CIASSMA	ite
Date :	
Page:	

	CIASSMATE Date: Page:
	Name: Tejas V. Redkan
	Roll No: PC-44
	PRN: 1032210937
	Batch: C2
	I was a series for the series and the series are the series and the series are the series and the series are th
	OEC Lab Assignment - 1
(11	a vive quissite extinctive and a vive to
*	Puroblem Statement:
	and the first of the contract
	Data Handling, locate any open source data: Load
	data into data grame. Perform bata forame
. 5-4	operations. Perform basic statistical operations
1.100	title mean, median, wandard deviation.
4.10.	Alle San Carrier Company of the second secon
*	objectives:
4	The first of the f
~	To explouse various data isources & data irrepositories
۵	To exploure the operations on a dataset file using
	data frame with basic whatistical operations in Python.
	so to the state of
*	Execution:
	stall hear . As it for the said the sai
48 - M 1	purguam implementation of output study.
1	The state of the s
×	conclusion:
	the contract of the second
	Basic operations were performed on the cov data
39	file using Python.
	
A Comment	

*	FAOLS
<u>Q1)</u>	state the wignificance of handling missing
	state the vigniticance of handling missing values in a data set.
. Ans	Handling missing values is crucial your veveral
	veasons:
	1) Haintaining data integrity: Hissing values can
	introduce reveral errours & inconsistencies in
	data analysis. Ignowing them can lead to
	incourrect our binsod viesults, making it
	essential to add sess them peroposely.
	2) Pureventing Biased Analysis: Tyl missing values
	are not handled, the data analysis may become
	I based towards the observationals with complete
	data, potentially skewing the viesults.
	3) Imperoving Model Performance: When building
	madrine reading models, missing values can
	cause issues during training & prediction.
	tranding them appropriately can visualt in
	better model performance & generalization.
	4) Meeting Analysis veguivement: Some statistical
	techniques & machine learning algorithms
	require complete datacets. Addressing missing
	values is necessary to meet the precisequisites of
	these methods.
	5) Enhancing data visavisualizations: Missing values
	can disvupt data visualizations, making it
	difficult do create meaning ful charts & graphs.

1	
6)	Facilitating Data Shaving: If the dataset is
	to be whaved, addressing missing values makes
	it movre accessible & usable you other vieseauchers,
	analysts our stakeholders.
(02)	Explain the centural tendency measures with
	examples.
Ans	Centeral tendency measures are estatistical
	values that oropresent the center our average
	of a dataset, punviding insights into where the
	data tends to duster.
1.7.	There are three primary measures: mean,
	median, & mode.
¥	
1)	Mean (Average):
	The mean is the most common measure of
	centural tendency of is calculated by summing all
	the values in a dataset & dividing by the total
	number of values.
1	and the second of the second o
. 1	Fournula: Mean = (Sum of values)/(Number of values)
, ·	Example: Mean of a scoure: [85,92,78,88,95]
	mean = (85+92+ 78+88 +95)/5
	= 87.6
	the first and all are named to be and the
<u>2)</u>	Median:
6. Or	and booking the many the year wish
_ , , , , , ,	It is the middle value when a datacret is is overled



in ascending our descending auder. Ty there is an even number of values, the median is average of the middle-two values.

Example: In above dutaset: [85,892,78,88,95] when sowled in ascending ourder-[78,885,88,92,95]

Median = 88

3) Mode:

The mode is the value that occurs most frequently in a dataset. A dataset can have one mode (unimodal) multiple modes our no modely all values are unique.

Example: In the dataset: [82, 78, 95, 82, 83,

82 is the mode.

values in a daterset.

Handling missing values is an essential cater in

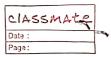
data pure puro cessing to ensure accurate & vieliable analysis. Some common methods to handle missing values are:

1) Deletion:

List wise deletion: In this method, entire vous containing missing values.

Pairwise deletion: In this method, missing values are ignored your specific calculations.

Date:	CIASSMATE Date: Page:
	when analyzing agis of waviables.
	when analyzing pair of variables. 2) Hean, median, our mode imputation: Hissing
	2) Mean, median, our mode reparce
	values can be veplaced using mean (you
	continuous data), median (your ourdinal data)
	de mode (your categorical data). Imputation
	perovides a wimple way of filling missing
	3) Interpolation: Interpolation methods, such as
	3) Interpresatation: Interpolation interpolation,
<u></u>	linear vegression our apline interpolation, estimate missing values based on values of
	d'adjaunt data points.
	a dajaunt aux poins.
Bli	Gratain dillerant types of data types.
O4 Ans	
VALD	as they provide apecitications as to how data is
	storred, manipulates & interpretted.
	1) Numeric Data Types:
	1) Integer lint): Represents whole numbers
	2 Floating point (yloat): Represents number with
	decimal point.
v	3 Double precision (double) : Higher precision de
	larger veange of
	values
	4) Long (long): Repuresents large integers (
	ex: 1234567890L)
	2) Text Data Types:
,	JI



	1) Sturing (when our text): Repuresents isoquence
1 2	(characters): Represents a wingle characters.
1	
8 .	3) Boolean Data Type:
	(1) Boolean (bool): Repuresents binary values, typically Turne our False.
1	
	4) Date & Time Data Types:
	① Date: Repuresents a specific data (formatted) ② Time: Repuresents specific time (ex:14:30:00) ③ Datetime: Repuresents date & time together.
	5) Categornical Data Types:
200	D Enumerations (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of predefined named values. Descriptions (enums). Represents a uset of