# *Distributed Systems*

Assignment 01

## TASK #1:

I solved the 1st task by using 2 methods in the class AverageSalaryPerJob, which are: mapper and reducer.

a. **Mapper**: I used the mapper to map/ extract the job titles and salaries line by line from the ds_salaries text file (which I specified in the CLI). Then, I yielded the job titles and the salaries as key-value pairs.

b. **Reducer**: I used the reducer to reduce/ calculate the average salary for each job by its title. I used (statistics.mean(salaries)) to calculate the average salary for each job. Lastly, I yielded each job and its average salary.
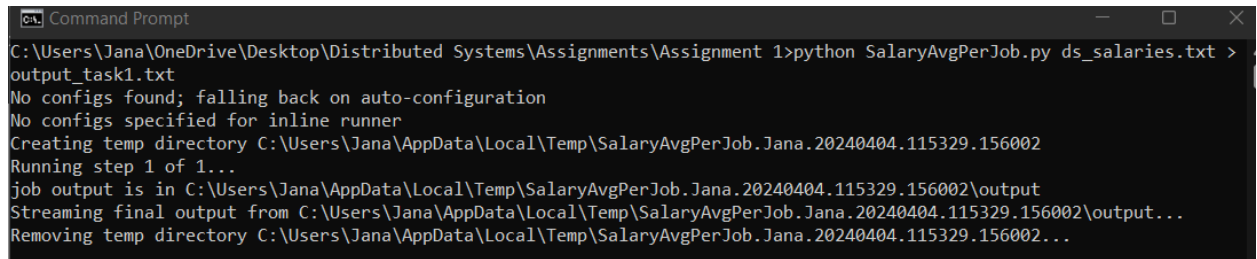
- **The code file**:

```python
# Libraries:
from mrjob.job import MRJob
import statistics
import pandas as pd

# Class implementation:
class AverageSalaryPerJob(MRJob):

    def mapper(self, _, line): # The mapper method is used to extract the job title and salary from the file that i'll specify in the command line.
        year, job_title, salary, company_location = line.split(',') # Splitting the lines by the commas.
        # Yielding the job title and salary as a key-value pair <job_title, salary>:
        yield job_title.strip(), float(salary) #.strip() removes leading and trailing whitespaces

    def reducer(self, job_title, salaries): # The reducer method is used to calculate the average salary for each job title.
        avg_salary = statistics.mean(salaries) # Calculating the average salary for each job title.
        # Yielding each job title and its average salary
        yield job_title, avg_salary

if __name__ == '__main__': # To run the MRJob class
    AverageSalaryPerJob.run()
    # Now, i'll extract from the output_task1.txt file the highest paying job:
    df = pd.read_csv('output_task1.txt', delimiter='\t', header=None, names=['Job_Title', 'Average_Salary']) # Reading the output file
    highest_paying_job = df[df['Average_Salary'] == df['Average_Salary'].max()] # Extracting the highest paying job
    print("Highest Paying Job:") # Printing the highest paying job
    print(highest_paying_job)
```

And then I Just ran the file by using the Command Line Prompt.

I used the command below to output the solution in an output file called: output_task1.txt

