
Group A

Assignment No: 3

Contents for Theory:

1. Summary statistics
 2. Types of Variables
 3. Summary statistics of income grouped by the age groups
 4. Display basic statistical details on the iris dataset.
-

1. Summary statistics:

- **What is Statistics?**

Statistics is the science of collecting data and analysing them to infer proportions (sample) that are representative of the population. In other words, statistics is interpreting data in order to make predictions for the population.

Branches of Statistics:

There are two branches of Statistics.

DESCRIPTIVE STATISTICS : Descriptive Statistics is a statistics or a measure that describes the data.

INFERENCE STATISTICS : Using a random sample of data taken from a population to describe and make inferences about the population is called Inferential Statistics.

Descriptive Statistics

Descriptive Statistics is summarising the data at hand through certain numbers like mean, median etc. so as to make the understanding of the data easier. It does not involve any generalisation or inference beyond what is available. This means that the descriptive statistics are just the representation of the data (sample) available and not based on any theory of probability.

Commonly Used Measures

1. Measures of Central Tendency
2. Measures of Dispersion (or Variability)

- **Measures of Central Tendency**

A Measure of Central Tendency is a one number summary of the data that typically describes the centre of the data. This one number summary is of three types.

- a. **Mean :** Mean is defined as the ratio of the sum of all the observations in the data to the total number of observations. This is also known as Average. Thus mean is a number around which the entire data set is spread.

Consider the following data points.

17, 16, 21, 18, 15, 17, 21, 19, 11, 23

$$\text{Mean} = \frac{17 + 16 + 21 + 18 + 15 + 17 + 21 + 19 + 11 + 23}{10} = \frac{178}{10} = 17.8$$

- b. **Median :** Median is the point which divides the entire data into two equal halves. One-half of the data is less than the median, and the other half is greater than the same. Median is calculated by first arranging the data in either ascending or descending order.

- If the number of observations is odd, the median is given by the middle observation in the sorted form.
- If the number of observations are even, median is given by the mean of the two middle observations in the sorted form.

An important point to note is that the order of the data (ascending or descending) does not affect the median.

To calculate Median, let's arrange the data in ascending order.

11, 15, 16, 17, 17, 18, 19, 21, 21, 23

Since the number of observations is even (10), median is given by the average of the two middle observations (5th and 6th here).

$$\text{Median} = \frac{5^{\text{th}} \text{ Obs} + 6^{\text{th}} \text{ Obs}}{2} = \frac{17 + 18}{2} = 17.5$$

c. **Mode** : Mode is the number which has the maximum frequency in the entire data set, or in other words, mode is the number that appears the maximum number of times. A data can have one or more than one mode.

- If there is only one number that appears the maximum number of times, the data has one mode, and is called Uni-modal.
- If there are two numbers that appear the maximum number of times, the data has two modes, and is called Bi-modal.
- If there are more than two numbers that appear the maximum number of times, the data has more than two modes, and is called Multi-modal.

Consider the following data points.

17, 16, 21, 18, 15, 17, 21, 19, 11, 23

Mode is given by the number that occurs the maximum number of times. Here, 17 and 21 both occur twice. Hence, this is a Bimodal data and the modes are 17 and 21.

- **Measures of Dispersion (or Variability)**

Measures of Dispersion describes the spread of the data around the central value (or the Measures of Central Tendency)

1. **Absolute Deviation from Mean** — The Absolute Deviation from Mean, also called Mean Absolute Deviation (MAD), describes the variation in the data set, in the sense that it tells the average absolute distance of each data point in the set. It is calculated as

$$\text{Mean Absolute Deviation} = \frac{1}{N} \sum_{i=1}^N |X_i - \bar{X}|$$

2. **Variance** — Variance measures how far are data points spread out from the mean. A high variance indicates that data points are spread widely and a small variance indicates that the data points are closer to the mean of the data set. It is calculated as

$$\text{Variance} = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2$$

3. **Standard Deviation** — The square root of Variance is called the Standard Deviation. It is calculated as

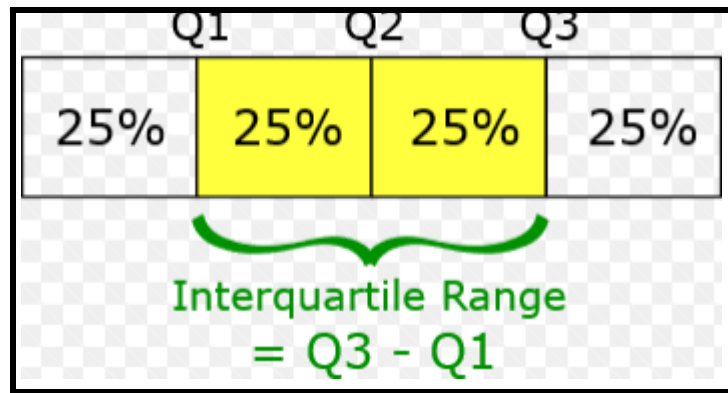
$$\text{Std Deviation} = \sqrt{\text{Variance}} = \sqrt{\frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2}$$

