

STATISTICS & PROBABILITY

Topic: Descriptive Statistics and Inferential Statistics

source code <https://github.com/nitishbuzzpro/Statistics-and-Hypothesis-Tetsing--Titanic-Dataset---Data-Science.git> (<https://github.com/nitishbuzzpro/Statistics-and-Hypothesis-Tetsing--Titanic-Dataset---Data-Science.git>)

```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5 from scipy import stats
```

Task 1

```
In [2]: 1 # Tasks:
        2 # Load the Titanic dataset using Python and perform
        3 # the following descriptive statistics:
        4 # Calculate the mean, median, and mode of the age
        5 # of passengers.
        6 # Calculate the range, variance, and standard
        7 # deviation of the fare paid by passengers.
        8 # Calculate the correlation coefficient between the
        9 # age and fare columns.
       10 # Compute the quartiles (25th, 50th, and 75th
       11 # percentiles) of the age and fare columns.
       12 # Generate a histogram and box plot for the age distribution of passeng
       13 # Explain Box plot, and How Box plot could be useful in a Real-world Sc
```

```
In [3]: 1 # Load the Titanic dataset using Python and perform the following descr
2 df = pd.read_excel('../input/train-dataset/train.xlsx')
3 df.head()
```

```
Out[3]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500

```
In [4]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
#   Column          Non-Null Count  Dtype
---  -
0   PassengerId     891 non-null   int64
1   Survived        891 non-null   int64
2   Pclass          891 non-null   int64
3   Name            891 non-null   object
4   Sex             891 non-null   object
5   Age            714 non-null   float64
6   SibSp           891 non-null   int64
7   Parch          891 non-null   int64
8   Ticket          891 non-null   object
9   Fare           891 non-null   float64
10  Cabin          204 non-null   object
11  Embarked       889 non-null   object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

```
In [5]: 1 # Calculate the mean, median, and mode of the age of passengers.
2 print('mean:', df['Age'].mean())
3 print('median:', df['Age'].median())
4 print('mode', df['Age'].mode()[0])
```

```
mean: 29.69911764705882
median: 28.0
mode 24.0
```

```
In [6]: 1 # Calculate the range, variance, and standard deviation of the fare pai
2 print('range: ',(df['Fare'].max(),df['Fare'].min()))
3 print('variance: ', df['Fare'].var())
4 print('Std Deviation', df['Fare'].std())
```

```
range: (512.3292, 0.0)
variance: 2469.436845743116
Std Deviation 49.6934285971809
```

```
In [7]: 1 # Calculate the correlation coefficient between the age and fare column
2 print('correlation coefficient between the age and fare columns',df['Ag
```

```
correlation coefficient between the age and fare columns 0.096066691769038
9
```

```
In [8]: 1 # Compute the quartiles (25th, 50th, and 75th percentiles) of the age a
2 print('quartiles (25th, 50th, and 75th percentiles) of the age: \n',df[
3 print('\n')
4 print('quartiles (25th, 50th, and 75th percentiles) of the age: \n',df[
```

```
quartiles (25th, 50th, and 75th percentiles) of the age:
0.25    20.125
0.50    28.000
0.75    38.000
Name: Age, dtype: float64
```

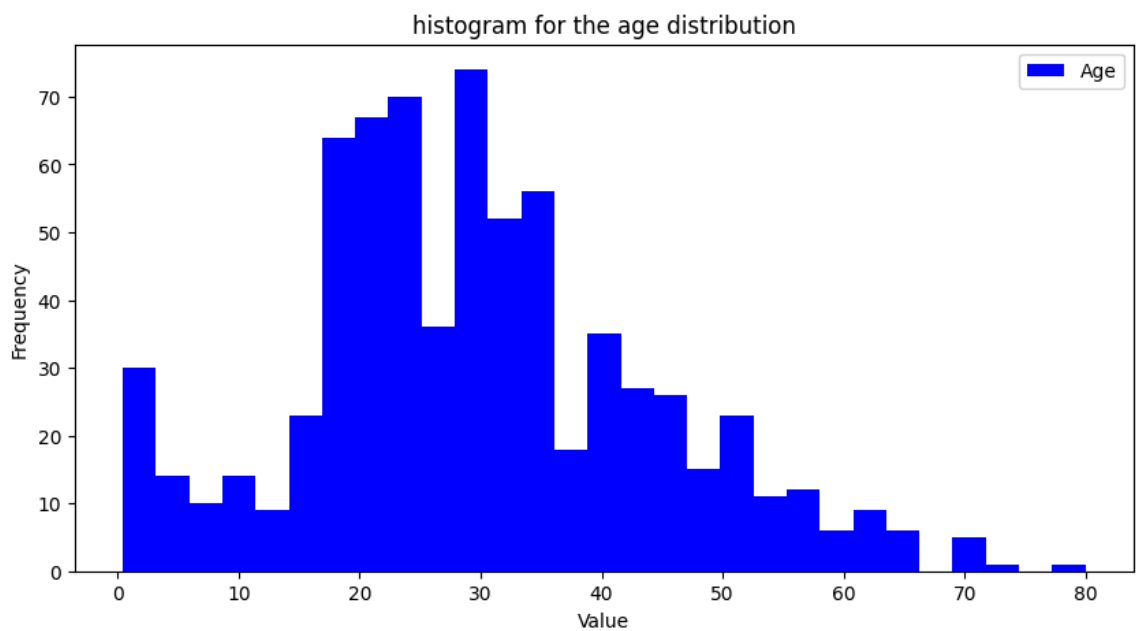
```
quartiles (25th, 50th, and 75th percentiles) of the age:
0.25     7.9104
0.50    14.4542
0.75    31.0000
Name: Fare, dtype: float64
```

```

In [9]: 1 # Generate a histogram and box plot for the age distribution of passeng
2 plt.figure(figsize=(10,5))
3 plt.hist(
4     df['Age'],
5     bins=int(np.sqrt(len(df['Age']))),
6     range=None,
7     density=False,
8     weights=None,
9     cumulative=False,
10    bottom=None,
11    histtype='stepfilled',
12    align='mid',
13    orientation='vertical',
14    rwidth=None,
15    log=False,
16    color='blue',
17    label='Age',
18    stacked=False,
19 )
20 plt.title('histogram for the age distribution')
21 plt.xlabel('Value')
22 plt.ylabel('Frequency')
23 plt.legend()

