Documentation-

The original code has only file_manager.py, I have added the following python scripts to address issues and added new operations

```
file_converter.py
file_copy_move_rm.py
file_parquet_reader.py
file_reader.py
file_writer.py
file_xlsx.py
gz_zip_reader.py
```

Testing Results:

Under Testing results, the following open issues in GitHub are resolved-

1. StopIteration Exception while reading a JSON or Parquet file

This has been handled by adding an exception, which will print 'End of the file reached' at EOF. Here is the link to the code-

https://github.com/ankan-mazumdar/file-manager/blob/Bg-Data/file_manager/file_reader.py

```
class Reader:
    def __init__(self, file_handler):
        self.file_handler = file_handler

    def read(self):
        while True:
        line = self.file_handler.readline()
        if not line:
            print("End of file reached")
            break
        yield line
```

```
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\env> & C:/Python/python.exe "C:\Users\Ankan Mazumdar\Downloads\Project_BigData CSP_554\file-manager-master\file_manager\file_reader.py"

Enter the file path: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\zip_test.json
{'col1': 1, 'col2': 'A'}
{'col1': 2, 'col2': 'B'}
{'col1': 3, 'col2': 'C'}

End of file reached

S.:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file_manager-master\data\zip_test.json
```

2. Parquet read file failed due to unknown encoding

For handling the encoding-decoding failure error, I initially tried setting the encoding standard value as 'utf-8', 'Latin-1' etc., however, none of them fixed the issue. Finally, it was a simple pandas read_parquet function that helped here. Please find the code link-

https://github.com/ankan-mazumdar/file-manager/blob/Bg-Data/file manager/file parguet reader.py

Please find the test results screenshots-

```
BigDataCSP_554\file-manager-master\data\employees.parquet^C ile-manager-master\data>
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data> & C:/Python/python.exe "C:\Users\Ankan Maz
Enter the path of the Parquet file: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\employees.
    EMPLOYEE_ID FIRST_NAME LAST_NAME
                                              EMAIL PHONE_NUMBER ...
                                                                            JOB_ID_SALARY_COMMISSION_PCT_MANAGER_ID_DEPARTMENT_
ID
            198
                      Donald
                                OConnell DOCONNEL 650.507.9833 ...
                                                                           SH CLERK 2600
                                                                                                                   124
50
            199
                     Douglas
                                    Grant
                                             DGRANT 650.507.9844 ...
                                                                           SH CLERK 2600
                                                                                                                   124
50
            200
                    Jennifer
                                   Whalen JWHALEN 515.123.4444 ...
                                                                            AD_ASST 4400
                                                                                                                   101
                     Michael
                              Hartstein MHARTSTE 515.123.5555 ...
                                                                            MK_MAN 13000
 PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\env> & C:/Python/python.exe "C:\Users\Ankan Mazumdar\Downloads\Project_BigData
 Enter the path of the Parquet file: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\sample_data.parquet
     name age
Alice 30
  Charlie
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\env>
```

PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\env> & C:/Python/python.exe "C:\Users\Ankan Mazumdar\Downloads\Project_BigData CSP_554\file-manager-master\file_manager\file_parquet_reader.py"										
Enter the path of the Parquet file: C:\Users\Ankan Max										
	url_length	hostname_length	count-	count@		count-www	count-digits	count-letters	count_dir	Lab
el https://www.google.com	22	14	9			1		17		
0 https://www.google.com 0	22	14	О	в		1	0	1/	0	
1 https://www.webfx.com/digital-marketing/	37	20	0	а		0	9	30	3	
0		20	v	· ·		Ü	· ·	50		
2 https://www.facebook.com	24	16	0	0		1	0	19	9	
0										
3 https://www.baidu.com	21	13	0	0		1	0	16	0	
0										
4 https://www.wikipedia.org	25	17	0	0		1	0	20	0	
0										
5 http://faboleena.com/js/infortis/jquery/plugin	159	13	0	0		0	21	118	12	
1	447	43					20	400	40	
6 http://faboleena.com/js/infortis/jquery/plugin	147	13	0	О		0	20	109	12	
1 7 https://www.google.co.in	24	16	0	a		1	9	18	0	
0	24	10					U	10	· ·	
8 https://www.qq.com	18	10	0	9		1	9	13	0	
0										
9 https://www.amazon.com	22	14	0	0		1	0	17	0	
0										
[10 rows x 11 columns]			_							
PS C:\Users\Ankan Mazumdar\Downloads\Project BigDataCSP 554\file-manager-master\env>										

3. The gzip file read operation Error: Raised OSError on reading a gzip file due to an invalid argument (Invalid characters in the file path).

The error occurs as the encoding can't recognize the characters, on researching, I found that it is Latin-1 characters rather than utf-8, hence I kept encoding ='Latin-1' in the solution code that worked perfectly. Here's the link to the code-

https://github.com/ankan-mazumdar/file-manager/blob/Bg-Data/file manager/gz zip reader.pv

Below is the testing result-

```
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\env> & C:\Python/python.exe "C:\Users\Ankan Mazumdar\Downloads\Project_BigData CSP_554\file-manager-master\file_manager\gz_zip_reader.py"

Choose an action (read, write, convert): read

Enter the file name (include the extension): C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\zip_test.csv.gz

coll,col2

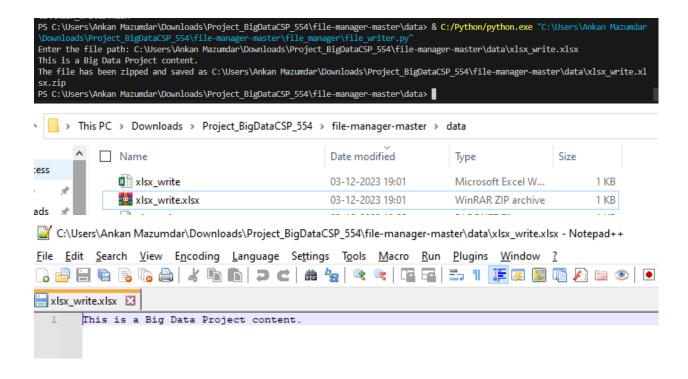
1,A

2,B

3,C
```

4. Next, as proposed I have added the functionality for read-write operations for xlsx and zip files-

```
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\env> & C:/Python/python.exe "C:\Users\Ankan Mazumdar
\Downloads\Project BigDataCSP 554\file-manager-master\file manager\file xlsx.py
Enter the file name (include the extension): C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\data
Choose an action (read, write, convert): read
Content of XLSX file 'C:\Users\Ankan Mazumdar\Downloads\Project BigDataCSP 554\file-manager-master\data\data2.xlsx':
                                                url url_length hostname_length ... count-letters count_dir label
                             https://www.google.com
                                                                                                              a
            https://www.webfx.com/digital-marketing/
                                                             37
                                                                                                  30
                           https://www.facebook.com
                                                             24
                                                                                                  19
                              https://www.baidu.com
                                                             21
                                                                                                  16
                                                                                                              0
                                                                                                                     0
                          https://www.wikipedia.org
                                                             25
                                                                                                  20
                                                                                                              0
                                                                                                                     0
  http://faboleena.com/js/infortis/jquery/plugin...
                                                             159
                                                                                                 118
                                                                                                             12
  http://faboleena.com/js/infortis/jquery/plugin...
                                                            147
                                                                                                 109
                                                                                                             12
                           https://www.google.co.in
                                                                                                  18
                                                                                                              a
                                                             24
                                 https://www.qq.com
                                                             18
                                                                              10
                                                                                                  13
                                                                                                              0
                                                                                                                     0
                             https://www.amazon.com
                                                             22
[10 rows x 11 columns]
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\env>
```



5. Converter- This script can do all-round file format conversion from csv to xml, json and parquet. Parquet to csv json xml and further in this fashion. Code link-

https://github.com/ankan-mazumdar/file-manager/blob/Bg-Data/file manager/file converter.py

Testing results-Csv to Parquet conversion-

```
S C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master> & C:/Python/python.exe "C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSF
594\file-manager-master\file_manager\file_converter.py
Enter the source file path: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\employees.csv
 SV file 'C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\employees.csv' converted to Parquet yees.parquet
file 'C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\employees.parquet' successfully.
BigDataCSP 554\file-manager-master\data\employees.parquet^C ile-manager-master\data>
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data> & C:/Python/python.exe "C:\Users\Ankan Maz
 umdar\Downloads\Project BigDataCSP 554\file-manager-master\file manager\file parquet reader.
Enter the path of the Parquet file: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\employees.
    EMPLOYEE_ID FIRST_NAME LAST_NAME
                                                EMAIL PHONE_NUMBER ...
                                                                                JOB_ID SALARY COMMISSION_PCT MANAGER_ID DEPARTMENT_
ID
                       Donald
                                OConnell DOCONNEL 650.507.9833 ...
                                                                              SH_CLERK 2600
50
            199
                     Douglas
                                   Grant DGRANT 650.507.9844 ...
                                                                              SH CLERK 2600
                                                                                                                       124
50
                                Whalen JWHALEN 515.123.4444 ...
            200
                     Jennifer
                                                                              AD ASST 4400
                                                                                                                       101
10
             201
                     Michael Hartstein MHARTSTE 515.123.5555 ...
                                                                               MK_MAN 13000
                                                                                                                       100
```

Json to parquet -

```
rs\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\file_manager> & C:/Python/python.exe
\\file-manager-master\file_manager\file_converter.py"
Enter the source file path: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\userdata1.json
Enter the destination Parquet file path: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\userdata1.parquet
DSON file 'C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\userdata1.json' converted to Parquet file 'C:\Users\Ankan Mazumdar\Downloads\Pro
ect_BigDataCSP_554\file-manager-master\data\userdata1.parquet' successfully.
'S C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\file_manager>
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\file manager> & C:/Python/python.exe "C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_55
4\file-manager-master\file_manager\file_parquet_reader.py"
Enter the path of the Parquet file: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\userdata1.parquet
                                                                                                                                                 salary
                            id first_name last_name
                                                                                                                                                                               title commen
                                                             ajordan@@com.com Female ... 6759521864920116 Indonesia 3/8/1971 49<u>7</u>56.53
         1454486129000 1 Amanda Jordan
                                                                                                                                                                   Internal Auditor 1F+
         1454519043000 2
                                  Albert Freeman
                                                                afreeman1@is.gd Male ...
                                                                                                                         Canada 1/16/1968 150280.17
                                                                                                                                                                      Accountant IV
         1454461771000 3
                                   Evelyn Morgan emorgan2@altervista.org Female ... 6767119071901597 Russia 2/1/1960 144972.51 Structural Engineer
         1454459781000 4
                                                               driley3@gmpg.org Female ... 3576031598965625
                                                                                                                            China 4/8/1997 90263.05 Senior Cost Accountant
                                   Denise
                                                                                         ... 5602256255204850 South Africa
         1454495459000 996
                                   Dennis
                                                           996
         1454519813000 997
                                   Gloria Hamilton ghamiltonro@rambler.ru Female ...
                                                                                                                            China 4/22/1975 83183.54 VP Product Management
                                                           nmorrisrp@ask.com ... 3553564071014997
997
         1454475740000 998
                                    Nancy
                                            Morris
                                                                                                                            Sweden 5/1/1979 NaN
                                                                                                                                                                  Junior Executive
                                    Annie Daniels adanielsrq@squidoo.com Female ... 30424803513734
                                                                                                                            China 10/9/1991 18433.85
998
         1454467292000 999
                                                                                                                                                                             Editor
         1454493138000 1000
                                                            jmeyerrr@flavors.me Female ... 374288099198540
                                                                                                                            China
                                                                                                                                                222561.13
                                    Julie
                                               Meyer
[1000 rows x 13 columns]
```

```
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\file_manager\& C:\Python/python.exe "C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_555 4\file-manager-master\file_manager\file_converter.py"

Enter the source file path: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\userdatal.parquet

Enter the destination file path: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\userdatal.parquet

Enter the destination file path: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\userdatal.json

Choose conversion type (parquet_to_csv, parquet_to_json) parquet_to_json

Parquet file 'C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\userdatal.parquet' converted to JSON file 'C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\file manager-master\file manager-
```

6. Copy move, create remove operations- As name suggests this script basically can copy , move, make, and remove files from any mentioned path. Here is the code linkhttps://github.com/ankan-mazumdar/file-manager/blob/Bg-Data/file_manager/file_copy_m_ove_rm.pv

```
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data> & C:/Python/python.exe "C:\Users\Ankan Maz
umdar\Downloads\Project_BigDataCSP_554\file-manager-master\file_manager\file_copy_move_rm.py
Choose an action (copy, move, remove, create, or exit): move
Enter the source file/folder path: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\xlsx_write.
xlsx
Enter the destination file/folder path: C:\Users\Ankan Mazumdar\Downloads\Project BigDataCSP 554\file-manager-master\xlsx write.
xlsx
File moved successfully from C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data\xlsx_write.xlsx t
o C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\xlsx_write.xlsx
Do you want to do anything else? (yes/no): no
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\data> cd ...
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master> dir
   Directory: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master
Mode
                     LastWriteTime
                                            Length Name
             03-12-2023
                              21:55
                                                   data
             18-11-2023
                             96:38
              03-12-2023
                                                    file manager
                             13:17
d----
             18-11-2023
                             19:03
                                                   test2
             28-04-2020
                                              2043 .gitignore
                             06:55
              18-11-2023
                             07:05
                                                22 abcd
              03-12-2023
                             19:55
                                            192504 data.zip
              18-11-2023
                             07:19
                                                 0 file_
                                              1527 gzip error.txt
390 README.md
              18-11-2023
                              16:51
              28-04-2020
                              06:55
              05-11-2023
                              13:10
                                                 87 requirements.txt
                                              39437 Screenshot 2023-11-18 165235.png
              18-11-2023
                              16:53
              28-04-2020
                                               677 setup.py
35 xlsx_write.xlsx
              03-12-2023
```

```
PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master> & C:/Python/python.exe "C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\file_manager\file_copy_move_rm.py"

Choose an action (copy, move, remove, create, or exit): remove

Enter the source file/folder path: C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\xlsx_write.xlsx

File C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master\xlsx_write.xlsx removed successfully

Do you want to do anything else? (yes/no): no

PS C:\Users\Ankan Mazumdar\Downloads\Project_BigDataCSP_554\file-manager-master>
```

- 7. Uploaded the scripts over EMR cluster and tested using python and Spark-.
- a. Csv file -

```
[hadoop8ip-172-31-46-19 file-manager-master]$ python file_manager/file_reader.py
Enter the file path: /home/hadoop/file-manager-master/data/data2.csv
OrderedDict([['url', 'https://www.ogo]e.com), ('urllength', '22'), ('hostname_length', '14'), ('count-', '0'), ('count-ktps', '1'), ('count-www', '1'), ('count-digits', '0'), ('count-letters', '17'), ('count-dir', '0'), ('label', '0')]
OrderedDict([['url', 'https://www.wefsebook.com/) ('urllength', '22'), ('hostname_length', '20'), ('count-', '0'), ('count-', '0')
```

b. Json file-

```
Enter the file path: /home/hadoop/file-manager-master/data/zip_test.json {'col1': 1, 'col2': 'A'} {'col1': 2, 'col2': 'B'} {'col1': 3, 'col2': 'C'} End of file reached [hadoop@ip-172-31-46-19 file_manager]$ |
```

c. Parquet file reading-

d. XIsx file-

e. Compressed file-

8. Using Spark-submit command-

Issue -Spark-submit to run Python script gets stuck and never proceeds when there is user interaction/inputs are required at run time

```
[hadoop@ip-172-31-46-19 file_manager]$ spark-submit fle_xlsx.py
Enter the file name (include the extension): /home/hadoop/file-manager-master/data/userdata1.xlsx
```

https://stackoverflow.com/questions/40910869/python-script-hangs-on-input-method-when-running-spark.

When we submit an application using spark-submit, we don't interact with Python code, but with Java one, which doesn't expect any input from stdin.

If we want to make it work we have to skip spark-submit and execute this as a Python script. We could alternatively achieve this by using pyspark or hardcode the filename in the script

9. Read the file using Pyspark

10. Writing parquet file and zipping it-

```
>>> import zipfile
>>> mode = 'w'
>>> # Ask the user to input the file path
>>> file_path = input("Enter the file path: ")
Enter the file path: /home/hadoop/file-manager-master/data/sample_test.parquet
>>> with open(file_path, "w") as file:
... file.write("This is a Big Data Project content.")
...
35
>>> with open(file_path, "r") as file:
... print(file.read())
...
This is a Big Data Project content.
>>> zip_file_path = file_path + ".zip"
>>> zip_file_zipFile(zip_file_path, mode='w').write(file_path)
>>> print(f'The file has been zipped and saved as {zip_file_path}')
The file has been zipped and saved as /home/hadoop/file-manager-master/data/sample_test.parquet.zip
```

The file is created -

```
[hadoop@ip-172-31-46-19 data]$ ls -l /home/hadoop/file-manager-master/data/sample_test.parquet.zip -rw-rw-r-- 1 hadoop hadoop 245 Dec 4 03:06 /home/hadoop/file-manager-master/data/sample_test.parquet.zip
```

11.Xlsx file in spark

```
[hadoop@ip-172-31-46-19 file_manager]$ spark-submit fle_xlsx.py

Content of XLSX file '/home/hadoop/file-manager-master/data/userdata1.xlsx':
    registration_dttm id ... title comments

0 2016-02-03 07:55:29 1 ... Internal Auditor 100.0

1 2016-02-03 17:04:03 2 ... Accountant IV NaN

2 2016-02-03 01:09:31 3 ... Structural Engineer NaN

3 2016-02-03 00:36:21 4 ... Senior Cost Accountant NaN

4 2016-02-03 05:05:31 5 ... NaN NaN

5 2016-02-03 07:22:34 6 ... Account Executive NaN

6 2016-02-03 08:33:08 7 ... Senior Financial Analyst NaN

7 2016-02-03 06:47:06 8 ... Web Developer IV NaN

8 2016-02-03 03:52:53 9 ... Software Test Engineer I 100.0

9 2016-02-03 18:29:47 10 ... Health Coach IV NaN

[10 rows x 13 columns]

23/12/04 03:31:13 INFO ShutdownHookManager: Shutdown hook called

23/12/04 03:31:13 INFO ShutdownHookManager: Deleting directory /mnt/tmp/spark-aadb23ce-44af-42c9-bc08-12bfd518cfd1
```