

set

Manjari (Bk.), Hadapsar, Pune-412307.

Accredited by NAAC

ASSIGNMENT No 3

Title: Descriptive statistics - Measure of central Tendency of variability.

Aim: Descriptive statistics-Measure of central tendency of variability.

Perform the following operations on any open source dataset. (eg. data, csv).

1. Petvide summary statistics for a dataset with numeric variables grouped by one of the quantitative variable. For example, if your categorial variable is age, groups of quantitative variable in income, then provide summary statistic of income grouped by the age groups. Create a list that contains a numeric value for each response to the categorial variable.

2. Weite a python program to display some basic statistical details like percentile, mean, standard deviation etc. of the species of 'Tris-setosa', 'Iris-versicolor' of Tris-setosa', dataset.

Provide the codes with ofp of explain everything that you do in this step.

- P.T.O.





set

Manjari (Bk.), Hadapsar, Pune-412307.

Accredited by NAAC

Objectives: students should be able to perform the statistical operations using python on any open sourcedataset.

Requirement: 1. Basic of python programming.

2. Concept of Data preprocessing, Data formatting, Data
Normalization & cleaning.

Theory: 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 of central tendency f
these measures can be broken down further into the
measures of central tendency f the measures of
dispersion.

Measures of central tendency include mean, median of mode, while the measures of variability include standard deviation, variance of the interquartile sange. In this guide, you will learn how to compute these measures of descriptive statistics of use them to interpret the data

- P.T.O.



Suns Sun

sel

Manjari (Bk.), Hadapsar, Punc-412307.

Accredited by NAAC

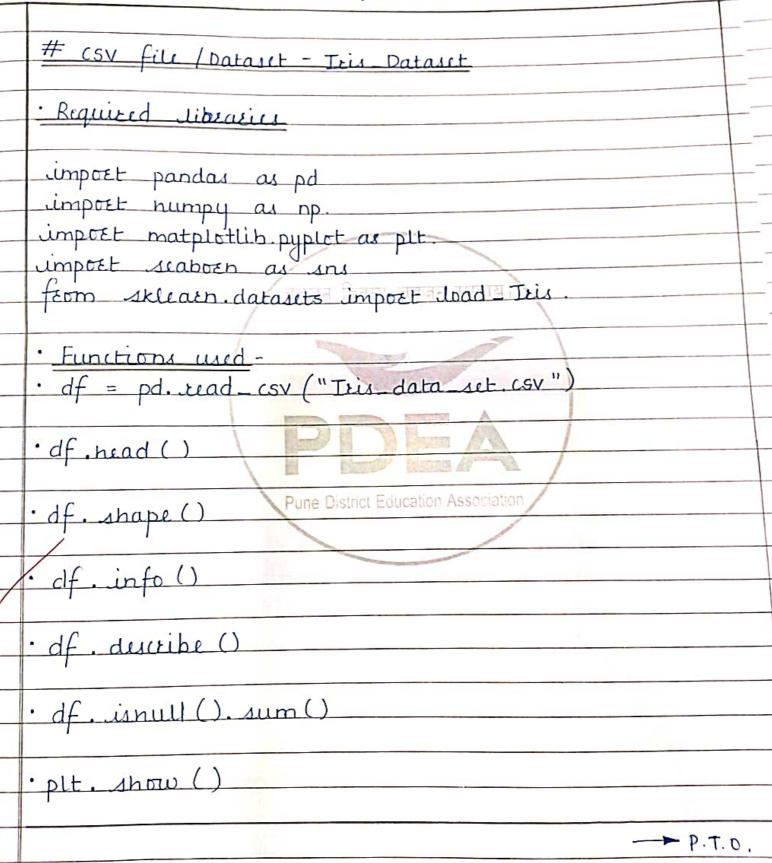
We will cover the topics given below :-1. Mean 2. Median 1. Standard deviation S. Variance 6. Interquartile 1. Skewness. We will begin by loading the dataset to be used un this guide. Data: In this guide, we will be using fictitions data of Joan applications containing 600 observations of 10 variables, 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 graduate. ("Yes") OE ("Not") 1. Income: Annual Income of the applicant (in USD) 8. Loan_amount: Loan amount (in USD) for which the application was submitted 6. Term_Monthy: Tenuce of the Joan (in months) 7. Credit - score; whether the applicants credit score was good (" satisfactory") or Not ("Not satisfactory") 8. Age: The applicants age (in years) I. sex: whether the applicant is female or male. 10. approval-status: whether the Joan application was approved ("Ya") or not ("No").





Manjari (Bk.), Hadapsar, Pune-412307.

Accredited by NAAC



set





Manjari (Bk.), Hadapsar, Pune-412307.

Accredited by NAAC



set

_	Trecredited by NAAC
	· mean = geouped_df.mean()
	· median = geouped_df. median ()
	· min = geouped - df. min ()
	· max = geouped _ df. max ()
	· std = geouped_df. std () F. std () F. std + tolal
	· df. skew ()
	Do all operations on each column. Also Draw
	boxplot for each column
	Conclusion: In this guide, you have learned about the fundamentals of the most widely used descriptive
	the fundamentals of the most widely used descriptive
	statistics of their calculations with python we
	lovered the following topics in this guide.
1	1. Mean 2. Median 3. Mode.
	4. Standard Deviation S. Variance 6. Interquartile range
	7. Skewness.
	It is important to understand the usage of these
	statistics & which one to use, depending on the
	problem statement of the data.
	·





Manjari (Bk.), Hadapsar, Pune-412307.

Accredited by NAAC

15et

	To learn more about data preparation of building machine models using python's 'scikit-learn' library, please refer to the following guides- 1. scikit machine learning. 2. Ensemble modeling with scikit-learn.	
1	Which teems include under Measures of central tendency? Mean, Median, Mode.	
2	Which terms include under Measures of variability? Standard deviation, variance & interquartile range.	
3	Describe Tris dataset. Tris is a collection of instruments, materials, stimuli, data & data coding & analysis tools used for research into languages, signed language learning, etc. It contains four features (length & width of Lepals & petals) of 50 samples of three species of Tris (Tris setosa, Tris virginica & Tris versicolor). Tris dataset contains fire columns such as Petal length, Petal width, Sepal length, Sepal Width & species Type. & sows being the samples.	
4.	Explain Mean Mean represents the arithmetic average of the data. P.T.O.	1
		_





150

-- P.T.O.

Manjari (Bk.), Hadapsar, Pune-412307.

Accredited by NAAC Describe Measures of Central Tendency. It describe the center of the data of often supresented by mean, median, mode. Explain Median Median represents the 50th percentile of the middle value of the data, that separates the distribution into two halves. Explain Mode It represents the most frequent value of a variable in the data Explain Standard Deviation. It 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 & vice versa. Explain Variance It is square of the standard deviation of the covariance of the random variable with itself.



150

Manjari (Bk.), Hadapsar, Pune-412307.

-	Accredited by NAAC
10.	Explain Interquartile Range (IQR). Measure of statistical dispersion. It is calculated as the difference between upper quartile (75th percentile) f the lower quartile (25th percentile). IQR is very important measure for identifying outliers & could be visualize using boxplot.
12.	Explain term outliers. Outliers are always given wrong direction for your expected results. Outliers always talk about extremities. Too small or too large.
13.	Explain skewness It is used to measure of symmetry or lack of symmetry. The skewness value can be positive, negative or undefined. In a perfectly symmetrical distribution, the nacan, median & mode will all have the same value.