Page No.:

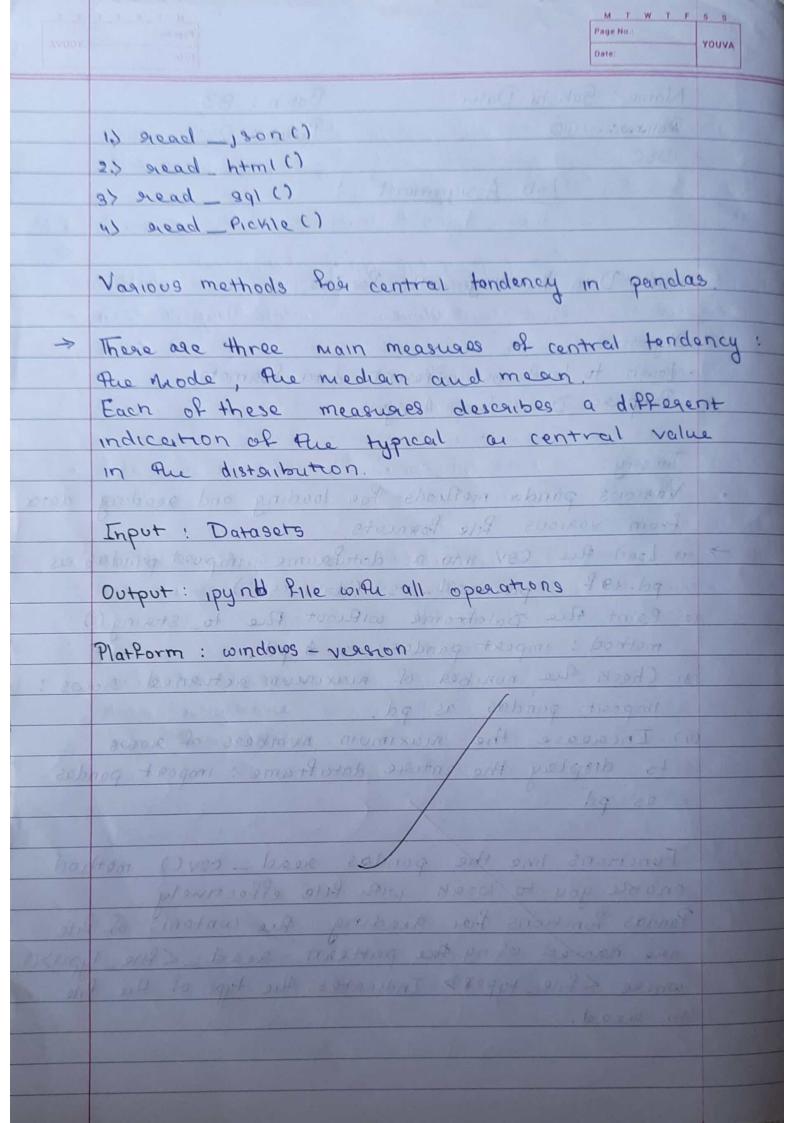
Page No.:

VOUVA

Name: Gakshi Dalvi Batch: B3 Roll no.: 60 Panel : 2 DEC Lab Assignment: 1 (2) per house to Aim : Data handling wife Pandas Objective: 101 for Edward minus south all and · Leasn to load use different data formats · Prepare Data for Analysis of boot to whom havenous is whodget by the to normanism Theory: · Various panda methods for loading and reading data from various file formats. Hospital + tuga -> 11) Load Que CSV into a databaame: imposit pandas as pd.d9f = pd. soad _ csv ('data.csv') (2) Paint the DataFrame without the to_ staing () method: imposit pandas as & pd. (3) Check the number of maximum geturned nows: imposit pandas as pd. (u) Increase the maximum numbers of nows to display the entire dataframe: impost pandas as pd. Functions like the pandas read _ csv() method enable you to work with file effectively. Pandas functions for sending the contents of file one harried using the pertern send _ < file types >0),

where < Prie type > Indicates the type of the file

to read.



of FAQS togularous with events of through and whole

What are different types data file format?

There are there types of data file format:

- File based data format! This type of data format includes either one can file on more than one file.

 These files are then around in any of the arbitrary folders. In most of the cases it uses the single file only for ex Dan. But then there cases, which even includes atleast these files. The filename extension of all these these files is different from each other
- Disectory Based Data format: In this type of data francat, whether Areae is one file or There is more than one file they are all stored in the parent facilities in a particular manner. There are some cases where the arguinement of an additional folder is there in the file tree in some other location so that it can be accessed easily. It is a possibility that data source is the directory itself. There are many files present in the directory, which are represented at the available data's layers.
 - Database connection: In one respect, the database connections are quite similar to the above mentioned data formats that are file and directory-based data format. For interpreting, for map server, they give the geographic coordinate

	data. One needs to access the coordinates Alherde
	The map server that'se creating vector datasals.
	a se accomenta land otal saget those 27 h as took 10
Q2	Compage various central tendency measures
=>	
100	Mode ich for egyt ent Median & whole beend Hean
	most frequent dater 100. Availe Prot divides 0. x = 2x
	point Profession of la sea sanked data pts 301 for seath N
1 1	180 to 2000 suff for into halves: 50% of permost stable
	Proces And Mod 18 rollanger Prantit, espire meesure.
	II roke south tarolto 50% Smaller a device 202 do
	mode exist as a doter 11. may not exist as we may not
	point la Pa date point in Rie exist as ex
	set elatapt in set
	unaffected by to influenced by affected by
ero,	
	extreme values position of items, extreme
ant Pl	position of items extreme
DIP!	rextreme values an apost aposition of items of extreme
9 ROJ	destar nous have a seed a seed as a
9 ROJ	destar nous have a seed a seed as a
or Pl	dater hay have ment a sent a sent and a sent a sent and
9 RO)	dester nous have no sent some destance of thems of extreme destance have the sent sent sent sent sent sent sent sen
9 RO)	dester nous have no sent some destance of thems of extreme destance have the sent sent sent sent sent sent sent sen
enspl 9 ken 10.00 10	extreme values and position of items of extreme I besent a 110 as a part of the form of items of extreme Useful for qualitative and a sent as a part of the part I more from I value man to the part of the part Useful for qualitative and to the part of the part I more from I value man to the part of the part Useful form of the part of t
enspl 9 ken 10.00 10	extreme values and position of items of extreme The sect a 110 as a part of the last and values Useful for qualitative as sector as a part of the sector dates. I may have of sector as a part of the sector mode Phan I value may be to be proposed as a part of the sector when so self present sector as a part of the sector.
and property	extreme values a position of items of extreme where the last part of the last and a position of items of extreme where the last part of the last and a position of items of extreme where the last part of the last and a position of items of extreme where the last part of the last and a position of items of extreme where the last part of the
ensel en	extreme values and appell appoint of items of extreme whose to the service and allow through a values about the service and t
order	extreme values a position of items of extreme where the last part of the last and a position of items of extreme where the last part of the last and a position of items of extreme where the last part of the last and a position of items of extreme where the last part of the last and a position of items of extreme where the last part of the
ensper 9 Reg 10.00 1	extreme values an appel aposition of Hems of extreme who so to 110 also want allow not free of the so values was unaw harden through a values about the so was ful for qualitative and a season as a marked a season and a season and a season as a
ensel en	extreme relies and appeal appeals of items of extreme who works 110 appears to and allow the foreign a relies about relies (seful for qualitative and a relies about appearance appearance) dates. I may have or a rest at a per dark about a per dark a per dark about a per dark abou

a c 7 Y W 1 B