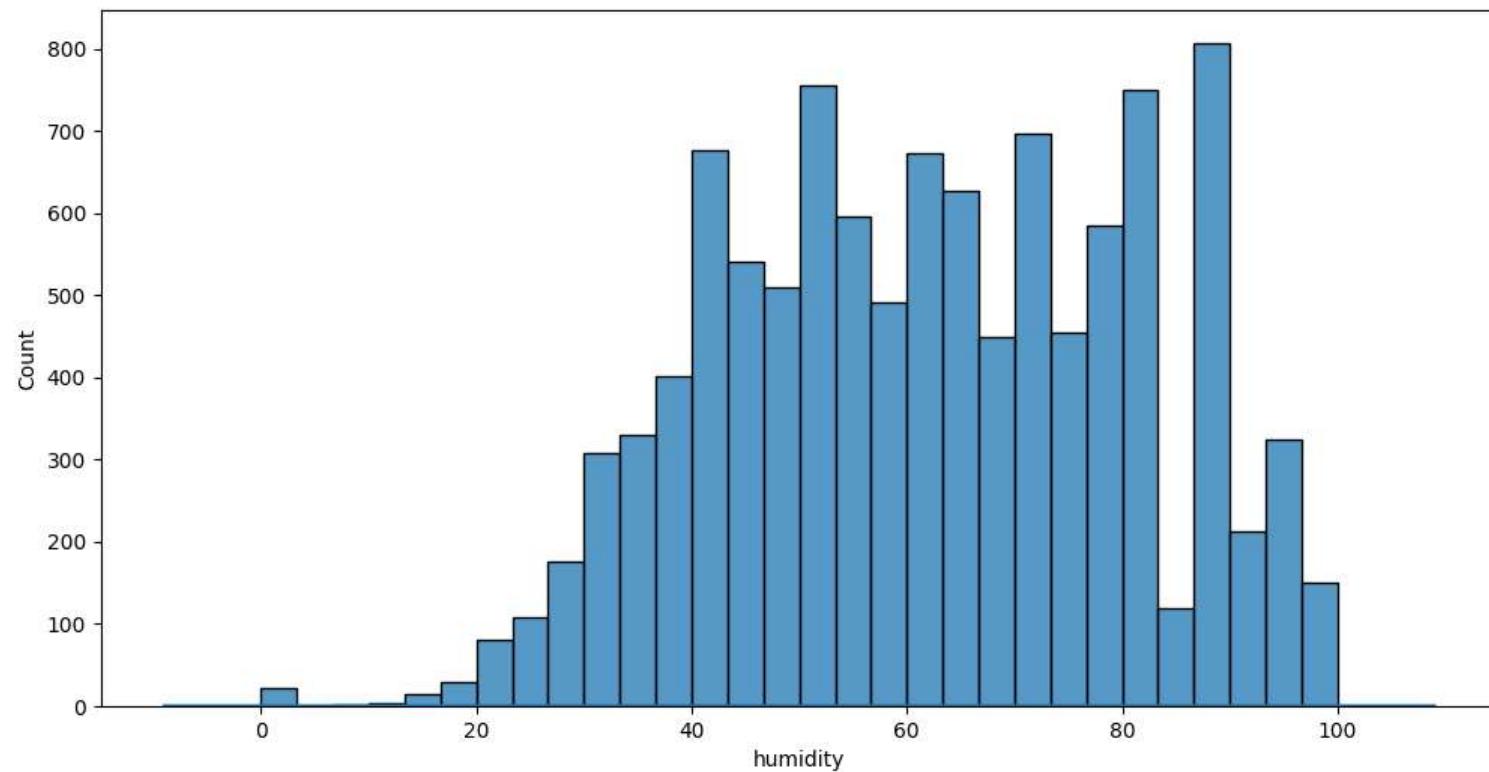
1 **### atemp distribution can be broadly considered as normally distributed**

In [61]:

```
1 plt.figure(figsize=(12,6))
2 sns.histplot(data=df,x='humidity',bins=30)
3 plt.show()
```

In [61]:

```
1 plt.figure(figsize=(12,6))
2 sns.histplot(data=df,x='humidity',bins=30)
3 plt.show()
```

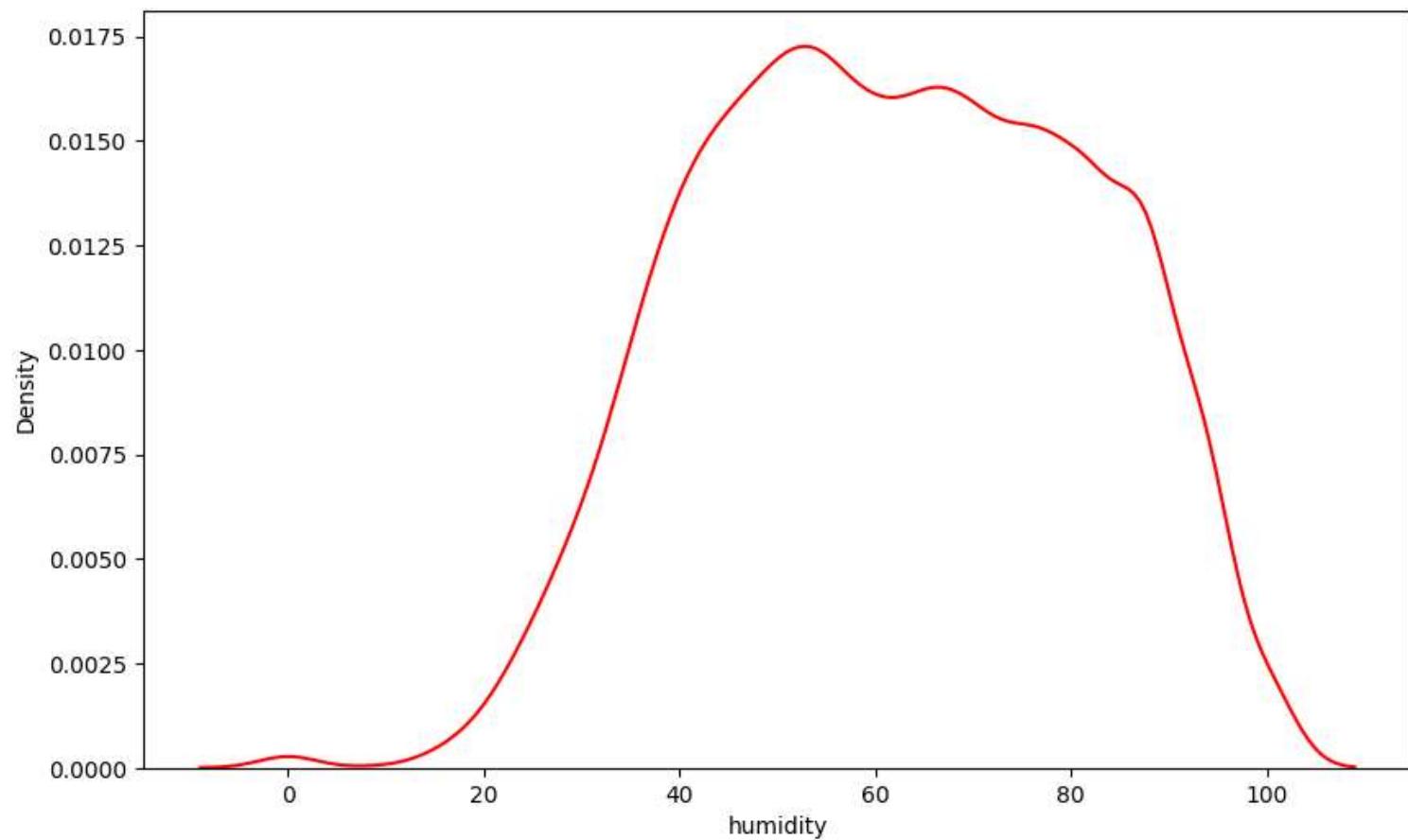


In [97]:

```
1 plt.figure(figsize=(10,6))
2 sns.kdeplot(x=df['humidity'],c='red')
3 plt.show()
```

In [97]:

```
1 plt.figure(figsize=(10,6))
2 sns.kdeplot(x=df['humidity'],c='red')
3 plt.show()
```



```
1 ### humidity distribution can be broadly considered as normally distributed
```

In []:

```
1
```

```
1 ### tried to see whether there are any fluctuations in bike usage over a period of time;
2 ### even though plot is very cluttered, increasing trend can be seen in usage of bike over a
### period of time;
3 ### winter ( around Dec-Jan) shows dip in usage of bikes
```

