For hadoop.3x version

http://localhost:9870

http://localhost:8088/cluster

Place the mapper file ,reducer file and hadoop dtearming jar file in Documents,Create input foler in hadoop and place the wordcount.txt file on it.

hadoop/Documents\$ give below comments to run

hadoop@Ubuntu:~/Documents\$ hadoop jar hadoop-streaming-2.7.3.jar -input /home/hadoop/input/word_count_data.txt -output /home/hadoop/output -mapper mapper.py reducer reducer.py

To check the output folder part-oooo file is created or not hadoop@Ubuntu:~/Documents\$ hadoop fs -ls /home/hadoop/output

hadoop@Ubuntu:~/Documents\$ hadoop fs -ls /home/hadoop/output

Found 2 items

```
-rw-r--r-- 1 hadoop supergroup 0 2024-08-03 08:59 /home/hadoop/output/_SUCCESS -rw-r--r-- 1 hadoop supergroup 592 2024-08-03 08:59 /home/hadoop/output/part-00000
```

hadoop@Ubuntu:~/Documents\$ hdfs dfs -cat /user/hadoop/output/part-00000 cat: `/user/hadoop/output/part-00000': No such file or directory

Verify the output

```
hadoop@Ubuntu:~/Documents$ hdfs dfs -cat /home/hadoop/output/part-00000 2,000 1
ChatGPT 1
Did 1
Roman 2
```

Romans 1
Some 1
Sure! 1
This 1
a 3
actually 1
ancient 1 and
3
ash 1 ash,

.

```
because
buildings
      1
called 1
concrete 2 concrete, 1
concrete.
             1
construction 1 durable
1 for 1 form 2 from 1
gets
       1 harbors
has
       1 have 1
impressive
incredibly
       1 know
is
1 lime, 1
longevity
       1 made
1 many 1
mineral
mixture
       1
modern
       1 of
       3
outlasted
       1 over
       2 partly
1 reacts 1
reinforces
       1
seawater
       1
seawater,
       1
showcasing
stronger
structures?
survived
techniques
       1
       2
that
       5
the
their
      1
```

```
time. 1 to
tobermorite,
              1
used
      1
volcanic
              2
       1
was
which 1
with
years, 1 you
       1
path of hadoop input file
/home/hadoop/input/word count data.txt
path of hadoop output file /home/hadoop/output
Commands on hadoop to check the input and output file
1. List Contents of a Hadoop Directory
To list the contents of a directory in HDFS, use the -ls option with the hadoop fs command.
List the Input Folder
bash
hadoop fs -ls /user/hduser/input
List the Output Folder
```

2. View Detailed Information

hadoop fs -ls /user/hduser/output

bash

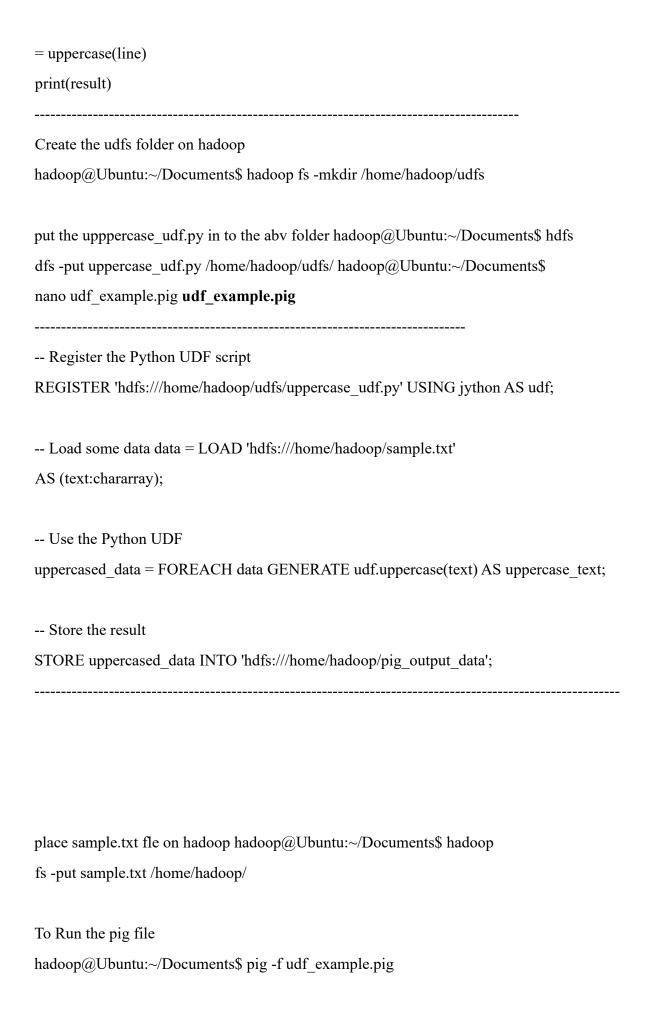
| group, size, and modification date. |
|--|
| 3. View File Contents |
| |
| To view the contents of a file, use the -cat option. For example: |
| View a File in the Input Folder |
| • |
| bash |
| |
| hadoop fs -cat /user/hduser/input/filename.txt |
| nadoop is -cat/usei/indusei/input/mename.txt |
| Replace filename.txt with the actual name of the file you want to view. |
| View a File in the Output Folder |
| |
| If your output folder contains multiple files (e.g., part-r-00000), you can view one of the files: |
| If your output folder contains multiple files (e.g., part-1-00000), you can view one of the files. |
| 1 1 |
| bash |
| |
| hadoop fs -cat /user/hduser/output/part-r-00000 |
| |
| 4. Check for Folder Existence |
| |
| To check if a folder exists in HDFS, you can use the -test command with the -d option: |
| Check if Input Folder Exists |
| |
| bash |
| |
| hadoop fs -test -d /user/hduser/input && echo "Input folder exists" echo "Input folder does not |
| exist" |
| |
| Check if Output Folder Exists |
| |
| bash |
| |

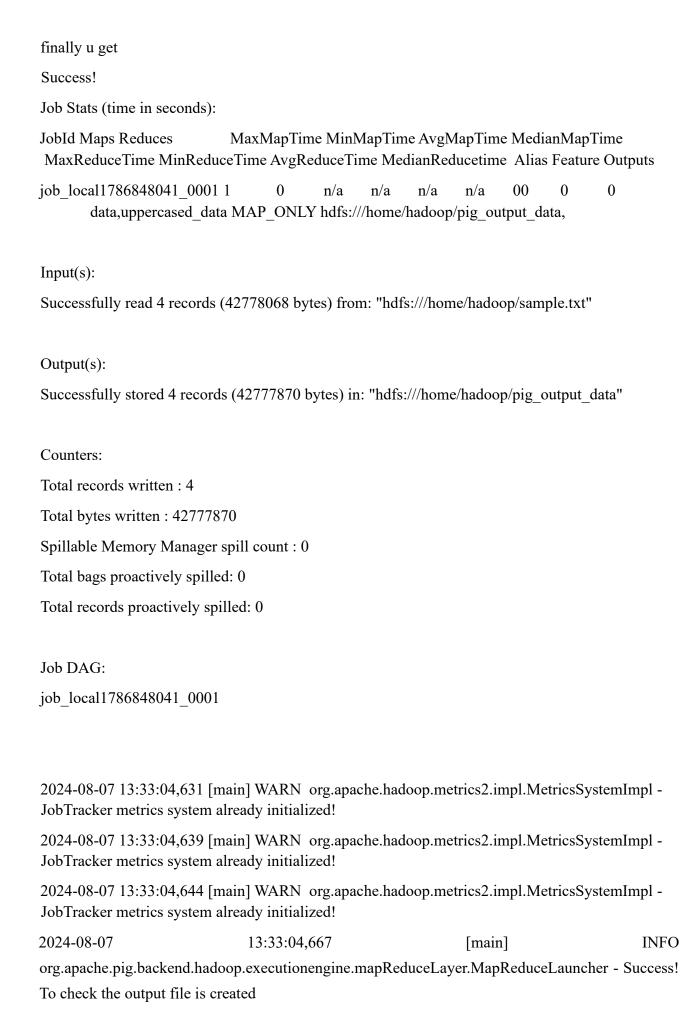
The -ls command provides a detailed listing of files and directories, including permissions, owner,

| hadoop fs -test -d /user/hduser/output && echo "Output folder exists" \parallel echo "Output folder does not exist" |
|---|
| 5. Additional Commands |
| Make a Directory: |
| bash |
| hadoop fs -mkdir /user/hduser/new_directory |
| Delete a File or Directory: |
| bash |
| hadoop fs -rm /user/hduser/input/filename.txt |
| hadoop fs -rm -r /user/hduser/output |
| Get File Status: |
| bash |
| hadoop fs -stat /user/hduser/input/filename.txt |
| Example Usage |
| Assuming you have the following setup: |
| Input Folder: /user/hduser/input |
| Output Folder: /user/hduser/output |

| You would list the contents like this: |
|--|
| bash |
| hadoop fs -ls /user/hduser/input |
| And to check the output: |
| bash |
| hadoop fs -ls /user/hduser/output |
| bash |
| hadoop fs -cat /user/hduser/output/part-r-00000 |
| PIG UDF PROGRAM |
| To check the pig program |
| hadoop@Ubuntu:~/Documents\$ nano sample.txt |
| Paste the below content to sample.txt |
| 1,John |
| 2,Jane |
| 3,Joe |
| 4,Emma |
| hadoop@Ubuntu:~/Documents\$ hadoop fs -put sample.txt /home/hadoop/piginput/ |
| hadoop@Ubuntu:~/Documents\$ nano demo_pig.pig |
| paste the below the content to demo_pig.pig |

| Load the data from HDFS | | | | | |
|---|--|--|--|--|--|
| data = LOAD '/home/hadoop/piginput/sample.txt' USING PigStorage(',') AS (id:int> | | | | | |
| | | | | | |
| Dump the data to check if it was loaded correctly | | | | | |
| DUMP data; | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| hadoop@Ubuntu:~/Documents\$ pig demo_pig.pig | | | | | |
| 2024 09 07 12 12 09 701 [] DIFO | | | | | |
| 2024-08-07 12:13:08,791 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process: 1 | | | | | |
| (1,John) | | | | | |
| (2,Jane) | | | | | |
| (3,Joe) | | | | | |
| (4,Emma) | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| By using these commands, you can manage and inspect files and directories in your Hadoop setup. | | | | | |
| upup | | | | | |
| To Run pig basic program and uf program | | | | | |
| uppercase_udf.py | | | | | |
| | | | | | |
| - def uppercase(text): return text.upper() | | | | | |
| | | | | | |
| ifname == | | | | | |
| "main": import sys | | | | | |
| for line in sys.stdin: | | | | | |
| line = line.strip() result | | | | | |





```
Found 2 items
If you need to examine the files in the output folder, use:
To view the output
hadoop@Ubuntu:~/Documents$ hdfs dfs -cat /home/hadoop/pig output data/part-m-00000
1,JOHN
2,JANE
3,JOE
4,EMMA
Create json file on bash & save as emp.json
nano emp.json; Paste the below content on it
ſ
  {"name": "John Doe", "age": 30, "department": "HR", "salary": 50000},
  {"name": "Jane Smith", "age": 25, "department": "IT", "salary": 60000},
  {"name": "Alice Johnson", "age": 35, "department": "Finance", "salary": 70000},
  {"name": "Bob Brown", "age": 28, "department": "Marketing", "salary": 55000},
{"name": "Charlie Black", "age": 45, "department": "IT", "salary": 80000}
]
Check json is readable or any error by giving install
jq by sudo apt-get install jq hadoop@Ubuntu:~$
jq.emp.json
ſ
  "name": "John Doe",
  "age": 30,
  "department": "HR",
  "salary": 50000
 },
 {
```

hadoop@Ubuntu:~/Documents\$ hdfs dfs -ls /home/hadoop/pig output data

```
"name": "Jane Smith",
  "age": 25,
  "department": "IT",
  "salary": 60000
 },
  "name": "Alice Johnson",
  "age": 35,
  "department": "Finance",
  "salary": 70000
 },
  "name": "Bob Brown",
  "age": 28,
  "department": "Marketing",
  "salary": 55000
 },
  "name": "Charlie Black",
  "age": 45,
  "department": "IT",
  "salary": 80000
]
```

bash: put the employees.json local directory to home/hadoop directory

Example

Suppose the original employees relation has the following data:

name age department salary

John Doe 30 HR 50000 Jane Smith 25 IT 60000 Alice Johnson 35 Finance 70000 Bob Brown 28 Marketing 55000 Charlie Black 45IT 80000

After executing:

pig shell: Load the json file by giving following command

grunt>-- Load the data employees = LOAD '/home/hadoop/emp.json' USING

JsonLoader('name:chararray,age:int,department:chararray,salary:float'); grunt>projected

= FOREACH employees GENERATE name, salary;

DUMP projected;

The projected relation will look like:

| name | salary |
|---------------|---------|
| John Doe | 50000 |
| Jane Smith | 60000 |
| Alice Johnson | n 70000 |
| Bob Brown | 55000 |
| Charlie Black | 80000 |

.....

Assume your employees dataset looks like this:

name age department salary John Doe 30 HR 50000 Jane Smith 25 IT 60000

Alice Johnson 35 Finance 70000

Bob Brown 28 Marketing 55000

Charlie Black 45 IT 80000

1. Aggregation

Aggregate the total salary:

pig

-- Load the data

employees = LOAD '/home/hadoop/employees.json' USING JsonLoader('name:chararray,age:int,department:chararray,salary:float');

```
-- Aggregate: Calculate the total salary
total salary = FOREACH (GROUP employees ALL) GENERATE SUM(employees.salary) AS
total salary;
DUMP total salary;
Output:
scss
(315000.0)
2. Skip
Skip the first 2 records:
pig
-- Load the data
employees = LOAD '/home/hadoop/employees.json' USING
JsonLoader('name:chararray,age:int,department:chararray,salary:float');
-- Skip the first 2 records skipped employees = LIMIT employees 1000000; -- Use
LIMIT to handle skipping
DUMP skipped employees;
Output:
             department
                           salary
name age
                    Finance
Alice Johnson 35
                                   70000
Bob Brown
                    Marketing
                                   55000
```

Charlie Black 45

IT

80000

Note: The LIMIT command should be used with an appropriate number, as Pig does not directly support skipping a specific number of records. 3. Limit Limit the results to the top 3 records: pig -- Load the data employees = LOAD '/home/hadoop/employees.json' USING JsonLoader('name:chararray,age:int,department:chararray,salary:float'); -- Limit: Get the top 3 highest earners top 3 employees = LIMIT employees 3; DUMP top_3_employees; Output: name age department salary Charlie Black 45 IT 80000 Alice Johnson 35 Finance 70000 Jane Smith 25 IT 60000 4. Count Count the number of employees: pig -- Load the data

-- Count the number of employees

employees = LOAD '/home/hadoop/employees.json' USING

JsonLoader('name:chararray,age:int,department:chararray,salary:float');

| <pre>employee_count = FOREACH (GROUP employees ALL) GENERATE COUNT(employees) AS total_count;</pre> | | | | | | | |
|---|---------|----------|-----------|---|--|--|--|
| DUMP emplo | yee_co | unt; | | | | | |
| Output: | | | | | | | |
| scss | | | | | | | |
| (5) | | | | | | | |
| 5. Remove | | | | | | | |
| Remove employees from a specific department, e.g., "IT": | | | | | | | |
| pig | | | | | | | |
| | LOAD ' | | _ | employees.json' USING department:chararray,salary:float'); | | | |
| Remove em | ployees | s from t | he 'IT' d | lepartment | | | |
| | | | | byees BY department != 'IT'; | | | |
| DUMP filtered | d_empl | oyees; | | | | | |
| Output: | | | | | | | |
| name age | departi | ment | salary | | | | |
| John Doe | 30 | HR | 50000 | | | | |
| Alice Johnson | 35 | Financ | e | 70000 | | | |
| Bob Brown | 28 | Market | ting | 55000 | | | |

import Json file and do projetion, aggregation, limit, count, skip and remove using python and hdfs

Steps to be followed:

Install pandas and hdfs using pip.

- Optionally install pyarrow or hdfs3 if needed based on your specific requirements.
- Verify the installation to ensure everything is set up correctly.

| Required Packages |
|--|
| pandas: Purpose: Provides data structures and functions to efficiently manipulate and analyze data. Installation: Use pip to install pandas. |
| bash |
| pip install pandas |
| hdfs: |
| Purpose: Provides a Python interface to interact with HDFS. Installation: Use pip to install hdfs. |
| bash |
| pip install hdfs |
| Additional Considerations |
| While the script should work with just the above packages, here are some additional considerations: |
| pyarrow (Optional but useful): |

| datasets or different file formats, pyarrow can be useful. |
|---|
| Installation: Use pip to install pyarrow. |
| bash |
| pip install pyarrow |
| hdfs3 (Alternative to hdfs): |
| Purpose: Another Python library for interacting with HDFS. It's an alternative to the hdfs package and might be preferred in some scenarios. Installation: Use pip to install hdfs3. |
| bash |
| pip install hdfs3 |
| Verifying Package Installation |
| After installing the required packages, you can verify that they are correctly installed and accessible in your Python environment: |
| python |
| import pandas as pd from hdfs import InsecureClient |
| # Check pandas version print("Pandas version:", pdversion) |