- **The command Line:**



- **The output_task1.txt file:**

```
salaryAvgPerJob.py        output_task1.txt  ×        salaryAvgPerCountry.py

output_task1.txt
 1    "3D Computer Vision Researcher" 5409.0
 2    "AI Scientist"  66135.57142857143
 3    "Analytics Engineer"    175000.0
 4    "Applied Data Scientist"    175655.0
 5    "Applied Machine Learning Scientist"    142068.75
 6    "BI Data Analyst"  74755.16666666667
 7    "Big Data Architect"    99703.0
 8    "Big Data Engineer" 51974.0
 9    "Business Data Analyst" 76691.2
10    "Cloud Data Engineer"   124647.0
11    "Computer Vision Engineer"  44419.333333333336
12    "Computer Vision Software Engineer" 105248.66666666667
13    "Data Analyst"  92893.06185567011
14    "Data Analytics Engineer"   64799.25
15    "Data Analytics Lead"   405000.0
16    "Data Analytics Manager"    127134.28571428571
17    "Data Architect"    177873.9090909091
18    "Data Engineer" 112725.0
19    "Data Engineering Manager"  123227.2
20    "Data Science Consultant"   69420.71428571429
21    "Data Science Engineer" 75803.33333333333
22    "Data Science Manager"  158328.5
23    "Data Scientist"    108187.83216783217
24    "Data Specialist"   165000.0
25    "Director of Data Engineering"  156738.0
26    "Director of Data Science"  195074.0
27    "ETL Developer" 54957.0
28    "Finance Data Analyst"  61896.0
29    "Financial Data Analyst"    275000.0
30    "Head of Data Science"  146718.75
31    "Head of Data"  160162.6
32    "Head of Machine Learning"  79039.0
33    "Lead Data Analyst" 92203.0
34    "Lead Data Engineer"    139724.5
35    "Lead Data Scientist"   115190.0
36    "Lead Machine Learning Engineer"    87932.0
37    "ML Engineer"   117504.0
38    "Machine Learning Developer"    85860.66666666667
39    "Machine Learning Engineer" 104880.14634146342
40    "Machine Learning Infrastructure Engineer"  101145.0
```

```
39    "Machine Learning Engineer" 104880.14634146342
40    "Machine Learning Infrastructure Engineer"  101145.0
41    "Machine Learning Manager"  117104.0
42    "Machine Learning Scientist"    158412.5
43    "Marketing Data Analyst"    88654.0
44    "NLP Engineer"  37236.0
45    "Principal Data Analyst"    122500.0
46    "Principal Data Engineer"   328333.3333333333
47    "Principal Data Scientist"  215242.42857142858
48    "Product Data Analyst"  13036.0
49    "Research Scientist"    109019.5
50    "Staff Data Scientist"  105000.0
51    Highest Paying Job:
52            |   |   |   Job_Title  Average_Salary
53    14  Data Analytics Lead       405000.0
54
```

So, we deduce that the highest paying job is "Data Analytics Lead" with average salary 405000.0

# TASK #2:

I solved the 2nd task by using 2 methods in the class AverageSalaryPerCountry, which are: mapper and reducer.

a. **Mapper**: I used the mapper to map/ extract the countries and salaries line by line from the ds_salaries text file (which I specified in the CLI). I extracted the countries from the company_location column: company_location.strip().split()[-1]. Then, I yielded the countries and the salaries as key-value pairs.