4. **Range** — Range is the difference between the Maximum value and the Minimum value in the data set. It is given as

$$\text{Range} = \text{Maximum} - \text{Minimum}$$

5. **Quartiles** — Quartiles are the points in the data set that divides the data set into four equal parts. Q1, Q2 and Q3 are the first, second and third quartile of the data set.

- 25% of the data points lie below Q1 and 75% lie above it.
- 50% of the data points lie below Q2 and 50% lie above it. Q2 is nothing but Median.
- 75% of the data points lie below Q3 and 25% lie above it.



6. **Skewness** — The measure of asymmetry in a probability distribution is defined by Skewness. It can either be positive, negative or undefined.

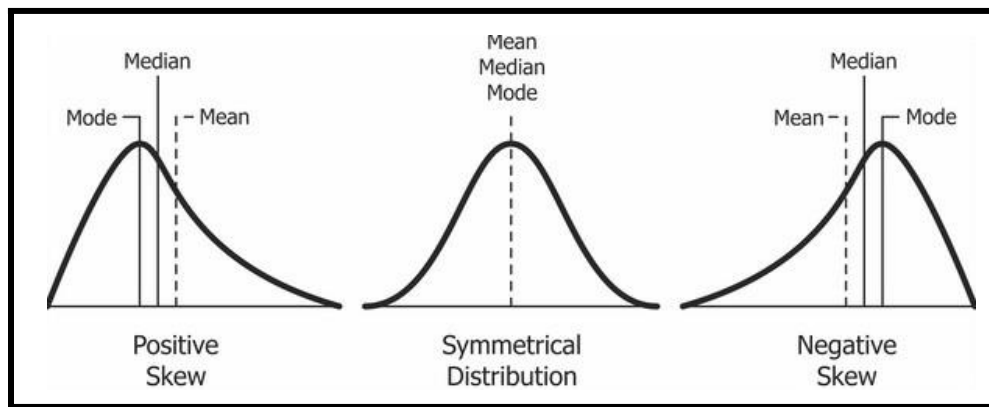
$$Skewness = \frac{3 (Mean - Median)}{Std \ Deviation}$$

Positive Skew — This is the case when the tail on the right side of the curve is bigger than that on the left side. For these distributions, mean is greater than the mode.

Negative Skew — This is the case when the tail on the left side of the curve is bigger than that on the right side. For these distributions, mean is smaller than the mode.

The most commonly used method of calculating Skewness is

If the skewness is zero, the distribution is symmetrical. If it is negative, the distribution is Negatively Skewed and if it is positive, it is Positively Skewed.



Python Code:

1. Mean

To find mean of all columns

Syntax:

```
df.mean()
```

Output:

```
CustomerID      100.50
Age              38.85
Annual Income (k$)  60.56
Spending Score (1-100)  50.20
dtype: float64
```

To find mean of specific column

Syntax:

```
df.loc[:, 'Age'].mean()
```

Output:

```
38.85
```

To find mean row wise

Syntax:

```
df.mean(axis=1)[0:4]
```

Output:

```
0    18.50
1    29.75
2    11.25
3    30.00
dtype: float64
```

2. Median

To find median of all columns

Syntax:

```
df.median()
```

Output:

```
CustomerID      100.5
Age              36.0
Annual Income (k$)  61.5
Spending Score (1-100)  50.0
dtype: float64
```

To find median of specific column

Syntax:

```
df.loc[:, 'Age'].median()
```

Output:

```
36.0
```

To find median row wise

Syntax:

```
df.median(axis=1)[0:4]
```

Output:

```
0    17.0
1    18.0
2    11.0
3    19.5
dtype: float64
```

3. Mode

To find mode of all columns

Syntax:

```
df.mode()
```

Output:

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Female	32.0	54.0	42.0
1	2	NaN	NaN	78.0	NaN
2	3	NaN	NaN	NaN	NaN
3	4	NaN	NaN	NaN	NaN
4	5	NaN	NaN	NaN	NaN
...
195	196	NaN	NaN	NaN	NaN
196	197	NaN	NaN	NaN	NaN
197	198	NaN	NaN	NaN	NaN
198	199	NaN	NaN	NaN	NaN
199	200	NaN	NaN	NaN	NaN

200 rows x 5 columns

In the Genre Column mode is Female, for column Age mode is 32 etc. If a particular column does not have mode all the values will be displayed in the column.