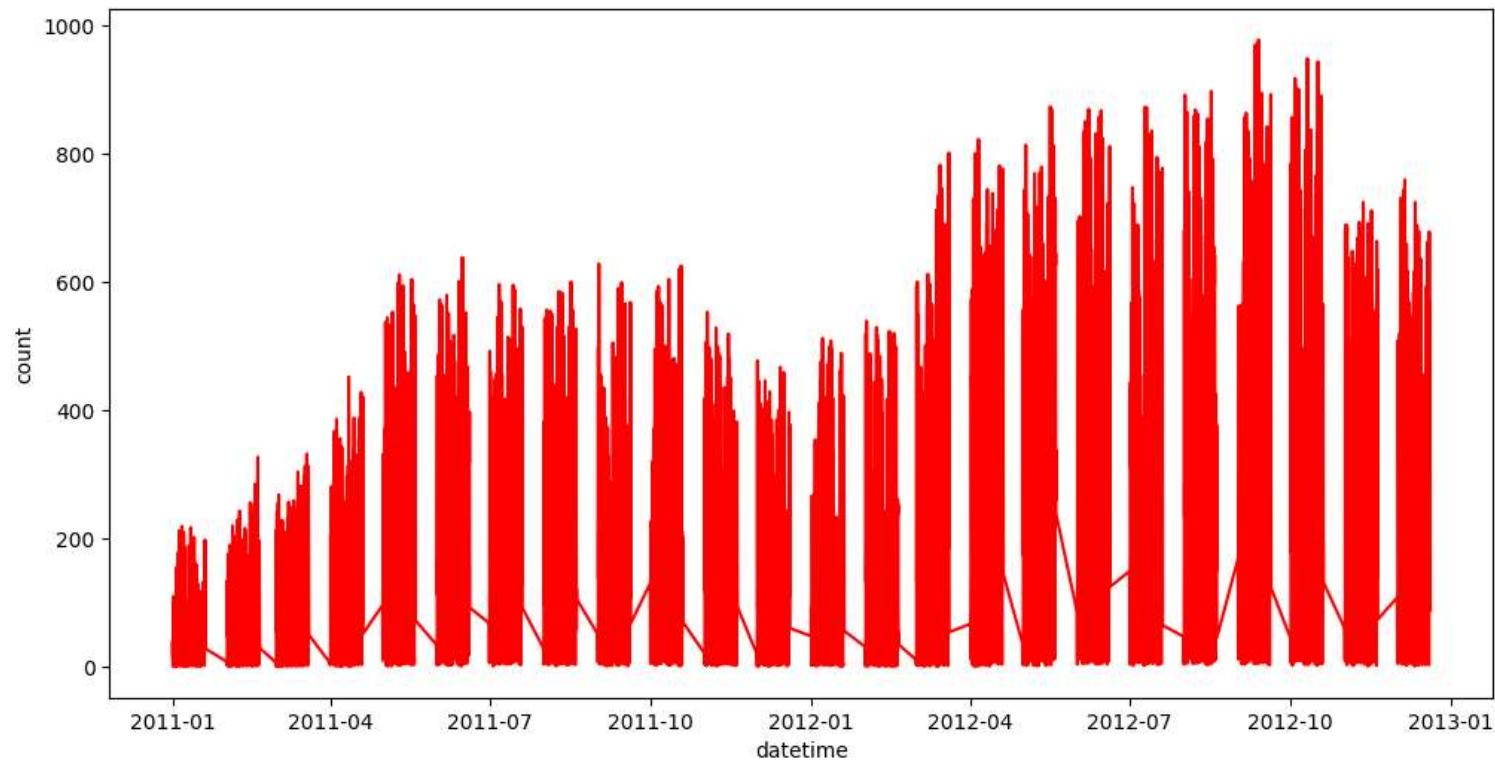
In [67]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df, x='date_time', y='count', hue='humidity')
```

3 **### winter (around Dec-Jan) shows dip in usage of bikes**

In [67]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df, x='datetime',y='count',c='red')
3 plt.show()
```



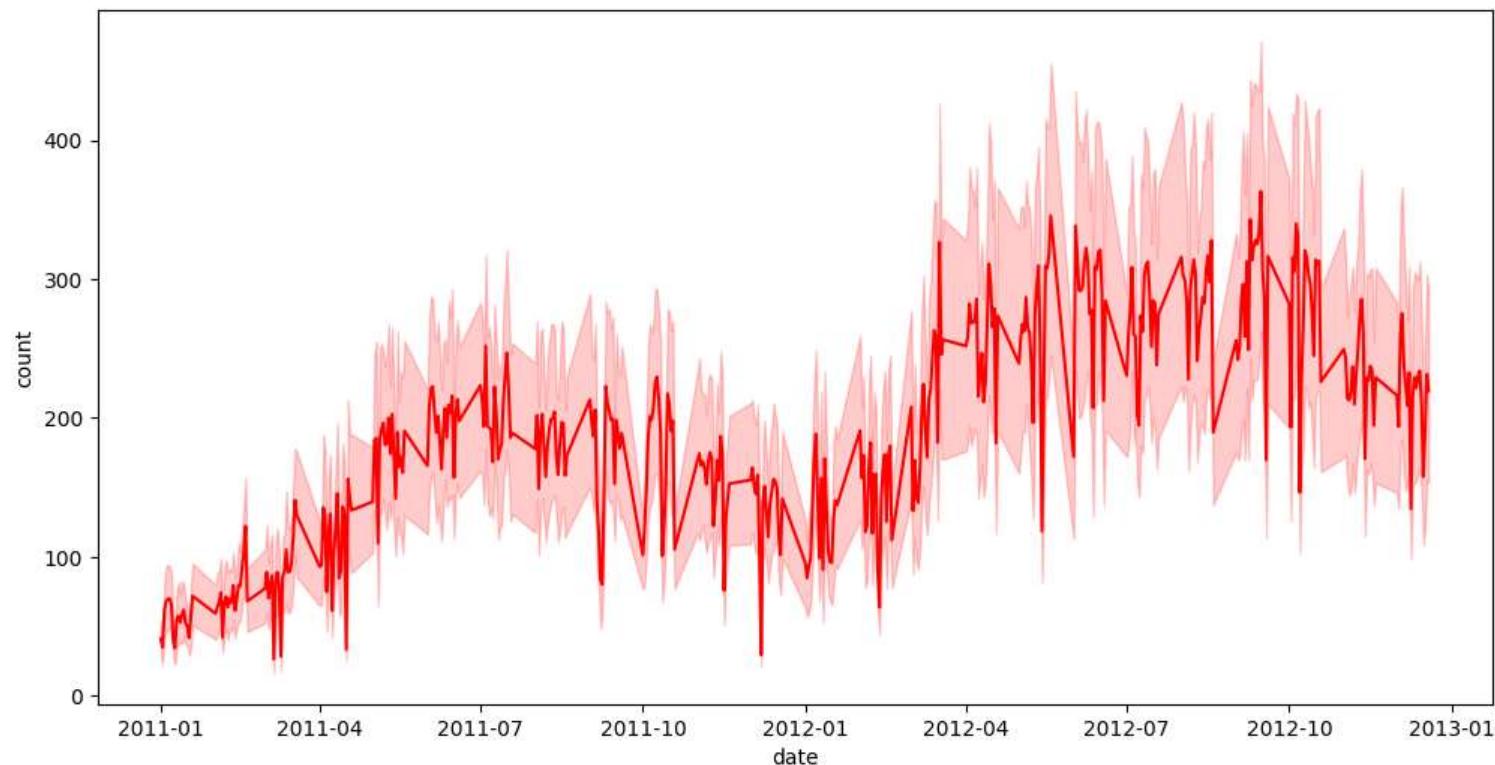
1 **### granularity is increased in subsequent steps in next few charts;**
2 **### also month wise data has been plotted to see the trends more clearly**

In [98]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df, x='date',y='count',c='red')
3 plt.show()
```

In [98]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df, x='date',y='count',c='red')
3 plt.show()
```

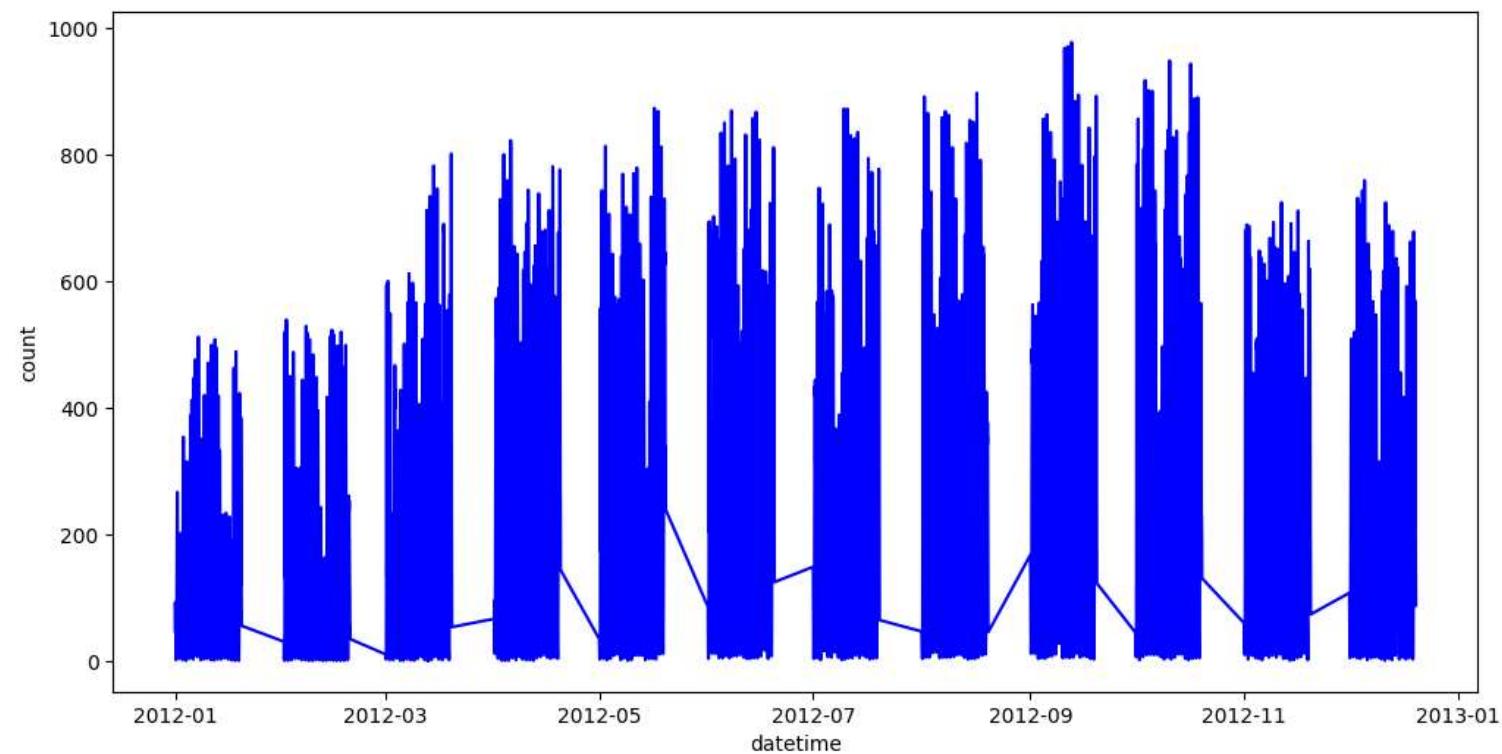


In [73]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df1, x='datetime',y='count',c='blue')
3 plt.show()
```

In [73]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df1, x='datetime',y='count',c='blue')
3 plt.show()
```