```

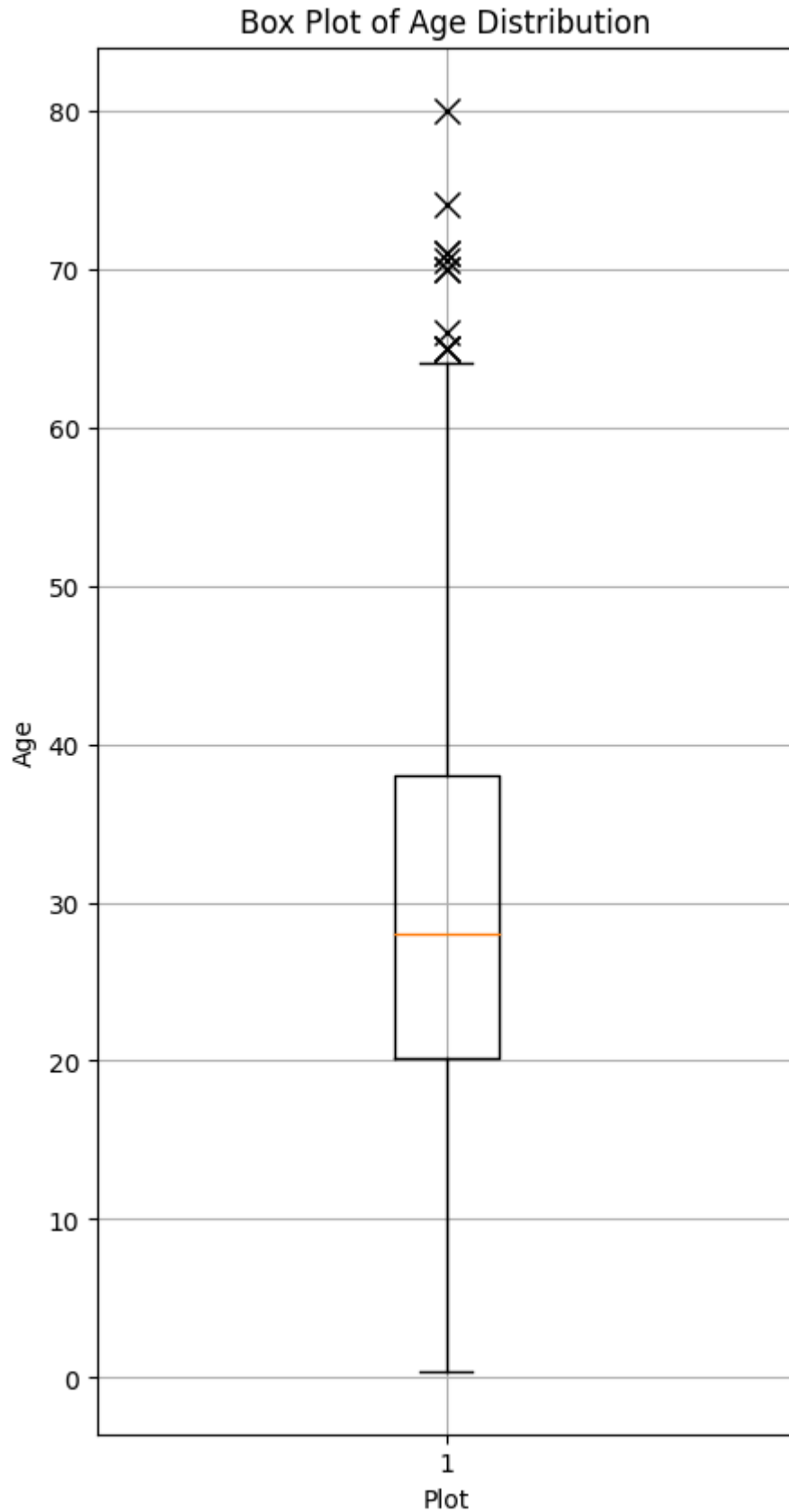
Out[9]: <matplotlib.legend.Legend at 0x7c3e5b93fd00>



```

In [10]: 1 # box plot for the age distribution of passengers
2 plt.figure(figsize=(5, 10))
3 plt.boxplot(df['Age'].dropna(),flierprops=dict(marker='x', color='red'),
4 plt.title('Box Plot of Age Distribution')
5 plt.ylabel('Age')
6 plt.xlabel('Plot')
7 plt.grid(True) # Add grid lines for better readability
8 plt.show()

