

Write-up	Correctness of Program	Documentation of Program	Viva	Timely Completion	Total	Dated Sign of Subject Teacher
4	4	4	4	4	20	

Expected Date of Completion:.....

Actual Date of Completion:.....

Group A

Assignment No: 3

Title of the Assignment: Descriptive Statistics - Measures of Central Tendency and variability Perform the following operations on any open source dataset (e.g., data.csv)

1. Provide summary statistics (mean, median, minimum, maximum, standard deviation) for a dataset (age, income etc.) with numeric variables grouped by one of the qualitative (categorical) variable. 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. 2. Write a Python program to display some basic statistical details like percentile, mean, standard deviation etc. of the species of 'Iris-setosa', 'Iris-versicolor' and 'Iris-versicolor' of iris.csv dataset. Provide the codes with outputs and explain everything that you do in this step.

Objective of the Assignment: Students should be able to data analysis using liner regression using Python for any open source dataset

Prerequisite:

1. Basic of Python Programming

2. Concept of mean, median, minimum, maximum, standard deviation

Introduction:

Descriptive Statistics is the building block of data science. Advanced analytics is often incomplete without analyzing descriptive statistics of the key metrics. In simple terms, descriptive statistics can be defined as the measures that summarize a given data, and these measures can be broken down further into the measures of central tendency and the measures of dispersion. Measures of central tendency include mean, median, and the mode, while the measures of variability include standard deviation, variance, and the interquartile range. In this guide, you will learn how to compute these measures of descriptive statistics and use them to interpret the data.

We will cover the topics given below:

1. Mean
2. Median
3. Mode
4. Standard Deviation
5. Variance
6. Interquartile Range
7. Skewness

We will begin by loading the dataset to be used in this guide.

Data:

In this guide, we will be using fictitious data of loan applicants containing 600 observations and 10 variables, as described below:

1. Marital_status: Whether the applicant is married ("Yes") or not ("No").
2. Dependents: Number of dependents of the applicant.
3. Is_graduate: Whether the applicant is a graduate ("Yes") or not ("No").
4. Income: Annual Income of the applicant (in USD).
5. Loan_amount: Loan amount (in USD) for which the application was submitted.
6. Term_months: Tenure of the loan (in months).
7. Credit_score: Whether the applicant's credit score was good ("Satisfactory") or not ("Not_satisfactory").
8. Age: The applicant's age in years.
9. Sex: Whether the applicant is female (F) or male (M).
10. approval_status: Whether the loan application was approved ("Yes") or not ("No").

Let's start by loading the required libraries and the data.

```
1import pandas as pd
2import numpy as np
3import statistics as st
4
```

```
5# Load the data
6df = pd.read_csv("data_desc.csv")
7print(df.shape)
8print(df.info())
```

python

Output:

```
1 (600, 10)
2
3 <class 'pandas.core.frame.DataFrame'>
4 RangeIndex: 600 entries, 0 to 599
5 Data columns (total 10 columns):
6 Marital_status    600 non-null object
7 Dependents        600 non-null int64
8 Is_graduate       600 non-null object
9 Income            600 non-null int64
10 Loan_amount      600 non-null int64
11 Term_months      600 non-null int64
12 Credit_score     600 non-null object
13 approval_status  600 non-null object
14 Age              600 non-null int64
15 Sex              600 non-null object
16 dtypes: int64(5), object(5)
17 memory usage: 47.0+ KB
18 None
```

Five of the variables are categorical (labelled as 'object') while the remaining five are numerical (labelled as 'int').

Measures of Central Tendency:

Measures of central tendency describe the center of the data, and are often represented by the mean, the median, and the mode.

