# Assignment of MRjob and Pig script

Assignment:

**ASSIGNMENT using MRJOB AND PIGLATIN SCRIPTS**

Create two data files with names and datasets as shown below and then implement the following queries using mrjob package in python for map reduce programming and using pig also for the same.Compare the two processing types.

File consisting of employee details (id , name, salary and rating as fields) and expenses details with fields(id and expenses)

**employee.txt**

```
101,Abhay,20000,1
102,Shiv,10000,2
103,Aarav,11000,3
104,Anubhav,5000,4
105,Palash,2500,5
106,Aman,25000,1
107,Sahil,17500,2
108,Ram,14000,3
109,Karan,1000,4
110,Priya,2000,5
111,Tushar,500,1
112,Ajay,5000,2
113,Jay,1000,1
114,Maddy,2000,2
```

**expenses.txt**

```
101     200
102     100
110     400
114     200
119     200
105     100
101     100
104     300
102     400
```

1. Top 5 employees (employee id and employee name) with highest rating. (In case two employees have same rating, employee with name coming first in dictionary should get preference)
2. Top 3 employees (employee id and employee name) with highest salary, whose employee id is an odd number. (In case two employees have same salary, employee with name coming first in dictionary should get preference)
3. Employee (employee id and employee name) with maximum expense (In case two employees have same expense, employee with name coming first in dictionary should get preference)
4. List of employees (employee id and employee name) having entries in expenses.txt.
5. List of employees (employee id and employee name) having no entry in expenses.txt.

By: Chayan Gulati 18csu054

Q1.

```
from mrjob.job import MRJob
from mrjob.step import MRStep
import re

word = re.compile(r"[\w*]+")
class MRHighestRating(MRJob):
    def steps(self):
        return [
            MRStep(mapper=self.mapper,
                    reducer=self.reducer),
            MRStep(reducer = self.reduce1)
        ]

    def mapper(self,key, line):
        (id, name, salary, rating) = line.split(',')
        yield [int(id), name], int(rating)

    def reducer(self, id, rating):
        yield None, (rating, id)

    def reduce1(self, _, pair):
        for v,k in sorted(pair, reverse=True)[:5]:
            yield v, k


if __name__ == '__main__':
    MRHighestRating.run()
```

o/p:



```
(base) C:\Users\HP\Desktop\Assignment_1>python one.py employee.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory C:\Users\HP\AppData\Local\Temp\one.HP.20201010.092214.215512
Running step 1 of 2...
Running step 2 of 2...
job output is in C:\Users\HP\AppData\Local\Temp\one.HP.20201010.092214.215512\output
Streaming final output from C:\Users\HP\AppData\Local\Temp\one.HP.20201010.092214.215512\output...
["5"]    ["105","Palash"]
["5"]    ["110","Priya"]
["4"]    ["104","Anubhav"]
["4"]    ["109","Karan"]
["3"]    ["103","Aarav"]
Removing temp directory C:\Users\HP\AppData\Local\Temp\one.HP.20201010.092214.215512...
```
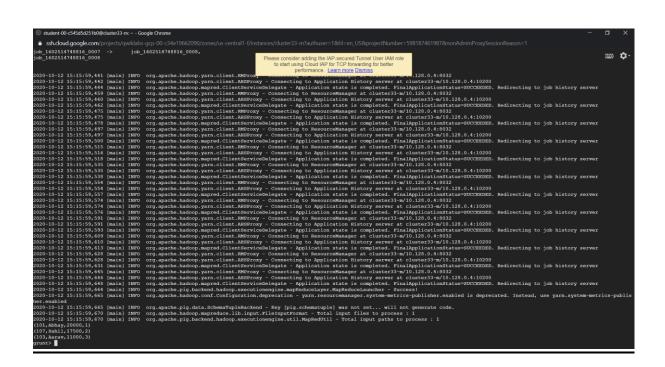
Using Pig:

```
2020-10-12 15:11:59,458 [main] INFO  org.apache.hadoop.yarn.client.RMProx                          128.0.4:8032
2020-10-12 15:11:59,458 [main] INFO  org.apache.hadoop.yarn.client.AHSPro         Please consider adding the IAP-secured Tunnel User IAM role        uster33-m/10.128.0.4:10200
2020-10-12 15:11:59,461 [main] INFO  org.apache.hadoop.mapred.ClientServi          to start using Cloud IAP for TCP forwarding for better           pplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,478 [main] INFO  org.apache.hadoop.yarn.client.RMProxy                   performance. Learn more Dismiss                   to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,479 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,481 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,498 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,498 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,515 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,540 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,541 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,546 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,573 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,573 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,576 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,595 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,595 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,598 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,615 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,616 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,619 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,647 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,648 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,650 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,667 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,667 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,670 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,692 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,693 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,696 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,713 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,713 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,716 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,742 [main] INFO  org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster33-m/10.128.0.4:8032
2020-10-12 15:11:59,743 [main] INFO  org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster33-m/10.128.0.4:10200
2020-10-12 15:11:59,750 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-10-12 15:11:59,774 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2020-10-12 15:11:59,779 [main] INFO  org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2020-10-12 15:11:59,779 [main] INFO  org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2020-10-12 15:11:59,784 [main] INFO  org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2020-10-12 15:11:59,785 [main] INFO  org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(110,Priya)
(105,Palash)
(109,Karan)
(104,Anubhav)
(108,Ram)
grunt>
```