```



```
In [11]: 1 # Explain Box plot, and How Box plot could be useful in a Real-world Sc
2
3
4 # Answer:- A box plot, also known as a whisker plot, is a graphical rep
5 # of a dataset that displays its distribution through five main summary
6 # statistics: minimum, first quartile (Q1), median (Q2), third quartile
7 # and maximum. It's a powerful tool for visualizing the spread and skew
8 # of the data, identifying outliers, and comparing different datasets.
9
10 # Benefits:
11 # Identifying Variability: By comparing the box plots of different batc
12
13 # Detecting Outliers: Outliers in the box plot can indicate defective c
14
15 # Comparing Batches: Multiple box plots allow for direct comparison of
16
17 # Ensuring Consistency: Regularly creating box plots for production bat
18
19 # Additional Applications:
20 # Healthcare: Monitoring patient vital signs (e.g., blood pressure) to
21 # Finance: Analyzing the distribution of daily stock returns to assess
22 # Education: Comparing student test scores across different classes or
```

```
In [ ]: 1
```

Task 2

```
In [12]: 1 # Analyze the categorical variables using the following descriptive sta
2
3 # Calculate the frequency distribution of the passenger classes (Pclass
4 # Calculate the percentage distribution of survival (Survived column).
5 # Determine the mode and percentage distribution of the passenger title
```

```
In [13]: 1 # Calculate the frequency distribution of the passenger classes (Pclass
2 print('frequency distribution of the passenger classes (Pclass column):
```

frequency distribution of the passenger classes (Pclass column):

Pclass

1 216

2 184

3 491

Name: count, dtype: int64

```
In [14]: 1 # Calculate the percentage distribution of survival (Survived column).
2 print('the percentage distribution of survival (Survived column):\n',df
```

the percentage distribution of survival (Survived column):

Survived

0 61.616162

1 38.383838

Name: proportion, dtype: float64

```
In [15]: 1 # Determine the mode and percentage distribution of the passenger title
2 df['title'] = df['Name'].apply(lambda name : name.split(',')[1].split('
3 print('the mode of the passenger titles (Name column):.\n', df['title']
4 print('\n')
5 print('the percentage distribution of the passenger titles (Name colum
```

```
the mode of the passenger titles (Name column):.
Mr
```

```
the percentage distribution of the passenger titles (Name column):.
title
```

```
Mr          58.024691
Miss        20.426487
Mrs         14.029181
Master      4.489338
Dr          0.785634
Rev         0.673401
Mlle        0.224467
Major       0.224467
Col         0.224467
the Countess 0.112233
Capt       0.112233
Ms          0.112233
Sir         0.112233
Lady        0.112233
Mme         0.112233
Don         0.112233
Jonkheer    0.112233
Name: proportion, dtype: float64
```

```
In [ ]: 1
```

Task 3

```
In [16]: 1 # Analyze the following statistical distributions:
2
3 # Determine the distribution of the passenger ages and visualize it usi
4 # Determine the distribution of the fares paid by passengers and visual
```

```

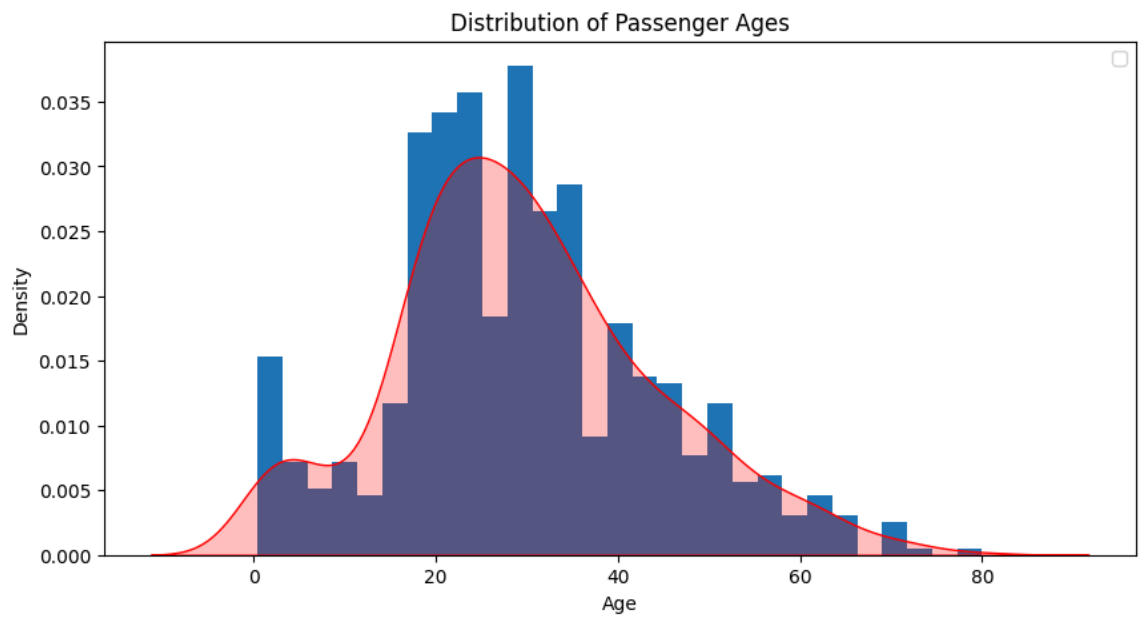
In [17]: 1 # Determine the distribution of the passenger ages and visualize it usi
2 plt.figure(figsize=(10,5))
3 plt.hist(
4     df['Age'],
5     bins=int(np.sqrt(len(df['Age']))),
6     range=None,
7     density=True,
8     weights=None,
9     cumulative=False,
10    bottom=None,
11    histtype='bar',
12    align='mid',
13    orientation='vertical',
14    rwidth=None,
15    log=False,
16    color=None,
17    label=None,
18    stacked=False,
19 )
20 sns.kdeplot(
21     data=df['Age'],
22     hue=None,
23     weights=None,
24     palette=None,
25     hue_order=None,
26     hue_norm=None,
27     color='Red',
28     fill=True,
29     multiple='layer',
30     common_norm=True,
31     common_grid=True,
32     cumulative=False,
33     bw_method='scott',
34     bw_adjust=1,
35     warn_singular=True,
36     log_scale=None,
37     levels=10,
38     thresh=0.05,
39     gridsize=200,
40     cut=3,
41     clip=None,
42     legend=True,
43     cbar=False,
44     cbar_ax=None,
45     cbar_kws=None,
46     ax=None,
47 )
48
49 plt.xlabel('Age')
50 plt.ylabel('Density')
51 plt.title('Distribution of Passenger Ages')
52 plt.legend()

