```
1  import numpy as np # linear algebra
2  import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
3  from scipy import stats
4  import matplotlib.pyplot as plt
5  import seaborn as sns
```

```
1  csv_path = "yulu dataset.txt"
2  df = pd.read_csv(csv_path, delimiter=",")
3  df.head()
```

|   | 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 | 0.0 | 3 | 13 | 16 |
| 1 | 2011-01-01 01:00:00 | 1 | 0 | 0 | 1 | 9.02 | 13.635 | 80 | 0.0 | 8 | 32 | 40 |
| 2 | 2011-01-01 02:00:00 | 1 | 0 | 0 | 1 | 9.02 | 13.635 | 80 | 0.0 | 5 | 27 | 32 |
| 3 | 2011-01-01 03:00:00 | 1 | 0 | 0 | 1 | 9.84 | 14.395 | 75 | 0.0 | 3 | 10 | 13 |
| 4 | 2011-01-01 04:00:00 | 1 | 0 | 0 | 1 | 9.84 | 14.395 | 75 | 0.0 | 0 | 1 | 1 |

```
1  # no of rows amd columns in dataset
2  print(f"# rows: {df.shape[0]} \n# columns: {df.shape[1]}")
```

```
# rows: 10886
# columns: 12
```

```
1  df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10886 entries, 0 to 10885
Data columns (total 12 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   datetime    10886 non-null  object
 1   season      10886 non-null  int64
