

Loading the Data

```
import pandas as pd
import numpy as np
import seaborn as sns
```

```
ufo_df = pd.read_csv('ufo_data.csv')
ufo_df.head()
```

	Number	City	Colors Reported	Shape Reported	State	Time
0	1234560	Ithaca	NaN	TRIANGLE	NY	6/1/1930 22:00
1	543210	Willingboro	green	NaN	NJ	6/30/1930 20:00
2	876540	Holyoke	NaN	OVAL	CO	2/15/1931 14:00
3	34560	Abilene	blue	DISK	KS	1/12/1900 0:00
4	76540	New York Worlds Fair	NaN	LIGHT	NY	4/18/1933 19:00

```
#Printing the size of the data
ufo_df.shape
```

```
(15, 6)
```

```
ufo_df.dtypes
```

```
Number      float64
City         object
Colors Reported  object
Shape Reported  object
State         object
Time         object
dtype: object
```

Checking for null values Column wise

```
ufo_df.isnull().any()
```

```
Number      True
City         False
Colors Reported  True
Shape Reported  True
State        True
Time         False
dtype: bool
```

```
ufo_df.isnull()
```

	Number	City	Colors Reported	Shape Reported	State	Time
0	False	False	True	False	False	False
1	False	False	False	True	False	False
2	False	False	True	False	False	False
3	False	False	False	False	False	False
4	False	False	True	False	False	False
5	True	False	True	True	False	False
6	False	False	False	False	False	False
7	False	False	True	False	False	False
8	True	False	False	False	False	False
9	False	False	False	False	True	False
10	False	False	True	False	False	False
11	True	False	True	False	False	False
12	False	False	False	False	False	False
13	False	False	True	False	False	False
14	False	False	True	False	False	False

Printing the number of null Values in each column.

```
print('Total Missing values in the Data :- ',ufo_df.isnull().sum())

Total Missing values in the Data :- Number      3
                                City         0
                                Colors Reported  9
                                Shape Reported  2
                                State         1
                                Time         0
                                dtype: int64
```

Printing the total number of null Values in Dataset

```
print('Total Missing values in the Data :- ',ufo_df.isnull().sum().sum())

Total Missing values in the Data :- 15
```

Missing Value Treatment

1)Dropping Null/Missing Values

```
#Shape of the Data before Dropping missing values
ufo_df.shape
```

```
(15, 6)
```

```
# Dropping the missing or null values
data1 = ufo_df.dropna()
```

```
data1
```

	Number	City	Colors Reported	Shape Reported	State	Time
3	34560	Abilene	blue	DISK	KS	1/12/1900 0:00
6	876540	Crater Lake	yellow	CIRCLE	CA	6/15/1935 0:00
12	65430	Belton	red	SPHERE	Y	6/30/1939 20:00

```
data1.shape
```

```
(3, 6)
```

```
print('Loss of the Data : ',round((ufo_df.shape[0]-data1.shape[0])*100/ufo_df.shape[0],2),'%')
```

```
Loss of the Data : 80.0 %
```

We are losing nearly 80 % data by dropping the missing value rows which would be a loss for the prediction

Not a good technique for the Data Analytics

2) Filling Missing Values

Filling Missing Values function is used to replace the null values with a test statistic like mean, median or mode of the particular feature

We can also specify a forward-fill or back-fill to propagate the next values backward or previous value forward.

```
data_new = ufo_df
```

```
data_new.isnull().any()
```

```
Number      True
City         False
Colors Reported  True
Shape Reported  True
State        True
Time         False
dtype: bool
```

Using Back-fill or Forward-fill to propagate next or previous values respectively.

```
## Before applying fill function on the data.
ufo_df.head()
```

	Number	City	Colors Reported	Shape Reported	State	Time
0	1234560	Ithaca	NaN	TRIANGLE	NY	6/1/1930 22:00
1	543210	Willingboro	green	NaN	NJ	6/30/1930 20:00
2	876540	Holyoke	NaN	OVAL	CO	2/15/1931 14:00
3	34560	Abilene	blue	DISK	KS	1/12/1900 0:00
4	76540	New York Worlds Fair	NaN	LIGHT	NY	4/18/1933 19:00

```
## After backward-fill
data_new = data_new.fillna(method='bfill')
#for forward-fill
data_new.fillna(method='ffill')
```