```

/opt/conda/lib/python3.10/site-packages/seaborn/_oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.

with pd.option_context('mode.use_inf_as_na', True):

Out[17]: <matplotlib.legend.Legend at 0x7c3e5b999060>



```

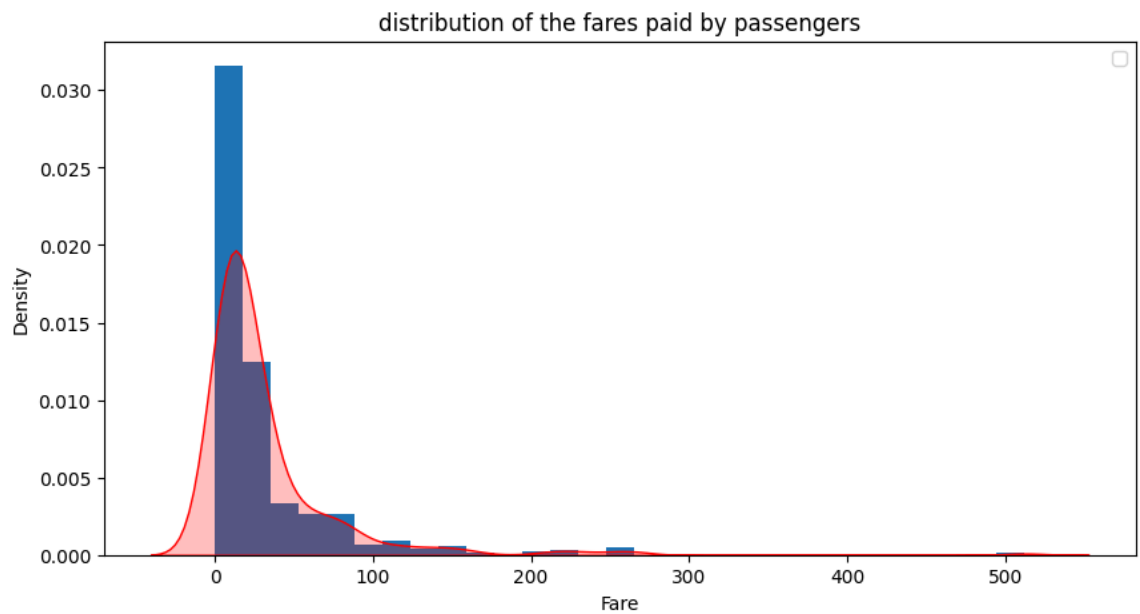
In [18]: 1 # Determine the distribution of the fares paid by passengers and visual
2 plt.figure(figsize=(10,5))
3 plt.hist(
4     df['Fare'],
5     bins=int(np.sqrt(len(df['Fare']))),
6     range=None,
7     density=True,
8     weights=None,
9     cumulative=False,
10    bottom=None,
11    histtype='bar',
12    align='mid',
13    orientation='vertical',
14    rwidth=None,
15    log=False,
16    color=None,
17    label=None,
18    stacked=False,
19 )
20 sns.kdeplot(
21     data=df['Fare'],
22     hue=None,
23     weights=None,
24     palette=None,
25     hue_order=None,
26     hue_norm=None,
27     color='Red',
28     fill=True,
29     multiple='layer',
30     common_norm=True,
31     common_grid=True,
32     cumulative=False,
33     bw_method='silverman',
34     bw_adjust=1,
35     warn_singular=True,
36     log_scale=None,
37     levels=10,
38     thresh=0.05,
39     gridsize=200,
40     cut=3,
41     clip=None,
42     legend=True,
43     cbar=False,
44     cbar_ax=None,
45     cbar_kws=None,
46     ax=None,
47 )
48
49 plt.xlabel('Fare')
50 plt.ylabel('Density')
51 plt.title('distribution of the fares paid by passengers')
52 plt.legend()