b. **Reducer**: I used the reducer to reduce/ calculate the average salary for each country. I used (statistics.mean(salaries)) to calculate the average salary for each country. Lastly, I yielded each country and its average salary.
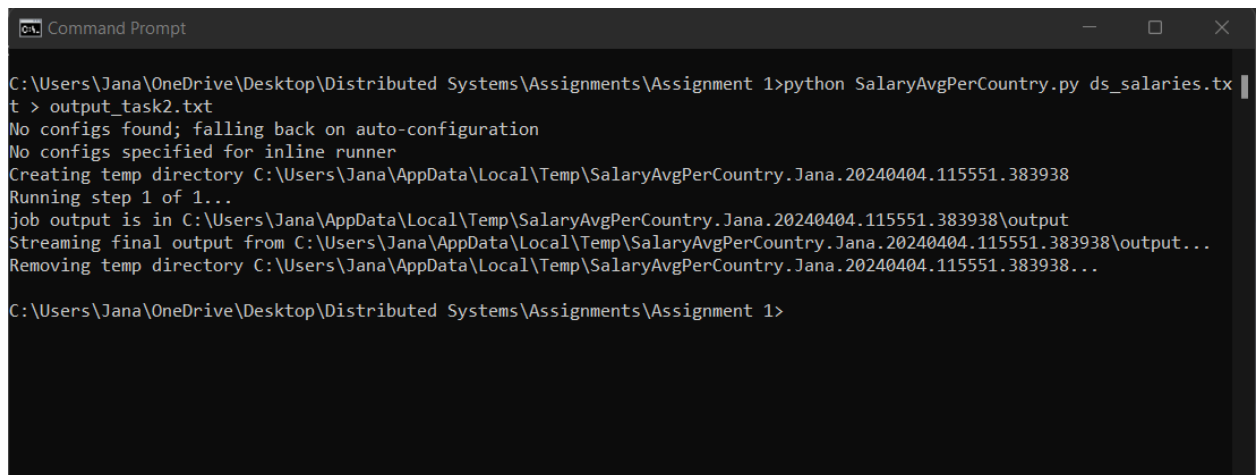
- **The code file:**

```python
# Libraries:
from mrjob.job import MRJob
import statistics
import pandas as pd

#Class implementation:
class AverageSalaryPerCountry(MRJob):

    def mapper(self, _, line): # The mapper method is used to extract the country and salary from the file that i'll specify in the command line.
        year, job_title, salary, company_location = line.split(',') # Splitting the lines by the commas.
        country = company_location.strip().split()[-1] # Extracting the country from the company location
        # Yielding the country and salary as a key-value pair <country, salary>:
        yield country, float(salary)

    def reducer(self, country, salaries): # The reducer method is used to calculate the average salary for each country.
        avg_salary = statistics.mean(salaries) # Calculating the average salary for each country.
        # Yielding each country and its average salary
        yield country, avg_salary

if __name__ == '__main__': # To run the MRJob class
    AverageSalaryPerCountry.run()
    # Now, i'll extract from the output_task1.txt file the highest paying country:
    df = pd.read_csv('output_task2.txt', delimiter='\t', header=None, names=['Country', 'Average_Salary']) # Reading the output file
    highest_paying_country = df[df['Average_Salary'] == df['Average_Salary'].max()] # Extracting the highest paying country
    print("Highest Paying Country:") # Printing the highest paying country
    print(highest_paying_country)
```
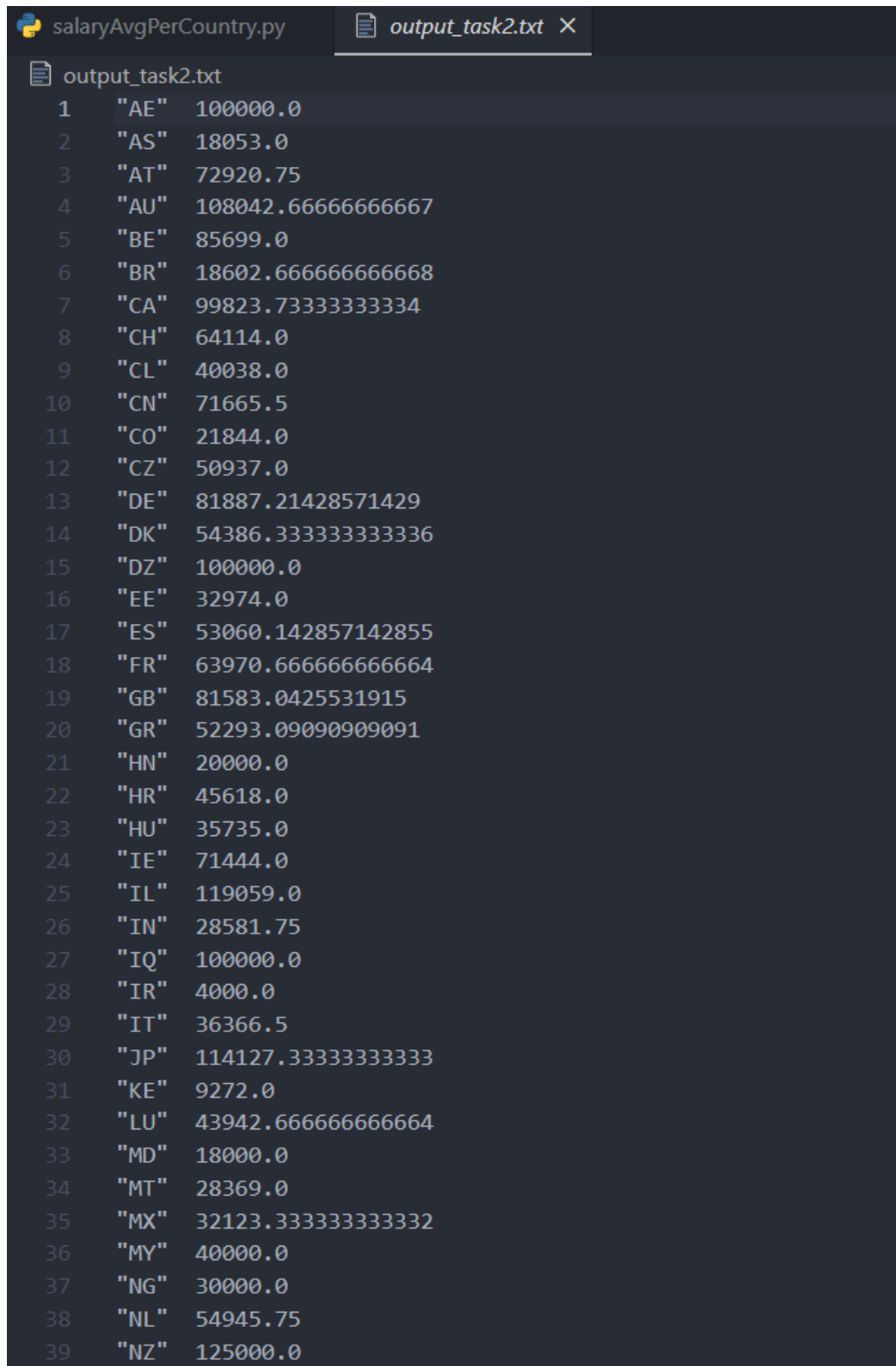
- **The CLI:**