 2   holiday     10886 non-null  int64
 3   workingday  10886 non-null  int64
 4   weather     10886 non-null  int64
 5   temp        10886 non-null  float64
 6   atemp       10886 non-null  float64
 7   humidity    10886 non-null  int64
 8   windspeed   10886 non-null  float64
 9   casual      10886 non-null  int64
 10  registered  10886 non-null  int64
 11  count       10886 non-null  int64
dtypes: float64(3), int64(8), object(1)
memory usage: 1020.7+ KB
```

```
1 df['datetime'] = pd.to_datetime(df['datetime'])
2
3 cat_cols= ['season', 'holiday', 'workingday', 'weather']
4 for col in cat_cols:
5     df[col] = df[col].astype('object')
```

```
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  object
 2   holiday     10886 non-null  object
 3   workingday  10886 non-null  object
 4   weather     10886 non-null  object
 5   temp        10886 non-null  float64
 6   atemp       10886 non-null  float64
 7   humidity    10886 non-null  int64
 8   windspeed   10886 non-null  float64
 9   casual      10886 non-null  int64
 10  registered  10886 non-null  int64
 11  count       10886 non-null  int64
dtypes: datetime64[ns](1), float64(3), int64(4), object(4)
memory usage: 1020.7+ KB
```

```
1 df.iloc[:, 1:].describe(include='all')
```

| | season | holiday | workingday | weather | temp | atemp | humidity | windspeed | casual | registered | 108 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 10886.0 | 10886.0 | 10886.0 | 10886.0 | 10886.00000 | 10886.000000 | 10886.000000 | 10886.000000 | 10886.000000 | 10886.000000 | 108 |
| unique | 4.0 | 2.0 | 2.0 | 4.0 | NaN | NaN | NaN | NaN | NaN | NaN | |
| top | 4.0 | 0.0 | 1.0 | 1.0 | NaN | NaN | NaN | NaN | NaN | NaN | |
| freq | 2734.0 | 10575.0 | 7412.0 | 7192.0 | NaN | NaN | NaN | NaN | NaN | NaN | |
| mean | NaN | NaN | NaN | NaN | 20.23086 | 23.655084 | 61.886460 | 12.799395 | 36.021955 | 155.552177 | 1 |
| std | NaN | NaN | NaN | NaN | 7.79159 | 8.474601 | 19.245033 | 8.164537 | 49.960477 | 151.039033 | 1 |
| min | NaN | NaN | NaN | NaN | 0.82000 | 0.760000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | |
| 25% | NaN | NaN | NaN | NaN | 13.94000 | 16.665000 | 47.000000 | 7.001500 | 4.000000 | 36.000000 | |
