

Amrita School of Engineering Chennai – 601 103, Tamil Nadu, India.

AMRITA SCHOOL OF ENGINEERING

Degree of Bachelor of Technology in **COMPUTER SCIENCE & ENGINEERING**

COGNIZANCE TASK-3 [Python - Medicore Lvl]

Submitted by:

GUNDE ELIYAJER
CSE-A
CH.EN.U4CSE20027

Question-1:

File Handling is one of the basic important task when it comes to building machine learning models or neural networks. Building a good model always starts with finding datasets and processing it, for which, file handling acts as a stepping stone.

Write a python program that reads the contents from the given file 'onelinefile.txt'. The file contains a single line which is of the format (int)(string)(float)(string) repeatedly.

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

ds = pd.read_csv("filename2.csv")
ds.head()

1Aaa3.5Maths2Bbb4.2Physics3Ccc7.62Chemistry4Ddd9.55Biology5Eee4.0Social6Fff7.6English76gg3.111Maths8Hhh9.99Physics9Iii1.23Civics
```

```
def split_file(filename2):
header=0
header line=""
file_count=0
for line in filename2:
    line=line.rstrip()
    a=line.split()
    if header==0:
        header line=line
        header+=1
    else:
        if a[-1] not in 1:
            1.append(a[-1])
            file count+=1
            if file count>1:
                dest.close()
            else:
                pass
            dest=open(a[-1], 'a')
            dest.write(header_line+"\n"+line+"\n")
        else:
            dest.write(line+"\n")
source.close()
dest.close()
```

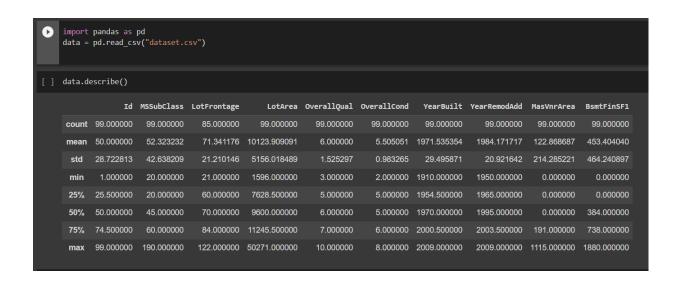
OUTPUT:-

```
1,Aaa,3.5,Maths
2,Bbb,4.2,Physics
3,Ccc,7.62,Chemistry
```

Question-2

Data formatting

Python libraries represent missing numbers as nan which is short for "not a number". Most libraries (including scikit-learn) will give you an error if you try to build a model using data with missing values. One of the common solution to get around this issue is to impute or fill in the missing value with a number or value of same format. From the given dataset, find the missing values(Nan/NA/-/Nil) and change those values into an appropriate number.



```
[ ] data.info()
    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 99 entries, 0 to 98
    Data columns (total 36 columns):
                        Non-Null Count Dtype
         Column
     #
         Td
                        99 non-null
                                        int64
     0
         MSSubClass
                        99 non-null
                                        int64
     1
     2
         MSZoning
                        99 non-null
                                        object
                        85 non-null
                                        float64
     3
         LotFrontage
         LotArea
                                        int64
     4
                        99 non-null
     5
         Street
                        99 non-null
                                        object
         Alley
                        6 non-null
                                        object
     6
     7
         LotShape
                        99 non-null
                                        object
         LandContour
                        99 non-null
     8
                                        object
     9
         Utilities
                        99 non-null
                                        object
         LotConfig
                        99 non-null
                                        object
     10
     11
         LandSlope
                        99 non-null
                                        object
     12
         Neighborhood
                        99 non-null
                                        object
         Condition1
     13
                        99 non-null
                                        object
         Condition2
                        99 non-null
     14
                                        object
     15
         BldgType
                        99 non-null
                                        object
         HouseStyle
     16
                        99 non-null
                                        object
         OverallQual
     17
                        99 non-null
                                        int64
         OverallCond
     18
                        99 non-null
                                        int64
                        99 non-null
         YearBuilt
     19
                                        int64
     20
         YearRemodAdd
                        99 non-null
                                        int64
     21
         RoofStyle
                        99 non-null
                                        object
         RoofMat1
     22
                        99 non-null
                                        object
     23
         Exterior1st
                        99 non-null
                                        object
         Exterior2nd
     24
                        99 non-null
                                        object
```

```
[ ] data.duplicated(keep="first")
          False
    0
    1
          False
          False
          False
    4
          False
    94
          False
         False
    95
          False
    96
    97
         False
         False
    98
    Length: 99, dtype: bool
[ ] data_without_missing_values = data.dropna(axis=1)
[ ] cols_with_missing = [col for col in data.columns
                                     if data[col].isnull().any()]
    redued original data = data.drop(cols with missing, axis=1)
    reduced test data = data.drop(cols with missing, axis=1)
```

	one_hot_encoded_training_predictors = pd.get_dummies(train_predictors) one_hot_encoded_training_predictors[:5]												
	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	BsmtFinSF1	MSZoning_C (all)	MSZoning_FV	MSZoning_RL		ExterQual_Gd Ex	
	8450			2003	2003	196	706						
1	9600		8	1976	1976		978						
2	11250			2001	2002	162	486						
3	9550			1915	1970		216						
4	14260			2000	2000	350	655						
5 ro	ws × 72 co	olumns											
4.	+												

```
cols_with_missing = [col for col in X_train.columns if X_train[col].isnull().any()]
reduced_X_train = X_train.drop(cols_with_missing, axis=1)
reduced_X_test = X_test.drop(cols_with_missing, axis=1)
print("Mean Absolute Error after dropping columns with missing values:")
print(score_dataset(reduced_X_train, reduced_X_test, y_train, y_test))
Mean Absolute Error after dropping columns with missing values:
```

26.65

Question-3

Read the file 'about.txt' and find the words with atleast 6 letters and the most frequently used word.

Contents of the file 'about.txt':

Python has tools for almost every aspect of scientific computing. The Bank of America uses Python to crunch its financial data and Facebook looks upon the Python library Pandas for its data analysis. While there are many libraries available to perform data analysis in Python, here are a few: NumPy, SciPy, Pandas and Matplotlib.



```
[1] count = 0;
    word = "";
    maxCount = 0;
    words = [];
    file = open("about.csv", "r")
    for line in file:
        string = line.lower().replace(',','').replace('.','').split(" ");
        for s in string:
            if len(s)>=6: #len greater than 6
                words.append(s);
    for i in range(0, len(words)):
        count = 1;
        for j in range(i+1, len(words)):
            if(words[i] == words[j]):
                count = count + 1;
        if(count > maxCount):
            maxCount = count; #frequency
            word = words[i];
    print("Most repeated word: " + word);
    print("Frequency: " ,maxCount);
    file.close();
```

Output:-

```
Most repeated word: python
Frequency: 4
```

-----THE END-----