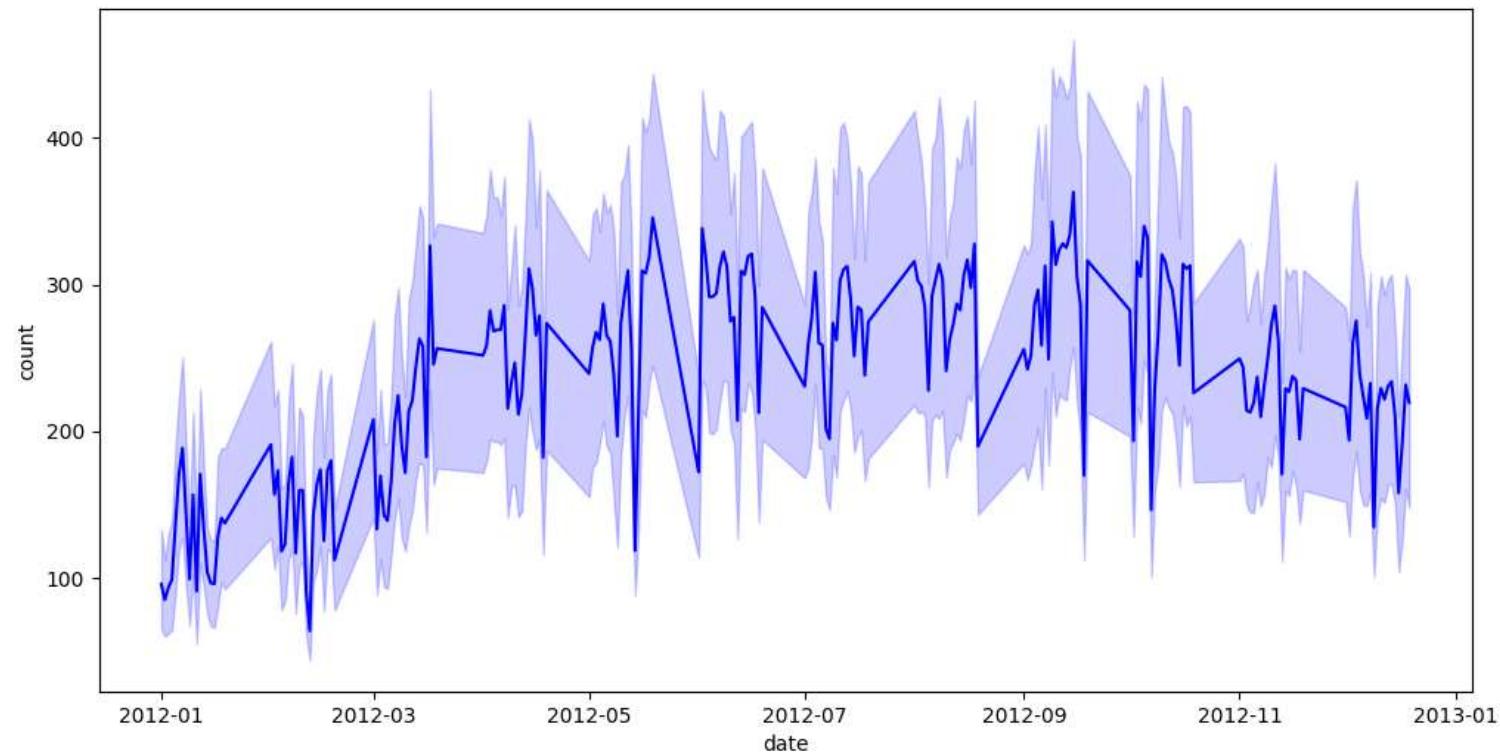


In [93]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df1, x='date',y='count',c='blue')
3 plt.show()
```

In [93]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df1, x='date',y='count',c='blue')
3 plt.show()
```

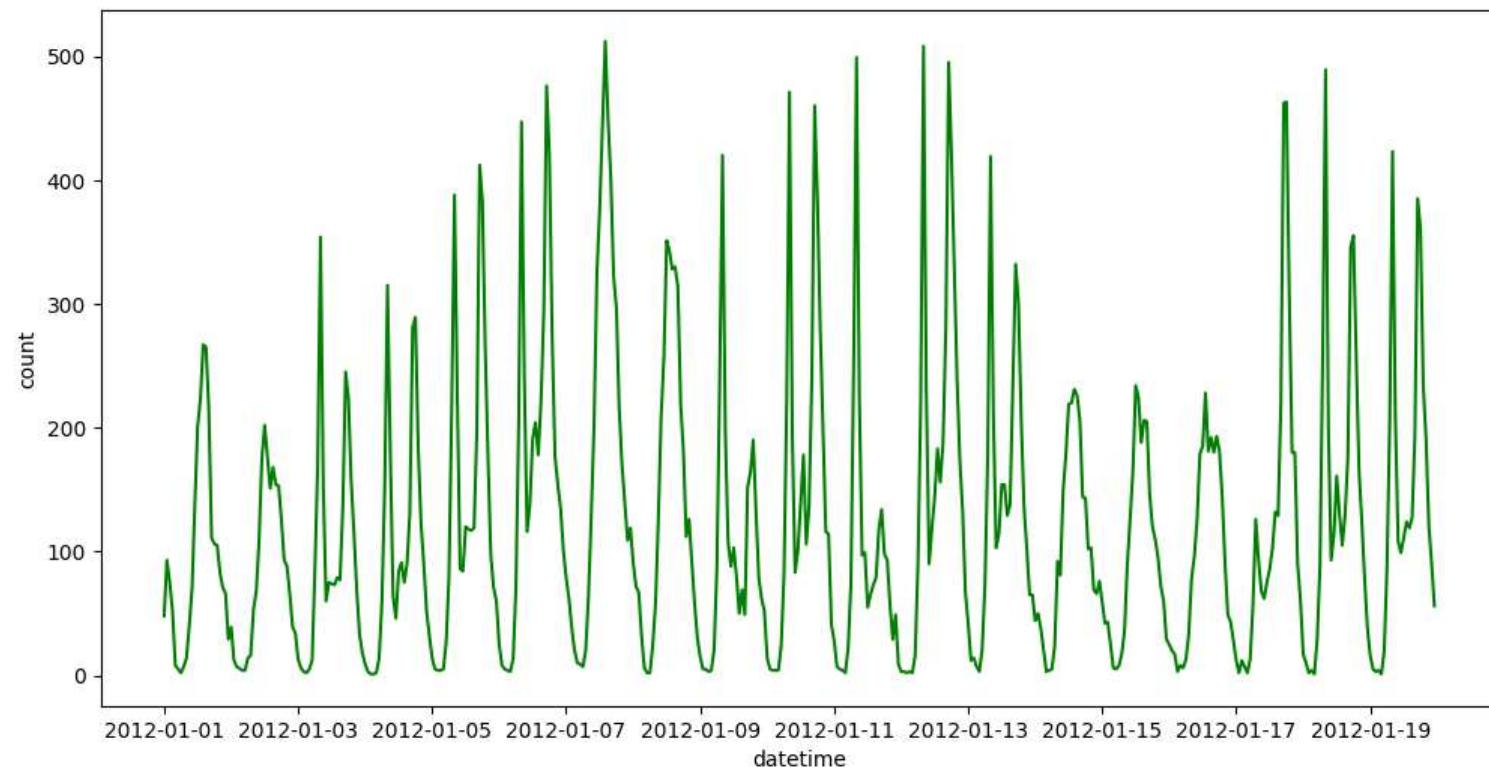


In [89]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df11, x='datetime',y='count',c='green')
3 plt.show()
```

In [89]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df11, x='datetime',y='count',c='green')
3 plt.show()
```