```

/opt/conda/lib/python3.10/site-packages/seaborn/_oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.

with pd.option_context('mode.use_inf_as_na', True):

Out[18]: <matplotlib.legend.Legend at 0x7c3e5b7bc5e0>



In []:

1

Task 4

In [19]:

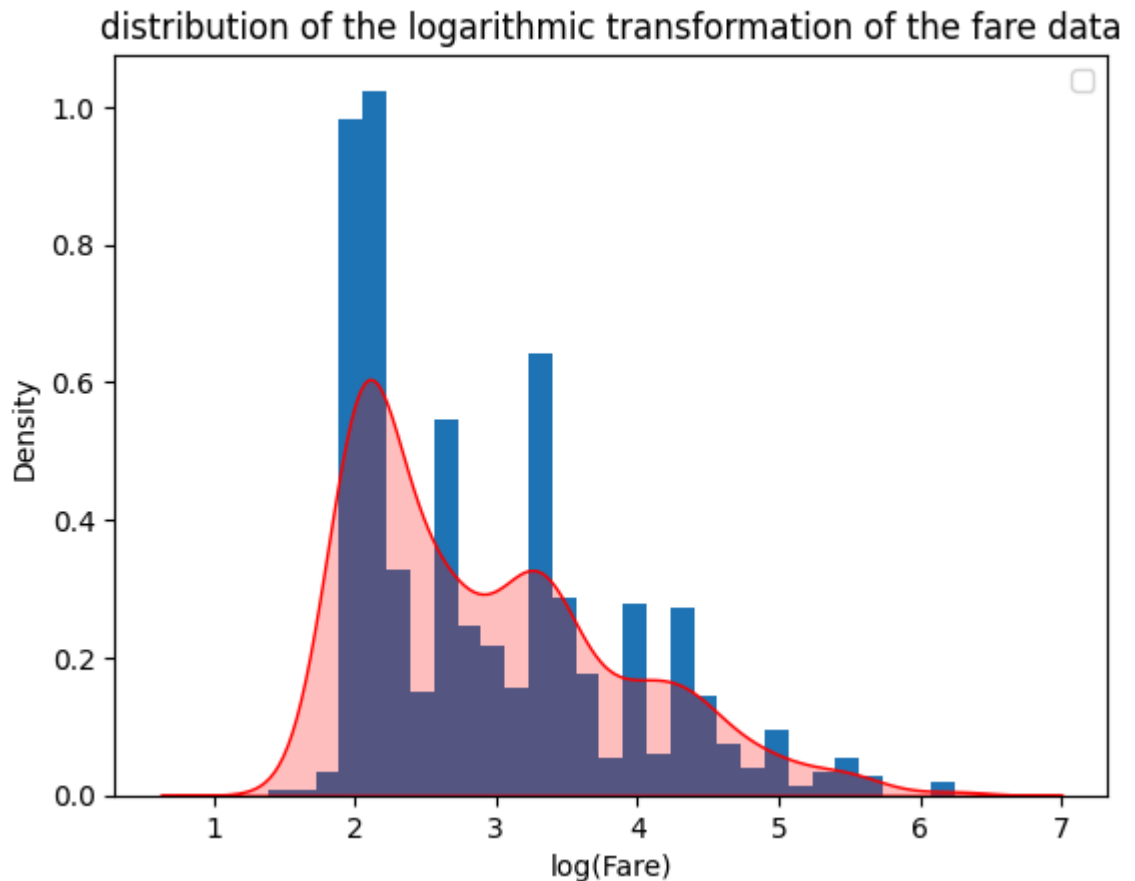
```
1 # Perform statistical transformations on the data:
2
3 # Apply logarithmic transformation to the fare data
4 # and visualize the transformed distribution using a
5 # histogram and a kernel density plot.
6
7 # Apply square root transformation to the age data
8 # and visualize the transformed distribution using a
9 # histogram and a kernel density plot
```


In [20]:

```
1  # Apply logarithmic transformation to the fare data
2  # and visualize the transformed distribution using a
3  # histogram and a kernel density plot.
4
5  log_fare = np.log(df['Fare'])
6  log_fare = log_fare[log_fare>0]
7
8  plt.hist(
9      log_fare,
10     bins=int(np.sqrt(len(log_fare))),
11     range=None,
12     density=True,
13     weights=None,
14     cumulative=False,
15     bottom=None,
16     histtype='bar',
17     align='mid',
18     orientation='vertical',
19     rwidth=None,
20     log=False,
21     color=None,
22     label=None,
23     stacked=False,
24 )
25 sns.kdeplot(
26     data=log_fare,
27     hue=None,
28     weights=None,
29     palette=None,
30     hue_order=None,
31     hue_norm=None,
32     color='Red',
33     fill=True,
34     multiple='layer',
35     common_norm=True,
36     common_grid=True,
37     cumulative=False,
38     bw_method='silverman',
39     bw_adjust=1,
40     warn_singular=True,
41     log_scale=None,
42     levels=10,
43     thresh=0.05,
44     gridsize=200,
45     cut=3,
46     clip=None,
47     legend=True,
48     cbar=False,
49     cbar_ax=None,
50     cbar_kws=None,
51     ax=None,
52 )
53
54 plt.xlabel('log(Fare)')
55 plt.ylabel('Density')
56 plt.title('distribution of the logarithmic transformation of the fare d
57 plt.legend()
```

```
/opt/conda/lib/python3.10/site-packages/pandas/core/arraylike.py:399: RuntimeWarning: divide by zero encountered in log
  result = getattr(ufunc, method)(*inputs, **kwargs)
/opt/conda/lib/python3.10/site-packages/seaborn/_oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
```