| 50% | NaN | NaN | NaN | NaN | 20.50000 | 24.240000 | 62.000000 | 12.998000 | 17.000000 | 118.000000 | 1 |
| 75% | NaN | NaN | NaN | NaN | 26.24000 | 31.060000 | 77.000000 | 16.997900 | 49.000000 | 222.000000 | 2 |
| max | NaN | NaN | NaN | NaN | 41.00000 | 45.455000 | 100.000000 | 56.996900 | 367.000000 | 886.000000 | 9 |

```
1 # detecting missing values in the dataset
2 df.isnull().sum()
```

```
datetime      0
season        0
holiday       0
workingday    0
weather       0
temp          0
atemp         0
humidity      0
windspeed     0
casual        0
registered    0
count         0
dtype: int64
```

```
1 # minimum datetime and maximum datetime
2 print(df['datetime'].min(), df['datetime'].max())
3 # number of unique values in each categorical columns
4 df[cat_cols].melt().groupby(['variable', 'value'])[['value']].count()
```
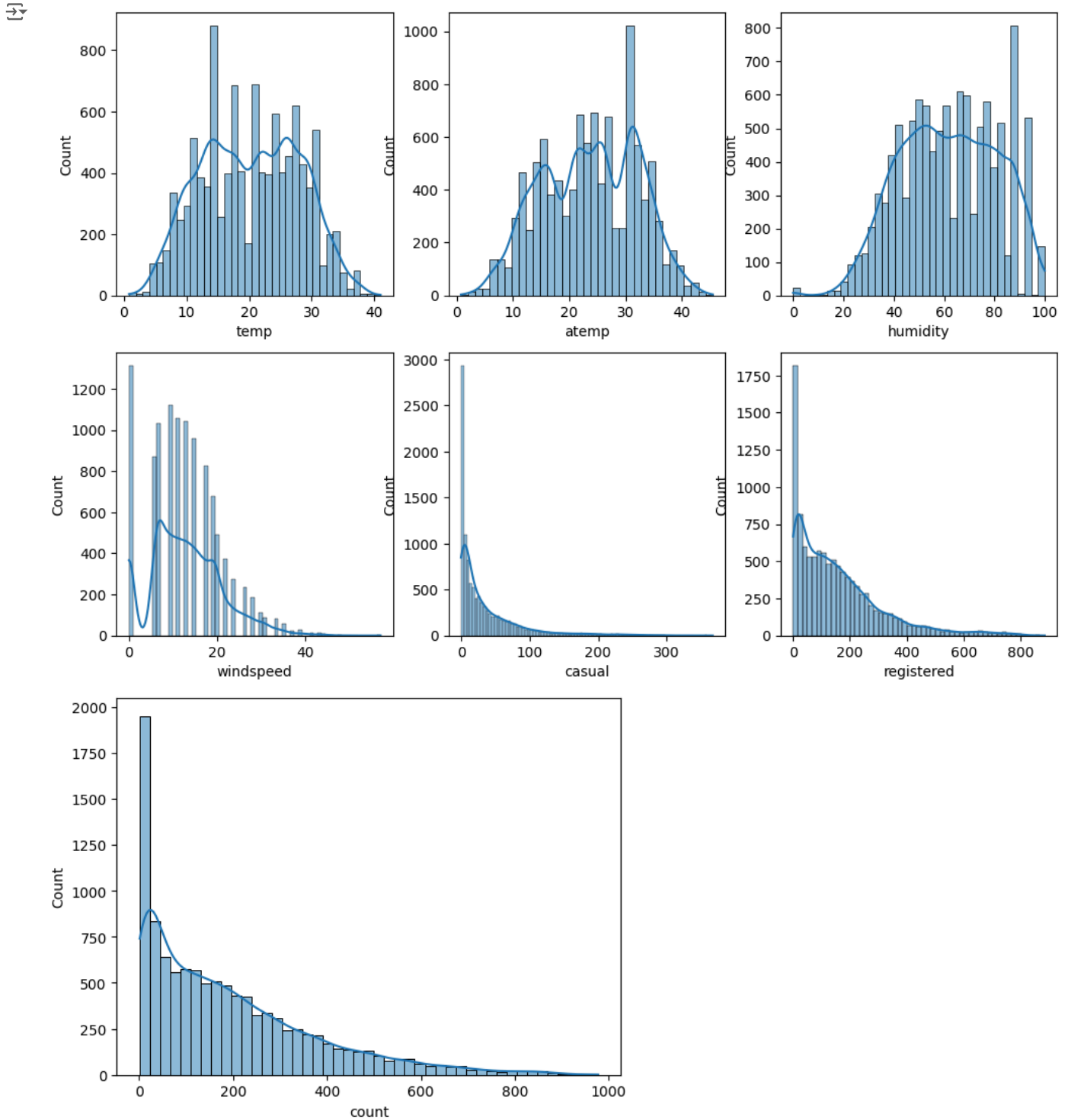
```
2011-01-01 00:00:00 2012-12-19 23:00:00
```

| | | value |
|---|---|---|
| **variable** | **value** | |
| holiday | 0 | 10575 |
| | 1 | 311 |
| season | 1 | 2686 |
| | 2 | 2733 |
| | 3 | 2733 |
| | 4 | 2734 |