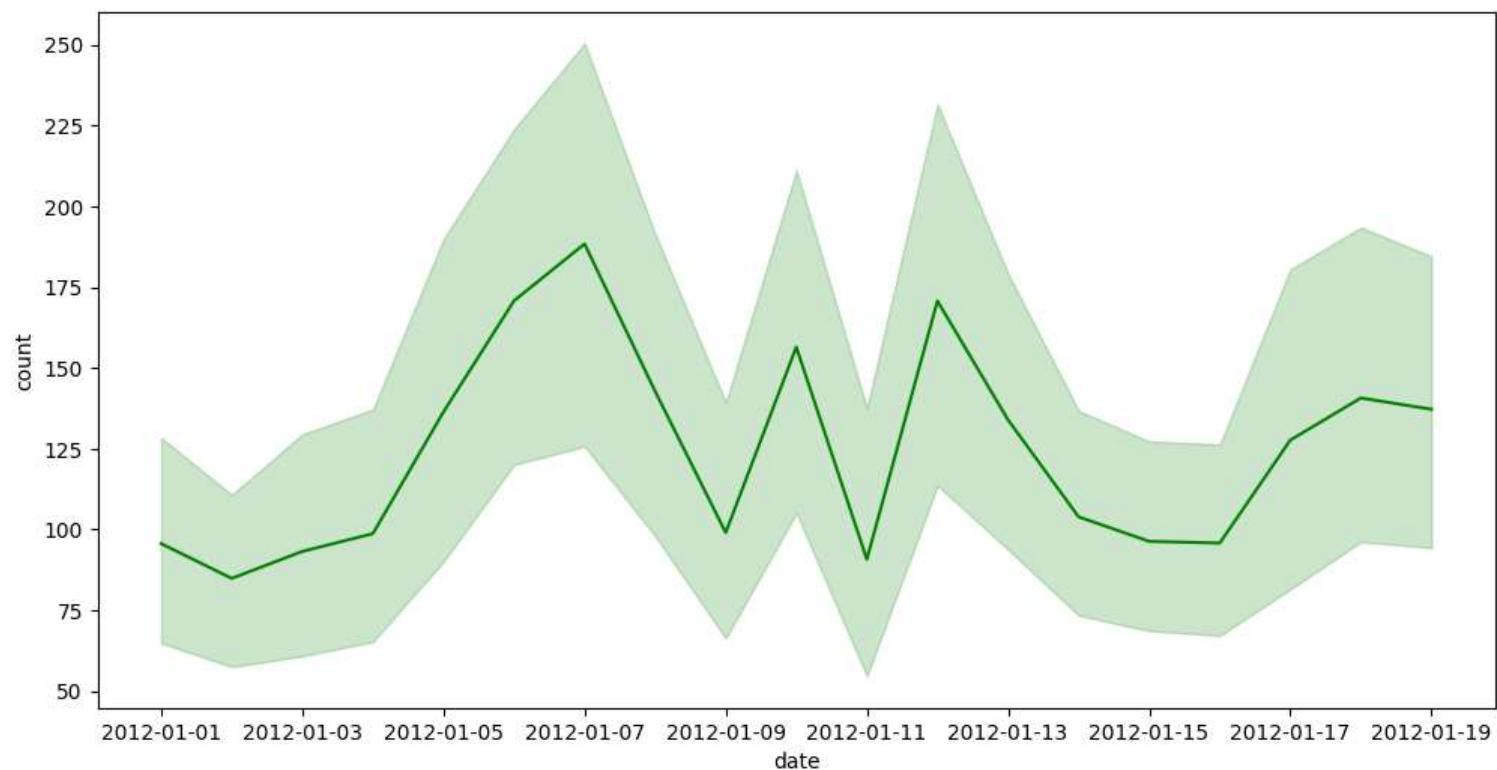
1 **### when I take datetime (timestamp) on x axis, chart is more cluttered... (above chart)**
2 **### when I take date on x axis, chart is more readable (below chart)**

In [88]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df11, x='date',y='count',c='green')
3 plt.show()
```

In [88]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df11, x='date',y='count',c='green')
3 plt.show()
```

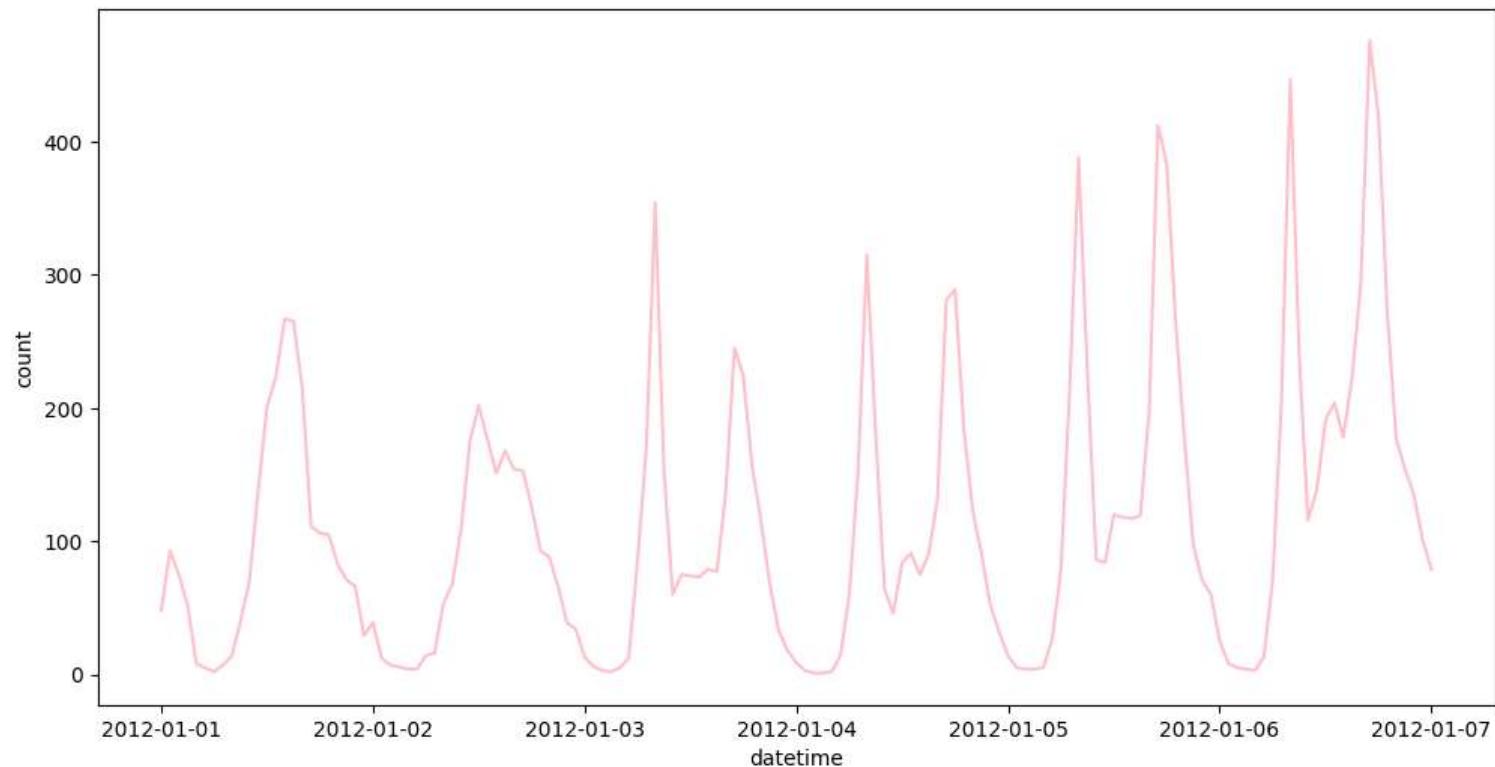


In [91]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df111, x='datetime',y='count',c='pink')
3 plt.show()
```

In [91]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df111, x='datetime',y='count',c='pink')
3 plt.show()
```

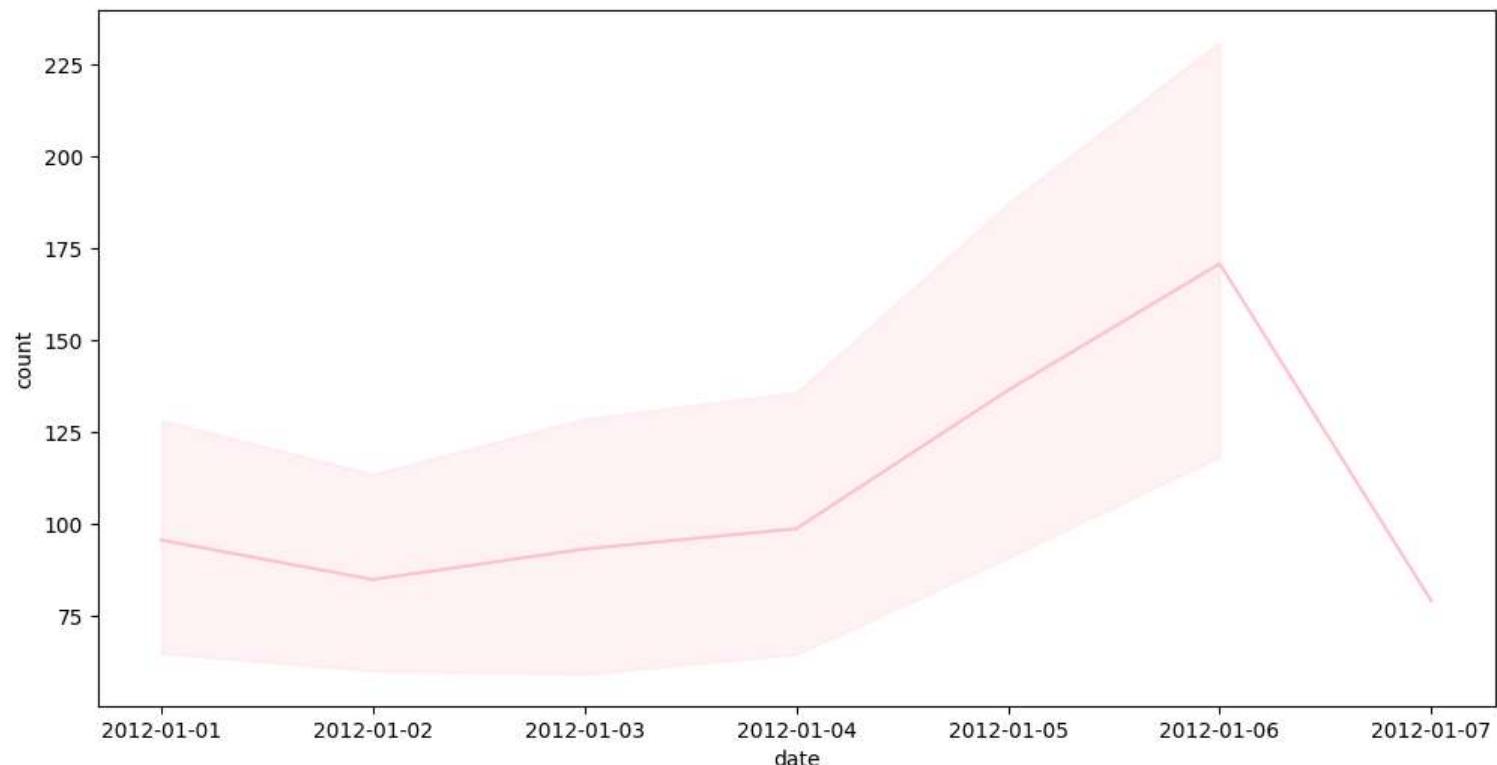


In [90]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df111, x='date',y='count',c='pink')
3 plt.show()
```

In [90]:

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(data=df111, x='date',y='count',c='pink')
3 plt.show()
```



