



Assignment of

Title- Data Wrangling

problem statement

perform the following operations using python or any open source data set (e.g. data.csv)

- 1) import all the required libraries
- 2) locate an open source data from web.
- 3) load data set into pandas data frame.
- 4) Data processing: checks for missing values in the data using pandas isnull(), describe() etc.
- 5) Data formatting and data Normalization: summarise the types of variables by checking the data types. (i.e. character, numeric, integer, logical)
the variables in the data set, IF the variable used are not in the correct data type, apply further type conversion.
- 6) Turn categorical variables into quantitative variables in python in addition to the codes and outputs explain every operation that you do in the above steps and explain everything that you do to import / read / scrap the data set.

Learning objectives

one should be able to make raw data useful by applying scientific data processing libraries.



Learning outcomes

understand the importance of data wrangling and role of scientific data processing libraries.

H/w & S/w Requirement

8 GB RAM

Windows 10

Jupyter Notebook

python 3.0

Theory

Data Wrangling

Data Wrangling, sometimes referred to as data munging, is the process of transforming and mapping data from one raw data into another format with the intent of making it more appropriate and valuable for variety of downstream purpose such as analysis.

- The goal of the data wrangling is to assure quality and useful data.
- Data Analysis typically spend the majority of their time in the process of data wrangling compared to the actual analysis of data.
- the process of data wrangling may include further munging, data visualization, data aggregation, training and statistical model as well as many other uses.
- Data wrangling typically follows a set of general steps which begin with extracting the data in raw form from the data source, munging the raw data is forcing the data into classified data structure.

methods and functions used.

(i) `pandas.read.csv()`

- This method is used to read a comma separated values file into Data frame.
- Also supports optionally iterating or breaking of the file into channels.

2. Data Formatting

- Data formatting is the process of transforming data into common format, which helps users to perform comparisons.
- An example of not formatted data is the following where the same entity is referred in the same column with different values, such as Newyork & NY.
- We should also make sure that every column is assigned to the correct data type.

Data Normalization

Normalization is a technique often applied as a part of data preparation of me. the goal of normalization is to change the values of numeric columns in the dataset to a common scale.

without distorting difference in the ranges values.

- Normalization is also required for some algorithms to model the data correctly.

2) data frame.head(limit)

returns the first 5 rows of the data frame, to override the default, you may insert value between the braces to change the number of rows.

Ex. head(10) will return 10 rows.

3) data frame.shape

returns a tuple representing the dimensions.

for ex. an output (100,14) represents 100 rows & 14 columns.

4) data frame.describe()

provides descriptive statistics that summarize the central tendency, dispersion and shape of a dataset.

5) value count

this method returns count for each unique value in the column you selected.

6) is NULL

returns data frame with value 'true' where it finds null value and false where it encounters any value.

7) df['col'].unique()

returns the unique values in column 'col'



PICT, PUNE

page No. 5

8) `isnull().sum()`

returns count of null values in the column.

9) `dataframe.dtypes`

this attribute returns the dtype in the data frame
it returns the series with the data type of each column

10) `pandas.map()`

it is used to map values from two series having
one column same, for mapping two series, the
last column of the first series should be same as
index column of the second series, also the values
should be unique.

packages/ module/ libraries used.

1) `pandas`:

- `pandas` is software library written for the python programming language for data manipulation and analysis.
- In particular it offers data structure and operations for manipulating numerical tables and time series it is a free software library released under the three clause BSD license.
- The name is derived from the term 'panel data' an econometrics terms for data set that includes observations are multiple time period for the same individuals.
- It allows importing various file formats such as CSV, JSON, SQL datasets & queries.



PICT, PUNE

page No. 6

2) Numpy

- Numpy is a python library adding support for large, multi dimensional arrays and matrices along with large collections of high level mathematical functions to operate on these arrays.
- The ancestor of numpy, numeric was originally created by several other developers.
- In 2005, Travis Oliphant created numpy by incorporating features of numeric arrays into numeric with extensive modifications.
- Numpy is open-source software and has many contributors. It is a numfocus fiscally sponsored project.

Conclusion

Data wrangling is one of the most important technique to turn raw data into useful asset where pandas perform an very important role of converting the raw data into productive data set that can be utilized for different purpose. Hence we performed different operation using these libraries.

Observations.

- The Melbourne Lassing market contained all the null values for 'car', 'Building Area', 'YearBuil', 'CouncilArea', columns for dealing with it.
- We have removed the rows and added it with mean values.
- Type & method columns have been ^{converted} from categorical variable to quantitative variable.