

Practical Lab Assignment 3

Pre-processing of missing values (using mean value of each attribute class):

Problem Statement: Replace the missing values for given automobile dataset “imports-85.data” with mean value of each attribute class. (Consider no. of doors as the class attribute - 6th attribute)

Create a Menu Driven Python program.

Dataset:

<https://github.com/nyuvis/datasets/blob/master/auto/imports-85.data>

Dataset Information :

<https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.names>

Implementation:

```
#import the Libraries
import pandas as pd
import numpy as np
from google.colab import files

#import the dataset
ds = pd.read_csv("https://raw.githubusercontent.com/nyuvis/datasets/master/auto/imports-85.data")

#replace the NaN value for "?" value
ds.replace("?", np.nan, inplace=True)

# convert the object class into float object
ds['normalized-losses'] = ds['normalized-losses'].astype(float)
ds['bore'] = ds['bore'].astype(float)
ds['stroke'] = ds['stroke'].astype(float)
ds['horsepower'] = ds['horsepower'].astype(float)
ds['peak-rpm'] = ds['peak-rpm'].astype(float)
ds['price'] = ds['price'].astype(float)

#creating the menu driven program
print("=====MENU=====")
print("1. Missing Values using MEAN")
print("2. Missing Values using MEDIAN")
print("3. Missing Values using MODE")
option = int(input("Enter Choice :- "))

if option == 1:

    #taking the individual mean by taking num-of-doors as class
    mean_t_nod = ds[ds['num-of-doors'] == 'two']['normalized-losses'].mean()
    mean_f_nod = ds[ds['num-of-doors'] == 'four']['normalized-losses'].mean()

    mean_t_bore = ds[ds['num-of-doors'] == 'two']['bore'].mean()
    mean_f_bore = ds[ds['num-of-doors'] == 'four']['bore'].mean()

    mean_t_stroke = ds[ds['num-of-doors'] == 'two']['stroke'].mean()
    mean_f_stroke = ds[ds['num-of-doors'] == 'four']['stroke'].mean()

    mean_t_horsepower = ds[ds['num-of-doors'] == 'two']['horsepower'].mean()
```

```
mean_f_horsepower = ds[ds['num-of-doors'] == 'four']['horsepower'].mean()

mean_t_peak_rpm = ds[ds['num-of-doors'] == 'two']['peak-rpm'].mean()
mean_f_peak_rpm = ds[ds['num-of-doors'] == 'four']['peak-rpm'].mean()

mean_t_price = ds[ds['num-of-doors'] == 'two']['price'].mean()
mean_f_price = ds[ds['num-of-doors'] == 'four']['price'].mean()

#filling all the missing value by mean value of each attribute class
ds['normalized-losses'][ds['num-of-doors'] == 'two'] = ds['normalized-
losses'][ds['num-of-doors'] == 'two'].fillna(mean_t_nod)
ds['normalized-losses'][ds['num-of-doors'] == 'four'] = ds['normalized-
losses'][ds['num-of-doors'] == 'four'].fillna(mean_f_nod)

ds['bore'][ds['num-of-doors'] == 'two'] = ds['bore'][ds['num-of-
doors'] == 'two'].fillna(mean_t_bore)
ds['bore'][ds['num-of-doors'] == 'four'] = ds['bore'][ds['num-of-
doors'] == 'four'].fillna(mean_f_bore)

ds['stroke'][ds['num-of-doors'] == 'two'] = ds['stroke'][ds['num-of-
doors'] == 'two'].fillna(mean_t_stroke)
ds['stroke'][ds['num-of-doors'] == 'four'] = ds['stroke'][ds['num-of-
doors'] == 'four'].fillna(mean_f_stroke)

ds['horsepower'][ds['num-of-doors'] == 'two'] = ds['horsepower'][ds['num-of-
doors'] == 'two'].fillna(mean_t_horsepower)