Out[20]: <matplotlib.legend.Legend at 0x7c3e5965e230>



In [21]:

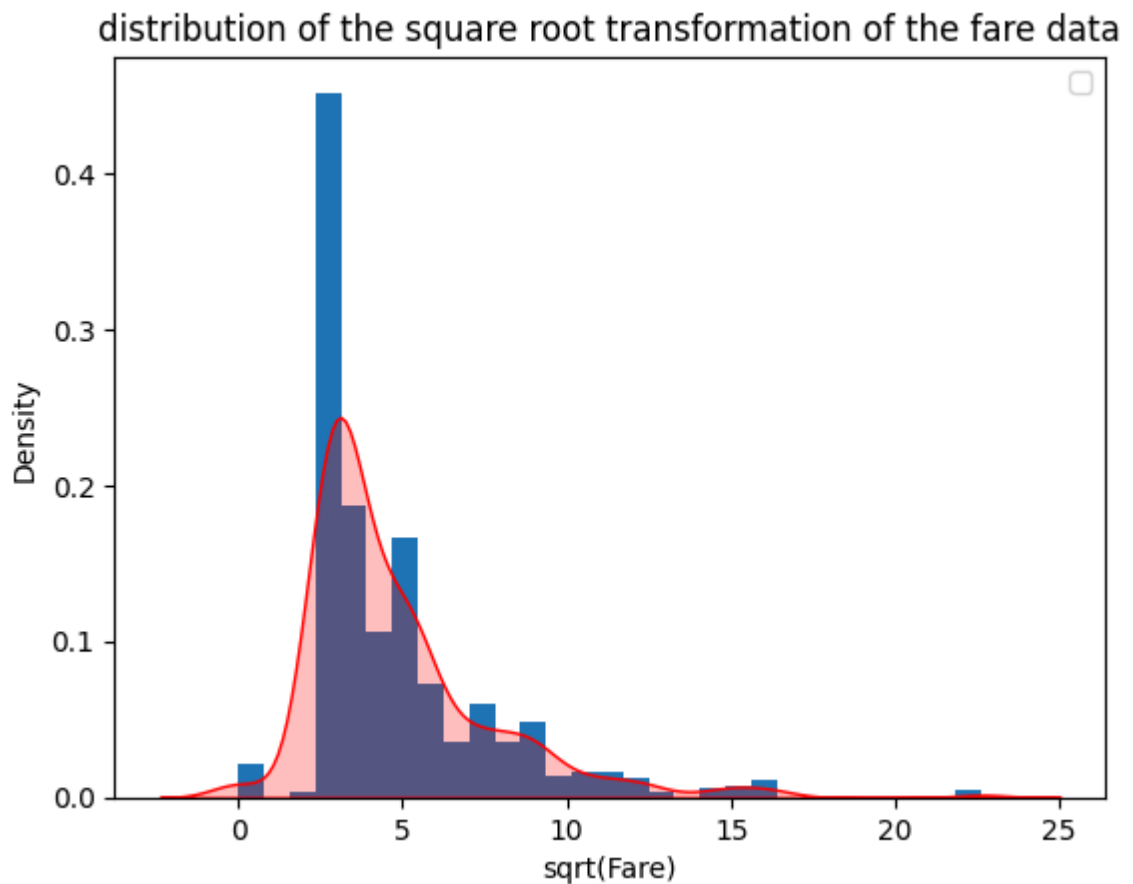
```
1  # Apply square root transformation to the age data
2  # and visualize the transformed distribution using a
3  # histogram and a kernel density plot
4
5  sqroot_fare = np.sqrt(df['Fare'])
6
7  plt.hist(
8      sqroot_fare,
9      bins=int(np.sqrt(len(sqroot_fare))),
10     range=None,
11     density=True,
12     weights=None,
13     cumulative=False,
14     bottom=None,
15     histtype='bar',
16     align='mid',
17     orientation='vertical',
18     rwidth=None,
19     log=False,
20     color=None,
21     label=None,
22     stacked=False,
23 )
24 sns.kdeplot(
25     data=sqroot_fare,
26     hue=None,
27     weights=None,
28     palette='viridis',
29     hue_order=None,
30     hue_norm=None,
31     color='Red',
32     fill=True,
33     multiple='layer',
34     common_norm=True,
35     common_grid=True,
36     cumulative=False,
37     bw_method='silverman',
38     bw_adjust=1,
39     warn_singular=True,
40     log_scale=None,
41     levels=10,
42     thresh=0.05,
43     gridsize=200,
44     cut=3,
45     clip=None,
46     legend=True,
47     cbar=False,
48     cbar_ax=None,
49     cbar_kws=None,
50     ax=None,
51 )
52
53 plt.xlabel('sqrt(Fare)')
54 plt.ylabel('Density')
55 plt.title('distribution of the square root transformation of the fare d
56 plt.legend()
```

```

/tmp/ipykernel_33/2867368794.py:24: UserWarning: Ignoring `palette` because no `hue` variable has been assigned.
  sns.kdeplot(
/opt/conda/lib/python3.10/site-packages/seaborn/_oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):