Purpose: If you're working with Apache Arrow or need additional features for handling large

```
# Test HDFS client connection client =
InsecureClient('http://localhost:9870', user='hadoop')
print("HDFS status:", client.status('/'))
If you run this script and see the version of pandas and a status message from HDFS without any
errors, the packages are installed correctly.
Create process data.py file
from hdfs import
InsecureClient import pandas
as pd import json
# Connect to HDFS hdfs client =
InsecureClient('http://localhost:9870', user='hdfs')
# Read JSON data from HDFS try: with
hdfs client.read('/home/hadoop/emp.json', encoding='utf-8') as reader:
    json data = reader.read() # Read the raw data as a
string
           if not json data.strip(): # Check if data is empty
raise ValueError("The JSON file is empty.")
     print(f"Raw JSON Data: {json data[:1000]}") # Print first 1000 characters for
debugging
                data = json.loads(json data) # Load the JSON data except
json.JSONDecodeError as e: print(f"JSON Decode Error: {e}") exit(1) except Exception
as e:
  print(f"Error reading or parsing JSON data: {e}")
exit(1)
# Convert JSON data to DataFrame
try:
  df = pd.DataFrame(data)
except ValueError as e:
```

```
print(f"Error converting JSON data to DataFrame: {e}")
exit(1)
# Projection: Select only 'name' and 'salary' columns
projected df = df[['name', 'salary']]
# Aggregation: Calculate total salary
total salary = df['salary'].sum()
# Count: Number of employees earning more than 50000
high earners count = df[df['salary'] > 50000].shape[0]
# Limit: Get the top 5 highest earners
top 5 earners = df.nlargest(5, 'salary')
# Skip: Skip the first 2 employees
skipped df = df.iloc[2:]
# Remove: Remove employees from a specific department
filtered df = df[df['department'] != 'IT']
# Save the filtered result back to HDFS
filtered json = filtered df.to json(orient='records')
try:
  with hdfs_client.write('/home/hadoop/filtered_employees.json', encoding='utf-8',
overwrite=True) as writer:
     writer.write(filtered json)
  print("Filtered JSON file saved successfully.")
except Exception as e:
  print(f"Error saving filtered JSON data: {e}")
exit(1)
```

```
# Print results
print(f"Projection: Select only name and salary columns")
print(f"{projected df}")
print(f"Aggregation: Calculate total salary")
print(f"Total Salary: {total salary}")
print(f"\n")
print(f"# Count: Number of employees earning more than 50000")
print(f"Number of High Earners (>50000):
{high_earners_count}") print(f"\n") print(f"limit Top 5 highest
salary")
print(f"Top 5 Earners: \n{top_5_earners}")
print(f"\n")
print(f"Skipped DataFrame (First 2 rows skipped): \n{skipped df}")
print(f'' \setminus n'')
print(f"Filtered DataFrame (Sales department removed): \n{filtered df}")
run the file by bash: python3
process data.py
```

output

```
o@ubuntu:~/Ex6$ hdfs dfs -chmod 777 /ex6
o@ubuntu:~/Ex6$ python3 process_data.py
 Raw JSON Data: [
 Raw JSON Data: [
{"name": "John Doe", "age": 30, "department": "HR", "salary": 50000},
{"name": "Jane Smith", "age": 25, "department": "II", "salary": 60000},
{"name": "Alice Johnson", "age": 35, "department": "Finance", "salary": 70000},
{"name": "Bob Brown", "age": 28, "department": "Marketing", "salary": 55000},
{"name": "Charlie Black", "age": 45, "department": "II", "salary": 80000}
Filtered JSON file saved successfully.
Projection: Select only name and salary columns
name salary
0 John Doe 50000
1 Jane Smith 60000
2 Alice Johnson 70000
3 Bob Brown 55000
4 Charlie Black 80000
Angregation: Calculate total salary
Aggregation: Calculate total salary
Total Salary: 315000
# Count: Number of employees earning more than 50000 Number of High Earners (>50000): 4
Limit: Top 5 highest salary
Top 5 Earners:
                                                age department salary
45 IT 80000
                             name
        Charlie Black
        Alice Johnson
Jane Smith
Bob Brown
                                                                 Finance
IT
                                                                                             70000
60000
2
1
3
0
                                                25 IT
28 Marketing
30 HR
                      John Doe
                                                                                              50000
```

```
Limit: Top 5 highest salary
Top 5 Earners:

a mame age department salary
4 Charlie Black 45 IT 80000
2 Alice Johnson 35 Finance 70000
1 Jane Smith 25 IT 60000
3 Bob Brown 28 Marketing 55000
0 John Doe 30 HR 50000

Skipped DataFrame (First 2 rows skipped):

name age department salary
2 Alice Johnson 35 Finance 70000
3 Bob Brown 28 Marketing 55000
4 Charlie Black 45 IT 80000

Filtered DataFrame (IT department removed):

name age department salary
0 John Doe 30 HR 50000
2 Alice Johnson 35 Finance 70000
3 Bob Brown 28 Marketing 55000
4 Alice Johnson 35 Finance 70000
3 Bob Brown 28 Marketing 55000
4 Alice Johnson 35 Finance 70000
3 Bob Brown 28 Marketing 55000
```