To find the mode of a specific column.

Syntax:

```
df.loc[:, 'Age'].mode()
```

Output:

```
32
```

4. Minimum

To find minimum of all columns

Syntax:

```
df.min()
```

Output:

```
CustomerID      1
Genre          Female
Age            18
Annual Income (k$)  15
Spending Score (1-100)  1
dtype: object
```

To find minimum of Specific column

Syntax:

```
df.loc[:, 'Age'].min(skipna = False)
```


Output:

18

5. Maximum

To find Maximum of all columns

Syntax:

```
df.max()
```

Output:

CustomerID	200
Genre	Male
Age	70
Annual Income (k\$)	137
Spending Score (1-100)	99
dtype: object	

To find Maximum of Specific column

Syntax:

```
df.loc[:, 'Age'].max(skipna = False)
```

Output:

70

6. Standard Deviation

To find Standard Deviation of all columns

Syntax:

```
df.std()
```

Output:

CustomerID	57.879185
Age	13.969007
Annual Income (k\$)	26.264721
Spending Score (1-100)	25.823522
dtype: float64	

To find Standard Deviation of specific column

Syntax:

```
df.loc[:, 'Age'].std()
```

Output:

13.969007331558883

To find Standard Deviation row wise

Syntax:

```
df.std(axis=1)[0:4]
```

Output:

0	15.695010
1	35.074920
2	8.057088
3	32.300671
dtype: float64	

2. Types of Variables:

A variable is a characteristic that can be measured and that can assume different values. Height, age, income, province or country of birth, grades obtained at school and type of housing are all examples of variables.

Variables may be classified into two main categories:

- Categorical and
- Numeric.

Each category is then classified in two subcategories: nominal or ordinal for categorical variables, discrete or continuous for numeric variables.

- **Categorical variables**

A categorical variable (also called qualitative variable) refers to a characteristic that can't be quantifiable.

Categorical variables can be either nominal or ordinal.

- **Nominal Variable**

A nominal variable is one that describes a name, label or category without natural order. In the given table, the variable “mode of transportation for travel to work” is also nominal.

Method of travel to work for Canadians	
Mode of transportation for travel to work	Number of people
Car, truck, van as driver	9,929,470
Car, truck, van as passenger	923,975
Public transit	1,406,585
Walked	881,085
Bicycle	162,910
Other methods	146,835

- **Ordinal Variable**

An ordinal variable is a variable whose values are defined by an order relation between the different categories. In the following table, the variable “behaviour” is ordinal because the category “Excellent” is better than the category “Very good,” which is better than the category “Good,” etc. There is some natural ordering, but it is limited since we do not know by how much “Excellent” behaviour is better than “Very good” behaviour.

Student behaviour ranking	
Behaviour	Number of students
Excellent	5
Very good	12
Good	10
Bad	2
Very bad	1

- **Numerical Variables**

A numeric variable (also called quantitative variable) is a quantifiable characteristic whose values are numbers (except numbers which are codes standing up for categories). Numeric variables may be either continuous or discrete.

- **Continuous variables**

A variable is said to be continuous if it can assume an infinite number of real values within a given interval.

For instance, consider the height of a student. The height can't take any values. It can't be negative and it can't be higher than three metres. But between 0 and 3, the number of possible values is theoretically infinite. A student may be 1.6321748755 ... metres tall.

- **Discrete variables**

As opposed to a continuous variable, a discrete variable can assume only a finite number of real values within a given interval.

An example of a discrete variable would be the score given by a judge to a gymnast in competition: the range is 0 to 10 and the score is always given to one decimal (e.g. a score of 8.5)

3. Summary statistics of income grouped by the age groups

Problem Statement: For example, if your categorical variable is age groups and quantitative variable is income, then provide summary statistics of income grouped by the age groups. Create a list that contains a numeric value for each response to the categorical variable.