ds['horsepower'][ds['num-of-doors'] == 'four'] = ds['horsepower'][ds['num-of-
doors'] == 'four'].fillna(mean_f_horsepower)

ds['peak-rpm'][ds['num-of-doors'] == 'two'] = ds['peak-rpm'][ds['num-of-
doors'] == 'two'].fillna(mean_t_peak_rpm)
ds['peak-rpm'][ds['num-of-doors'] == 'four'] = ds['peak-rpm'][ds['num-of-
doors'] == 'four'].fillna(mean_f_peak_rpm)

ds['price'][ds['num-of-doors'] == 'two'] = ds['price'][ds['num-of-
doors'] == 'two'].fillna(mean_t_price)
ds['price'][ds['num-of-doors'] == 'four'] = ds['price'][ds['num-of-
doors'] == 'four'].fillna(mean_f_price)

#this converts the dataset into CSV files and then download it.
ds.to_csv('mean.csv')
files.download('mean.csv')
```

```
elif option == 2:

    #taking the individual median by taking num-of-doors as class
    median_t_nod = ds[ds['num-of-doors'] == 'two']['normalized-losses'].median()
    median_f_nod = ds[ds['num-of-doors'] == 'four']['normalized-losses'].median()

    median_t_bore = ds[ds['num-of-doors'] == 'two']['bore'].median()
    median_f_bore = ds[ds['num-of-doors'] == 'four']['bore'].median()

    median_t_stroke = ds[ds['num-of-doors'] == 'two']['stroke'].median()
    median_f_stroke = ds[ds['num-of-doors'] == 'four']['stroke'].median()

    median_t_horsepower = ds[ds['num-of-doors'] == 'two']['horsepower'].median()
    median_f_horsepower = ds[ds['num-of-doors'] == 'four']['horsepower'].median()

    median_t_peak_rpm = ds[ds['num-of-doors'] == 'two']['peak-rpm'].median()
    median_f_peak_rpm = ds[ds['num-of-doors'] == 'four']['peak-rpm'].median()

    median_t_price = ds[ds['num-of-doors'] == 'two']['price'].median()
    median_f_price = ds[ds['num-of-doors'] == 'four']['price'].median()

    #filling all the missing value by median value of each attribute class
    ds['normalized-losses'][ds['num-of-doors'] == 'two'] = ds['normalized-losses'][ds['num-of-doors'] == 'two'].fillna(median_t_nod)
    ds['normalized-losses'][ds['num-of-doors'] == 'four'] = ds['normalized-losses'][ds['num-of-doors'] == 'four'].fillna(median_f_nod)

    ds['bore'][ds['num-of-doors'] == 'two'] = ds['bore'][ds['num-of-doors'] == 'two'].fillna(median_t_bore)
    ds['bore'][ds['num-of-doors'] == 'four'] = ds['bore'][ds['num-of-doors'] == 'four'].fillna(median_f_bore)

    ds['stroke'][ds['num-of-doors'] == 'two'] = ds['stroke'][ds['num-of-doors'] == 'two'].fillna(median_t_stroke)
    ds['stroke'][ds['num-of-doors'] == 'four'] = ds['stroke'][ds['num-of-doors'] == 'four'].fillna(median_f_stroke)

    ds['horsepower'][ds['num-of-doors'] == 'two'] = ds['horsepower'][ds['num-of-doors'] == 'two'].fillna(median_t_horsepower)
    ds['horsepower'][ds['num-of-doors'] == 'four'] = ds['horsepower'][ds['num-of-doors'] == 'four'].fillna(median_f_horsepower)
```

```
ds['peak-rpm'][ds['num-of-doors'] == 'two'] = ds['peak-rpm'][ds['num-of-doors'] == 'two'].fillna(median_t_peak_rpm)
ds['peak-rpm'][ds['num-of-doors'] == 'four'] = ds['peak-rpm'][ds['num-of-doors'] == 'four'].fillna(median_f_peak_rpm)

ds['price'][ds['num-of-doors'] == 'two'] = ds['price'][ds['num-of-doors'] == 'two'].fillna(median_t_price)
ds['price'][ds['num-of-doors'] == 'four'] = ds['price'][ds['num-of-doors'] == 'four'].fillna(median_f_price)

#this converts the dataset into CSV files and then download it.
ds.to_csv('median.csv')
files.download('median.csv')

elif option == 3:

    #taking the individual mode by taking num-of-doors as class
    mode_t_nod = ds[ds['num-of-doors'] == 'two']['normalized-losses'].mode()[0]
    mode_f_nod = ds[ds['num-of-doors'] == 'four']['normalized-losses'].mode()[0]

    mode_t_bore = ds[ds['num-of-doors'] == 'two']['bore'].mode()[0]
    mode_f_bore = ds[ds['num-of-doors'] == 'four']['bore'].mode()[0]

    mode_t_stroke = ds[ds['num-of-doors'] == 'two']['stroke'].mode()[0]
    mode_f_stroke = ds[ds['num-of-doors'] == 'four']['stroke'].mode()[0]

    mode_t_horsepower = ds[ds['num-of-doors'] == 'two']['horsepower'].mode()[0]
    mode_f_horsepower = ds[ds['num-of-doors'] == 'four']['horsepower'].mode()[0]

    mode_t_peak_rpm = ds[ds['num-of-doors'] == 'two']['peak-rpm'].mode()[0]
    mode_f_peak_rpm = ds[ds['num-of-doors'] == 'four']['peak-rpm'].mode()[0]

    mode_t_price = ds[ds['num-of-doors'] == 'two']['price'].mode()[0]
    mode_f_price = ds[ds['num-of-doors'] == 'four']['price'].mode()[0]

    #filling all the missing value by mode value of each attribute class
    ds['normalized-losses'][ds['num-of-doors'] == 'two'] = ds['normalized-losses'][ds['num-of-doors'] == 'two'].fillna(mode_t_nod)
    ds['normalized-losses'][ds['num-of-doors'] == 'four'] = ds['normalized-losses'][ds['num-of-doors'] == 'four'].fillna(mode_f_nod)

    ds['bore'][ds['num-of-doors'] == 'two'] = ds['bore'][ds['num-of-doors'] == 'two'].fillna(mode_t_bore)
```

```
ds['bore'][ds['num-of-doors'] == 'four'] = ds['bore'][ds['num-of-doors'] == 'four'].fillna(mode_f_bore)

ds['stroke'][ds['num-of-doors'] == 'two'] = ds['stroke'][ds['num-of-doors'] == 'two'].fillna(mode_t_stroke)
ds['stroke'][ds['num-of-doors'] == 'four'] = ds['stroke'][ds['num-of-doors'] == 'four'].fillna(mode_f_stroke)

ds['horsepower'][ds['num-of-doors'] == 'two'] = ds['horsepower'][ds['num-of-doors'] == 'two'].fillna(mode_t_horsepower)
ds['horsepower'][ds['num-of-doors'] == 'four'] = ds['horsepower'][ds['num-of-doors'] == 'four'].fillna(mode_f_horsepower)

ds['peak-rpm'][ds['num-of-doors'] == 'two'] = ds['peak-rpm'][ds['num-of-doors'] == 'two'].fillna(mode_t_peak_rpm)
ds['peak-rpm'][ds['num-of-doors'] == 'four'] = ds['peak-rpm'][ds['num-of-doors'] == 'four'].fillna(mode_f_peak_rpm)

ds['price'][ds['num-of-doors'] == 'two'] = ds['price'][ds['num-of-doors'] == 'two'].fillna(mode_t_price)
ds['price'][ds['num-of-doors'] == 'four'] = ds['price'][ds['num-of-doors'] == 'four'].fillna(mode_f_price)

#this converts the dataset into CSV files and then download it.
ds.to_csv('mode.csv')
files.download('mode.csv')
```

Output:

```
=====MENU=====
1. Missing Values using MEAN
2. Missing Values using MEDIAN
3. Missing Values using MODE
Enter Choice :- 
```

For Value 1 (Mean):

Whole Processed Dataset:

https://github.com/darshanjoshi16/DataMiningPracticals/blob/main/mean_after.csv

For Value 2 (Median):

Whole Processed Dataset:

https://github.com/darshanjoshi16/DataMiningPracticals/blob/main/median_after.csv

For Value 3 (Mode):

Whole Processed Dataset:

https://github.com/darshanjoshi16/DataMiningPracticals/blob/main/mode_after.csv