	Number	City	Colors Reported	Shape Reported	State	Time
0	1234560	Ithaca	green	TRIANGLE	NY	6/1/1930 22:00
1	543210	Willingboro	green	OVAL	NJ	6/30/1930 20:00
2	876540	Holyoke	blue	OVAL	CO	2/15/1931 14:00
3	34560	Abilene	blue	DISK	KS	1/12/1900 0:00
4	76540	New York Worlds Fair	yellow	LIGHT	NY	4/18/1933 19:00
5	876540	Valley City	yellow	CIRCLE	N	9/15/1934 15:30
6	876540	Crater Lake	yellow	CIRCLE	CA	6/15/1935 0:00
7	67890	Alma	green	DISK	MI	7/15/1936 0:00
8	45670	Eklutna	green	12	AK	10/15/1936 17:00
9	45670	Hubbard	red	CYLINDER	CA	6/15/1937 0:00
10	76540	Fontana	red	LIGHT	CA	8/15/1937 21:00
11	65430	Waterloo	red	FIREBALL	AL	6/1/1939 20:00
12	65430	Belton	red	SPHERE	Y	6/30/1939 20:00
13	8760	Keokuk	NaN	@	IA	7/7/1939 2:00
14	54320	Ludington	NaN	DISK	MI	6/1/1941 13:00

```
# After forward-fill
data_new = data_new.fillna(method='ffill')
data_new
```

	Number	City	Colors Reported	Shape Reported	State	Time
0	1234560	Ithaca	green	TRIANGLE	NY	6/1/1930 22:00
1	543210	Willingboro	green	OVAL	NJ	6/30/1930 20:00
2	876540	Holyoke	blue	OVAL	CO	2/15/1931 14:00
3	34560	Abilene	blue	DISK	KS	1/12/1900 0:00
4	76540	New York Worlds Fair	yellow	LIGHT	NY	4/18/1933 19:00
5	876540	Valley City	yellow	CIRCLE	N	9/15/1934 15:30
6	876540	Crater Lake	yellow	CIRCLE	CA	6/15/1935 0:00
7	67890	Alma	green	DISK	MI	7/15/1936 0:00
8	45670	Eklutna	green	12	AK	10/15/1936 17:00
9	45670	Hubbard	red	CYLINDER	CA	6/15/1937 0:00
10	76540	Fontana	red	LIGHT	CA	8/15/1937 21:00
11	65430	Waterloo	red	FIREBALL	AL	6/1/1939 20:00
12	65430	Belton	red	SPHERE	Y	6/30/1939 20:00
13	8760	Keokuk	red	@	IA	7/7/1939 2:00
14	54320	Ludington	red	DISK	MI	6/1/1941 13:00

What Happens if for a missing value the Next value or previous value is also a Nan

If a previous or next value isn't available or rather if it is also a NaN value, then, the NaN remains even after back-filling or forward-filling.

```
data_new.isnull().any()
```

```
Number      False
City         False
Colors Reported  False
Shape Reported  False
State        False
Time         False
dtype: bool
```

Here all columns are still having Missing values

Filling the missing value with a statistic

```
data_new = ufo_df
```

```
# 'Number' is a column name for our data
mean_value = data_new['Number'].mean()
print('Mean Value: ',mean_value)
#this will replace all NaN values with the mean of the non null values
data_new['Number']=data_new['Number'].fillna(mean_value)
data_new.head()
```

```
Mean Value: 33004.666666666664
```

	Number	City	Colors Reported	Shape Reported	State	Time
0	1234560	Ithaca	NaN	TRIANGLE	NY	6/1/1930 22:00
1	543210	Willingboro	green	NaN	NJ	6/30/1930 20:00
2	876540	Holyoke	NaN	OVAL	CO	2/15/1931 14:00
3	34560	Abilene	blue	DISK	KS	1/12/1900 0:00
4	76540	New York Worlds Fair	NaN	LIGHT	NY	4/18/1933 19:00

```
data_new
```