| weather | 1 | 7192 |
| | 2 | 2834 |
| | 3 | 859 |
| | 4 | 1 |
| workingday | 0 | 3474 |
| | 1 | 7412 |

```
 1 # understanding the distribution for numerical variables
 2 num_cols = ['temp', 'atemp', 'humidity', 'windspeed', 'casual', 'registered','count']
 3
 4 fig, axis = plt.subplots(nrows=2, ncols=3, figsize=(12, 8))
 5
 6 index = 0
 7 for row in range(2):
 8     for col in range(3):
 9         sns.histplot(df[num_cols[index]], ax=axis[row, col], kde=True)
10         index += 1
11
12 plt.show()
13 sns.histplot(df[num_cols[-1]], kde=True)
14 plt.show()
```
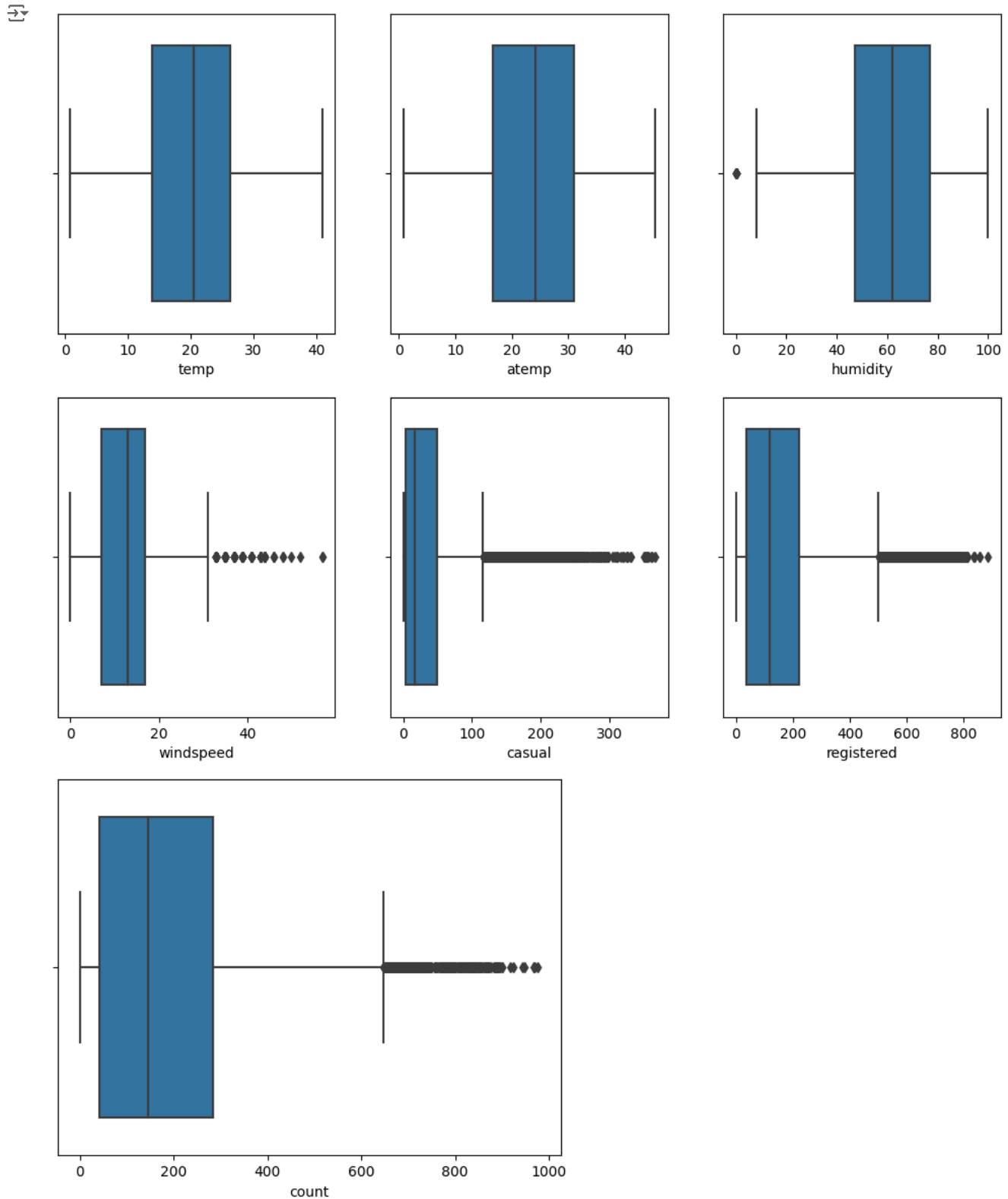
```
1 # plotting box plots to detect outliers in the data
2 fig, axis = plt.subplots(nrows=2, ncols=3, figsize=(12, 9))
3
4 index = 0
5 for row in range(2):
6     for col in range(3):
7         sns.boxplot(x=df[num_cols[index]], ax=axis[row, col])
8         index += 1
9
10 plt.show()
11 sns.boxplot(x=df[num_cols[-1]])
12 plt.show()
```
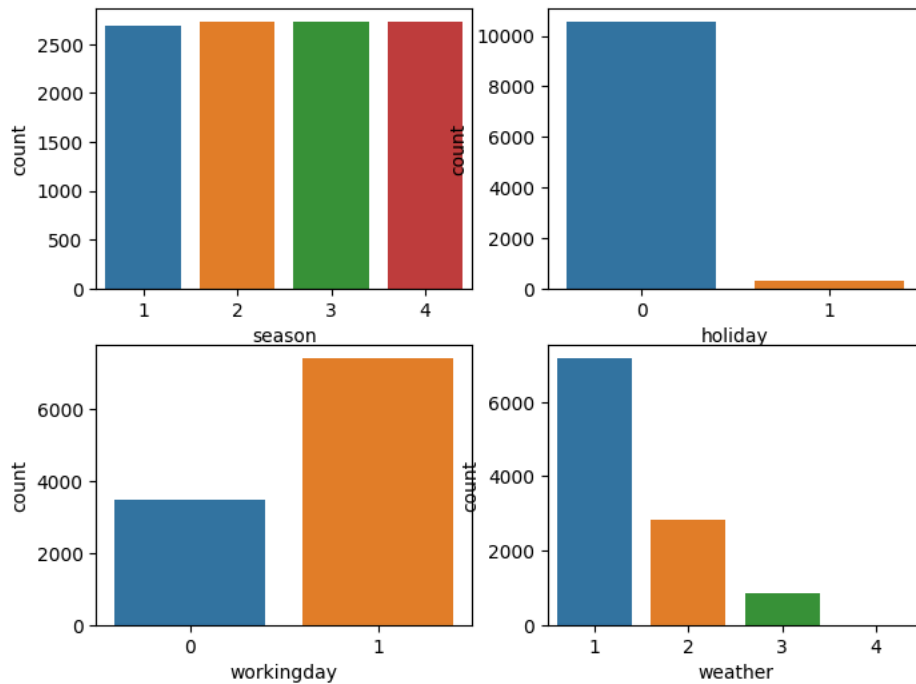
```
 1 # countplot of each categorical column
 2 fig, axis = plt.subplots(nrows=2, ncols=2, figsize=(8, 6))
 3