In [107]:

```
1 df_original=df_original.drop(columns='datetime')
```

In [108]:

```
1 df_original
```

Out[108]:

In [108]: 1 df_original

Out[108]:

	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
0	1	0	0	1	9.84	14.395	81	0.0000	3	13	16
1	1	0	0	1	9.02	13.635	80	0.0000	8	32	40
2	1	0	0	1	9.02	13.635	80	0.0000	5	27	32
3	1	0	0	1	9.84	14.395	75	0.0000	3	10	13
4	1	0	0	1	9.84	14.395	75	0.0000	0	1	1
...
10881	4	0	1	1	15.58	19.695	50	26.0027	7	329	336
10882	4	0	1	1	14.76	17.425	57	15.0013	10	231	241
10883	4	0	1	1	13.94	15.910	61	15.0013	4	164	168
10884	4	0	1	1	13.94	17.425	61	6.0032	12	117	129
10885	4	0	1	1	13.12	16.665	66	8.9981	4	84	88

10886 rows × 11 columns

1 **### correlation matrix is created to check the degree of correlation between features and target variable and ### also correlation among the features**

In [109]: 1 df_original.corr()

Out[109]:

season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
--------	---------	------------	---------	------	-------	----------	-----------	--------	------------	-------

In [109]: 1 df_original.corr()

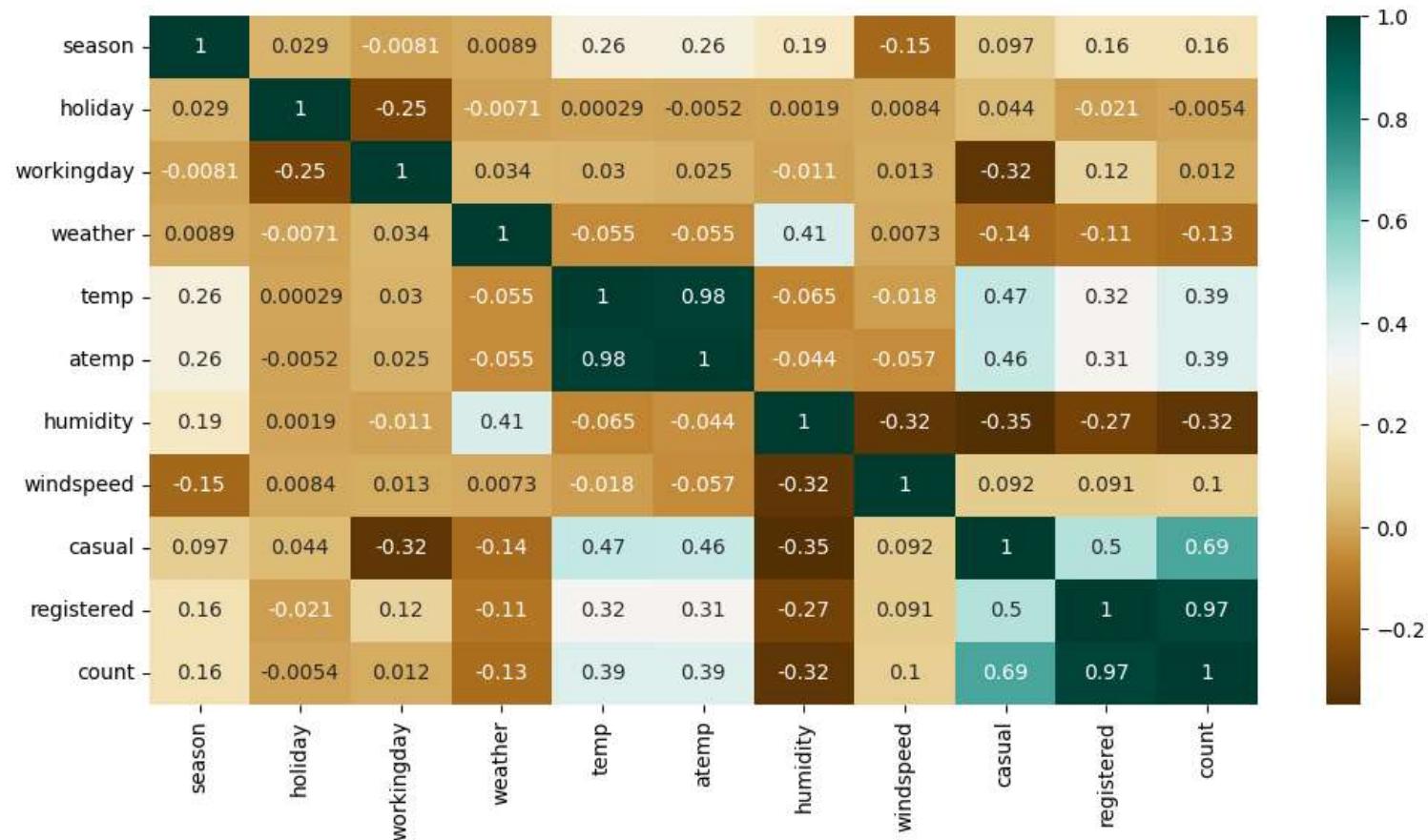
Out[109]:

	season	holiday	workingday	weather	temp	atemp	humidity	windspeed	casual	registered	count
season	1.000000	0.029368	-0.008126	0.008879	0.258689	0.264744	0.190610	-0.147121	0.096758	0.164011	0.163439
holiday	0.029368	1.000000	-0.250491	-0.007074	0.000295	-0.005215	0.001929	0.008409	0.043799	-0.020956	-0.005393
workingday	-0.008126	-0.250491	1.000000	0.033772	0.029966	0.024660	-0.010880	0.013373	-0.319111	0.119460	0.011594
weather	0.008879	-0.007074	0.033772	1.000000	-0.055035	-0.055376	0.406244	0.007261	-0.135918	-0.109340	-0.128655
temp	0.258689	0.000295	0.029966	-0.055035	1.000000	0.984948	-0.064949	-0.017852	0.467097	0.318571	0.394454
atemp	0.264744	-0.005215	0.024660	-0.055376	0.984948	1.000000	-0.043536	-0.057473	0.462067	0.314635	0.389784
humidity	0.190610	0.001929	-0.010880	0.406244	-0.064949	-0.043536	1.000000	-0.318607	-0.348187	-0.265458	-0.317371
windspeed	-0.147121	0.008409	0.013373	0.007261	-0.017852	-0.057473	-0.318607	1.000000	0.092276	0.091052	0.101369
casual	0.096758	0.043799	-0.319111	-0.135918	0.467097	0.462067	-0.348187	0.092276	1.000000	0.497250	0.690414
registered	0.164011	-0.020956	0.119460	-0.109340	0.318571	0.314635	-0.265458	0.091052	0.497250	1.000000	0.970948
count	0.163439	-0.005393	0.011594	-0.128655	0.394454	0.389784	-0.317371	0.101369	0.690414	0.970948	1.000000

In [111]: 1 plt.figure(figsize=(12,6))
2 sns.heatmap(df_original.corr(), annot=True, cmap='BrBG')

```
In [111]: 1 plt.figure(figsize=(12,6))
2 sns.heatmap(df_original.corr(), annot=True, cmap='BrBG')
```

Out[111]: <Axes: >



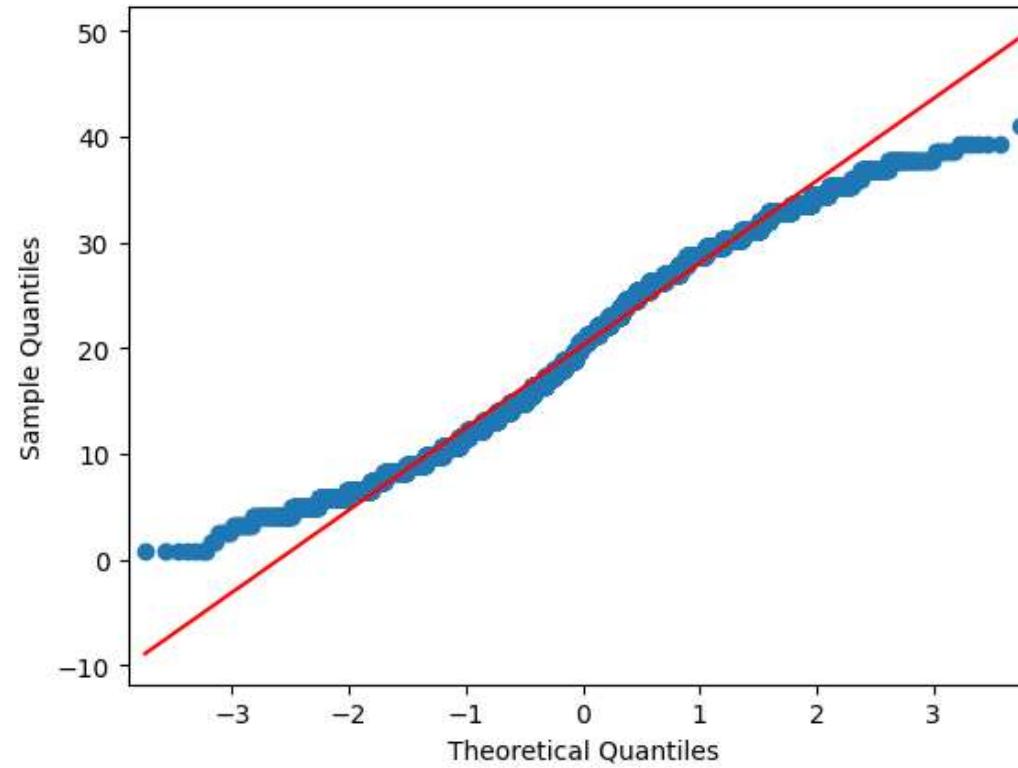
```
1 ##### qqplots have been plotted to check normality of feature data
2 ##### (I dont know why 2 charts are coming for each qq plot, even though they are same)
```

