

Assi 1 - Data Wrangling - I

csv file | Dataset - titanic_train.csv

Required libraries -

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

Functions used -

`df = pd.read_csv` - To load files
`df.head()` - It's showing top 5 results.
`df.tail()` - It's showing bottom 5 results.
`df.isnull().sum()` - Calculating the null values.
`sns.heatmap`
`df.describe()` - Get some initial statistics.
`df.info()` - Getting some info about dataset.
`df.dtypes` - Finding data types.
`df.shape` - finding Dimensions of the data frame.
(no. of rows & no. of columns)
`df.drop_duplicates()`
`df.columns`

Deal Questions with Answers.

① What is mean by Data Preprocessing?

→ ① checking of missing values using pandas `isnull()` function.

② By using `describe()` function to get some initial statistics.

③ Provide variable description.

④ Type of Variable

⑤ checking the dimensions of data frame.

② What is mean by Data Wrangling?

→ ① Data exploration - The data is studied, analyzed & understood by visualizing representations of data.

② Dealing with missing values

③ Reshaping data - Data is manipulated according to the requirement, where new data can be added or pre-existing data can be modified.

④ Filtering data - Some data sets exist with unwanted rows & columns which are required to be removed or filtered.

③ Importance of Data Wrangling.

→ - Improve data usability.

- Converts data into compatible format.

④ What is mean by data Wrangling process?

→ Cleaning, organising & enriching raw data so that it can be used for decision making process.

⑤ Uses of pandas library.

→ Panda is open source library.

Uses-

① Data cleaning

② Data fill

③ Data normalization

④ Merges & joins

⑤ Data visualization

⑥ Statistical analysis

⑦ Data inspection

⑧ Loading & saving data

⑥ Uses of numpy library.

→ numpy is the numerical python library.

Uses -

- ① Working with arrays
- ② Working in domain of linear algebra.
- ③ Fourier transform
- ④ Working with matrices.

It is also open source library.

- ⑤ Working in vector-vector multiplication

⑦ Uses of matplotlib library.

→ matplotlib is used for data visualization & graphical plotting library (histograms, scatter plots, bar chart etc)

⑧ Uses of seaborn library.

→ seaborn library used for making statistical graphics in python.

⑨ Difference between matplotlib & seaborn library

matplotlib -

- Matplotlib is a python library used to plot various graphs with the help of additional libraries like Numpy & Pandas.
- Matplotlib creates simple graphs, including histograms, bar graphs, pie charts, scatter plots, lines & other visual representation of data.
- Mainly used to plot 2D graphs of arrays.
- It uses syntax that is relatively complicated & extensive.
eg. `matplotlib.pyplot.bar(x-axis, y-axis)` is the syntax for bar graph.

Seaborn -

- Seaborn is also a pandas library that utilizes Matplotlib, Pandas & Numpy to plot graphs.
- It is a superset of Matplotlib library.
- It has relatively simple syntax
eg. `seaborn.barplot(x-axis, y-axis)` syntax for a bar graph.
- Seaborn is more comfortable with pandas data frames.
- Seaborn prevents overlapping with the help of default themes.

⑩

How to install any library in python program

`pip install package-name`

eg `pip install seaborn`