```

Out[21]: <matplotlib.legend.Legend at 0x7c3e596f89a0>



In []:

1

Task 5

In [22]:

```

1  # Perform Statistical sampling and calculate confidence
2  # intervals:
3
4  # a. Randomly sample 200 passengers' ages from the
5  # dataset and calculate the 95% confidence interval for
6  # the population mean age.
7
8  # b. Repeat step a multiple times (e.g., 1000 times) and
9  # calculate the confidence intervals each time.
10
11 # c. Calculate the proportion of confidence intervals
12 # that contain the true population mean.

```



```

In [23]: 1 # a. Randomly sample 200 passengers' ages from the
          2 # dataset and calculate the 95% confidence interval for
          3 # the population mean age.
          4
          5 sample_size = 200
          6 sample_passanger_age = np.random.choice(df['Age'].dropna(),size=sample_
          7 sample_mean = np.mean(sample_passanger_age)
          8
          9 sample_std = np.std(sample_passanger_age,ddof=1)
         10
         11 z_critical = stats.norm.ppf(0.975) # 0.975 for 95% confidence level, tw
         12
         13 standard_error = sample_std/ np.sqrt(sample_size)
         14
         15 confidence_interval = ((sample_mean - z_critical*standard_error),(sampl
         16 confidence_interval

```

Out[23]: (28.145661282601434, 32.272638717398564)

```

In [24]: 1 # b. Repeat step a multiple times (e.g., 1000 times) and
          2 # calculate the confidence intervals each time.
          3 intervals = []
          4 sample_size = 200
          5 for i in range(1000):
          6     sample_passanger_age = np.random.choice(df['Age'].dropna(),size=sam
          7     sample_mean = np.mean(sample_passanger_age)
          8
          9     sample_std = np.std(sample_passanger_age,ddof=1)
         10
         11     z_critical = stats.norm.ppf(0.975) # 0.975 for 95% confidence level
         12
         13     standard_error = sample_std/ np.sqrt(sample_size)
         14
         15     confidence_interval = ((sample_mean - z_critical*standard_error),(s
         16     print(i+1, confidence_interval)
         17     intervals.append( confidence_interval)

```

```

1 (28.609789308949175, 32.59021069105083)
2 (27.920187241428007, 31.561512758571997)
3 (29.014251272361328, 33.109948727638674)
4 (26.824110228893957, 30.93758977110604)
5 (26.2453433493244, 30.304656650675597)
6 (28.062349257411075, 32.261850742588926)
7 (27.447698956807624, 31.693201043192378)
8 (27.48440494085015, 31.503095059149846)
9 (28.819987443301024, 32.98831255669898)
10 (28.148161964726384, 32.22773803527362)
11 (27.30904037179329, 31.094359628206714)
12 (28.319443209419223, 32.19055679058077)
13 (27.8502636874476, 31.8730363125524)
14 (27.338091741862886, 31.215308258137117)
15 (26.3291645749079, 30.3058354250921)
16 (27.804042013740244, 32.06095798625976)
17 (27.777145723280544, 32.02955427671946)
18 (27.536339990857567, 31.897860009142438)
19 (26.877493903777193, 30.687506096222805)
20 (27.62086025005066, 31.70102064014124)

```

```
In [25]: 1 # c. Calculate the proportion of confidence intervals
2 # that contain the true population mean.
3
4 true_population_mean_age = df['Age'].mean()
5 count_mean = 0
6 for interval in range(len(intervals)):
7     if (true_population_mean_age >= intervals[0][0]) and (true_populati
8         count_mean += 1
9 print(f'the proportion of confidence intervals that contain the true po
```

the proportion of confidence intervals that contain the true population mean: 100.0 %