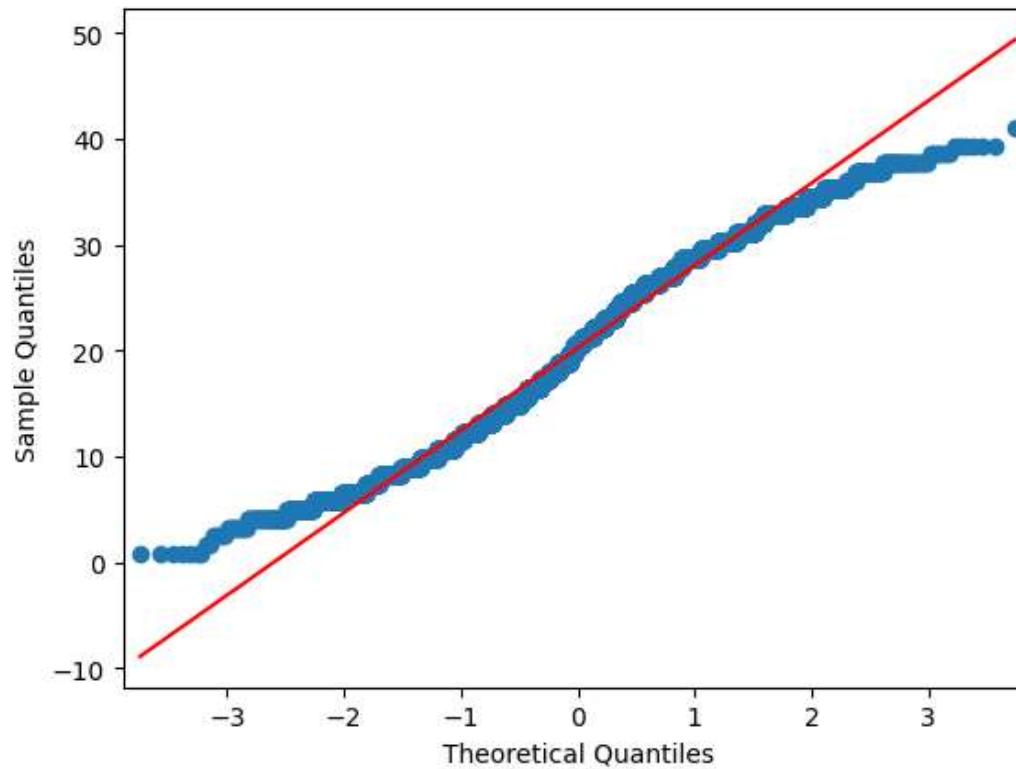
```
In [177]: 1 qqplot(df["temp"], line="s")
```

Out[177]:

```
In [177]: 1 qqplot(df[ "temp"],line="s")
```

Out[177]:





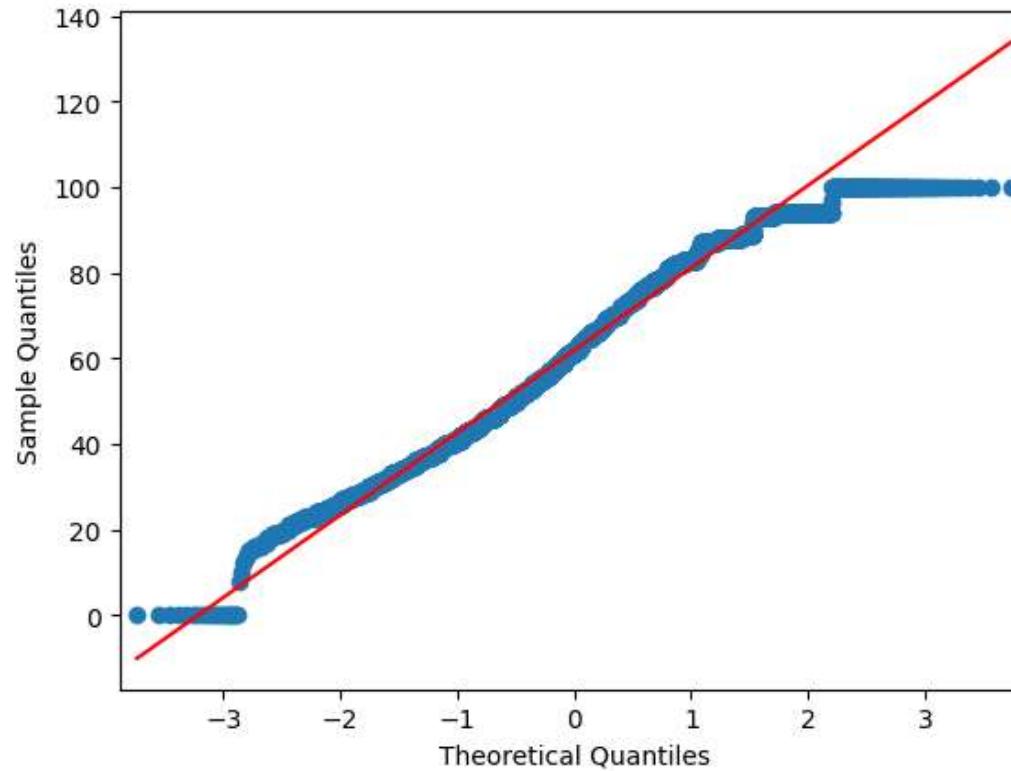
In [170]: 1 `qqplot(df["humidity"], line="s")`

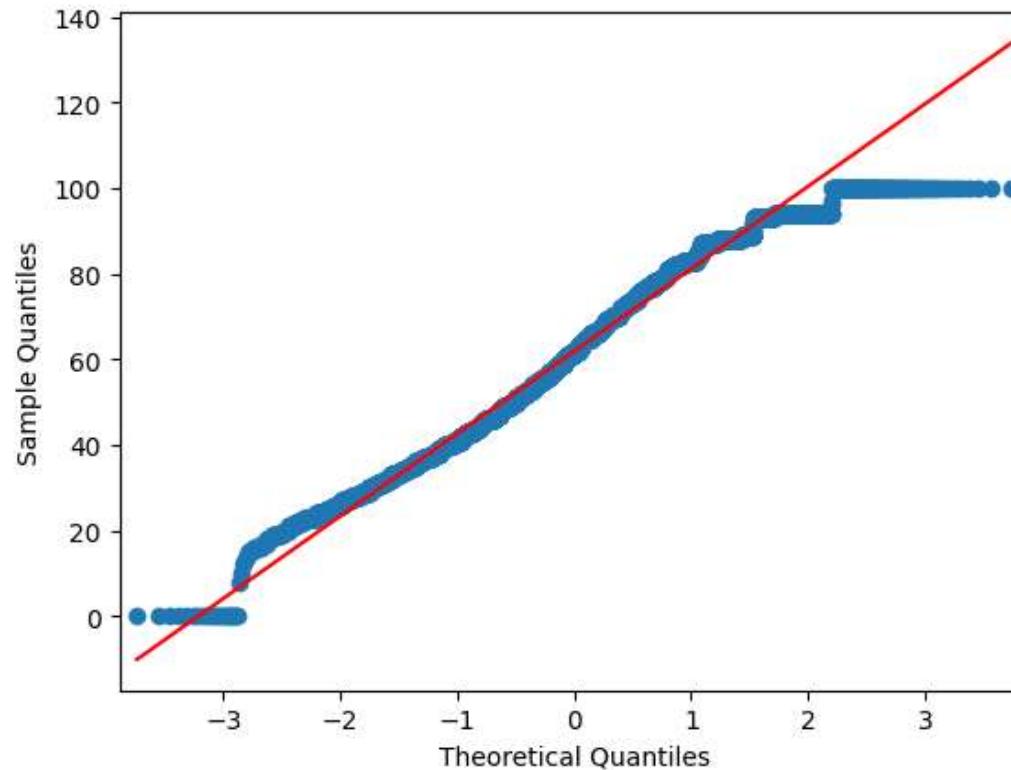
Out[170]:

140 ↴

```
In [170]: 1 qqplot(df["humidity"],line="s")
```

Out[170]:





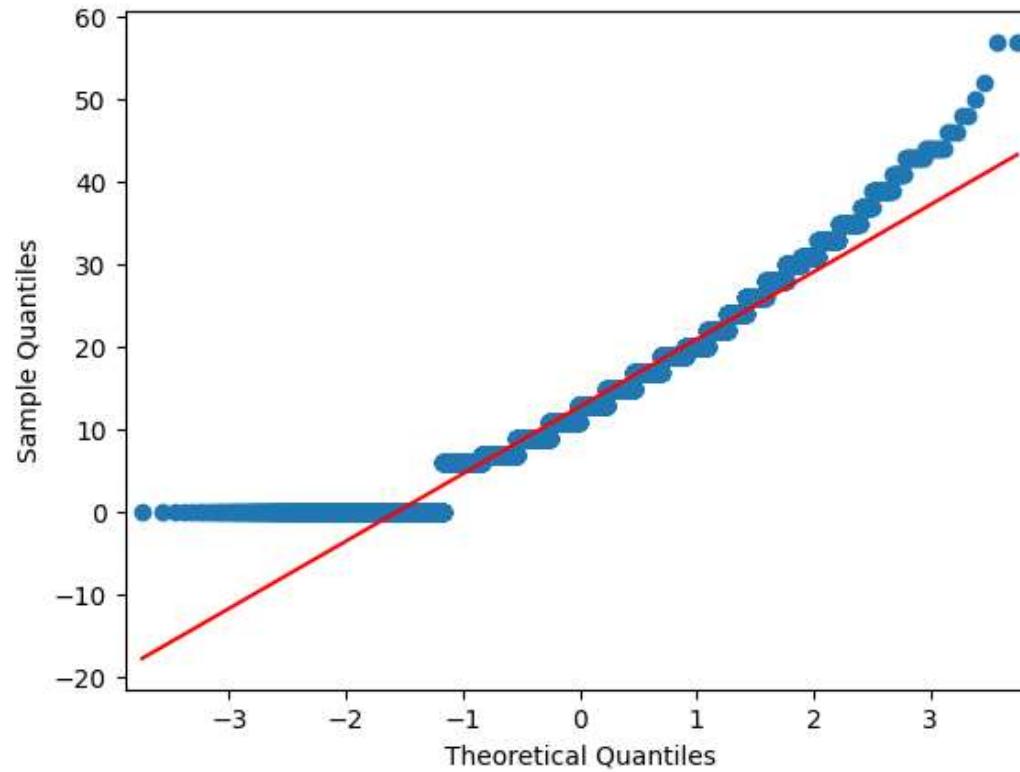
In [173]: 1 | `qqplot(df["windspeed"], line="s")`

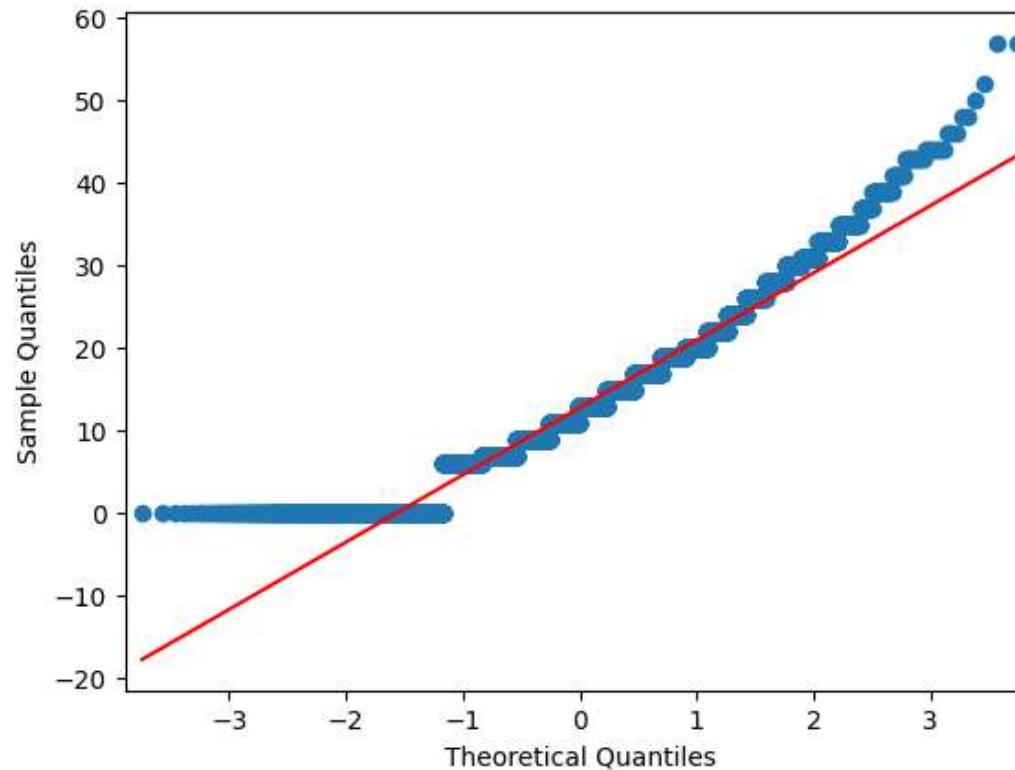
Out[173]:

60

In [173]: 1 | `qqplot(df["windspeed"],line="s")`

Out[173]:





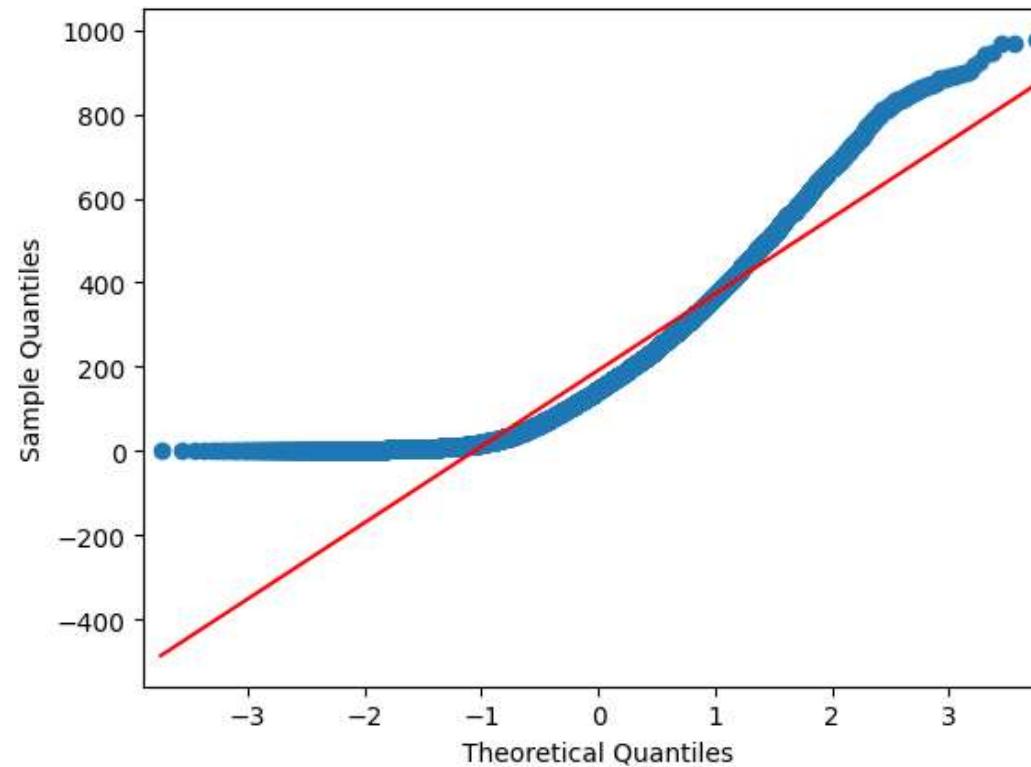
In [174]: 1 qqplot(df["count"], line="s")

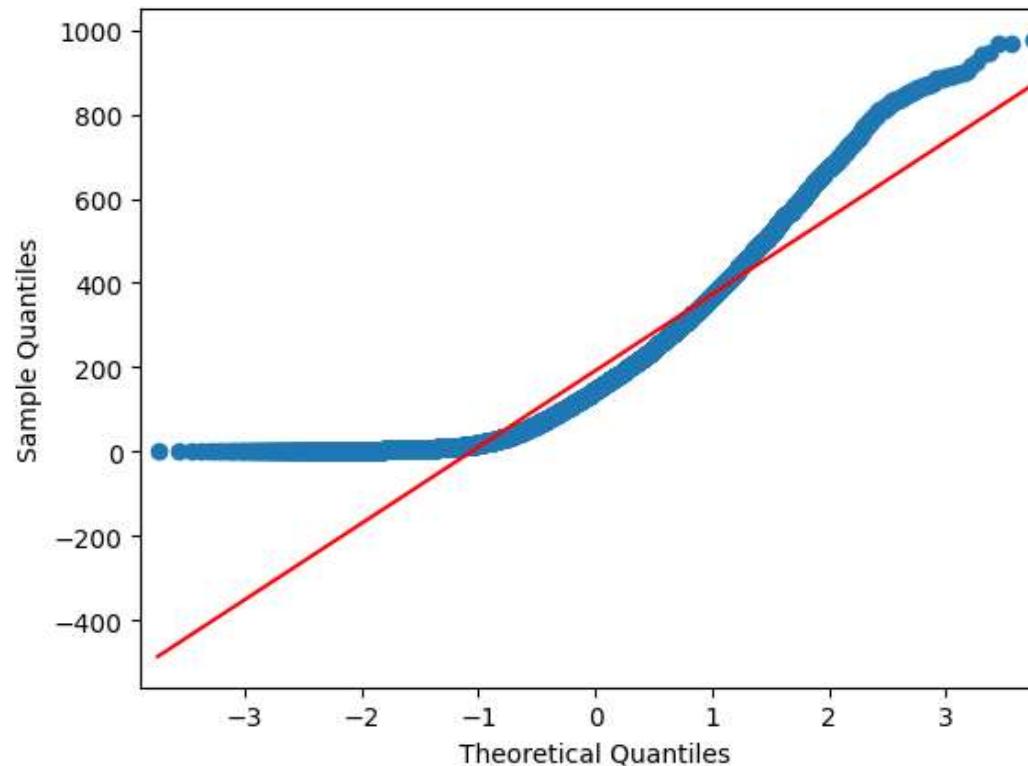
Out[174]:

1000

In [174]: 1 | `qqplot(df["count"],line="s")`

Out[174]:





In []: 1

```
1 #### temp, humidity, windspeed features and count (target variable) do not show normality
```

```
1 ##### from the correlation matrix we see working day feature has 0.012 correlation with count (bikes used or rented)
2 ##### also we can conduct a hypothesis test to check whether working day has any effect on bike usage
3 ##### before conducting 2 sample t test, we need to check whether variances are equal or not
```

```
In [120]: 1 data_group1 = df[df['workingday']==0]['count'].values
2 data_group2 = df[df['workingday']==1]['count'].values
3
4 var1=np.var(data_group1)
5 var2=np.var(data_group2)
```

```
In [121]: 1 print(var1)
```

30171.346098942427

```
In [121]: 1 print(var1)
```

```
30171.346098942427
```

```
In [122]: 1 print(var2)
```

```
34040.69710674686
```

```
In [126]: 1 var2/var1
```

```
Out[126]: 1.1282458858519429
```

```
In [113]: 1 data_group1
```

```
Out[113]: array([ 16,  40,  32, ..., 106,  89,  33], dtype=int64)
```

```
In [115]: 1 data_group1.shape
```

```
Out[115]: (3474,)
```

```
In [114]: 1 data_group2
```

```
Out[114]: array([ 5,  2,  1, ..., 168, 129,  88], dtype=int64)
```

```
In [116]: 1 data_group2.shape
```

```
Out[116]: (7412,)
```

```
In [119]: 1 # Ho : Var is Equal  
2 # Ha : Var is NOT Equal  
3 # alpha = 0.05  
4 levene(data_group1,data_group2)  
5  
6
```

```
Out[119]: LeveneResult(statistic=0.004972848886504472, pvalue=0.9437823280916695)
```

1 ##### according to levene test, since p value is more than 0.05, we fail to reject null hypothesis; variances are equal
2 ##### also the ratio var1/var2 less than 4, we can say variances are equal; we can proceed with ttest

```
In [127]: 1 # Ho: working day has no effect on number of bikes used  
2 # Ha: working day has effect on number of bikes used
```

3 ##### also the ratio var1/var2 less than 4, we can say variances are equal; we can proceed with ttest

In [127]:

```
1 # Ho: working day has no effect on number of bikes used  
2 # Ha: working day has effect on number of bikes used  
3 # alpha=0.05  
4  
5 ttest_ind(data_group1, data_group2)
```

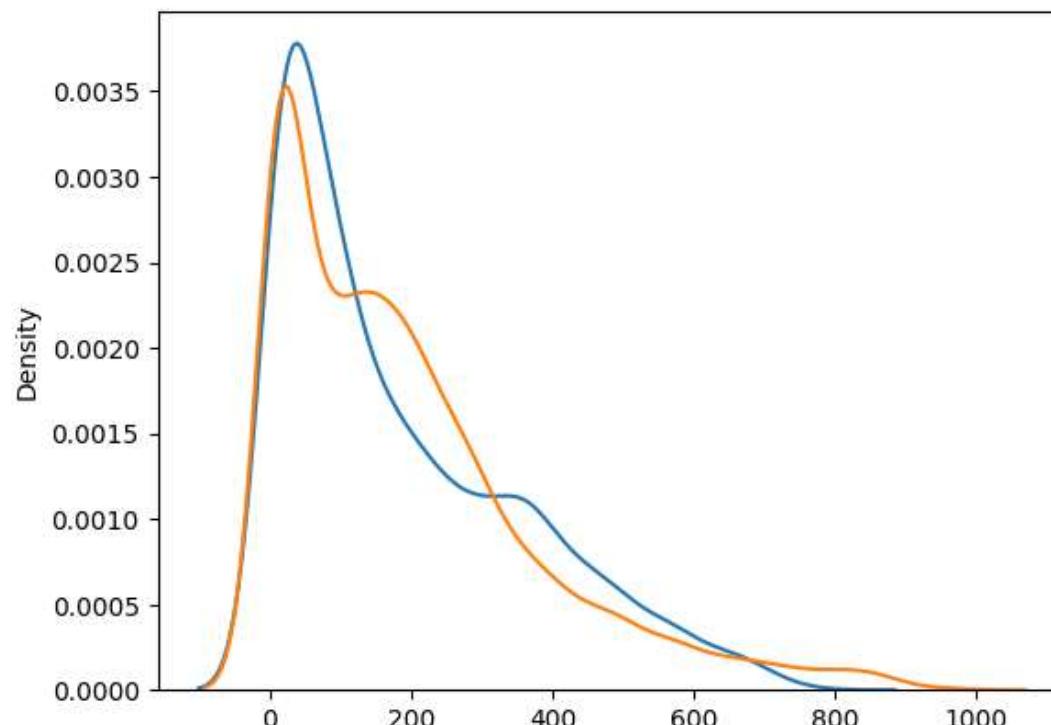
Out[127]: Ttest_indResult(statistic=-1.2096277376026694, pvalue=0.22644804226361348)

1 ##### p value is more than alpha, we fail to reject null hypothesis => working day feature has no effect on bikes used

In [125]:

```
1 sns.kdeplot(data_group1)  
2 sns.kdeplot(data_group2)
```

Out[125]: <Axes: ylabel='Density'>



2 ##### workingday=1

1 ##### kde plot also shows NOT MUCH BIG difference between the distributions of workingday=0 (not working day) and correlation matrix we see correlation coefficient = 0.0089 between weather and season

```

2 ##### workingday=1
1 ##### kde plot also shows NOT MUCH BIG difference between the distributions of workingday=0 (not working
1 day) and correlation matrix we see correlation coefficient = 0.0089 between weather and season
2
3 ##### for Testing whether weather and season are independent or ralated; we will do a chisquare test to test this

```

In []: █

```

1 # Ho: weather and season are independent
2 # Ha: weather and season are NOT independent
3 # alpha=0.05 (significance Level)

```

In [141]: █

```

1 alpha=0.05

```

In [131]: █

```

1 observed = pd.crosstab(df['season'], df['weather'])
2 print("observed values:")
3 observed

```

observed values:

Out[131]:

	weather	1	2	3	4
season					
1	1759	715	211	1	
2	1801	708	224	0	
3	1930	604	199	0	
4	1702	807	225	0	

In [136]: █

```

1 chi_ststistic, p_value, dof,exp_values = chi2_contingency(observed)

```

In [137]: █

```

1 chi_ststistic

```

Out[137]: 49.158655596893624

In [138]: █

```

1 p_value

```

Out[138]: 1.549925073686492e-07

In [139]: █

```

1 dof

```

Out[139]: 9

In [139]: 1 dof

Out[139]: 9

In [140]: 1 exp_values

Out[140]: array([[1.77454639e+03, 6.99258130e+02, 2.11948742e+02, 2.46738931e-01],
 [1.80559765e+03, 7.11493845e+02, 2.15657450e+02, 2.51056403e-01],
 [1.80559765e+03, 7.11493845e+02, 2.15657450e+02, 2.51056403e-01],
 [1.80625831e+03, 7.11754180e+02, 2.15736359e+02, 2.51148264e-01]])

In [142]: 1 if p_value < alpha:
 2 print("Reject Null Hypothesis")
 3 else:
 4 print("Fail to reject null hypothesis")

Reject Null Hypothesis

1 ##### weather and season are NOT independent ; they are dependent

1 ## testing whether different weathers have impact on bike usage

1 ### we will do ANOVA of all weathers with usage of bikes

In []: 1 # Ho: bike usage is independent of weathers
 2 # Ha: bike usage is dependent on atleast one weather
 3 # alpha=0.05 (significance Level)

In [149]: 1 w1 = df[df['weather']==1]['count'].values
 2 w2 = df[df['weather']==2]['count'].values
 3 w3 = df[df['weather']==3]['count'].values
 4 w4 = df[df['weather']==4]['count'].values

In [150]: 1 f_statistic, p_value = f_oneway(w1,w2,w3,w4)

In [151]: 1 print('f_statistic', f_statistic)

f_statistic 65.53024112793271

In [152]: 1 print('p_value', p_value)

p_value 5.482069475935669e-42

```
In [152]: t_statistic 65.53024112/932/1
1 print('p_value', p_value)

p_value 5.482069475935669e-42
```

```
In [153]: 1 if p_value < alpha:
2     print("Reject Null Hypothesis")
3 else:
4     print("Fail to reject null hypothesis")

Reject Null Hypothesis
```

```
1 ##### bike usage is dependent on atleast one weather
2 ##### we can also see a correlation of (-0.13) between weather feature and bike usage (count) from correlation
matrix
3
```

```
In [ ]: 1

1 ## testing whether different seasons have impact on bike usage

1 ### we will do ANOVA of all seasons with usage of bikes
```

```
In [ ]: 1 # Ho: bike usage is independent of seasons
2 # Ha: bike usage is dependent on atleast one season
3 # alpha=0.05 (significance Level)
```

```
In [ ]: 1
```

```
In [154]: 1 s1 = df[df['season']==1]['count'].values
2 s2 = df[df['season']==2]['count'].values
3 s3 = df[df['season']==3]['count'].values
4 s4 = df[df['season']==4]['count'].values
```

```
In [155]: 1 f_statistic, p_value = f_oneway(s1,s2,s3,s4)
```

```
In [156]: 1 print('f_statistic', f_statistic)
```

```
In [157]: 1 f_statistic 236.94671081032106
1 print('p_value', p_value)

p_value 6.164843386499654e-149
```

```
In [157]: f_statistic(p236.194671081032106)
```

```
p_value 6.164843386499654e-149
```

```
In [158]: 1 if p_value < alpha:  
2     print("Reject Null Hypothesis")  
3 else:  
4     print("Fail to reject null hypothesis")
```

```
Reject Null Hypothesis
```

```
1 ##### bike usage is dependent on atleast one season  
2 ##### we can also see a correlation of (0.16) between season feature and bike usage (count) from correlation  
3 matrix
```

```
In [ ]: 1
```