	Number	City	Colors Reported	Shape Reported	State	Time
0	123456000000	Ithaca	NaN	TRIANGLE	NY	6/1/1930 22:00
1	543210000000	Willingboro	green	NaN	NJ	6/30/1930 20:00
2	876540000000	Holyoke	NaN	OVAL	CO	2/15/1931 14:00
3	345600000000	Abilene	blue	DISK	KS	1/12/1900 0:00
4	765400000000	New York Worlds Fair	NaN	LIGHT	NY	4/18/1933 19:00
5	33004.666667	Valley City	green	NaN	N	9/15/1934 15:30
6	876540000000	Crater Lake	yellow	CIRCLE	CA	6/15/1935 0:00
7	678900000000	Alma	NaN	DISK	MI	7/15/1936 0:00
8	33004.666667	Eklutna	green	12	AK	10/15/1936 17:00
9	456700000000	Hubbard	red	CYLINDER	NaN	6/15/1937 0:00
10	765400000000	Fontana	NaN	LIGHT	CA	8/15/1937 21:00
11	33004.666667	Waterloo	NaN	FIREBALL	AL	6/1/1939 20:00
12	654300000000	Belton	red	SPHERE	Y	6/30/1939 20:00
13	876000000000	Keokuk	NaN	@	IA	7/7/1939 2:00
14	543200000000	Ludington	NaN	DISK	MI	6/1/1941 13:00

```
data_new.isnull().any()
```

```
Number      False
City         False
Colors Reported  False
Shape Reported  True
State         False
Time         False
dtype: bool
```

Observations

By mean value imputation for numerical data type and mode value imputation for categorical data ,we are able to completely fill the null values, but it is an effective method.

Let's find out with an example

Let x be list of age of people who are attending a Gym in an Office

Data without Outlier's

```
x = [20,25,26,28,29,30,35,38,40,45,30]
```

```
mean(x) = 31
```

Data with missing values and a single Outlier

```
x = [20, ,26,28,29,30,35,38,40,45,30]
```

```
mean(x) = 38
```

We can observe how much the mean value is deviated due to single outlier. So mean is ineffective in presence of outlier's

So, if outlier's are present in the data then median is the best out of box method

Which method would work best suitable for Imputing the missing values in Class Labels

Mode method of imputation would work but by that we may introduce bias into the model to avoid this, we generally move those datapoints into test set and try to predict with model

Outlier Treatment

Outlier Detection

```
import pandas as pd
import numpy as np
import seaborn as sns

data = sns.load_dataset("tips")
data.head()
```

