

▼ Data Visualization I

1. Use the inbuilt dataset 'titanic'. The dataset contains 891 rows and contains information about the passengers who boarded the unfortunate Titanic ship. Use the Seaborn library to see if we can find any patterns in the data.

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

df=sns.load_dataset('titanic')

df.head(3)

df.shape

df.info()

df.describe()

df.isnull().sum()

df.drop(columns=['embark_town', 'deck'],axis=1,inplace=True)

df.isnull().sum()

df.impute values in age column

df.check distribution of age column

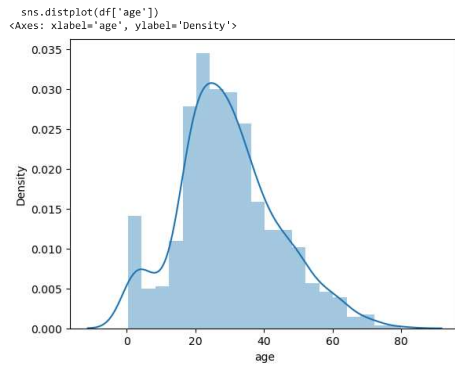
sns.distplot(df['age'])
```

C:\Users\admin\AppData\Local\Temp\ipykernel\_2984\3234920688.py:1: UserWarning:

'distplot' is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de4447ed2974457ad6372750bbe5751>



```
df['age'].skew()
```

```
0.38910778230082704
```

#age column has normal data distribution

```
df['age'].mean()
```

```
29.69911764705882
```

```
df['age'].fillna(df['age'].mean(),inplace=True)
```

```
df.isnull().sum()
```

```
survived    0
pclass      0
sex          0
age          0
sibsp       0
parach      0
fare        0
embarked     2
class       0
who         0
adult_male  0
alive       0
alone       0
dtype: int64
```

✓ impute values of embarked column

```
df['embarked'].nunique()
```

```
3
```

```
df['embarked'].unique()
```

```
array(['S', 'C', 'Q', nan], dtype=object)
```

```
df['embarked'].mode()
```

```
0    S
Name: embarked, dtype: object
```

```
df['embarked'].mode()[0]
```

```
'S'
```

```
df['embarked'].fillna(df['embarked'].mode()[0],inplace=True)
```

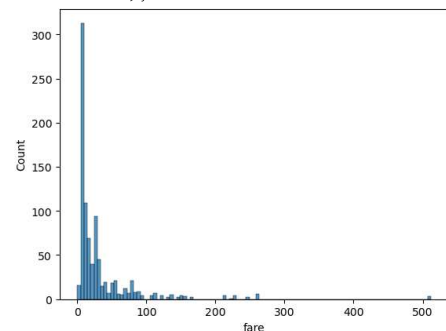
```
df.isnull().sum()
```

```
survived    0
pclass      0
sex          0
age          0
sibsp       0
parach      0
fare        0
embarked     0
class       0
who         0
adult_male  0
alive       0
alone       0
dtype: int64
```

✓ 2. Write a code to check how the price of the ticket (column name: 'fare') for each passenger is distributed by plotting a histogram.

```
sns.histplot(df['fare'])
```

<Axes: xlabel='fare', ylabel='Count'>



```
df['fare'].skew()
4.787316519674893
```

```
# highly right distributed
```

```
df.nunique()
survived      2
pclass       3
sex           2
age          89
sibsp        7
parch        7
fare        248
embarked      3
class         3
who           3
adult_male    2
alive         2
alone         2
dtype: int64
```

out of 13 columns age and fare are quantitative columns remaining are qualitative

✓ Quantitative= distplot, histplot, kdeplot, boxplot

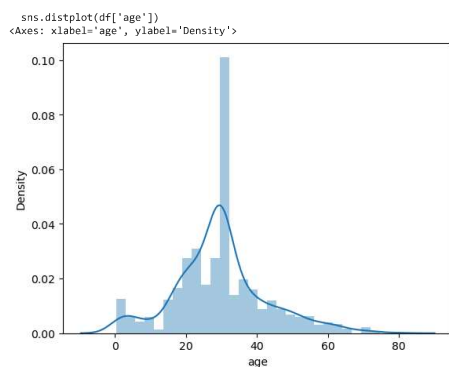
Qualitative= countplot, pie chart

```
sns.distplot(df['age'])
```

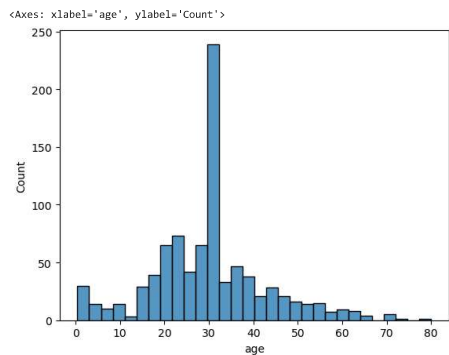
```
C:\Users\admin\AppData\Local\Temp\ipykernel_2984\3234928688.py:1: UserWarning:
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with
similar flexibility) or `histplot` (an axes-level function for histograms).