Mean: Mean represents the arithmetic average of the data. The line of code below prints the mean of the numerical variables in the data. From the output, we can infer that the average age of the applicant is 49 years, the average annual income is USD 705,541, and the average tenure of loans is 183 months. The command `df.mean(axis = 0)` will also give the same output.

```
1df.mean()
```

python

Sandip Institute of Engineering and Management, Nashik

Output:

```
1 Dependents      0.748333
2 Income          705541.333333
3 Loan_amount     323793.666667
4 Term_months     183.350000
5 Age             49.450000
6 dtype: float64
```

It is also possible to calculate the mean of a particular variable in a data, as shown below, where we calculate the mean of the variables 'Age' and 'Income'.

```
1 print(df.loc[:, 'Age'].mean())
2 print(df.loc[:, 'Income'].mean())
```

python

Output:

```
1 49.45
2 705541.33
```

In the previous sections, we computed the column-wise mean. It is also possible to calculate the mean of the rows by specifying the (axis = 1) argument. The code below calculates the mean of the first five rows.

```
1 df.mean(axis = 1)[0:5]
```

python

Output:

```
1 0    70096.0
2 1    161274.0
3 2    125113.4
4 3    119853.8
5 4    120653.8
6 dtype: float64
```

Median

In simple terms, median represents the 50th percentile, or the middle value of the data, that separates the distribution into two halves. The line of code below prints the median of the numerical variables in the data. The command `df.median(axis = 0)` will also give the same output.

```
1 df.median()
```

python

Output:

```
1 Dependents      0.0
2 Income          508350.0
```

```

3 Loan_amount    76000.0
4 Term_months    192.0
5 Age            51.0
6 dtype: float64

```

From the output, we can infer that the median age of the applicants is 51 years, the median annual income is USD 508,350, and the median tenure of loans is 192 months. There is a difference between the mean and the median values of these variables, which is because of the distribution of the data. We will learn more about this in the subsequent sections.

It is also possible to calculate the median of a particular variable in a data, as shown in the *first two lines of code* below. We can also calculate the median of the rows by specifying the (`axis = 1`) argument. The *third line* below calculates the median of the first five rows.

```

1#to calculate a median of a particular column
2print(df.loc[:, 'Age'].median())
3print(df.loc[:, 'Income'].median())
4
5df.median(axis = 1)[0:5]

```

python
Output:

```

1 51.0
2 508350.0
3
4 0    102.0
5 1    192.0
6 2    192.0
7 3    192.0
8 4    192.0
9 dtype: float64

```

Mode: Mode represents the most frequent value of a variable in the data. This is the only central tendency measure that can be used with categorical variables, unlike the mean and the median which can be used only with quantitative data.

The line of code below prints the mode of all the variables in the data. The `.mode()` function returns the most common value or most repeated value of a variable. The command `df.mode(axis = 0)` will also give the same output.

```
1df.mode()
```

python
Output:

```

1|      | Marital_status | Dependents | Is_graduate | Income | Loan_amount |
Term_months | Credit_score | approval_status | Age | Sex |

```

2 ---	-----	-----	-----	-----	-----	-----
---	-----	-----	---	---		
3 0	Yes	0	Yes	333300	70000	192.0

| Satisfacto

The interpretation of the mode is simple. The output above shows that most of the applicants are married, as depicted by the 'Marital_status' value of "Yes". Similar interpretation could be done for the other categorical variables like 'Sex' and 'Credit-Score'. For numerical variables, the mode value represents the value that occurs most frequently. For example, the mode value of 55 for the variable 'Age' means that the highest number (or frequency) of applicants are 55 years old.

Measures of Dispersion:

In the previous sections, we have discussed the various measures of central tendency. However, as we have seen in the data, the values of these measures differ for many variables. This is because of the extent to which a distribution is stretched or squeezed. In statistics, this is measured by dispersion which is also referred to as variability, scatter, or spread. The most popular measures of dispersion are standard deviation, variance, and the interquartile range.