	total_bill	tip	sex	smoker	day	time	size
0	16.99	1.01	Female	No	Sun	Dinner	2
1	10.34	1.66	Male	No	Sun	Dinner	3
2	21.01	3.50	Male	No	Sun	Dinner	3
3	23.68	3.31	Male	No	Sun	Dinner	2
4	24.59	3.61	Female	No	Sun	Dinner	4

```
data.dtypes
```

```
total_bill    float64
tip           float64
sex           object
smoker        object
day           object
time         object
size          int64
dtype: object
```

```
#Separating the continues variables
import matplotlib.pyplot as plt
import seaborn as sns

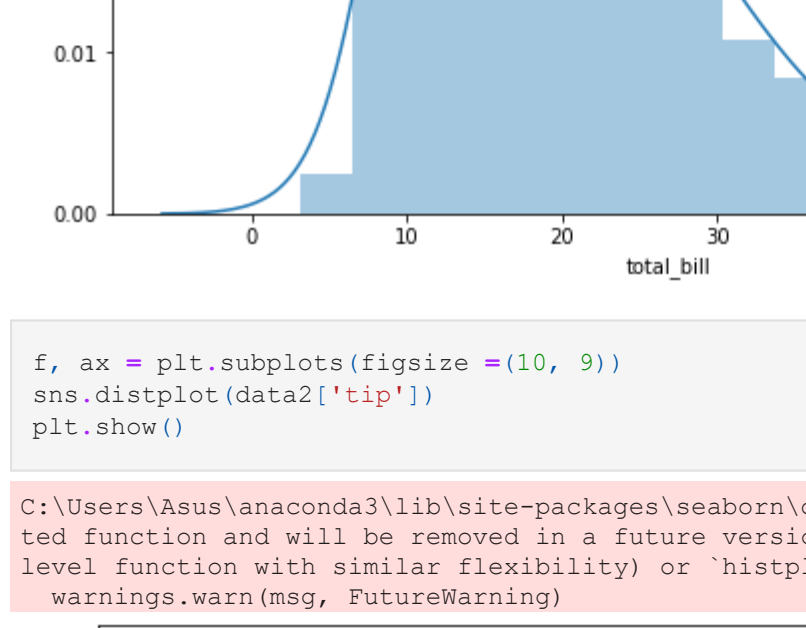
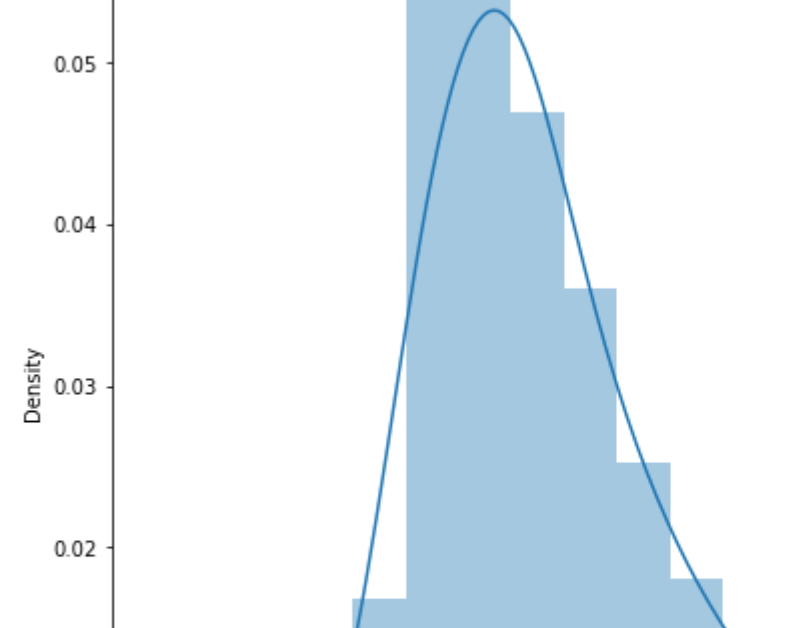
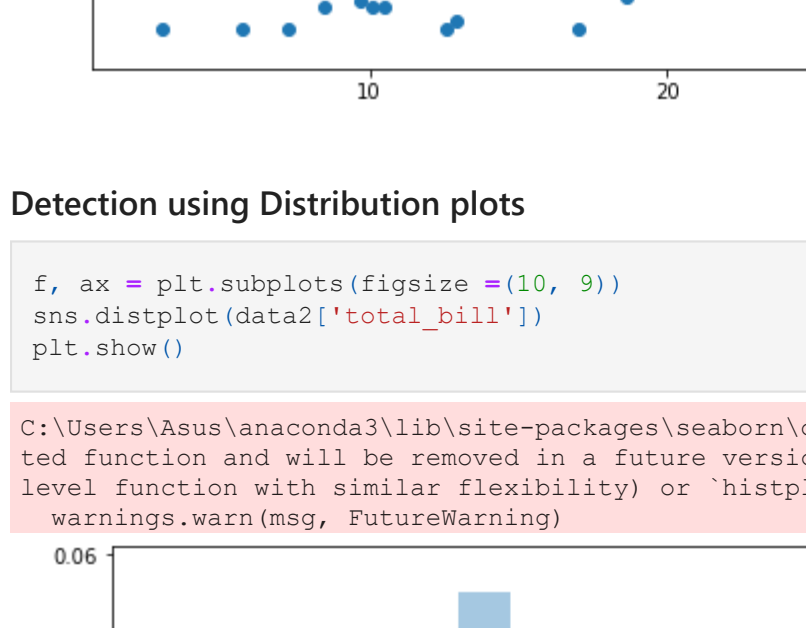
data2 = data.select_dtypes(exclude = 'object')
data2.head()
```

```
total_bill    size
```

	total_bill	size
0	16.99	101
1	10.34	1.66
2	21.01	3.50
3	23.68	3.1
4	24.59	3.61

Detection Using Box Plot

```
for i in data2.columns:
    print('\n')
    print('Feature :- ',i)
    sns.boxplot(x = 'size',y = i,data = data2)
    plt.show()
```

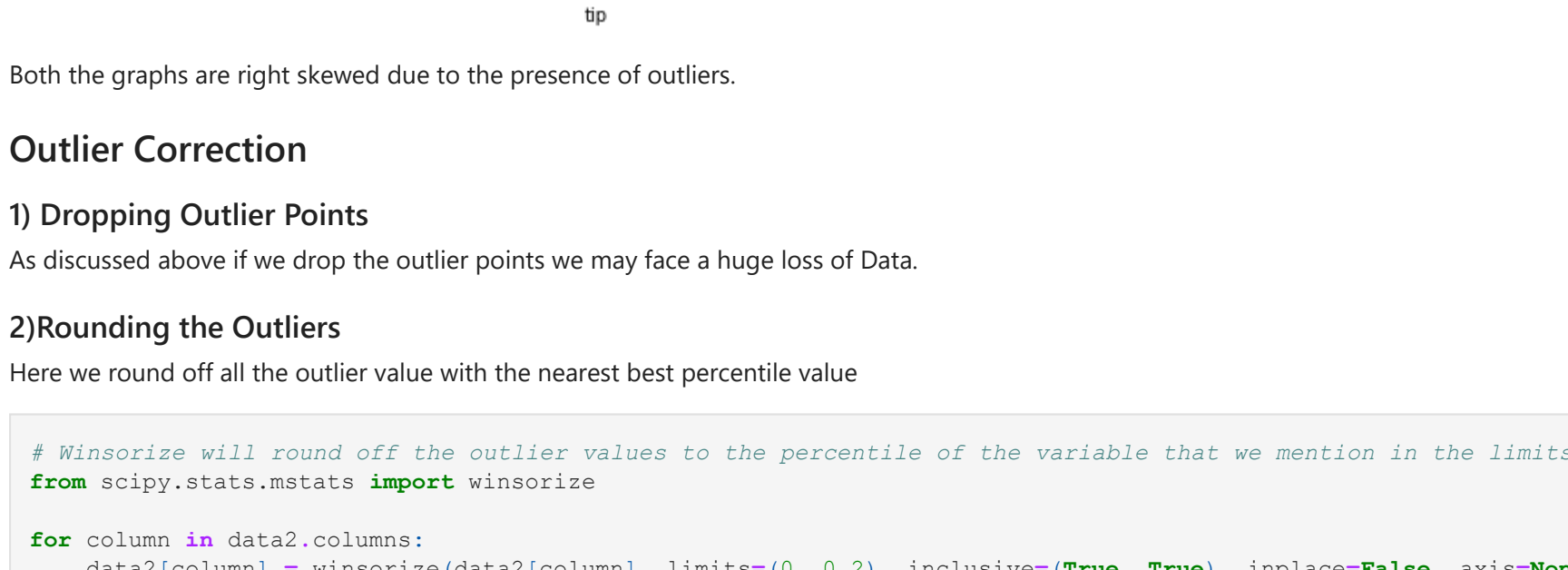


The dots in the box plots of the total bill and tips represent the distribution of the outliers class wise in the data.

The presence of the outlier can severely affect the model performance and may lead to improper conclusions. For example you want to calculate the total turn over of your restaurant for the next month financial planning. If you calculate based on outlier data you entire planning may collapse.

Detection using Scatter Plot

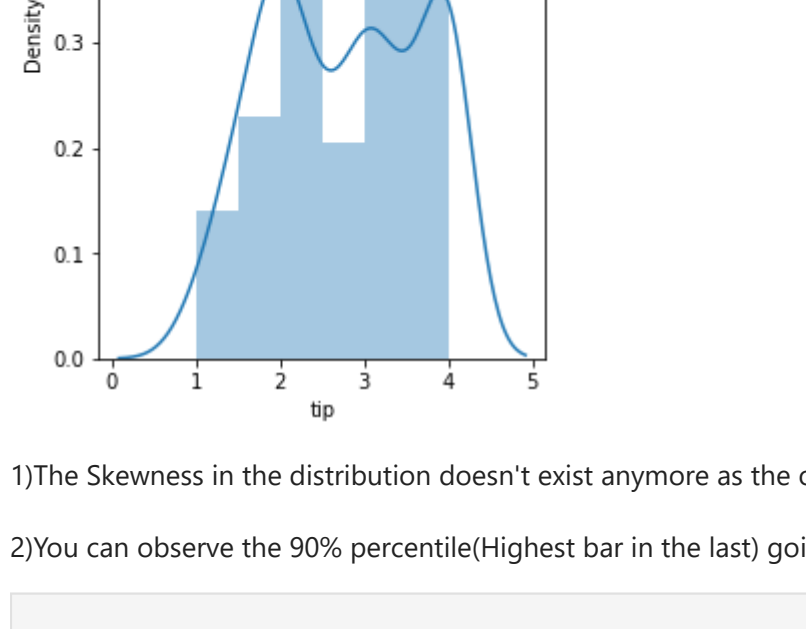
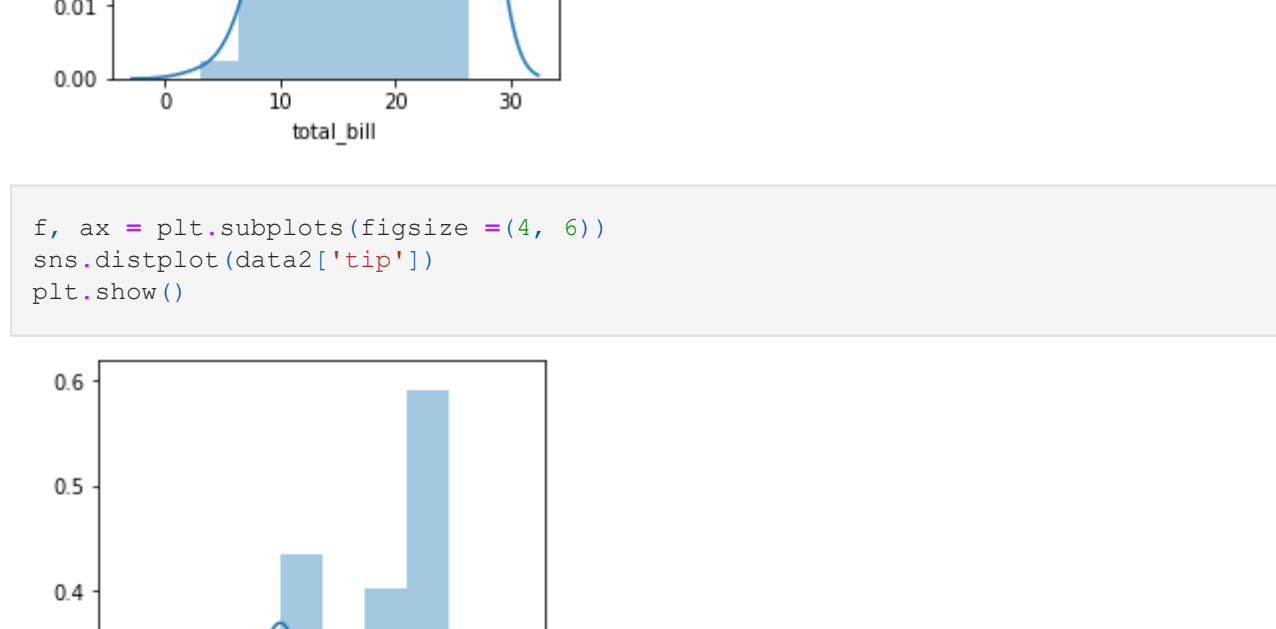
```
fig, ax = plt.subplots(figsize=(14,8))
ax.scatter(data2['total_bill'], data2['tip'])
ax.set_xlabel('Total_bill')
ax.set_ylabel('tip received')
plt.show()
```



Detection using Distribution plots

```
f, ax = plt.subplots(figsize=(10, 9))
sns.distplot(data2['total_bill'])
plt.show()
```

C:\Users\Anus\anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)



1)The Skewness in the distribution doesn't exist anymore as the outliers has been rounded off.

2)You can observe the 90% percentile(Highest bar in the last) going high than before as all outliers are rounded to it.

```
In [ ]:
```