## Q2.

emphighsalary - Notepad

File   Edit   Format   View   Help

```python
from mrjob.job import MRJob
from mrjob.step import MRStep
import re

WORD_REGEXP = re.compile(r"[\w*]+")
class MRHighestSalary(MRJob):
    def steps(self):
        return [
            MRStep(mapper=self.mapper,
                    reducer=self.reducer),
            MRStep(reducer = self.reduce1)
        ]


    def mapper(self,key, line):
        (id, name, salary, rating) = line.split(',')
        yield [int(id), name], int(salary)


    def reducer(self, id, salary):
        yield None, (max(salary), id)


    def reduce1(self, _, pair):
        for k, v in sorted(pair, reverse=True)[:5]:
            if(int(v[0])%2 != 0):
                yield k, v


if __name__ == '__main__':
    MRHighestSalary.run()
```

Ln 1, Col 1          100%   Windows (CRLF)   UTF-8

## o/p:

Command Prompt

```
Microsoft Windows [Version 10.0.18362.1082]
(c) 2019 Microsoft Corporation. All rights reserved.

C:\Users\HP>python C:\Users\HP\Desktop\ML\emphighsalary.py C:\Users\HP\Desktop\ML\employee.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory C:\Users\HP\AppData\Local\Temp\emphighsalary.HP.20201012.170426.506518
Running step 1 of 2...
Running step 2 of 2...
job output is in C:\Users\HP\AppData\Local\Temp\emphighsalary.HP.20201012.170426.506518\output
Streaming final output from C:\Users\HP\AppData\Local\Temp\emphighsalary.HP.20201012.170426.506518\output...
20000   [101, "Abhay"]
17500   [107, "Sahil"]
11000   [103, "Aarav"]
Removing temp directory C:\Users\HP\AppData\Local\Temp\emphighsalary.HP.20201012.170426.506518...

C:\Users\HP>
```

## Using Pig:

## Q3:

```
empmaxexp - Notepad
File Edit Format View Help

from mrjob.step import MRStep
from mrjob.job import MRJob

class MaxExpense(MRJob):

    def configure_args(self):
        super(MaxExpense, self).configure_args()
        self.add_file_arg('--items', help='Path to expenses.txt')

    def steps(self):
        return [
            MRStep(mapper=self.mapper,
                    reducer_init=self.reducer_init,
                    reducer=self.reducer),
            MRStep(reducer=self.reducer1)
        ]

    def mapper(self, key, value):
        (id, name, salary, rating) = value.split(',')
        yield [id,name], 1

    def reducer_init(self):
        self.expense = {}

        with open("expenses.txt") as f:
            for line in f:
                fields = line.split('\t')
                self.expense[fields[0]] = fields[1]

    def reducer(self, key, values):
        if key[0] in self.expense:
            yield None, (self.expense[key[0]], key)

    def reducer1(self, key, values):
        for count, key in sorted(values, reverse=True):
            yield ((count), key)

if __name__ == '__main__':
    MaxExpense.run()
```

## o/p:

```
C:\Users\HP>python C:\Users\HP\Desktop\ML\empmaxexp.py --item=C:\Users\HP\Desktop\ML\expenses.txt C:\Users\HP\Desktop\ML\employee.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory C:\Users\HP\AppData\Local\Temp\empmaxexp.HP.20201012.170606.462308
Running step 1 of 2...
Running step 2 of 2...
job output is in C:\Users\HP\AppData\Local\Temp\empmaxexp.HP.20201012.170606.462308\output
Streaming final output from C:\Users\HP\AppData\Local\Temp\empmaxexp.HP.20201012.170606.462308\output...
"400\n"  ["110", "Priya"]
"400"    ["102", "Shiv"]
"300\n"  ["104", "Anubhav"]
"200\n"  ["114", "Maddy"]
"100\n"  ["105", "Palash"]
"100\n"  ["101", "Abhay"]
Removing temp directory C:\Users\HP\AppData\Local\Temp\empmaxexp.HP.20201012.170606.462308...

C:\Users\HP>
```