 4 index = 0
 5 for row in range(2):
 6     for col in range(2):
 7         sns.countplot(data=df, x=cat_cols[index], ax=axis[row, col])
 8         index += 1
 9
10 plt.show()
```
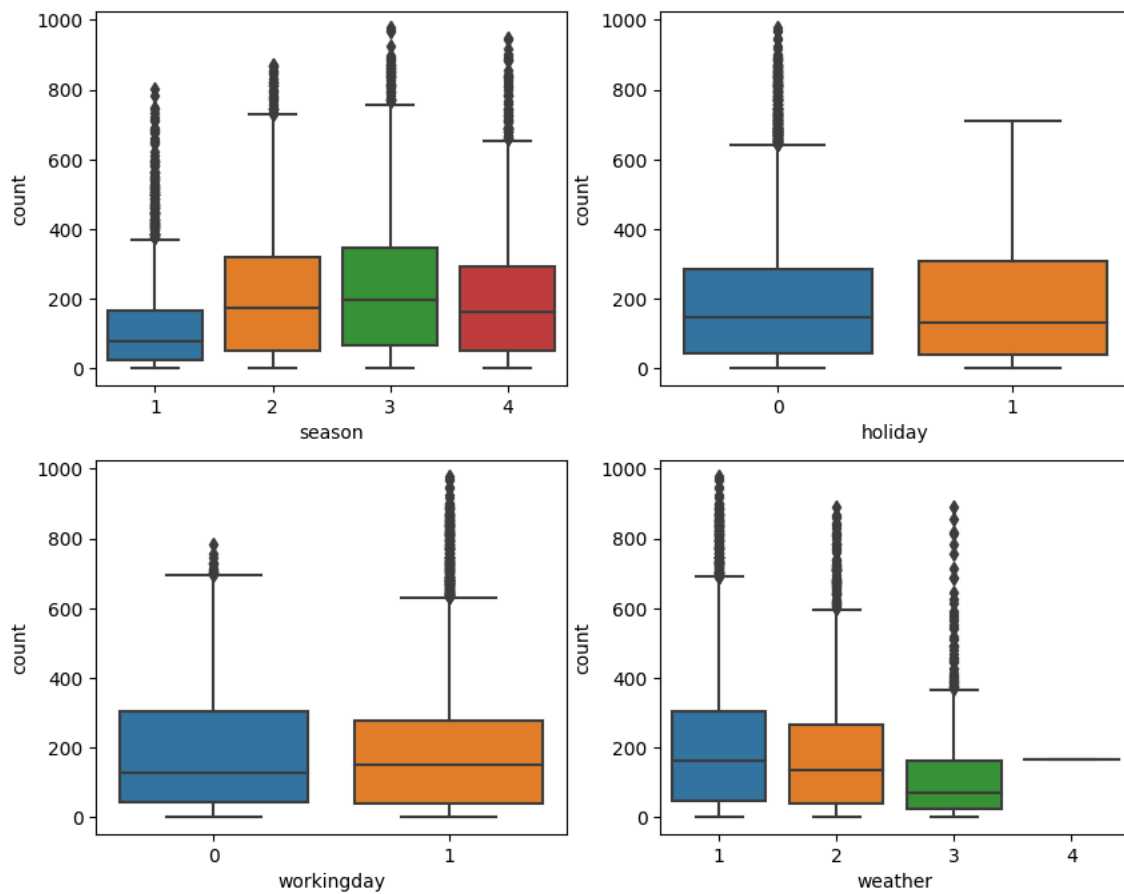


```
 1 # plotting categorical variables againt count using boxplots
 2 fig, axis = plt.subplots(nrows=2, ncols=2, figsize=(10, 8))
 3
 4 index = 0
 5 for row in range(2):
 6     for col in range(2):
 7         sns.boxplot(data=df, x=cat_cols[index], y='count', ax=axis[row, col])
 8         index += 1
 9
10 plt.show()
```
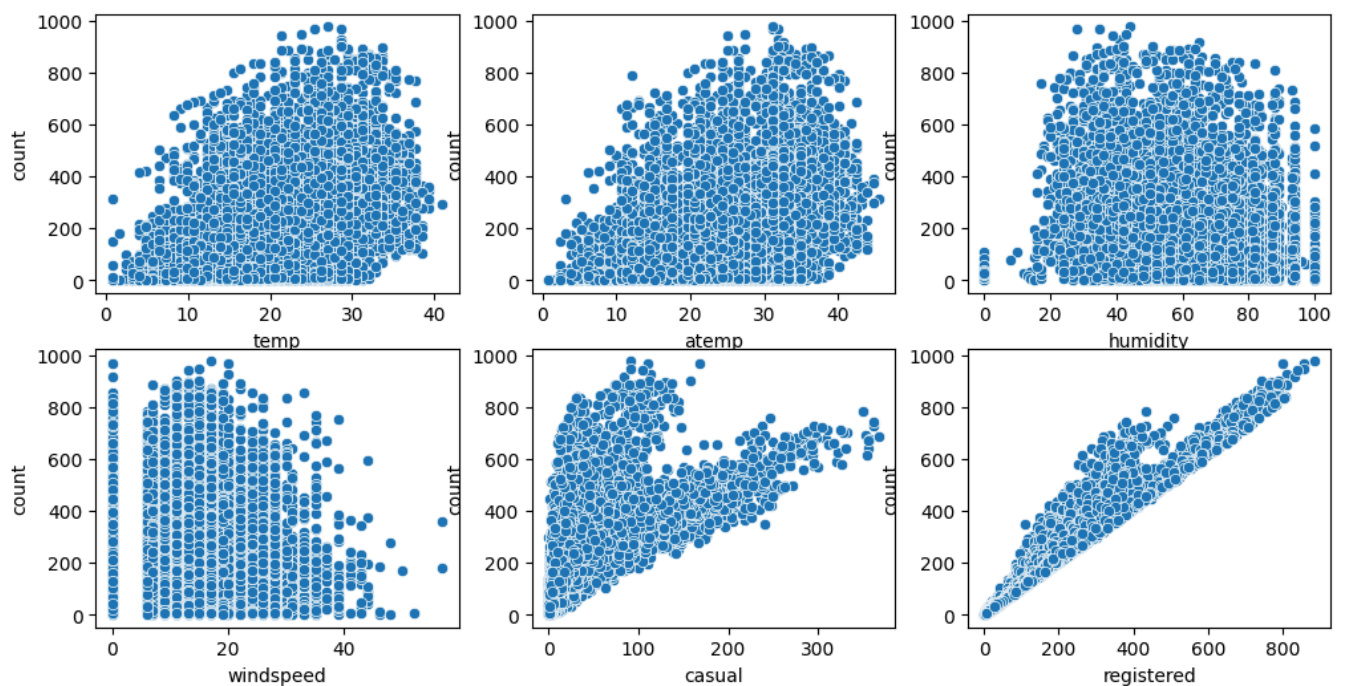
```
1  # plotting numerical variables againt count using scatterplot
2  fig, axis = plt.subplots(nrows=2, ncols=3, figsize=(12, 6))
3
4  index = 0
5  for row in range(2):
6      for col in range(3):
7          sns.scatterplot(data=df, x=num_cols[index], y='count', ax=axis[row, col])
8          index += 1
9
10 plt.show()
```
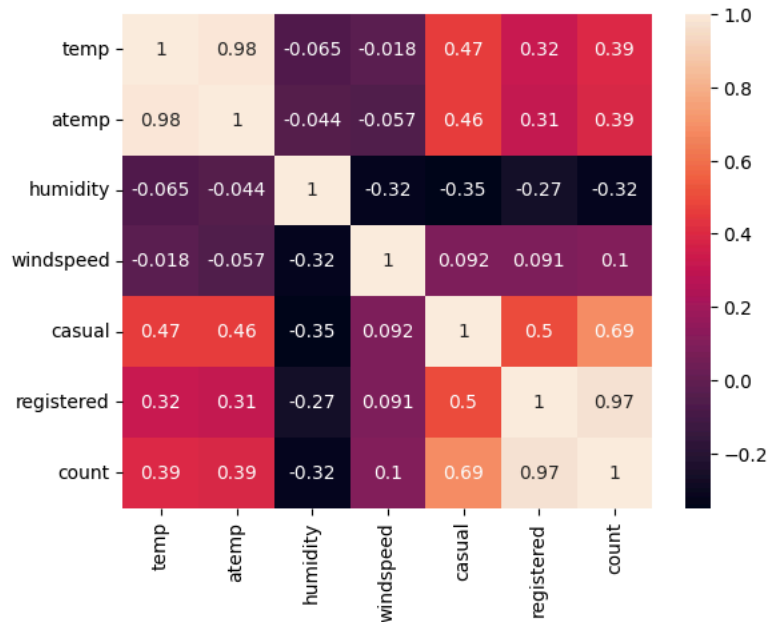
```
1 # understanding the correlation between count and numerical variables
2 df.corr()['count']
3 sns.heatmap(df.corr(), annot=True)
4 plt.show()
```

<ipython-input-22-b0729b22659f>:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In
    df.corr()['count']
<ipython-input-22-b0729b22659f>:3: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In
    sns.heatmap(df.corr(), annot=True)



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

Observed values:

| 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 |

```
1 val = stats.chi2_contingency(data_table)
2 print(val)
3 expected_values = val[3]
4 print(expected_values)
5 nrows, ncols = 4, 4
6 dof = (nrows-1)*(ncols-1)
7 print("degrees of freedom: ", dof)
8 alpha = 0.05
9
10
11 chi_sqr = sum([(o-e)**2/e for o, e in zip(data_table.values, expected_values)])
12 chi_sqr_statistic = chi_sqr[0] + chi_sqr[1]
13 print("chi-square test statistic: ", chi_sqr_statistic)
14
15 critical_val = stats.chi2.ppf(q=1-alpha, df=dof)
16 print(f"critical value: {critical_val}")
17
18 p_val = 1-stats.chi2.cdf(x=chi_sqr_statistic, df=dof)
19 print(f"p-value: {p_val}")
20
21 if p_val <= alpha:
22     print("\nSince p-value is less than the alpha 0.05, We reject the Null Hypothesis. Meaning that\