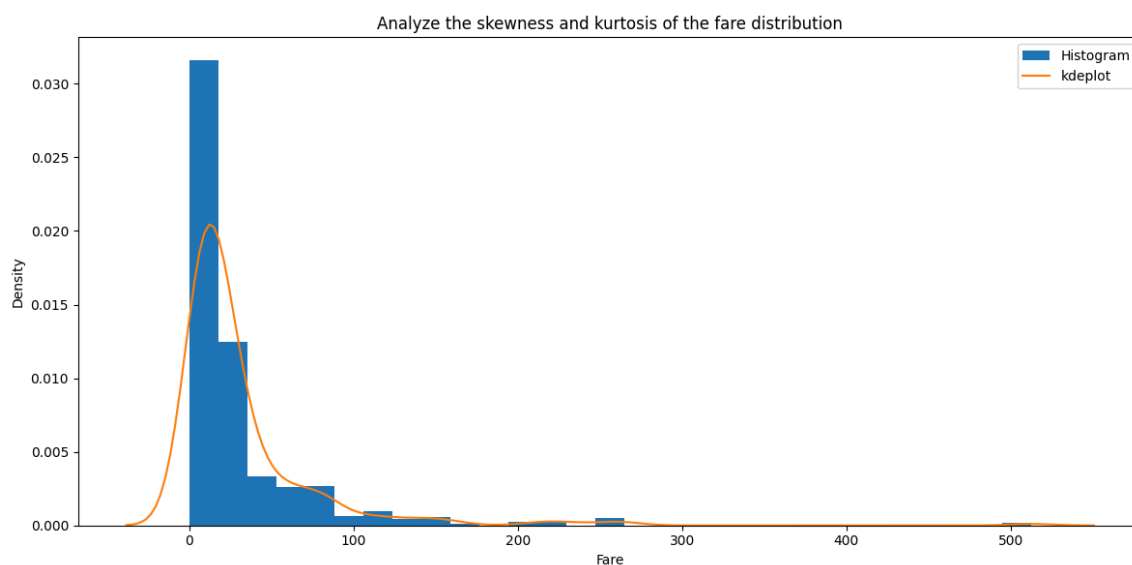
```
In [ ]: 1
```

Task 5

```
In [26]: 1 # Analyze the skewness and kurtosis of the fare
2 # distribution:
3
4 # Calculate the skewness and kurtosis of the fare
5 # data.
6
7 # Interpret the results and discuss the shape of the
8 # distribution
```

```
In [32]: 1 # Analyze the skewness and kurtosis of the fare
2 # distribution:
3
4 plt.figure(figsize=(12, 6))
5 plt.hist(df['Fare'], bins=int(np.sqrt(len(df['Fare']))), histtype='bar', d
6 sns.kdeplot(df['Fare'], bw_method='scott', label='kdeplot')
7
8 plt.xlabel('Fare')
9 plt.ylabel('Density')
10 plt.title('Analyze the skewness and kurtosis of the fare distribution')
11 plt.legend()
12
13 plt.tight_layout()
14 plt.show()
15
16
17 # Shape of the Distribution:
18
19 # The distribution is heavily skewed to the right, meaning most of the
20 # The high kurtosis indicates that the distribution has a sharp peak an
21
```

/opt/conda/lib/python3.10/site-packages/seaborn/_oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.
with pd.option_context('mode.use_inf_as_na', True):



```
In [28]: 1 # Calculate the skewness and kurtosis of the fare data.
2 fare_skewness = stats.skew(df['Fare'].dropna())
3 fare_kurtosis = stats.kurtosis(df['Fare'].dropna())
4
5 print(f'skewness of the fare distribution: {fare_skewness}')
6 print(f'kurtosis of the fare distribution: {fare_kurtosis}')
```

skewness of the fare distribution: 4.7792532923723545
kurtosis of the fare distribution: 33.20428925264474

```
In [29]: 1 # Interpret the results and discuss the shape of the distribution
2
3 # Skewness : Skewness measures the asymmetry of the distribution around
4
5 # Positive Skewness: A skewness value of 4.779 indicates a highly posit
6 # Implication: This suggests that most passengers paid relatively low f
7
8 # Kurtosis : Kurtosis measures the "tailedness" of the distribution.
9
10 # High Kurtosis (Leptokurtic): A kurtosis value of 33.204 is very high,
11 # Implication: The fare data has many extreme values, indicating a dist
```

```
In [ ]: 1
```

Task 6

```
In [33]: 1 # Calculate confidence intervals for population parameters:
2
3 # Calculate a 95% confidence interval for the population mean fare.
4
5 # Calculate a 95% confidence interval for the population proportion of
```

```
In [50]: 1 # Calculate a 95% confidence interval for the population mean fare.
2
3 population_mean = df['Fare'].mean() #population mean
4 population_std = df['Fare'].std() #population std
5
6 z_score = stats.norm.ppf(0.975) #for 95% confidence interval
7
8 standard_error = population_std/np.sqrt(len(df['Fare']))
9
10 confidence_interval = (population_mean - z_score*standard_error, populat
11 print("95% confidence interval for the population mean fare: ",confiden
12
13 count_fare=0
14 for fare in df['Fare']:
15     if (fare >= confidence_interval[0]) and (fare <= confidence_interva
16         count_fare+=1
17
18 print("Proportion of fares lie in the above mentioned confidence interv
19
```

95% confidence interval for the population mean fare: (28.94127463271884, 35.46714130443043)

Proportion of fares lie in the above mentioned confidence interval: 0.06

```
In [56]: 1 # Calculate a 95% confidence interval for the population proportion of
2
3 sample_proportion = df['Survived'].mean()
4
5 z_score = stats.norm.ppf(0.975) #for 95% confidence interval
6
7 standard_error = np.sqrt((sample_proportion * (1-sample_proportion))/le
8
9 margin_of_error = z_score * standard_error
10
11 confidence_interval = (sample_proportion - margin_of_error,sample_propo
12
13 sample_proportion
14
15 print('95% confidence interval for the population proportion of survivo
```

95% confidence interval for the population proportion of survivors (Survived column): (0.3519060427032577, 0.4157707249735099)