Using Pig:

Q4:

```
empentries - Notepad                                                                    —    □    ×
File Edit Format View Help
|
from mrjob.step import MRStep
from mrjob.job import MRJob

class Entries(MRJob):
    fields = []
    def configure_args(self):
        super(Entries, self).configure_args()
        self.add_file_arg('--items', help='Path to expenses.txt')

    def steps(self):
        return [
            MRStep(mapper=self.mapper,
                    reducer_init=self.reducer_init,
                    reducer=self.reducer),
            MRStep(reducer=self.reducer1)
        ]

    def mapper(self, key, value):
        (id, name, salary, rating) = value.split(',')
        yield [id,name], 1

    def reducer_init(self):
        self.expense = {}

        with open("expenses.txt") as f:
            for line in f:
                fields = line.split('\t')
                self.expense[fields[0]] = fields[1]

    def reducer(self, key, values):
        if key[0] in self.expense:
            yield None, key

    def reducer1(self, key, values):
        for key in sorted(values, reverse=True):
            yield key

if __name__ == '__main__':
    Entries.run()
```

o/p:

```
C:\Users\HP>python C:\Users\HP\Desktop\ML\empentries.py --item=C:\Users\HP\Desktop\ML\expenses.txt C:\Users\HP\Desktop\ML\employee.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory C:\Users\HP\AppData\Local\Temp\empentries.HP.20201012.170633.972624
Running step 1 of 2...
Running step 2 of 2...
job output is in C:\Users\HP\AppData\Local\Temp\empentries.HP.20201012.170633.972624\output
Streaming final output from C:\Users\HP\AppData\Local\Temp\empentries.HP.20201012.170633.972624\output...
"114"   "Maddy"
"110"   "Priya"
"105"   "Palash"
"104"   "Anubhav"
"102"   "Shiv"
"101"   "Abhay"
Removing temp directory C:\Users\HP\AppData\Local\Temp\empentries.HP.20201012.170633.972624...

C:\Users\HP>
```

## Using Pig:

## Q5:

```
empnoentries - Notepad
File Edit Format View Help

from mrjob.step import MRStep
from mrjob.job import MRJob

class NoEntries(MRJob):
    fields = []
    def configure_args(self):
        super(NoEntries, self).configure_args()
        self.add_file_arg('--items', help='Path to expenses.txt')

    def steps(self):
        return [
            MRStep(mapper=self.mapper,
                    reducer_init=self.reducer_init,
                    reducer=self.reducer),
            MRStep(reducer=self.reducer1)
        ]

    def mapper(self, key, value):
        (id, name, salary, rating) = value.split(',')
        yield [id,name], 1

    def reducer_init(self):
        self.expense = {}

        with open("expenses.txt") as f:
            for line in f:
                fields = line.split('\t')
                self.expense[fields[0]] = fields[1]

    def reducer(self, key, values):
        if key[0] not in self.expense:
            yield None, key

    def reducer1(self, key, values):
        for key in sorted(values, reverse=True):
            yield key

if __name__ == '__main__':
    NoEntries.run()
```

## o/p:

```
C:\Users\HP>python C:\Users\HP\Desktop\ML\empnoentries.py --item=C:\Users\HP\Desktop\ML\expenses.txt C:\Users\HP\Desktop\ML\employee.txt
No configs found; falling back on auto-configuration
No configs specified for inline runner
Creating temp directory C:\Users\HP\AppData\Local\Temp\empnoentries.HP.20201012.170653.945823
Running step 1 of 2...
Running step 2 of 2...
job output is in C:\Users\HP\AppData\Local\Temp\empnoentries.HP.20201012.170653.945823\output
Streaming final output from C:\Users\HP\AppData\Local\Temp\empnoentries.HP.20201012.170653.945823\output...
"113"   "Jay"
"112"   "Ajay"
"111"   "Tushar"
"109"   "Karan"
"108"   "Ram"
"107"   "Sahil"
"106"   "Aman"
"103"   "Aarav"
Removing temp directory C:\Users\HP\AppData\Local\Temp\empnoentries.HP.20201012.170653.945823...

C:\Users\HP>
```

## Using Pig: