Title of the Assignment: Descriptive Stastics- Measures of central tendency and variability.

perform the following operations on any open source dataset (e.g data, csv)

1. Provide summary stastics (mean, median, minimum, maximum, standard deviation) for a dataset (age, income etc) with numeric variables grouped by of the qualitative (categorical) variables.
for example. if your categorical variable is age groups quantitative variable is income. then provide summary statics of income grouped by the age groups.

2. write a python program to display some basic steastical details like percentile, mean standard deviation etc.

objective of the assignment: students should be able to perform the steastical operations using Python on any open source dataset.

prerequisite.
1. Basic of python programming
2. concept of stastistics such as mean, median, minimum, maximum, standard deviation etc

concepts for Theory:
1. summary stastistics
2. Types of variables.

3. Summary Statistics of income grouped by 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.

Branches of ~~statistics~~: statistics:
There are two branches of statistics.
DISCRIPTIVE STATISTICS: Descriptive statistics is statistics or measure that describes the data.

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 of inference beyond what is available.

Commonly used Measures
1. Measures of central Tendency
2. Measures of Dispersion (or variability)

- Measures of central Tendency
A measure of central tendency is one number summary of the data typically describes the centre of the data

a] Mean: Mean is defined at the ratio of the sum of all the observations, in the data to the total number of observation.

b] Median: Median is the point which devider the entire data into two equal halves. one-half of the data is less than the median is given b and other half is greater than the same.

c] Mode: Mode is the number - which has the maximum frequency in the entire data set or in the other words, mode is the number that appears the maximum number of the times.

• If there is only one number that appears the maximum number of times. the data her onemode and is called unimodel.

• Measures of Dispension (or variability)
measures of Dispession describes the spread of the data around the central value for measures of central Tendency)

1. Absolute Deviation from Mean - The absolute Deviation from mean, also called mean absolute deviation (MAD), describes the variation in dataset.

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

2) variance- variance measures how far are data points spread out from the mean

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

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) skewness- The measures of a symmetry in a probability distribution is defined by skewness. It can either be positive negative or undefined

$$\text{Skewness} = \frac{3(\text{mean - median})}{\text{std Deviation}}$$

2. Types of variables:

A variable is a characteristics that can be measured and the can assume different values. Height, Age, income, province of 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

- Categorical Variable-
A categorical variable (also called qualitative variables) refers to characteristics that can't be quantifiable.

- Norminal Variable-

    A norminal variable is one that describes a name, label or category without natural order in the given table.

- Ordinal Variable -

    An ordinal variable is a variable whose values are defined by an order relation between the different categories.

- Numerical variables-

    A numeric variable (also called quantitative variable) is a quantifiable characteristics whose values are numbers.

- Continuous variables-

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

    for instance, consider the height of student, The height can't take any values.

- Discreate variables:

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

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 income grouped by the age groups. Create a list that contains a numerical value for each response to the Categorial variable.

Categorial Variable = Genre
Quantitative variable: Age

Conclusion:

Descriptive statistics summarises or describes the characteristics of data set. Descriptive stastics consist of two basic categories of measures
- measures of central tendency and
- measures of variability (or spread)

Measures of central tendency describe the centre of 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, varience, minimum and maximum variables.