Categorical Variable: Genre

Quantitative Variable : Age

Syntax:

```
df.groupby(['Genre'])['Age'].mean()
```

Output:

Genre	
Female	38.098214
Male	39.806818
Name: Age, dtype: float64	

Categorical Variable: Genre

Quantitative Variable : Income

Syntax:

```
df_u=df.rename(columns= {'Annual Income  
k$'): 'Income'}, inplace=False)
```

```
(df_u.groupby(['Genre']).Income.mean())
```

Output:

```
Genre
Female    59.250000
Male      62.227273
Name: Income, dtype: float64
```

To create a list that contains a numeric value for each response to the categorical variable.

```
from sklearn import preprocessing
enc = preprocessing.OneHotEncoder()
enc_df = pd.DataFrame(enc.fit_transform(df[['Genre']]).toarray())
enc_df
```

	0	1
0	0.0	1.0
1	0.0	1.0
2	1.0	0.0
3	1.0	0.0
4	1.0	0.0

To concat numerical list to dataframe

```
df_encode = df_u.join(enc_df)
df_encode
```

	CustomerID	Genre	Age	Income	Spending Score (1-100)	0	1
0	1	Male	19	15	39	0.0	1.0
1	2	Male	21	15	81	0.0	1.0
2	3	Female	20	16	6	1.0	0.0
3	4	Female	23	16	77	1.0	0.0
4	5	Female	31	17	40	1.0	0.0
...
195	196	Female	35	120	79	1.0	0.0
196	197	Female	45	126	28	1.0	0.0
197	198	Male	32	126	74	0.0	1.0
198	199	Male	32	137	18	0.0	1.0
199	200	Male	30	137	83	0.0	1.0

200 rows x 7 columns

4. Display basic statistical details on the iris dataset.

Algorithm:

1. Import Pandas Library

2. The dataset is downloaded from UCI repository.

```
csv_url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data'
```

3. Assign Column names

```
col_names =
```

```
['Sepal_Length', 'Sepal_Width', 'Petal_Length', 'Petal_Width', 'Species']
```

4. Load Iris.csv into a Pandas data frame

```
iris = pd.read_csv(csv_url, names = col_names)
```

5. Load all rows with Iris-setosa species in variable irisSet

```
irisSet = (iris['Species'] == 'Iris-setosa')
```

6. To display basic statistical details like percentile, mean, standard deviation etc. for Iris-setosa use describe

```
print('Iris-setosa')
```

```
print(iris[irisSet].describe())
```

7. Load all rows with Iris-versicolor species in variable irisVer

```
irisVer = (iris['Species'] == 'Iris-versicolor')
```

8. To display basic statistical details like percentile, mean, standard deviation etc. for Iris-versicolor use describe

```
print('Iris-versicolor')

print(iris[irisVer].describe())
```

9. Load all rows with Iris-virginica species in variable irisVir

```
irisVir = (iris['Species']== 'Iris-virginica')
```

10. To display basic statistical details like percentile, mean, standard deviation etc. for Iris-virginica use describe

```
print('Iris-virginica')

print(iris[irisVir].describe())
```

Iris-setosa				
	Sepal_Length	Sepal_Width	Petal_Length	Petal_Width
count	50.00000	50.00000	50.00000	50.00000
mean	5.00600	3.41800	1.46400	0.24400
std	0.35249	0.38102	0.17351	0.10721
min	4.30000	2.30000	1.00000	0.10000
25%	4.80000	3.12500	1.40000	0.20000
50%	5.00000	3.40000	1.50000	0.20000
75%	5.20000	3.67500	1.57500	0.30000
max	5.80000	4.40000	1.90000	0.60000
Iris-versicolor				
	Sepal_Length	Sepal_Width	Petal_Length	Petal_Width
count	50.00000	50.00000	50.00000	50.00000
mean	5.93600	2.77000	4.26000	1.32600
std	0.51617	0.31379	0.46991	0.19775
min	4.90000	2.00000	3.00000	1.00000
25%	5.60000	2.52500	4.00000	1.20000
50%	5.90000	2.80000	4.35000	1.30000
75%	6.30000	3.00000	4.60000	1.50000
max	7.00000	3.40000	5.10000	1.80000
Iris-virginica				
	Sepal_Length	Sepal_Width	Petal_Length	Petal_Width
count	50.00000	50.00000	50.00000	50.00000
mean	6.58800	2.97400	5.55200	2.02600
std	0.63588	0.32249	0.55189	0.27465
min	4.90000	2.20000	4.50000	1.40000
25%	6.22500	2.80000	5.10000	1.80000
50%	6.50000	3.00000	5.55000	2.00000
75%	6.90000	3.17500	5.87500	2.30000
max	7.90000	3.80000	6.90000	2.50000

Conclusion:

Descriptive statistics summarises or describes the characteristics of a data set. Descriptive statistics consists of two basic categories of measures:

- measures of central tendency and
- measures of variability (or spread).

Measures of central tendency describe the centre of a data set. It includes the mean, median, and mode.

Measures of variability or spread describe the dispersion of data within the set and it includes standard deviation, variance, minimum and maximum variables.

Assignment Questions:

- 1. Explain Measures of Central Tendency with examples.**
- 2. What are the different types of variables? Explain with examples.**
- 3. Which method is used to statistic the dataframe? write the code.**