- **The output_task2.txt file:**

```
output_task2.txt
 1    "AE"   100000.0
 2    "AS"   18053.0
 3    "AT"   72920.75
 4    "AU"   108042.66666666667
 5    "BE"   85699.0
 6    "BR"   18602.666666666668
 7    "CA"   99823.73333333334
 8    "CH"   64114.0
 9    "CL"   40038.0
10    "CN"   71665.5
11    "CO"   21844.0
12    "CZ"   50937.0
13    "DE"   81887.21428571429
14    "DK"   54386.333333333336
15    "DZ"   100000.0
16    "EE"   32974.0
17    "ES"   53060.142857142855
18    "FR"   63970.666666666664
19    "GB"   81583.0425531915
20    "GR"   52293.09090909091
21    "HN"   20000.0
22    "HR"   45618.0
23    "HU"   35735.0
24    "IE"   71444.0
25    "IL"   119059.0
26    "IN"   28581.75
27    "IQ"   100000.0
28    "IR"   4000.0
29    "IT"   36366.5
30    "JP"   114127.33333333333
31    "KE"   9272.0
32    "LU"   43942.666666666664
33    "MD"   18000.0
34    "MT"   28369.0
35    "MX"   32123.333333333332
36    "MY"   40000.0
37    "NG"   30000.0
38    "NL"   54945.75
39    "NZ"   125000.0
```

```
39    "NZ"   125000.0
40    "PK"   13333.333333333334
41    "PL"   66082.5
42    "PT"   47793.75
43    "RO"   60000.0
44    "RU"   157500.0
45    "SG"   89294.0
46    "SI"   63831.0
47    "TR"   20096.666666666668
48    "UA"   13400.0
49    "US"   144055.26197183097
50    "VN"   4000.0
51    Highest Paying Country:
52       Country  Average_Salary
53    43      RU          157500.0
54
```

So, we deduce that the highest paying country is "RU" with average salary 157500.0