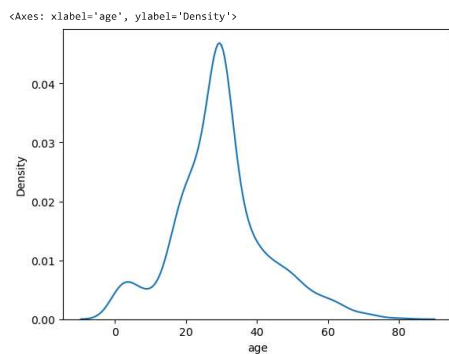
For a guide to updating your code to use the new functions, please see
https://rishi.github.io/mwaskom/de4d47ed2874457ade6372758bbe5751
```



```
sns.histplot(df['age'])
```

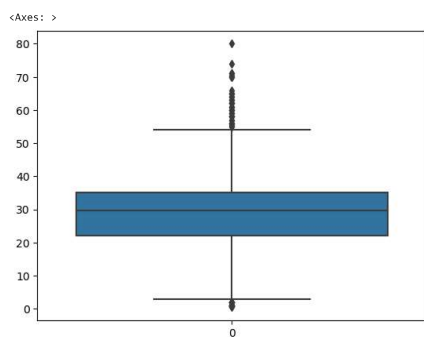


```
sns.kdeplot(df['age'])
```



✓ age column is normally distributed

```
sns.boxplot(df['age'])
```



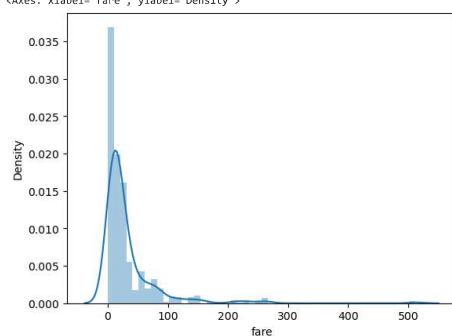
✓ there are outliers in age column

```
sns.distplot(df['fare'])
```

C:\Users\admin\AppData\Local\Temp\ipykernel\_2984\1195996103.py:1: UserWarning:  
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.  
Please adapt your code to use either `displot` (a figure-level function with  
similar flexibility) or `histplot` (an axes-level function for histograms).  
For a guide to updating your code to use the new functions, please see  
<https://gist.github.com/mwaskom/de44147ed2974457ad6372758bbe5751>

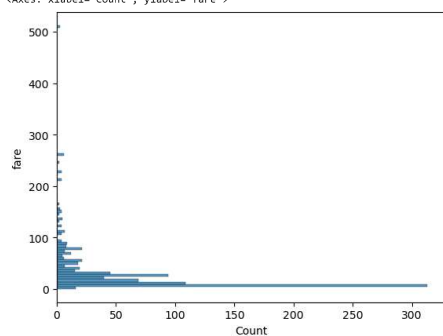
```
sns.distplot(df['fare'])
```

<Axes: xlabel='fare', ylabel='Density'>



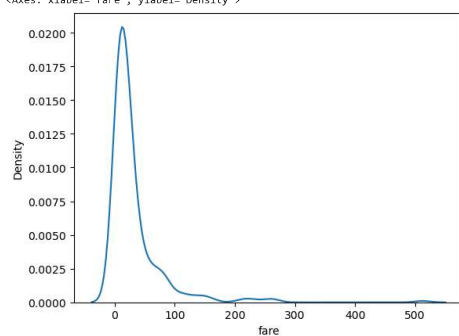
```
sns.histplot(y=df['fare'])
```

<Axes: xlabel='Count', ylabel='fare'>



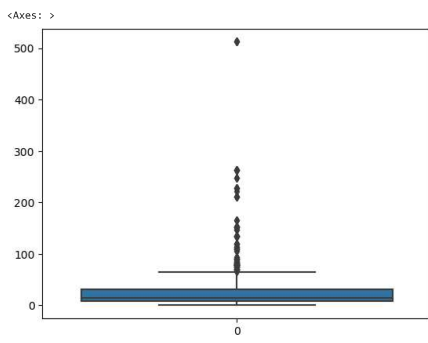
```
sns.kdeplot(df['fare'])
```

<Axes: xlabel='fare', ylabel='Density'>



✓ fare column is rightly distributed

```
sns.boxplot(df['fare'])
```



▼ fare has also outliers

```
sns.countplot(x=df['survived'])
```

<Axes: xlabel='survived', ylabel='count'>