23     Weather is dependent on the season.")
24 else:
25     print("Since p-value is greater than the alpha 0.05, We do not reject the Null Hypothesis")
```

```
Chi2ContingencyResult(statistic=49.158655596893624, pvalue=1.549925073686492e-07, dof=9, expected_freq=array([[1.7745463
        [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]]]))
[[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]]
degrees of freedom:  9
chi-square test statistic:  44.09441248632364
critical value: 16.918977604620448
p-value: 1.3560001579371317e-06

Since p-value is less than the alpha 0.05, We reject the Null Hypothesis. Meaning that    Weather is dependent on the se
```

```python
1 data_group1 = df[df['workingday']==0]['count'].values
2 data_group2 = df[df['workingday']==1]['count'].values
3
4 print(np.var(data_group1), np.var(data_group2))
5 np.var(data_group2)// np.var(data_group1)
```

```
30171.346098942427 34040.69710674686
1.0
```

```python
1 stats.ttest_ind(a=data_group1, b=data_group2, equal_var=True)
```

```
Ttest_indResult(statistic=-1.2096277376026694, pvalue=0.22644804226361348)
```

```python
 1 # defining the data groups for the ANOVA
 2 from statsmodels.graphics.gofplots import qqplot
 3 gp1 = df[df['weather']==1]['count'].values
 4 gp2 = df[df['weather']==2]['count'].values
 5 gp3 = df[df['weather']==3]['count'].values
 6 gp4 = df[df['weather']==4]['count'].values
 7
 8 gp5 = df[df['season']==1]['count'].values
 9 gp6 = df[df['season']==2]['count'].values
10 gp7 = df[df['season']==3]['count'].values
11 gp8 = df[df['season']==4]['count'].values
12 groups=[gp1,gp2,gp3,gp4,gp5,gp6,gp7,gp8]
13
14
```
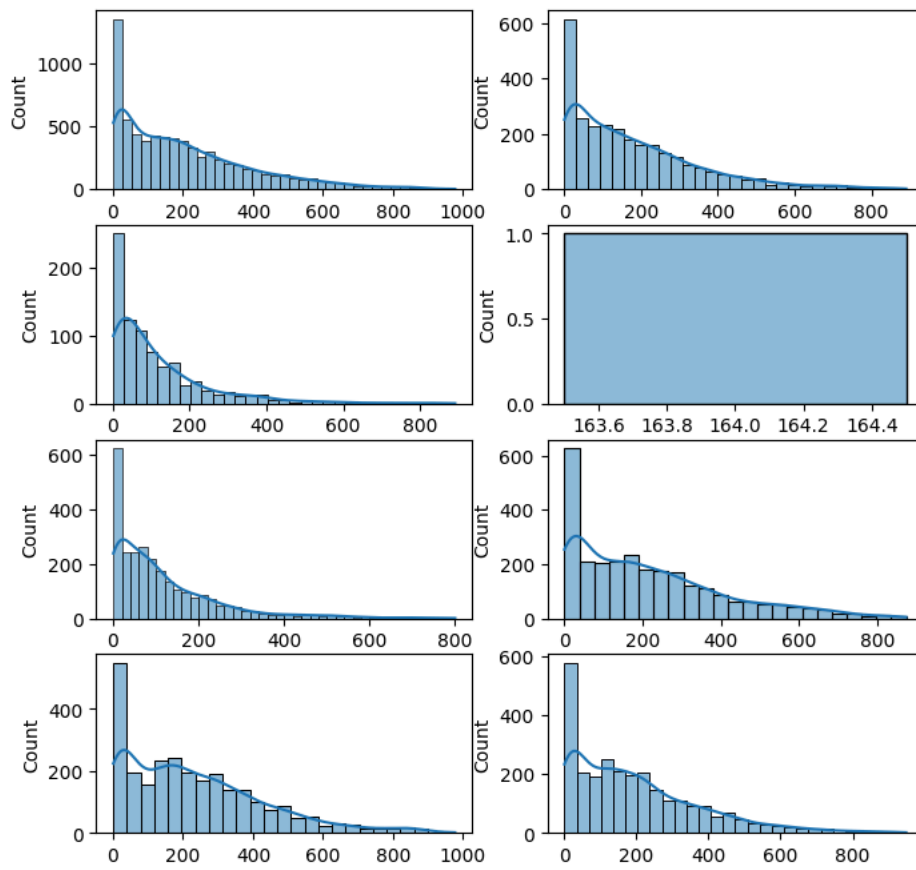
```python
1 fig, axis = plt.subplots(nrows=4, ncols=2, figsize=(8, 8))
2
3 index = 0
4 for row in range(4):
5     for col in range(2):
6         sns.histplot(groups[index], ax=axis[row, col], kde=True)
7         index += 1
8
9 plt.show()
```

```
1
2 index = 0
3 for row in range(4):
4     for col in range(2):
5         qqplot(groups[index], line="s")
6         index += 1
7
8 plt.show()
```