Standard Deviation:

Standard deviation is a measure that is used to quantify the amount of variation of a set of data values from its mean. A low standard deviation for a variable indicates that the data points tend to be close to its mean, and vice versa. The line of code below prints the standard deviation of all the numerical variables in the data.

```
1 df.std()
```

```
python
```

```
Output:
```

```
1 Dependents      1.026362
2 Income          711421.814154
3 Loan_amount     724293.480782
4 Term_months     31.933949
5 Age             14.728511
6 dtype: float64
```

While interpreting standard deviation values, it is important to understand them in conjunction with the mean. For example, in the above output, the standard deviation of the variable 'Income' is much higher than that of the variable 'Dependents'. However, the unit of these two variables is different and, therefore, comparing the dispersion of these two variables on the basis of standard deviation alone will be incorrect. This needs to be kept in mind.

It is also possible to calculate the standard deviation of a particular variable, as shown in the *first two lines of code* below. The *third line* calculates the standard deviation for the first five rows.

```
1 print(df.loc[:, 'Age'].std())
2 print(df.loc[:, 'Income'].std())
3
```

```
4#calculate the standard deviation of the first five rows
```

```
5df.std(axis = 1)[0:5]
```

python

Output:

```
1 14.728511412020659
2 711421.814154101
3
4 0    133651.842584
5 1    305660.733951
6 2    244137.726597
7 3    233466.205060
8 4    202769.786470
9 dtype: float64
```

Variance

Variance is another measure of dispersion. It is the square of the standard deviation and the covariance of the random variable with itself. The line of code below prints the variance of all the numerical variables in the dataset. The interpretation of the variance is similar to that of the standard deviation.

```
1df.var()
```

python

Output:

```
1 Dependents    1.053420e+00
2 Income        5.061210e+11
3 Loan_amount   5.246010e+11
4 Term_months   1.019777e+03
5 Age           2.169290e+02
6 dtype: float64
```

Interquartile Range (IQR)

The Interquartile Range (IQR) is a measure of statistical dispersion, and is calculated as the difference between the upper quartile (75th percentile) and the lower quartile (25th percentile). The IQR is also a very important measure for identifying outliers and could be visualized using a boxplot.

IQR can be calculated using the `iqr()` function. The *first line of code* below imports the 'iqr' function from the `scipy.stats` module, while the *second line* prints the IQR for the variable 'Age'.

```
1from scipy.stats import iqr
2iqr(df['Age'])
```

python
Output:

```
1 25.0
```

Skewness

Another useful statistic is skewness, which is the measure of the symmetry, or lack of it, for a real-valued random variable about its mean. The skewness value can be positive, negative, or undefined. In a perfectly symmetrical distribution, the mean, the median, and the mode will all have the same value. However, the variables in our data are not symmetrical, resulting in different values of the central tendency.

We can calculate the skewness of the numerical variables using the `skew()` function, as shown below.

```
1 print(df.skew())
```

python
Output:

```
1 Dependents    1.169632
2 Income        5.344587
3 Loan_amount   5.006374
4 Term_months   -2.471879
5 Age           -0.055537
6 dtype: float64
```

The skewness values can be interpreted in the following manner:

- Highly skewed distribution: If the skewness value is less than -1 or greater than $+1$.
- Moderately skewed distribution: If the skewness value is between -1 and $-\frac{1}{2}$ or between $+\frac{1}{2}$ and $+1$.
- Approximately symmetric distribution: If the skewness value is between $-\frac{1}{2}$ and $+\frac{1}{2}$.

Putting Everything Together

We have learned the measures of central tendency and dispersion, in the previous sections. It is important to analyse these individually, however, because there are certain useful functions in python that can be called upon to find these values. One such important function is the `.describe()` function that prints the summary statistic of the numerical variables. The line of code below performs this operation on the data.

```
1 df.describe()
```

python
Output:

```
1 |      | Dependents | Income | Loan_amount | Term_months | Age
2 |-----|-----|-----|-----|-----|-----
```



```

3| count | 600.000000 | 6.000000e+02 | 6.000000e+02 | 600.000000 | 600.000000
4| mean | 0.748333 | 7.055413e+05 | 3.237937e+05 | 183.350000 | 49.450000
5| std | 1.026362 | 7.114218e+05 | 7.242935e+05 | 31.933949 | 14.728511
6| min | 0.000000 | 3.000000e+04 | 1.090000e+04 | 18.000000 | 22.000000
7| 25% | 0.000000 | 3.849750e+05 | 6.100000e+04 | 192.000000 | 36.000000
8| 50% | 0.000000 | 5.083500e+05 | 7.600000e+04 | 192.000000 | 51.000000
9| 75% | 1.000000 | 7.661000e+05 | 1.302500e+05 | 192.000000 | 61.000000
10| max | 6.000000 | 8.444900e+06 | 7.780000e+06 | 252.000000 |
76.000000 |

```

The above output prints the important summary statistics of all the numerical variables like the mean, median (50%), minimum, and maximum values, along with the standard deviation. We can also calculate the IQR using the 25th and 75th percentile values.

However, the 'describe()' function only prints the statistics for the quantitative or numerical variable. In order to print the similar statistics for all the variables, an additional argument, include='all', needs to be added, as shown in the line of code below.

```
1 df.describe(include='all')
```

python

Output:

```

1|      | Marital_status | Dependents | Is_graduate | Income | Loan_amount |
2|-----|-----|-----|-----|-----|-----|
3| count | 600 | 600.000000 | 600 | 6.000000e+02 | 6.000000e+02 |
4| unique | 2 | NaN | 2 | NaN | NaN |
5| top | Yes | NaN | Yes | NaN | NaN |
6| freq | 391 | NaN | 470 | NaN | NaN |
7| mean | NaN | 0.748333 | NaN | 7.055413e+05 | 49.450000 |
3.237937e+05 | 183.350000 | NaN | NaN | 192.000000 | NaN |
8| std | NaN | 1.026362 | NaN | 7.114218e+05 | 7.242935e+05 |
9| min | NaN | 0.000000 | NaN | 3.000000e+04 | 1.090000e+04 |
10| 25% | NaN | 0.000000 | NaN | 3.849750e+05 | 6.100000e+04 |
6.100000e+04 | 192.000000 | NaN | NaN | 192.000000 | NaN |
11| 50% | NaN | 0.000000 | NaN | 5.083500e+05 | 7.600000e+04 |
7.600000e+04 | 192.000000 | NaN | NaN | 192.000000 | NaN |
12| 75% | NaN | 1.000000 | NaN | 7.661000e+05 | 1.302500e+05 |
1.302500e+05 | 192.000000 | NaN | NaN | 192.000000 | NaN |
13| max | NaN | 6.000000 | NaN | 8.444900e+06 | 7.780000e+06 |
7.780000e+06 | 252.000000 | NaN | NaN | 252.000000 | NaN |

```

Now we have the summary statistics for all the variables. For qualitative variables, we will not have the statistics such as the mean or the median, but we will have statistics like the frequency and the unique label.

Conclusion:

In this guide, you have learned about the fundamentals of the most widely used descriptive statistics and their calculations with Python. We covered the following topics in this guide:

1. Mean
2. Median
3. Mode
4. Standard Deviation
5. Variance
6. Interquartile Range
7. Skewness

It is important to understand the usage of these statistics and which one to use, depending on the problem statement and the data. To learn more about data preparation and building machine learning models using Python's 'scikit-learn' library, please refer to the following guides:

1. [Scikit Machine Learning](#)
2. [Linear, Lasso, and Ridge Regression with scikit-learn](#)
3. [Non-Linear Regression Trees with scikit-learn](#)
4. [Machine Learning with Neural Networks Using scikit-learn](#)
5. [Validating Machine Learning Models with scikit-learn](#)
6. [Ensemble Modeling with scikit-learn](#)
7. [Preparing Data for Modeling with scikit-learn](#)